# Physical Adversarial Attack Meets Computer Vision: A Decade Survey

Hui Wei, Hao Tang, Xuemei Jia, Zhixiang Wang,
Hanxun Yu, Zhubo Li, Shin'ichi Satoh, Luc Van Gool, and Zheng Wang

**Abstract**—Despite the impressive achievements of Deep Neural Networks (DNNs) in computer vision, their vulnerability to adversarial attacks remains a critical concern. Extensive research has demonstrated that incorporating sophisticated perturbations into input images can lead to a catastrophic degradation in DNNs' performance. This perplexing phenomenon not only exists in the digital space but also in the physical world. Consequently, it becomes imperative to evaluate the security of DNNs-based systems to ensure their safe deployment in real-world scenarios, particularly in security-sensitive applications. To facilitate a profound understanding of this topic, this paper presents a comprehensive overview of physical adversarial attacks. Firstly, we distill four general steps for launching physical adversarial attacks. Building upon this foundation, we uncover the pervasive role of artifacts carrying adversarial perturbations in the physical world. These artifacts influence each step. To denote them, we introduce a new term: adversarial medium. Then, we take the first step to systematically evaluate the performance of physical adversarial attacks, taking the adversarial medium as a first attempt. Our proposed evaluation metric, *hiPAA*, comprises six perspectives: *Effectiveness*, *Stealthiness*, *Robustness*, *Practicability*, *Aesthetics*, and *Economics*. We also provide comparative results across task categories, together with insightful observations and suggestions for future research directions.

**Index Terms**—Adversarial Attack, Physical World, Adversarial Medium, Computer Vision, Survey.

✦

## 1 INTRODUCTION

DEEP Neural Networks (DNNs) have achieved impressive results in a variety of fields: from computer vision [1] to natural language processing [2] to speech processing [3], and it is increasingly empowering many aspects of modern society. Nonetheless, Szegedy *et al.* [4] discovered in 2014 that adversarial samples can cause DNNs-based models to produce incorrect predictions, leading to a significant degradation in performance. This is a groundbreaking work that exposes the vulnerability of DNNs, casting a shadow over their reliability and security. Since then, researchers have conducted extensive explorations into adversarial samples, revealing their pervasive existence across a wide range of DNNs-based computer vision tasks [5], [6], [7], [8], [9], [10], [11], [12]. They designed adversarial clothing [13] to evade person detectors, adversarial eyeglasses [14] to deceive face recognizers, etc. Increasingly, these methods use a class of techniques called adversarial attacks.

Generally, adversarial attacks occur by adding imperceptible perturbations to input data (e.g., image, video) and

TABLE 1: Comparative analysis of current surveys on physical adversarial attack in computer vision.

| Survey | Physical Attack | Adversarial Medium | Evaluation | Number of Methods | Number of Tasks | Year |
|---|:---:|:---:|:---:|:---:|:---:|:---:|
| [15] | ✓ | ✗ | ✗ | 5 | 3 | 2018 |
| [16] | ✓ | ✗ | ✗ | 47 | 11 | 2022 |
| [17] | ✓ | ✗ | ✗ | 69 | 7 | 2022 |
| [18] | ✓ | ✗ | ✗ | 22 | 4 | 2023 |
| Ours | ✓ | ✓ | ✓ | 78 | 14 | 2023 |

fooling the trained DNNs-based models in the inference stage. Regarding the various domains, adversarial attacks can be categorized into two distinct classes: **(i) Digital Adversarial Attack**, which occurs in the digital space through the addition of subtle perturbations (e.g., style perturbations [19] and context-aware perturbations [20]). **(ii) Physical Adversarial Attack**, which occurs in the real world using tangible artifacts that contain adversarial perturbations (e.g., adversarial patches [21], [22] and adversarial stickers [23], [24]). Compared to the former, the latter pose an augmented threat to social security, raising significant apprehensions, particularly in safety-critical domains like autonomous driving, video surveillance, and facial biometric systems. To facilitate a profound understanding and provide in-depth insights into this topic, we present a comprehensive review of articles on physical adversarial attacks in computer vision tasks.

Though some existing surveys have also summarized the physical adversarial attack methods [15], [16], [17], [18], they primarily focus on listing and categorization, ignoring the evaluation and comparison (see TABLE 1). A unified evaluation criterion is still absent. This motivates us to take the first step to evaluate the performance of physical

- Hui Wei, Xuemei Jia, Hanxun Yu, and Zheng Wang are with the School of Computer Science, National Engineering Research Center for Multimedia Software, Wuhan University, Wuhan, P.R.China.
- Hao Tang is with the Department of Information Technology and Electrical Engineering, ETH Zurich, Zurich 8092, Switzerland.
- Zhubo Li is with the School of Cyber Science and Engineering, Wuhan University, Wuhan, P.R.China.
- Zhixiang Wang and Shin'ichi Satoh are with the Department of Information and Communication Engineering, Graduate School of Information Science and Technology, The University of Tokyo, Japan, and also with the Digital Content and Media Sciences Research Division, National Institute of Informatics, Japan.
- Luc Van Gool is with the Computer Vision Lab of ETH Zurich, 8092 Zürich, Switzerland, and also with KU Leuven, 3000 Leuven, Belgium, and INSAIT, Sofia.
- Zheng Wang is the corresponding author. E-mail: wangzwhu@whu.edu.cn
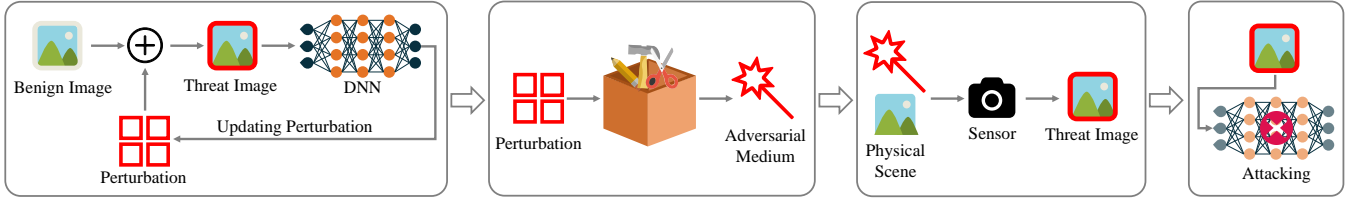
Fig. 1: The flow of designing a physical adversarial attack, including four main steps: 1) *Adversarial perturbation generation in the digital space*, 2) *Adversarial medium manufacturing in the physical space*, 3) *Threat image capturing*, and 4) *Attacking*.

adversarial attacks systematically.

To this end, we outline the four general steps (as shown in Fig. 1) required to build a physical adversarial attack:

1) Step 1: *Adversarial perturbation generation*. Generating perturbations in the digital domain based on given DNNs-based models, constrained by different attack forms and attack objectives.
2) Step 2: *Adversarial medium manufacturing*. Designing appropriate physical medium for carrying the perturbations in alignment with attack forms, and subsequently manufacturing them using suitable materials.
3) Step 3: *Threat image capturing*. Applying the adversarial medium in real-world scenarios to be captured by an imaging sensor, thereby generating threat images.
4) Step 4: *Attacking*. The captured threat images are fed to the DNNs-based model to initiate attacks.

Note that we introduce the concept of an "adversarial medium" to denote the tangible artifact responsible for carrying the adversarial perturbation in the real world. According to the four steps mentioned above, we discern the significant role of adversarial mediums in building a physical adversarial attack. They determine the form of perturbations (Step 1), impact manufacturing processes (Step 2), and hold relevance for real-world applications (Step 3). Therefore, we embrace an approach centered on adversarial mediums to examine the existing methods, systematically quantifying and evaluating them in the following six perspectives: *Effectiveness*, *Stealthiness*, *Robustness*, *Practicability*, *Aesthetics*, and *Economics*. Meanwhile, we introduce a comprehensive metric, the hexagonal indicator of Physical Adversarial Attack (*hiPAA*), and provide comparative results across task categories, along with insightful observations and suggestions for future research directions.

The major contributions of this work can be summarized as follows:

1) Through a comprehensive review of existing methodologies (see Tables 2, 3, 4), we abstract and summarize a general workflow for launching a physical adversarial attack, comprising four distinct steps (see Fig. 1).
2) Leveraging this general workflow, we discover that tangible artifacts carrying adversarial perturbations exert substantial influence over attacks, prompting the introduction of the new concept of adversarial medium to represent them.
3) As opposed to existing reviews [15], [16], [17], [18], we take the first step to systematically evaluate the performance of physical adversarial attack, taking adversarial medium as a first attempt. Our proposed evaluation metric, *hiPAA*, comprises six perspectives: *Effectiveness*,

*Stealthiness*, *Robustness*, *Practicability*, *Aesthetics*, and *Economics*.

4) We conduct comprehensive comparisons of existing physical adversarial attack methods, and discuss limitations, challenges, and potential directions from the standpoint of real-world applications to facilitate future research.

The remaining article is organized as follows. We first provide a brief introduction to the preliminaries in Section 2, which cover essential topics and concepts for the proper understanding of this work. Section 3 discusses the concept of adversarial mediums. Section 4 introduces the evaluation metric for conducting a comprehensive comparison. Then we present the recent advancements in physical adversarial attacks according to mainstream visual tasks and systematically evaluate them for each task in Section 5. Moreover, we provide discussions and future research opportunities in Section 6. Finally, we conclude this review in Section 7. We also provide a regularly updated project page on https://github.com/weihui1308/PAA.

## 2 PRELIMINARIES

In this section, we provide a concise introduction to problem formulation, key topics, and concepts, aiming to enhance comprehension of our work.

### 2.1 Computer Vision

Computer Vision (CV) aims to enable machines to perceive, observe, and understand the physical world like human eyes. An important milestone was reached when Krizhevsky *et al.* [25] proposed AlexNet, which secured victory in the ILSVRC (ImageNet Large Scale Visual Recognition Challenge), thus promoting the application of DNNs to address a wide variety of tasks. Up to the present, DNNs-based models have achieved impressive and competitive performances in CV tasks, including classification [26], [27], [28], segmentation [29], [30], [31], detection [32], [33], [34], image generation [35], [36], [37], and re-identification [38], [39], [40]. This survey concentrates on three prominent tasks: classification, detection, and re-identification, which are extensively utilized and encompass numerous attack scenarios. Given a DNNs-based model $f : X \rightarrow Y$ with pre-trained weights $\theta$, for the sake of conciseness, we formulate the DNNs-based CV models as follows:

$$\hat{y} = f(\theta, x), \qquad x \in X, y \in Y, \qquad (1)$$

where for any input data $x \in X$, the well-trained model $f(\cdot)$ is able to predict a $\hat{y}$ that closely approximates the corresponding ground truth $y \in Y$.
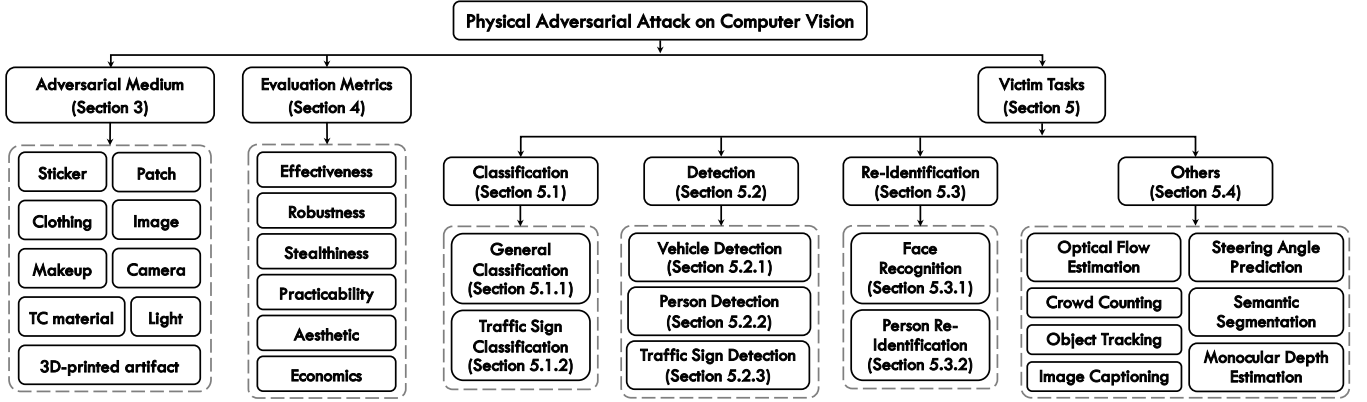
Fig. 2: A general overview of the scope in our survey.

## 2.2 Adversarial Attacks

### 2.2.1 Problem Formulation

Adversarial attacks refer to modifying input data by introducing perturbations that lead the model $f(\cdot)$ to make incorrect predictions $y'$. Note that the attacker's modification is limited to the input data. Exactly, the sample $x$, following the addition of perturbations $\delta$, is denoted as the adversarial sample $x'$. Mathematically, the joint representation is expressed as

$$\begin{cases} x' = x + \delta \\ y' = f(\theta, x') \end{cases} \quad s.t. \ \ y' \neq \hat{y}, \tag{2}$$

where perturbations $\delta$ are typically constrained by factors such as intensity, size, and the adversarial medium.

### 2.2.2 Distinguishing Adversarial Attacks from Backdoor Attacks and Poisoning Attacks.

Apart from adversarial attacks, there are two other widely used attack types: backdoor (a.k.a. trojan) attacks [41], [42], [43], [44] and poisoning attacks [45], [46], [47]. Observing the outcomes, these three categories of attacks do indeed appear analogous, as they all aim to misguide the model into producing incorrect predictions. However, their implementation methods diverge fundamentally. Herein, we elucidate these distinctions. For clarity, we introduce the DNNs-based model lifecycle, as shown in Fig. 3, comprising six discrete phases: data collection and preparation, model selection, model training, model testing, model deployment, and model updating.

**Poisoning Attack.** DNNs-based models are data-hungry, and achieving outstanding performance requires a substantial amount of data.Developers often utilize publicly available datasets or scrape data from online sources to train their models, which presents ample opportunities for poisoning attacks. In a poisoning attack, toxic samples are introduced into the model's training dataset. The inclusion of these samples during model training can result in diminished learning efficiency and, in certain instances, hinder convergence, ultimately causing disruption to the entire learning process [48]. Evidently, poisoning attacks occur during the data collection and preparation phase.

**Backdoor Attack.** Backdoor attacks have emerged as a significant research area in AI security in recent years.

Essentially, attackers manipulate neural network models by contaminating training data, altering model weights, or modifying the model architecture to make them produce specific outputs when presented with inputs containing specific triggers [49]. Backdoor attacks can be implemented at multiple stages in the model lifecycle, employing diverse methods like code poisoning [50], data poisoning [41], and control of the training process [51]. Backdoor attacks should be inconspicuous to users, resistant to removal, and should not disrupt the normal functioning of DNNs.

**Adversarial Attacks.** Compared to the two attack types mentioned above, adversarial attacks have the weakest underlying assumption: they solely modify input data, without any other alterations. This characteristic facilitates attackers in conducting practical applications more easily. Adversarial attacks occur during both the model testing and model deployment phases.

## 2.3 The Taxonomy of Adversarial Attacks

In general, adversarial attacks can be divided into three categories according to adversarial knowledge: white-box attacks, black-box attacks, and gray-box attacks.

**White-Box Attacks:** In a white-box attack, the attacker has access to both the data and model knowledge, including details like network structure, parameters, weights, and activation function types. Attackers leverage this information to assess the model's vulnerabilities and adapt their attack strategies accordingly. Therefore, this attack method is relatively easy to implement.

**Black-Box Attacks:** Black-box attacks mean that the structure and parameters of the target model are not accessible to the attacker, where the attack can only interact with the model to acquire the predictions of the corresponding samples. Employing paired data, i.e., samples and their predictions, the black-box attacker trains the substitute model to perform adversarial attacks, which resembles a real-world situation. Due to the transferability, the synthesized adversarial samples achieve great attack performance.

**Gray-Box Attacks:** Attacks under the gray-box setting regulate the attacker's interaction with the model outside and are aware of the model structure information. In other words, the attacker has query access to the model but not the model parameters. In this way, the attacker can use the already known structure information to construct a more accurate
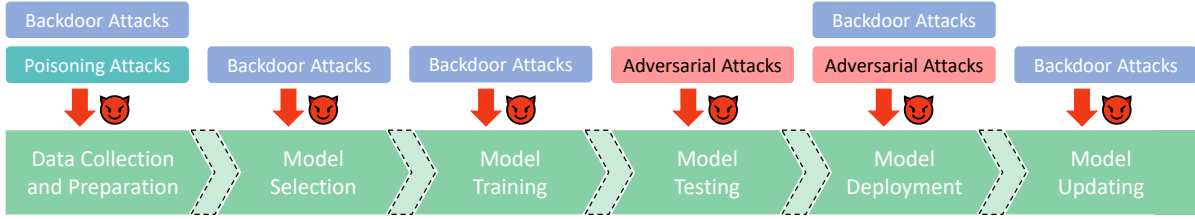
Fig. 3: An overview of the lifecycle in which the three types of attacks occur. Adversarial attacks occur only during the model deployment phase, without modifying the model and training data. Compared to backdoor attacks and poisoning attacks, adversarial attacks have weaker assumptions, focusing on the vulnerability of the model itself.

substitute classifier and then carry out the attack, which outperforms the black-box attack.

In addition, taking into consideration the adversarial goals, adversarial attacks can be divided into the following two categories:

**Targeted Attacks:** Targeted attacks aim to mislead DNNs-based models to the specified labels. For example, the attack recognizes the face of a state official as that of a specified criminal. The specified category labels make it more challenging to achieve targeted attacks with high success rates and considerable stealthiness and robustness.

**Untargeted Attacks:** The untargeted attacks mislead DNNs-based models to any wrong label. Designing an effective untargeted attack only requires making the model's predictions incorrect, without necessarily concerning what the incorrect results are.

### 2.4 Introduction to Physical Adversarial Attacks

Physical adversarial attacks refer to attacking DNNs-based models in the real world. As shown in Fig. 1, we summarize the general process of physical adversarial attack as four steps: 1) *Adversarial perturbation generation*, 2) *Adversarial medium manufacturing*, 3) *Threat image capturing*, and 4) *Attacking*.

Unlike digital adversarial attacks, which typically involve imperceptible perturbations to the human eye, as demonstrated in the One Pixel Attack [52], physical adversarial attacks demand more intense perturbations. This level of intensity is crucial for the perturbation to be detectable by sensors in real-world scenarios. Meanwhile, physical adversarial attacks face additional challenges due to of various physical constraints (e.g., spatial deformation) and environmental dynamics (e.g., lighting). These factors can degrade or even eliminate the effectiveness of adversarial perturbations.

It is worth noting that backdoor attacks can also be executed in the physical world [53]. For instance, Qi *et al.* [44] proposed a physically realizable backdoor attack algorithm against image classification tasks. However, it's important to clarify that this category is not within the scope of our discussion, as our primary focus is on physical adversarial attacks.

## 3 ADVERSARIAL MEDIUMS

We observe a commonality across all physical adversarial attack methods: the necessity of a physical entity to carry the specially designed perturbations. Therefore, we introduce the novel concept of adversarial medium for the first time to represent these entities. The *medium* is commonly used in physics. When a substance exists inside another substance, the latter is the medium of the former [54]. In most cases, properties such as form, density, and shape of the former would influence the properties of the latter. Analogously, for launching attacks in the physical space, the adversarial perturbation must have a carrier, i.e., the adversarial perturbation exists inside the carrier. Meanwhile, the carrier has an effect on the perturbation. Thus we define the carrier as the adversarial medium.

An adversarial medium is indispensable for physical adversarial attacks. Several milestone methods and their corresponding adversarial mediums are illustrated in Fig. 4. Additionally, as shown in Tables 2, 3, 4, we have compiled a comprehensive list of all the physical adversarial attack methods examined in this paper, organized by medium and chronological order. In this section, we discuss all existing adversarial mediums.

**Sticker** performs attacks by attaching to the surface of targeted objects, e.g., cars, traffic signs, and persons. Since it can be pasted over a large area, it is competitive in terms of effectiveness and robustness. But it is hard to hide. Stealthiness is usually improved by making the pattern of the sticker close to the natural textures.

**Patch,** which restricts perturbations to a small localized region without imposing any intensity limitations, is commonly used in a variety of attack tasks. Its popularity is primarily due to its user-friendly nature, often requiring nothing more than a straightforward printout. The patch and sticker share similar manufacturing processes. The key distinction lies in the fact that patches typically have regular shapes, while stickers are irregular and can adapt to non-rigid transformations.

**Clothing**, such as t-shirts, capes, or pants, can be utilized as carriers for adversarial perturbations in real-world scenarios. These perturbations can be integrated into wearable clothing, significantly enhancing their concealment and practicality. This approach is extensively employed to deceive person detectors and evade surveillance systems.

**Image** can also serve as a common adversarial medium. This form of attack disperses perturbations throughout the entire image, leading to incorrect predictions by DNN-based models. The primary challenge with this medium lies in effectively concealing the perturbation while maintaining attack efficacy.

**Light** can initiate a rapid and highly stealthy physical-world attack. This approach does not directly modify the object but instead leverages a lighting device (such as a laser pointer) or projector to project specific light onto the target object. However, the effectiveness of this method
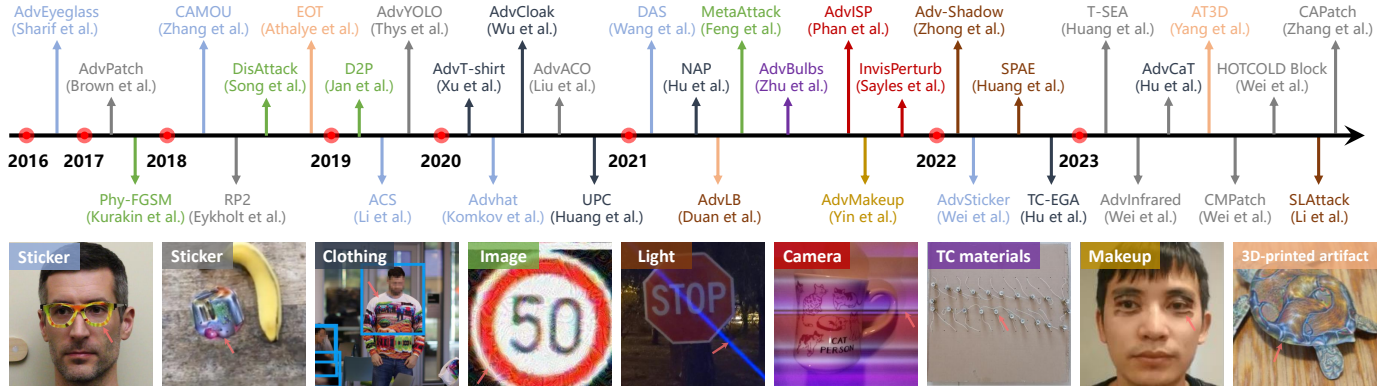
Fig. 4: Chronological overview of the milestone physical adversarial attach methods using different adversarial mediums. Colors indicate the category of adversarial mediums used.

diminishes in environments with strong lighting conditions, as the perturbations carried by the light are not visible.

**Camera** performs real-world attacks by modifying the camera instead of the targeted object. Two types of camera-based attacks have been developed: the camera's rolling shutter effect and the camera's image signal processing (ISP). Since the modification takes place within the camera, this method is highly stealthy.

**TC materials** is an abbreviation for temperature control materials. This medium is exclusively employed for physical adversarial attacks within the thermal infrared imaging modality. For instance, small glowing light bulbs appear as light spots shaped like Gaussian distributions when viewed through a thermal infrared camera. Exploiting this characteristic, attackers can create intricate patterns to disrupt DNNs-based models.

**Makeup** can successfully launch adversarial attacks on face recognition algorithms in real-world scenarios by employing special makeup techniques. The makeup is carefully designed to blend seamlessly with the human face, making it a highly stealthy adversarial medium.

**3D-printed artifact** utilize 3D-printing technology to create adversarial objects capable of deceiving classifiers in the physical world. This adversarial medium maintains robustness against physical constraints owing to its three-dimensional structure.

## 4 EVALUATION METRIC

Despite the widespread attention given to physical adversarial attacks, their performance has been subject to inconsistent and case-by-case evaluations. This phenomenon primarily stems from two factors:

1) *Challenges in Replication*: The creation of adversarial mediums, these physical entities, often incur high production costs and face various inherent factors during the manufacturing process. Factors such as the variability in printer models used for pattern printing and the difficulties in standardizing aspects like clothing style and color can influence the process. Furthermore, maintaining consistent real-world experimental conditions proves to be a daunting task. These challenges collectively hinder the replication of experiments conducted by others, making comparisons arduous.

2) *Varied Objectives*: Physical adversarial attack methods pursue diverse dimensions of enhancement. Some methods prioritize increasing attack effectiveness [13],

[21], [88], while others place emphasis on improving attack stealthiness [22], [59], [68]. A separate category strives to enhance robustness [91], [94], [102]. Consequently, a single metric struggles to impartially assess multiple methods.

These problems raise our question: *how does the actual performance of the physical adversarial attacks?* In this section, we take the first step in conducting a comprehensive assessment of physical adversarial attack methods, encompassing all the work in this field.

Since the diverse objectives across various works, we summarize six perspectives from existing literature, namely: *Effectiveness*, *Stealthiness*, *Robustness*, *Practicability*, *Aesthetics*, and *Economics*. Hingun *et al.* [119] introduced a large-scale realistic adversarial patch benchmark to assess the effectiveness of attacks. Parallel to this work, Wang *et al.* [120] evaluated the system-level effect of attack methods at autonomous driving. Li *et al.* [121] proposed an evaluation of visual naturalness These evaluation were conducted in settings with a *single perspective* and a *single task*. Towards a unified evaluation, we introduce the hexagonal indicator of Physical Adversarial Attack (*hiPAA*) to systematically quantify and compare attack methods across the aforementioned six perspectives. We define the *hiPAA* as the weighted sum of six components:

$$\begin{aligned} hiPAA = &\lambda_1 \cdot Eff. + \lambda_2 \cdot Rub. + \lambda_3 \cdot Ste. \\ &+ \lambda_4 \cdot Aes. + \lambda_5 \cdot Pra. + \lambda_6 \cdot Eco., \end{aligned} \tag{3}$$

where the values of $\lambda_1$ to $\lambda_6$ are assigned based on the importance of each component, and we have set them as {0.3, 0.2, 0.2, 0.1, 0.1, 0.1}, respectively.

### 4.1 Effectiveness

Attack effectiveness is employed to evaluate the influence a method can exert on the victim model. To evaluate the effectiveness of physical adversarial attacks, our primary concern lies in quantifying the degree of performance degradation induced by these attack methods. Note that while the attack success rate (ASR) serves as a prevalent metric, its applicability across all tasks is not universal. For instance, defining attack success in segmentation tasks can pose inherent challenges. Furthermore, to address variations across different tasks, we compute the percentage of performance degradation:

$$Eff. = 1 - Acc'/Acc, \tag{4}$$

TABLE 2: Physical adversarial attack methods that use **stickers** and **patches** as adversarial mediums. We list them by the adversarial medium and time order.

| Adversarial Medium | Description | | | Method | Victim Task | Venue | Year |
|---|---|---|---|---|---|---|---|
| | Manufacture | Instrument | Attack Type | | | | |
| Sticker | print, paste | eyeglasses | impersonation, dodging | AdvEyeglass [14] | Face Recognition | ACM CCS | 2016 |
| | cover | car body | hiding | CAMOU [23] | Vehicle Detection | ICLR | 2018 |
| | display | screen | hiding | InvisibleCloak [55] | Person Detection | UEMCON | 2018 |
| | paste | camera lens | misclassification | ACS [56] | General Classification | PMLR | 2019 |
| | print, paste | eyeglasses | impersonation, dodging | AdvEyeglass+ [57] | Face Recognition | TOPS | 2019 |
| | print, affix | hat | impersonation | Advhat [24] | Face Recognition | ICPR | 2020 |
| | print, paste | eyeglasses | misrecognition | CLBAAttack [58] | Face Recognition | BIOSIG | 2021 |
| | affix | car body | hiding | DAS [59] | Vehicle Detection | CVPR | 2021 |
| | affix | road marking | misdirection | AdvMarkings [60] | Lane Detection | USENIX | 2021 |
| | full cover | car body | hiding | FCA [61] | Vehicle Detection | AAAI | 2022 |
| | full cover | car body | hiding | DTA [62] | Vehicle Detection | CVPR | 2022 |
| | imprint | face mask | dodging | AdvMask [63] | Face Recognition | ECML PKDD | 2022 |
| | print, paste | face | impersonation | AdvSticker [64] | Face Recognition | TPAMI | 2022 |
| | cover | car body | hiding | CAC [65] | Vehicle Detection | IJCAI | 2022 |
| Patch | print, put | image patch | misclassification | AdvPatch [21] | General Classification | NIPS | 2017 |
| | print, paste | traffic sign | misclassification | RP$_2$ [66] | Sign Classification | CVPR | 2018 |
| | print, paste | traffic sign | misdetection | NestedAE [67] | Sign Detection | CCS | 2019 |
| | print, paste | traffic sign | misclassification | PS-GAN [68] | Sign Classification | AAAI | 2019 |
| | display | screen | lose track | PAT [69] | Object Tracking | ICCV | 2019 |
| | print | image patch | hiding | AdvYOLO [13] | Person Detection | CVPRW | 2019 |
| | print, paste | image patch | mismatching | AdvPattern [70] | Person Re-ID | ICCV | 2019 |
| | imprint | image patch | false estimation | FlowAttack [71] | Flow Estimation | ICCV | 2019 |
| | print, paste | image patch | misclassification | AdvACO [72] | General Classification | ECCV | 2020 |
| | paste | camera lens | misdetection | TransPatch [73] | Sign Detection | CVPR | 2021 |
| | print, paste | image patch | lose track | MTD [74] | Object Tracking | AAAI | 2021 |
| | print, paste | face | impersonation | TAP [75] | Face Recognition | CVPR | 2021 |
| | print, paste | image patch | misclassification | AdvACO+ [76] | General Classification | TIP | 2021 |
| | display | screen | misdetection | AITP [77] | Sign Detection | ACM AISec | 2022 |
| | print, paste | image patch | misclassification | CPAttack [78] | General Classification | NIPS | 2022 |
| | print, paste | image patch | misclassification | TnTAttack [79] | General Classification | TIFS | 2022 |
| | print, paste | image patch | false estimation | OAP [80] | Depth Estimation | ECCV | 2022 |
| | print | image patch | false segmentation | RWAEs [81] | Segmentation | WACV | 2022 |
| | print, paste | image patch | false estimation | PAP [82] | Crowd Counting | ACM CCS | 2022 |
| | print, put | image patch | misclassification | DAPatch [83] | General Classification | ECCV | 2022 |
| | print, paste | face | impersonation | SOPP [84] | Face Recognition | TPAMI | 2022 |
| | print, paste | car | hiding | AerialAttack [85] | Vehicle Detection | WACV | 2022 |
| | paste | aerogel patch | hiding | AdvInfrared [86] | Person Detection | CVPR | 2023 |
| | display | screen | hiding | T-SEA [87] | Person Detection | CVPR | 2023 |
| | paste | aerogel patch | hiding | CMPatch [88] | Person Detection | ICCV | 2023 |
| | print | image patch | misstatement | CAPatch [89] | Image Captioning | USENIX | 2023 |

where $Acc$ and $Acc'$ represent the model's accuracy without and with attacks, respectively.

## 4.2 Robustness

Attack robustness is employed to evaluate the method's ability to withstand continuous attacks. To evaluate the robustness of physical adversarial attack, we consider three evaluation scenarios: 1) Cross-model: whether the attack method remains effective when the model changes, 2) Cross-scenario: whether the attack method can consistently perform in various real-world scenarios, and 3) Transformation-resistant: whether the attack method can withstand various real-world transformations, including rotation, camera-to-object distances, view angles, *etc*. As shown in Table 7, we assess the areas in which the attack method provides countermeasures. A higher number of countermeasure strategies indicates greater attack robustness.

## 4.3 Stealthiness

Attack stealthiness is employed to evaluate whether the attack method is not easily detectable by human observers. In real-world attacks, stealthiness is crucial, as conspicuous attacks are easily detectable and can be thwarted by humans. Li *et al.* [121] introduced the Dual Prior Alignment (DPA) network to assess the visual naturalness of physical adversarial attacks. Stealthiness and naturalness are different terms but focus on the same attribute.

Since stealthiness is targeted at human observers, we conduct a user study in which participants rate the images using a 5-point Absolute Category Rating (ACR) scale [121].

## 4.4 Aesthetics

Aesthetics evaluates the social acceptability of the attack method. This metric is crucial, especially for wearable adversarial mediums [90], [93], [94]. Unusual patterns or appearances may lead users to reject them, whereas solutions that prioritize aesthetics are more likely to be accepted. Similar to stealthiness, we assess aesthetics through human ratings.

## 4.5 Practicability

Practicability evaluates the practicality and feasibility of using the attack method in real-world scenarios. It considers

TABLE 3: Physical adversarial attack methods that use **clothing**, **Images**, and **lights** as adversarial mediums. We list them by the adversarial medium and time order.

| Adversarial Medium | Description | | | Method | Victim Task | Venue | Year |
|---|---|---|---|---|---|---|---|
| | Manufacture | Instrument | Attack Type | | | | |
| Clothing | print | T-shirt | hiding | AdvT-shirt [90] | Person Detection | ECCV | 2020 |
| | print | pants, sweaters, mask | misdetection | UPC [91] | Person Detection | CVPR | 2020 |
| | print | sweatshirt | hiding | AdvCloak [92] | Person Detection | ECCV | 2020 |
| | print | sweatshirt | hiding | NAP [93] | Person Detection | ICCV | 2021 |
| | print | T-shirt | hiding | LAP [22] | Person Detection | ACM MM | 2021 |
| | print | dresses, T-shirts, skirts | hiding | TC-EGA [94] | Person Detection | CVPR | 2022 |
| | crop | aerogel material | hiding | InvisClothing [95] | Person Detection | CVPR | 2022 |
| | print | pants, sweatshirt | hiding | AdvCaT [96] | Person Detection | CVPR | 2023 |
| Image | print | picture | misclassification | Phy-FGSM [97] | General Classification | ICLR | 2017 |
| | print | picture | misdetection | ShapeShifter [98] | Sign Detection | ECML PKDD | 2018 |
| | print | picture | misdetection | RP$_2$+ [99] | Sign Detection | USENIX | 2018 |
| | print | picture | misclassification | D2P [100] | General Classification | AAAI | 2019 |
| | print | picture | misclassification | ABBA [101] | General Classification | NIPS | 2020 |
| | print | picture | false prediction | PhysGAN [102] | Steering Angle Prediction | CVPR | 2020 |
| | print | picture | misclassification | AdvCam [103] | General Classification | CVPR | 2020 |
| | print | picture | hiding | LPAttack [104] | Sign Detection | AAAI | 2020 |
| | print | picture | misclassification | MetaAttack [105] | General Classification | ICCV | 2021 |
| | print | picture | misclassification | Viewfool [106] | General Classification | NIPS | 2022 |
| | print | picture | misclassification | Meta-GAN [107] | General Classification | TIFS | 2023 |
| Light | project | projector | misclassification | PTAttack [108] | General Classification | AAAI | 2018 |
| | project | projector | misclassification | Poster [109] | General Classification | IEEE S&P | 2019 |
| | project | projector | impersonation | ALPA [110] | Face Recognition | CVPR | 2020 |
| | project | projector | hiding | SLAP [111] | Sign Detection | USENIX | 2021 |
| | project | projector | misclassification | OPAD [112] | General Classification | ICCV | 2021 |
| | project | laser pointer | misclassification | AdvLB [113] | General Classification | CVPR | 2021 |
| | project | shadow | misclassification | Adv-Shadow [114] | Sign Classification | CVPR | 2022 |
| | project | projector | misclassification | SPAA [115] | General Classification | IEEE VR | 2022 |
| | project | projector | hiding | AdvLight [116] | Vehicle Detection | ICASSP | 2023 |
| | project | laser pointer | hiding | AdvLS [117] | Sign Detection | PMLR | 2023 |
| | project | projector | impersonation | SLAttack [118] | Face Recognition | CVPR | 2023 |

factors such as the ease of implementation, availability of resources, and compatibility with existing systems. The higher the practicability, the more feasible and convenient it is to employ the adversarial medium in practical applications. Similar to stealthiness, we assess practicability through human ratings.

### 4.6 Economics

Economics evaluates the resource requirements and expenses associated with implementing and deploying the attack methods. Existing methods [122], [123] often provide descriptions of the costs associated with manufacturing adversarial mediums. We assess affordability based on these costs. For research that does not explicitly detail these expenses, we make estimations based on experience. Lower costs indicate greater cost-effectiveness.

## 5 VICTIM TASKS

In this section, we systematically review physical adversarial attacks according to their respective task categories. As shown in Fig. 2, to our knowledge, physical adversarial attacks currently involve 14 sub-tasks. We classify these tasks into four primary groups: classification, detection, re-identification, and others.

### 5.1 Attacks on Classification Tasks

Since classification is a fundamental task with significant downstream implications, numerous studies focus on attacking classifiers. Physical adversarial attacks on classification tasks primarily focus on general classification and

traffic sign classification. Please refer to Fig. 5 and Fig. 7 for illustrations.

#### 5.1.1 General Classification

Brown *et al.* [21] proposed the adversarial patch, a milestone in physical adversarial attack methods. The attacker places a printed circular image patch next to the "banana", effectively deceiving the classifier into recognizing the banana as a "toaster" with high confidence (see Fig. 5). Due to its simplicity in production, ease of deployment, and high attack potency, the adversarial patch immediately garnered widespread attention. Liu *et al.* [72], [76] developed class-agnostic universal adversarial patches with robust generalization capabilities for attacking classifiers in Automatic Check-out scenarios. Casper *et al.* [78] designed the "copy/paste" attacks, employing adversarial patches to investigate the reliability and interpretability of models. To unleash the patch-based attack potential, Chen *et al.* [83] designed a Deformable Adversarial Patch (DAPatch) to explore the impact of patch shapes. They simultaneously optimized shape and texture to enhance the attack performance. While these methods exhibit respectable attack performance, they are prone to visual abnormalities that render them conspicuously noticeable. To generate a natural-looking patch, Doan *et al.* [79] proposed searching for naturalistic patches with adversarial effects within the latent space $z$ of Generative Adversarial Networks (GANs) [130].

In contrast to patch-based methods that only add perturbations in a limited region, another category of methods adds perturbations across the entire image. Phy-FGSM [97]

TABLE 4: Physical adversarial attack methods that use niche adversarial mediums. We list them by the adversarial medium and time order.

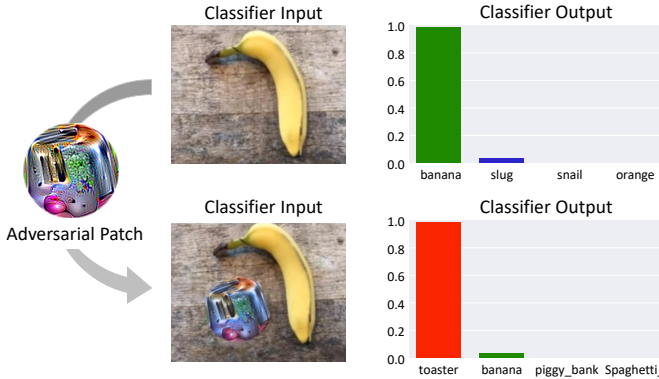| Adversarial Medium | Description | | | Method | Victim Task | Venue | Year |
|---|---|---|---|---|---|---|---|
| | Manufacture | Instrument | Attack Type | | | | |
| Camera | capture | rolling shutter effect | misclassification | InvisPerturb [124] | General Classification | CVPR | 2021 |
| | capture | image processing pipeline | misclassification | AdvISP [125] | General Classification | CVPR | 2021 |
| TC material | heating | cardboard | hiding | AdvBulbs [122] | Person Detection | AAAI | 2021 |
| | paste | warm/cool paste | hiding | HOTCOLD Block [123] | Person Detection | AAAI | 2023 |
| Makeup | cosmetics | face | impersonation | AdvMakeup [126] | Face Recognition | IJCAI | 2021 |
| 3D-printed artifact | 3D-print | tangible turtle | misclassification | EOT [127] | General Classification | PMLR | 2018 |
| | rendering | renderer | misclassification | 3DAttack [128] | General Classification | CVPR | 2019 |
| | 3D-print | tangible mesh | impersonation | AT3D [129] | Face Recognition | CVPR | 2023 |



Fig. 5: Display of the physical adversarial attack in general classification tasks. Initially, the classifier accurately labels the image as "banana". However, when an adversarial patch is placed adjacent to the banana, the classifier misclassifies the image as "toaster", despite the continued presence of the banana. Adapted from AdvPatch [21].

adds carefully designed perturbations to images from the ImageNet dataset [131], prints them out, and then captures them using a cell phone camera. The results demonstrate that the captured images can still successfully attack Inception v3 image classifier [132]. Given the generation of images in the digital domain and their subsequent application in the physical domain, a gap emerges between these two domains. D2P [100] leverages conditional GANs [133] to model the digital-to-physical transformation, aiming to alleviate the influence of this gap on attack performance. MetaAttack [105] and Meta-GAN [107] generate robust adversarial examples to maintain attacks in physical dynamics. They are class-agnostic and model-agnostic. These perturbations are subtle and constrained within a certain range $\varepsilon$, i.e., $\|\delta\| < \varepsilon$, but they are visible to the human eye, unlike digital adversarial attacks [52] that are imperceptible. To make these perturbations appear natural, ABBA [101] disguises the perturbations as motion blur, and AdvCam [103] conceals them within natural styles. In addition to adding perturbations for attacks, Viewfool [106] claims that changes in the viewpoint of the target object can also affect the classifier's predictions. It utilizes adversarial viewpoints to launch attacks and assess the classifier's robustness.

Inspired by One Pixel Attack [52], Nichols et al. [108] introduced light-based adversarial attacks on classifiers for the first time. This attack can be initiated simply by shining a specially designed demonstrating light on the target object. However, this method has only been validated on low-resolution image datasets CIFAR-10. Man et al. [109] manipulated the imaging of printed figures using the "flare effect" and the "blooming effect" to achieve attacks. Gnanasambandam et al. [112] modeled the light from a projector and used a physical projector to project the designed pattern onto the target object, achieving robust attacks. Duan et al. [113] proposed AdvLB, a method that enables the manipulation of the physical parameters of laser beams to execute adversarial attacks, which can be accomplished using hand-held laser pointers in real-world scenarios. To enhance the stealthiness of light-based attacks, Huang et al. [115] devised an optimization algorithm to balance adversarial loss and stealthiness loss.

As a key component of physical adversarial attack workflow, cameras are responsible for the transformation from the physical scene to a digital image. In contrast to altering objects in the physical scene, Li et al. [56] applied mainly translucent stickers directly onto the camera lens to achieve the goal of misleading classifiers. These specially crafted stickers introduce perturbations into the captured images. Subsequently, Sayles et al. [124] also focused on the camera, leveraging the Rolling Shutter Effect inherent in cameras, combined with adversarially illuminating a scene, to introduce an attack signal into the captured images. Phan et al. [125] designed adversarial attacks that are effective only under specific camera ISP parameter settings. These methods provide an alternative approach to implementing physical adversarial attacks by modifying certain camera attributes or settings.

Apart from utilizing 2D images, there are also methods involving 3D objects. EOT [127] generates robust adversarial examples that remain effective across an entire distribution of transformations simultaneously. It fabricates the first 3D adversarial object: a turtle classified as a rifle from multiple viewpoints. Zeng et al. [128] perturbed 3D physical properties under a renderer to fool 3D object classification and visual question-answering models.

The evaluation results of attack methods on general classification tasks are summarized in Table 5. From these results, we can conclude that the performance of these attack methods does not exhibit continuous improvement over the years. The absence of a comprehensive evaluation framework has led to the emergence of this phenomenon. For example, CPAttack [78] performs relatively well in terms of effectiveness and cost-effectiveness, but it neglects robustness and stealthiness.

(a) A Stop sign with black/white blocks is classified as Speed Limit 45.

(b) A Speed Limit 20 sign with an adversarial patch is classified as Slippery Road.

(c) A Speed Limit 25 sign with shadows is classified as Speed Limit 35.

Fig. 6: Display of the physical adversarial attack in traffic sign classification tasks. Adapted from RP$_2$ [66] (a), PS-GAN [68] (b), and Adv-Shadow [114] (c).

TABLE 5: Comparison of the *hiPPA* metric among attack methods for the general classification task. We highlight the minimum and maximum values using blue and red, respectively.

| Methods | Hexagonal Score | | | | | | hiPAA |
|---|---|---|---|---|---|---|---|
| | Eff. | Rob. | Ste. | Aes. | Pra. | Eco. | |
| AdvPatch [21] NIPS17 | 1.00 | 0.67 | 0.20 | 0.20 | 0.60 | 0.99 | 0.65 |
| Phy-FGSM [97] ICLR17 | 0.50 | 0.67 | 0.40 | 0.40 | 0.60 | 0.99 | 0.56 |
| PTAttack [108] AAAI18 | 0.78 | 0.50 | 0.60 | 0.60 | 0.80 | 0.95 | 0.69 |
| EOT [127] PMLR18 | 0.99 | 0.83 | 0.60 | 0.60 | 0.60 | 0.92 | 0.80 |
| 3DAttack [128] CVPR19 | 0.94 | 0.50 | 0.80 | 0.80 | 0.60 | 0.92 | 0.77 |
| Poster [109] S&P19 | 0.83 | 0.33 | 0.20 | 0.40 | 0.60 | 0.92 | 0.55 |
| ACS [56] PMLR19 | 0.49 | 0.33 | 0.80 | 0.80 | 0.80 | 0.99 | 0.63 |
| D2P [100] AAAI19 | 0.93 | 0.83 | 0.40 | 0.40 | 0.60 | 0.99 | 0.72 |
| AdvACO [72] ECCV20 | 0.44 | 0.83 | 0.80 | 0.80 | 0.60 | 0.99 | 0.70 |
| ABBA [101] NIPS20 | 0.85 | 0.50 | 0.60 | 0.60 | 0.60 | 0.99 | 0.69 |
| AdvCam [103] ICCV20 | 0.40 | 0.50 | 0.80 | 0.80 | 0.60 | 0.99 | 0.62 |
| MetaAttack [105] ICCV21 | 0.95 | 0.50 | 0.40 | 0.40 | 0.60 | 0.99 | 0.66 |
| AdvACO+ [76] TIP21 | 0.72 | 0.83 | 0.80 | 0.80 | 0.60 | 0.99 | 0.78 |
| InvisPerturb [124] CVPR21 | 0.94 | 0.67 | 0.80 | 0.40 | 0.40 | 0.00 | 0.66 |
| AdvISP [125] CVPR21 | 0.90 | 0.17 | 0.40 | 0.40 | 0.60 | 0.00 | 0.48 |
| OPAD [112] ICCV21 | 0.43 | 0.50 | 0.20 | 0.40 | 0.80 | 0.85 | 0.47 |
| AdvLB [113] CVPR21 | 0.88 | 0.67 | 0.80 | 0.40 | 0.80 | 0.92 | 0.77 |
| CPAttack [78] NIPS22 | 1.00 | 0.17 | 0.20 | 0.20 | 0.60 | 0.99 | 0.55 |
| TnTAttack [79] TIFS22 | 0.95 | 0.67 | 0.40 | 0.80 | 0.60 | 0.99 | 0.74 |
| DAPatch [83] ECCV22 | 0.44 | 0.67 | 0.20 | 0.20 | 0.60 | 0.99 | 0.49 |
| Viewfool [106] NIPS22 | 0.92 | 0.50 | 1.00 | 0.80 | 0.40 | 0.99 | 0.80 |
| SPAA [115] VR22 | 1.00 | 0.33 | 0.20 | 0.40 | 0.80 | 0.92 | 0.62 |
| Meta-GAN [107] TIFS23 | 0.95 | 0.67 | 0.40 | 0.40 | 0.60 | 0.99 | 0.70 |

### 5.1.2 Traffic Sign Classification

The classification of traffic signs plays a pivotal role in aiding autonomous driving systems to comprehend scenes and make informed decisions. Consequently, safety assessments in this domain have garnered considerable attention. Eykholt *et al.* [66] introduced Robust Physical Perturbations (RP$_2$), a method that misguides road sign classifiers by affixing black/white blocks onto signs. These blocks are visible but inconspicuous to human observers. Leveraging the generative capabilities of GAN models, Liu *et al.* [68] proposed a perceptual-sensitive generative adversarial network (PS-GAN), which generates adversarial patches that visually resemble the scrawls and patches typically found on signs in the real world, enhancing their stealthiness. Subsequently, building on potential real-world interferences, Zhong *et al.* [114] introduced Adv-Shadow to assess the impact of shadows cast on traffic sign detectors. Experimental results demonstrate that optimized shadows can effectively deceive detectors. Fig. 7 presents a comparison of RP$_2$, PS-GAN, and Adv-Shadow in terms of the six perspectives of the *hiPPA* metric.
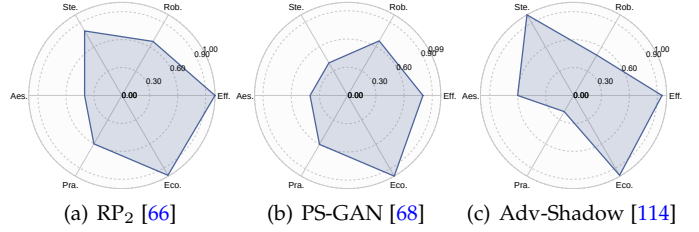


(a) RP$_2$ [66]

(b) PS-GAN [68]

(c) Adv-Shadow [114]

Fig. 7: Comparison of six perspectives of hiPAA across three physical adversarial attack methods on traffic sign classification task.

## 5.2 Attacks on Detection Tasks

Compared to classification tasks, detection tasks involve the prediction of additional information, including both the object's category and its spatial location. Therefore, physical adversarial attacks on detectors primarily focus on evading detection. As shown in Fig. 8, we display examples of attacking the vehicle, person, and sign detection tasks.

### 5.2.1 Vehicle Detection

DNNs-based models are extensively employed in autonomous driving systems for automatic vehicle detection. Adversarial attacks on vehicle detection are geared toward applying a distinctive pattern to a vehicle's exterior, thereby concealing it from detection.

To effectively camouflage a car in the real world, Zhang *et al.* [23] conducted experiments in a 3D space. This approach enabled them to simulate the intricate transformations induced by the physical environment comprehensively. Considering budget and time limitations, they utilized the photorealistic Unreal Engine 4 game engine[1] for their research. This engine provides a comprehensive set of configuration parameters, encompassing aspects such as camouflage resolution, patterns, 3D vehicle models, camera settings, environmental variables, and more. Similar to this work, Wu *et al.* [134] employed the open-source simulator CARLA [135]. They introduced an Enlarge-and-Repeat process and a discrete search method to craft physically adversarial textures. In addition, Duan *et al.* [65] utilized the Unity[2] and introduced a method called Coated Adversarial Camouflage (CAC).

The neural renderer is commonly used in 2D-to-3D transformation. One of the applications is to wrap the texture image to the 3D model, which then is rendered to the 2D image [136], [137], [138]. Thus utilizing the neural renderer to paint the adversarial stickers onto the vehicle surface is being pervasively used. Full-coverage Camouflage Attack (FCA) [61] tries rendering the non-planar texture over the full vehicle surface to overcome the partial occluded and long-distance issues. It bridges the gap between digital attacks and physical attacks via a differentiable neural renderer. Then, FCA introduces a transformation function to transfer the rendered camouflaged vehicle into a photo-realistic scenario. It outperforms other advanced attacks and achieves higher attack performance on both digital and physical attacks.

1. https://www.unrealengine.com/
2. https://unity.com/

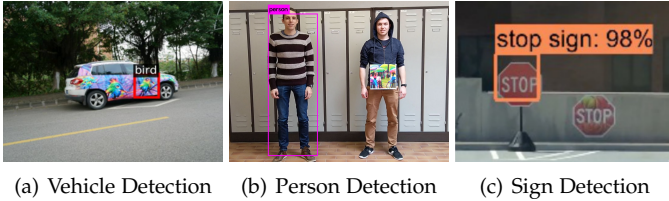(a) Vehicle Detection    (b) Person Detection    (c) Sign Detection

Fig. 8: Display of the physical adversarial attack in vehicle, person, and traffic sign detection tasks. The detectors fail to detect the perturbed target. Adapted from UPC [91] (a), AdvYOLO [13] (b), and ShapeShifter [98] (c).

However, existing neural renderers cannot fully represent various real-world transformations due to a lack of control of scene parameters compared to legacy photo-realistic renderers. Motivated by the challenge faced in prior works, Suryanto *et al.* [62] presented the Differentiable Transformation Attack (DTA), a framework for generating effective physical adversarial camouflage on 3D objects. It combines the advantages of a photo-realistic rendering engine with the differentiability of the novel rendering technique. In performance, DTA outperforms previous works in terms of effectiveness and transferability to other detection models.

The evaluation results of attack methods on vehicle detection tasks are summarized in Fig. 9. From these results, the following insights can be drawn. First, all attack methods do not perform ideally in terms of stealthiness. This is due to the fact that vehicles are large entities with a single color, and any external alterations are likely to be noticeable. Second, these methods do not excel in all six aspects; each has its own strengths. For instance, AdvLight [116] demonstrates strong robustness but fares poorly in terms of economics.

### 5.2.2 Person Detection

The objective of physical adversarial attacks on person detection is to conceal a person from detection models in the real world. Refer to Table 6 for the *hiPPA* evaluation, and Table 7 presents the robustness evaluation.

Yang *et al.* [55] were the first to propose attacking person detectors in the real world. The adversarial patches they generated caused the accuracy of the Tiny YOLO detector [139] to drop from 1.00 to 0.28. Thys *et al.* [13] designed a small (40cm × 40cm) adversarial patch that, when held by an attacker, can deceive the one-stage detector YOLOv2 [140]. Given an input image, the mainstream DNNs-based detectors have the ability to predict the position of the bounding boxes $\mathcal{V}_{pos}$, the object probability $\mathcal{V}_{obj}$ and the class score $\mathcal{V}_{cls}$. They minimized the $\mathcal{V}_{obj}$ and $\mathcal{V}_{cls}$ in the training phase to get the detector to ignore persons (target class). Meanwhile, the INRIAPerson dataset [141] provides a large number of person instances that support them to generate effective adversarial patches that can perform attacks in the real world successfully.

Not being satisfied with attacking detectors with printed cardboard, Xu *et al.* [90] crafted T-shirts with the generated adversarial patches. We called their method the AdvT-shirt. On the technical side, the AdvT-shirt develops a TPS-based transformer to model the temporal deformation of a T-shirt caused by pose changes of a moving person. Such non-rigid transformation ensures the attack effectiveness of the
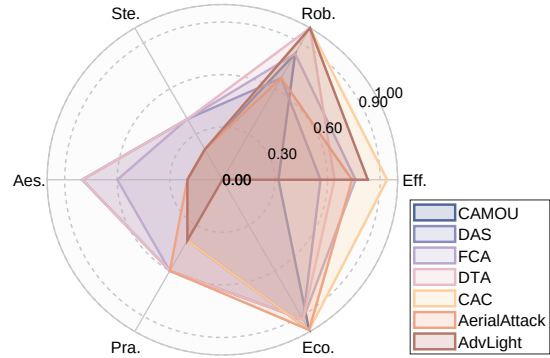


Fig. 9: Comparison of physical adversarial attack methods on vehicle detection tasks. These methods encompass CAMOU [23], DAS [59], FCA [61], DTA [62], CAC [65], AerialAttack [85], and AdvLight [116].

adversarial T-shirt in the physical world. Parallel to this work, Wu *et al.* [92] made a wearable invisibility cloak that, when placed over an object either digitally or physically, makes that object invisible to detectors. They quantify the success rate of attacks under various conditions and measure how algorithm and model choices impact success rates. Moreover, to push physical attacks to their limits with wearable adversarial clothing, they systematically quantify the success rate of attacks under complex fabric distortions and how well attacks on detectors transfer between models, classes, and datasets. To fairly evaluate the effectiveness of different physical attacks, Huang *et al.* [91] presented the first standardized dataset, AttackScenes, which simulates the real 3D world under controllable and reproducible settings to ensure that all experiments are conducted under fair comparisons for future research in this domain. In addition, they proposed the Universal Physical Camouflage (UPC) attack, which crafts adversarial patterns by jointly fooling the region proposal network, as well as misleading the classifier and the regressor to output errors.

Legitimate Adversarial Patches (LAP) [22] focuses on evading both human eyes and detection models in the physical world. To balance the attack effect and rationality of patches, LAP designed a two-stage training process. The first stage uses an original cartoon image as input and generates an initial patch. The second stage involves the input and the initial patch, which are the outputs from the first stage, to ultimately generate the adversarial patch in between. Generative adversarial networks (GANs) have the ability to efficiently generate desired samples [142]. Considering this, NAP [93] crafts adversarial patches for person detectors by leveraging the learned image manifold of BigGAN [143] and StyleGAN [144] pretrained on real-world images. Moreover, in this work, the MPII Human Pose dataset [145] is introduced to provide the diversity of training data. T-SEA [87] achieves high attack transferability through a series of strategies that involve self-ensembling the input data, the victim model, and the adversarial patch.

To enable multi-angle attacks, Hu *et al.* [94] proposed another generative method, named Toroidal-Cropping-based Expandable Generative Attack (TC-EGA). This method is designed to create adversarial textures with repetitive structures. TC-EGA aims to tackle the segment-missing problem. Unlike the prior studies, in the first stage, TC-EGA trains

TABLE 6: Comparison of the *hiPPA* metric among attack methods for both the person detection task and traffic sign detection task. We highlight the minimum and maximum values using blue and red, respectively.

| Methods | Hexagonal Score | | | | | | hiPAA |
|---|---|---|---|---|---|---|---|
| | Eff. | Rob. | Ste. | Aes. | Pra. | Eco. | |
| InvisibleCloak [55] UEMCON18 | 0.72 | 0.67 | 0.20 | 0.20 | 0.60 | 0.29 | 0.50 |
| AdvYOLO [13] CVPRW19 | 0.75 | 0.50 | 0.20 | 0.20 | 0.60 | 0.99 | 0.54 |
| AdvT-shirt [90] ECCV20 | 0.43 | 0.67 | 0.60 | 0.60 | 0.80 | 0.95 | 0.62 |
| UPC [91] CVPR20 | 0.93 | 0.67 | 0.20 | 0.20 | 0.80 | 0.91 | 0.64 |
| AdvCloak [92] ECCV20 | 0.50 | 0.83 | 0.60 | 0.60 | 0.80 | 0.95 | 0.67 |
| NAP [93] ICCV21 | 0.66 | 0.50 | 0.80 | 0.80 | 0.80 | 0.95 | 0.71 |
| LAP [22] ACM MM21 | 0.52 | 0.67 | 0.80 | 0.80 | 0.80 | 0.95 | 0.71 |
| AdvBulbs [122] AAAI21 | 0.65 | 0.83 | 0.20 | 0.40 | 0.60 | 0.99 | 0.60 |
| TC-EGA [94] CVPR22 | 0.65 | 1.00 | 0.40 | 0.20 | 0.80 | 0.95 | 0.67 |
| InvisClothing [95] CVPR22 | 0.88 | 1.00 | 0.20 | 0.20 | 0.20 | 0.00 | 0.54 |
| AdvInfrared [86] CVPR23 | 0.85 | 0.50 | 1.00 | 0.20 | 0.60 | 0.95 | 0.73 |
| T-SEA [87] CVPR23 | 0.99 | 0.67 | 0.20 | 0.20 | 0.60 | 0.29 | 0.58 |
| AdvCaT [96] CVPR23 | 0.85 | 1.00 | 0.80 | 0.80 | 0.80 | 0.91 | 0.87 |
| HOTCOLD Block [123] AAAI23 | 0.57 | 0.67 | 1.00 | 0.20 | 0.60 | 0.99 | 0.68 |
| CMPatch [88] ICCV23 | 0.62 | 0.67 | 1.00 | 0.20 | 0.80 | 0.95 | 0.72 |
| ShapeShifter [98] EP18 | 0.93 | 1.00 | 0.20 | 0.20 | 0.60 | 0.99 | 0.70 |
| RP$_2$+ [99] USENIX18 | 0.85 | 0.67 | 0.80 | 0.40 | 0.60 | 0.99 | 0.75 |
| NestedAE [67] CCS19 | 0.92 | 0.67 | 0.20 | 0.20 | 0.60 | 0.99 | 0.63 |
| LPAttack [104] AAAI20 | 0.80 | 1.00 | 0.20 | 0.40 | 0.40 | 0.99 | 0.66 |
| TransPatch [73] CVPR21 | 0.42 | 0.33 | 0.80 | 0.80 | 0.80 | 0.99 | 0.61 |
| AdvMarkings [60] USENIX21 | 1.00 | 0.83 | 1.00 | 0.60 | 0.80 | 0.99 | 0.91 |
| SLAP [111] USENIX21 | 0.99 | 1.00 | 0.20 | 0.20 | 0.60 | 0.92 | 0.71 |
| AITP [77] AISec22 | 0.90 | 0.83 | 0.20 | 0.40 | 0.60 | 0.99 | 0.68 |
| AdvLS [117] PMLR23 | 0.66 | 0.50 | 0.80 | 0.40 | 0.80 | 0.92 | 0.67 |

TABLE 7: Evaluation of the robustness dimension of the *hiPPA* metric on attacking person detection task.

| Cross -model | Cross-scenario | | Transformation | | | Method | Rob. |
|---|---|---|---|---|---|---|---|
| | Lig. | Bac. | Dis. | Ang. | Rot. | | |
| ✗ | ✓ | ✓ | ✓ | ✗ | ✓ | InvisibleCloak [55] | 0.67 |
| ✗ | ✓ | ✓ | ✗ | ✗ | ✓ | AdvYOLO [13] | 0.50 |
| ✗ | ✓ | ✓ | ✓ | ✗ | ✓ | AdvT-shirt [90] | 0.67 |
| ✓ | ✓ | ✓ | ✓ | ✗ | ✗ | UPC [91] | 0.67 |
| ✓ | ✓ | ✓ | ✓ | ✗ | ✓ | AdvCloak [92] | 0.83 |
| ✗ | ✗ | ✓ | ✗ | ✗ | ✓ | NAP [93] | 0.50 |
| ✗ | ✓ | ✓ | ✓ | ✗ | ✓ | LAP [22] | 0.67 |
| ✓ | ✓ | ✓ | ✓ | ✗ | ✓ | AdvBulbs [122] | 0.83 |
| ✓ | ✓ | ✓ | ✓ | ✓ | ✓ | TC-EGA [94] | 1.00 |
| ✓ | ✓ | ✓ | ✓ | ✓ | ✓ | InvisClothing [95] | 1.00 |
| ✗ | ✓ | ✓ | ✓ | ✗ | ✗ | AdvInfrared [86] | 0.50 |
| ✓ | ✓ | ✓ | ✗ | ✗ | ✓ | T-SEA [87] | 0.67 |
| ✓ | ✓ | ✓ | ✓ | ✓ | ✓ | AdvCaT [96] | 1.00 |
| ✓ | ✓ | ✓ | ✓ | ✗ | ✗ | HOTCOLD Block [123] | 0.67 |
| ✓ | ✓ | ✓ | ✓ | ✗ | ✗ | CMPatch [88] | 0.67 |

● *Lig.*, *Bac.*, *Dis.*, *Ang.*, and *Rot.* represent lighting, background, camera-to-object distances, view angles, and rotation respectively.

a fully convolutional network (FCN) [146], [147] as the generator to produce textures by sampling random latent variables as input. In the second stage, TC-EGA searches for the best local pattern of the latent variable with a cropping technique: Toroidal Cropping [148]. They have produced a variety of clothing items, such as T-shirts, skirts, and dresses, that have been printed with generated adversarial textures. These clothing items are effective for carrying out attacks when the wearer turns around or changes their posture. In the subsequent year, Hu *et al.* [96] improved the attack's stealthiness by refining the clothing texture based on the TC-EGA method. Specifically, they introduced camouflage textures with a high naturalness score, resulting in the production of natural-looking adversarial clothing.

The security of DNNs-based models has received substantial attention in the context of visible light, but its exploration in thermal infrared imaging remains incomplete. Additionally, thermal infrared detection systems hold significant relevance in various security-related domains, such as autonomous driving [149], night surveillance [150], and temperature measurement [151]. Thus, Zhu *et al.* [122] proposed a method to realize the adversarial infrared images in the real world. We named it the AdvBulbs. It is the first to realize physical attacks on the thermal infrared person detector. Their design belongs to the patch-based attack. Their experiments were conducted on the Teledyne FLIR ADAS Thermal dataset [152]. The following year, Zhu *et al.* [95] designed infrared invisible clothing based on a new material aerogel that successfully evades person detectors. Compared with the AdvBulbs, infrared invisible clothing hid from infrared detectors from multiple angles. Since then, this field has attracted attention. Wei *et al.* [86] also used aerogel material, with the difference that they designed irregularly shaped patches for the attack. Parallel to this work, Wei *et al.* [123] ingeniously employed anti-fever stickers and heating pads to create adversarial patches. Instead of

optimizing the texture and structural characteristics of the patches, their emphasis lies in studying the impact of patch size, shape, and placement on the attacks. Wei *et al.* [88] introduced the concept of cross-modal physical adversarial attacks, which remain effective under both thermal infrared and visible light imaging modalities.

### 5.2.3 Traffic Sign Detection

Traffic sign detection aims to recognize signage (e.g., stop signs and lane lines) on the roads in driving scenarios, which is widely used in autonomous driving. It is a security-critical domain because it assists the car in making driving decisions, and incorrect recognition could lead to driving violations or serious car accidents. In recent years, adversarial attacks against traffic sign detection have evolved to promote more robust detection algorithms. The bottom of Table 6 displays the comparative results of these methods.

The mainstream detectors prune the region proposals by using heuristics like non-maximum suppression (NMS) [153], [154]. The pruning operations are usually non-differentiable. However, generating adversarial perturbations pervasively requires calculating a backward gradient end to end. The non-differentiable operations make it hard to optimize the objective loss. To tackle this problem, Chen *et al.* [98] carefully studied the Faster R-CNN object detector [153] and successfully performed optimization-based attacks using gradient descent and backpropagation. Concretely, they ran the forward pass of the region proposal network and fixed the pruned region proposals as fixed constants to the second stage classification in each iteration. The RP$_2$ algorithm of [66] only focuses on attacking the traffic sign classification task. Following this line, Song *et al.* [99] extended the RP$_2$ to provide proof-of-concept attacks for object detection networks. They experimented with the YOLOv2 [140] and achieved 85.6% / 85.0% ASR in an indoor environment and 72.5% / 63.5% ASR in an outdoor environment.

Lane detection is important for autonomous driving because it supports steering decisions. Jing *et al.* [60] investigated the security of lane detection modules in real vehicles (Tesla Model S). We entitle their method as Adversarial Markings since they harnessed small markings (i.e., stickers) on the road surface to mislead the vehicles' visual

(a) Impersonation attack in face recognition task.

(b) Impersonation attack in person re-identification task.

Fig. 10: Display of the physical adversarial attack in re-identification (Re-ID) tasks. Adapted from AdvEyeglass [14] (a) and AdvPattern [70] (b).

TABLE 8: Comparison of the *hiPPA* metric among attack methods for both the face recognition task and person Re-ID task. We highlight the minimum and maximum values using blue and red, respectively.

| Methods | Hexagonal Score | | | | | | hiPAA |
|---|---|---|---|---|---|---|---|
| | Eff. | Rob. | Ste. | Aes. | Pra. | Eco. | |
| AdvEyeglass [14] CCS16 | 1.00 | 0.33 | 0.60 | 0.60 | 0.80 | 0.99 | 0.73 |
| AdvEyeglass+ [57] TOPS19 | 1.00 | 0.67 | 0.60 | 0.60 | 1.00 | 0.99 | 0.81 |
| Advhat [24] ICRP20 | 1.00 | 0.83 | 0.20 | 0.40 | 0.60 | 0.99 | 0.71 |
| ALPA [110] CVPR20 | 1.00 | 0.50 | 0.20 | 0.20 | 0.60 | 0.92 | 0.61 |
| CLBAAttack [58] BIOSIG21 | 0.95 | 0.33 | 0.20 | 0.20 | 0.60 | 0.99 | 0.57 |
| AdvMask [63] EP21 | 0.96 | 1.00 | 0.20 | 0.20 | 0.80 | 0.98 | 0.73 |
| AdvMakeup [126] IJCAI21 | 0.40 | 0.33 | 0.80 | 0.60 | 0.60 | 0.98 | 0.56 |
| TAP [75] CVPR21 | 1.00 | 0.17 | 0.40 | 0.60 | 0.60 | 0.99 | 0.63 |
| AdvSticker [64] TPAMI22 | 0.98 | 0.67 | 0.60 | 0.80 | 0.60 | 0.99 | 0.79 |
| SOPP [84] TPAMI22 | 0.96 | 0.83 | 0.40 | 0.60 | 0.60 | 0.99 | 0.75 |
| SLAttack [118] CVPR23 | 0.65 | 0.67 | 0.40 | 0.40 | 0.60 | 0.29 | 0.54 |
| AT3D [129] CVPR23 | 0.48 | 0.67 | 0.20 | 0.40 | 0.60 | 0.99 | 0.52 |
| AdvPattern [70] ICCV19 | 0.69 | 0.50 | 0.20 | 0.20 | 0.60 | 0.99 | 0.53 |

system. Extensive experiments show that Tesla Autopilot is vulnerable to Adversarial Markings in the physical world and follows the fake lane into oncoming traffic.

Zolfi *et al.* [73] designed a universal perturbation, called TransPatch, to fool the detector for all instances of a specific object class while maintaining the detection of other objects. TransPatch is a type of colored translucent sticker, which performs attacks by attaching this special sticker to the lens of the camera, resulting in disturbing the camera's imaging. In addition, Giulio *et al.* [111] proposed a light-based technique that allows attackers to realize physical attacks in the self-driving scenario, called Short-Lived Adversarial Perturbations (SLAP). Using a projector, SLAP shines a specific pattern on the Stop Sign causing YOLOv3 [32] and Mask-RCNN [30] detector to misdetect the targeted object. The experiment was carried out on a section of a private road in moving vehicle settings and SLAP obtained over 77% ASR.

### 5.3 Attacks on Re-Identification Tasks

In this section, we review physical adversarial attacks on re-identification (Re-ID) tasks. As shown in Fig. 10, we display examples of attacking the Re-ID tasks. Table 8 presents the comparative results of these methods based on the *hiPAA* metric.

#### 5.3.1 Face Recognition

Face Recognition Systems (FRS) are widely used in surveillance and access control [155], [156]. Therefore, it is valuable to explore the potential risks of FRS. Sharif *et al.* [14] developed a systematic method to attack the state-of-the-art face-recognition algorithm by printing a pair of eyeglass frames.

The person who wears the adversarial eyeglasses is able to evade being recognized or impersonate another individual. They demonstrate how an attacker that is unaware of the system's internals is able to achieve inconspicuous impersonation under a commercial FRS [157]. Meanwhile, In the dodging attack, an attacker can fool the most popular face-detection algorithm [158] into misidentifying as any other arbitrary face.

Pautov *et al.* [159] investigated the possibility of constructing a physical attack against ArcFace [160] by adversarial patches. They designed a cosine similarity loss that minimizes the similarity between the photo with patch and ground truth. The generated patch is gray, therefore it is readily printable. They tested the adversarial patch in three manners: one of the generated patches is an eyeglass, and two others are stickers on the nose and forehead. Numerical experiments showed that it is possible to efficiently attack ArcFace in the real world. Light-based attacks have been demonstrated to be feasible in classification tasks [112]. How about light-based attacks for face recognition systems? Nguyen *et al.* [110] designed a real-time adversarial light projection attack using an off-the-shelf camera-projector setup. And they choose to attack the state-of-the-art FRS, i.e., FaceNet [161] and SphereFace [162].

Advhat [24] implements an easily reproducible physical adversarial attack on the state-of-the-art public Face ID system [160], [163], [164]. In the digital space, Advhat uses Spatial Transformer Layer (STL) [165] to project the obtained sticker on the image of the face. In the physical space, Advhat launches attacks by wearing a hat with a special sticker on the forehead area, which significantly reduces the similarity to the ground truth class.

Yin *et al.* [126] analyzed the existing attacks against FRS and proposed the AdvMakeup, a unified adversarial face generation method. Adv-Makeup focuses on a common and practically implementable scenario: adding makeup to eye regions that shall mislead FRS models yet be visually unnoticeable (i.e., appearing as natural makeup). Concretely, AdvMakeup first introduces a makeup generation module, which can add natural eye shadow over the orbital region. Then, a task-driven fine-grained meta-learning adversarial attack strategy guarantees the attacking effectiveness of the generated makeup. Experimental results show that the Adv-Makeup' attack effectiveness is substantially higher than Advhat [24] and AdvEyeglass [14].

#### 5.3.2 Person Re-Identification

Person Re-ID is the task of identifying and tracking an individual of interest across multiple non-overlapping cam-

(a) Optical flow estimation      (b) Crowd counting      (c) Monocular depth estimation      (d) Semantic segmentation
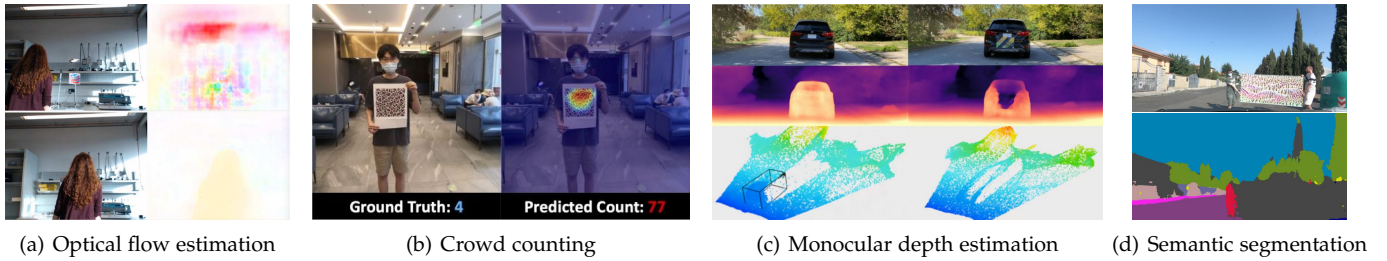
Fig. 11: Display of the physical adversarial attack in other tasks. Adapted from FlowAttack [71] (a), PAP [82] (b), OAP [80] (c), and RWAEs [81] (d).

eras [166]. This task plays an important role in surveillance and security applications. Wang *et al.* [70] were the first and only ones to propose a physical attack on the Re-ID model, known as AdvPattern. They accomplished evasion and impersonation attacks by formulating distinct optimization objectives. As shown in Fig. 10(b), AdvPattern employs adversarial patches featuring specially crafted patterns as the adversarial medium, which are affixed to a person's chest. The method degrades the rank-1 accuracy of person Re-ID models from 87.9% to 27.1% and under impersonation attack. This easily implementable approach exposes the vulnerability of the DNNs-based Re-ID system.

### 5.4 Attacks on Other Tasks

Besides the three aforementioned mainstream tasks, there are also physical adversarial attacks occurring in niche tasks. Below, we discuss attacks on seven tasks: optical flow estimation [71], steering angle prediction [102], crowd counting [82], semantic segmentation [81], object tracking [69], [74], monocular depth estimation [80], and image captioning [89]. Table 9 presents the comparative results of these methods based on the *hiPAA* metric.

**Optical Flow Estimation** (OFE) aims to measure the pixel 2D motion of an image sequence [167]. As shown in Fig. 11(a), Ranjan *et al.* [71] proposed FlowAttack to perturb the OFE models. FlowAttack utilizes the gradients from pretrained optical flow networks to update adversarial patches. Experimental results show that FlowAttack can cause large errors for encoder-decoder networks but not strongly affect spatial pyramid networks. This phenomenon demonstrates the correlation between network structure and vulnerability.

**Crowd Counting** aims to estimate the number of individuals within images or videos, with significant applications in public safety and traffic management [168]. Liu *et al.* [82] proposed a Perceptual Adversarial Patch (PAP) for attacking crowd-counting systems in the real world. PAP generates an adversarial patch by maximizing the model loss, leading the target victim model to overestimate the count by up to 100 on 80% of the samples (see Fig. 11(b)). To enhance robustness across multiple crowd-counting models with varying structures, PAP leverages the attention mechanism to capture scale and positional information.

**Monocular Depth Estimation** (MDE) aims to estimate the distance between the camera and a target object, which is crucial for autonomous driving [169]. Recently, Cheng *et al.* [80] developed an attack against MDE models. They generated a physical-object-oriented adversarial patch, called OAP, that can launch attacks in real-world driving scenarios (see Fig. 11(c)). OAP designs a rectangular patch

region optimization method to search for the optimal patch-pasting region. It achieves more than 6 meters mean depth estimation error and 93% ASR in downstream tasks. For Stealthiness, OAP alleviates the problem of unnaturalness from the patch's size and appearance. On the one hand, it minimizes the size while ensuring the attack. On the other hand, it designs the style transfer loss to incorporate the unobtrusive style. For Robustness, OAP applies EOT [127] and similar physical transformations (size, rotation, brightness, saturation, etc.) in the training stage.

**Semantic Segmentation** aims to classify each pixel into predefined categories without distinguishing between individual object instances [30]. Nesti *et al.* [81] crafted adversarial patches to perturb the semantic segmentation models. As shown in Fig. 11(d), they created a large adversarial patch, measuring 1m × 2m, which disrupts the predictions of segmentation models in the physical world. The adversarial patch is optimized using pixel-wise cross-entropy loss on the pre-trained ICNet [170]. Meanwhile, they built abundant and diverse scenes by the CARLA Simulator [135] for scene-specific attacks. Experimental results show that their attack method can reduce the baseline semantic segmentation model accuracy in the digital space, but in the real world, the attack is greatly downgraded.

**Steering Angle Prediction** assists autonomous driving systems in making informed decisions [171], [172], [173]. To ensure the safety and robustness of machine learning for autonomous driving, Kong *et al.* [102] introduced PhysGAN, a method that generates physically resilient adversarial examples to deceive autonomous steering systems. As shown in Fig. 12(a), by utilizing the discriminator within the GAN framework to assess the visual disparities between adversarial roadside signs and their original counterparts, Phys-GAN can generate realistic adversarial examples. Meanwhile, it can maintain attack effectiveness continuously across all frames throughout the entire trajectory.

**Object Tracking** aims to detect interesting moving objects and track them from frame to frame [174], which is an important task within the field of CV. Wiyatno *et al.* [69] proposed the first physical adversarial attack on this task. Specifically, they perform optimization to create a distinctive pattern, which is then presented on a large monitor as a background. When a person moves in front of the monitor, the tracker tends to prioritize locking onto the background and disregards the person. Subsequently, Ding *et al.* [74] proposed a patch-based attack method to launch universal physical attacks on single object tracking. As shown in Fig. 12(b), in the presence of the patch, the tracker neglects the originally tracked object. These explorations raise secu-

(a) Steering Angle Prediction

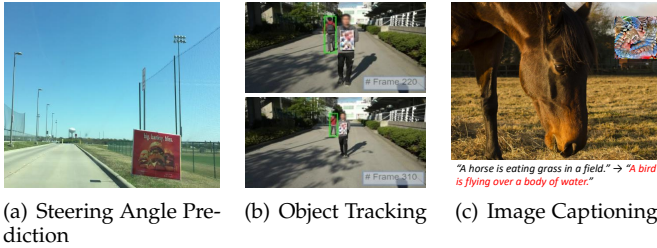(b) Object Tracking

(c) Image Captioning

Fig. 12: Display of the physical adversarial attack across three CV tasks: steering angle prediction, object tracking, and image captioning. Adapted from PhysGAN [102] (a), MTD [74] (b), and CAPatch [89] (c).

TABLE 9: Comparison of the hiPPA metric among attack methods for seven niche tasks.

| Methods | Hexagonal Score | | | | | | hiPAA |
|---|---|---|---|---|---|---|---|
| | Eff. | Rob. | Ste. | Aes. | Pra. | Eco. | |
| PAT [69] ICCV19 | 0.60 | 0.67 | 0.20 | 0.20 | 0.80 | 0.29 | 0.48 |
| MTD [74] AAAI21 | 0.74 | 0.67 | 0.20 | 0.40 | 0.60 | 0.99 | 0.60 |
| FlowAttack [71] ICCV19 | 1.00 | 0.67 | 0.20 | 0.20 | 0.60 | 0.99 | 0.65 |
| PhysGAN [102] CVPR20 | 1.00 | 0.67 | 0.60 | 0.80 | 0.60 | 0.95 | 0.79 |
| OAP [80] ECCV22 | 0.94 | 1.00 | 0.80 | 0.60 | 0.60 | 0.98 | 0.86 |
| RWAEs [81] WACV22 | 1.00 | 0.33 | 0.20 | 0.40 | 0.20 | 0.95 | 0.56 |
| PAP [82] CCS22 | 1.00 | 0.33 | 0.60 | 0.40 | 0.60 | 0.99 | 0.69 |
| CAPatch [89] USENIX23 | 1.00 | 1.00 | 0.20 | 0.20 | 0.60 | 0.99 | 0.67 |

rity concerns for real-world visual tracking.

**Image Captioning** focuses on generating a description of an image, which requires recognizing the important objects, their attributes, and their relationships in an image [175]. Inspired by adversarial patch attacks on CV models, Zhang *et al.* [89] designed CAPatch, an adversarial patch capable of inducing errors in final captions within real-world scenarios (see Fig. 12(c)). CAPatch had the capability to deceive image captioning systems, causing them to produce a specified caption or conceal certain keywords. In contrast to existing attack methods, this study represents the initial endeavor to employ an adversarial patch against multi-modal artificial intelligence systems.

# 6 DISCUSSION

During the development of this paper, it has been noticed physical adversarial attacks are diverse and threaten the security of many fields of human society. Despite the growth in published works over the past few years, there are still many potential risks to explore. In this section, we discuss the current challenges and opportunities in this field.

## 6.1 Current Challenges

### 6.1.1 Existing Domain Gaps

The workflow of physical adversarial attacks (see Fig. 1) reveals a process where attackers first design in the digital space, deploy in the physical space, and ultimately execute attacks in the digital domain. This workflow involves the transformation between the digital and physical domains. Current research has paid limited attention to addressing domain gaps. Jan *et al.* [100] designed a D2P network to model the transformation of images from the digital domain to the physical domain. Further exploration of methods to mitigate other domain gaps would be valuable.

### 6.1.2 Uncontrollable Evaluation Settings

Most existing works evaluate their physical adversarial attack methods in the real world using the adversarial mediums they manufacture. The real-world environment is dynamic, and the process of crafting adversarial mediums involves subjective factors, e.g., the material of the clothing, the quality of the printing, *etc*, all of which are uncontrollable. Future work with reliable and controllable evaluation setups is anticipated.

## 6.2 Future Work

### 6.2.1 New Adversarial Medium

In this survey, we define the adversarial medium as the object that carries the adversarial perturbations in the physical world. The adversarial medium plays a significant role in performing a physical attack. From the above discussion, we see that different tasks have different requirements for the adversarial medium, and the suitable adversarial medium can improve the performance of the attack in solving the trilemma, i.e., effectiveness, stealthiness, and stealthiness. The attacks using patches [21], light [110], camera ISP [125], makeup [126], 3D-printed object [127], clothing [90], *etc*, emerged in turn. Recently, the Laser Beam [113] and small lighting bulbs have been used to deceive the DNNs-based models, which inspire novel attack methods and expose the potential risks of these DNNs-based applications.

### 6.2.2 Transition from Digital Space to Physical Space

Adversarial perturbations are designed in the digital domain and carried by adversarial mediums in the physical domain. The transition from the digital domain to the physical domain introduces a gap. A typical example is the printing loss proposed by Sharif *et al.* [14], which specifically refers to the inability to accurately and reliably reproduce colors due to the smaller color space of printing devices compared to the RGB color space. They introduced the non-printability score (NPS) to address this issue. Image-to-image translation network can also be used to model this transformation [100]. However, it is still challenging when the fabrication deviates from the printing method, as is the case with light-based attack approaches [108], [109], [110]. The unknown gap in the transition needs to be explored and mitigated to the greatest extent possible.

### 6.2.3 Physical World Simulation

The fundamental characteristic of physical adversarial attacks lies in their feasibility in the real-world scenario. Precisely simulating the physical environment can bolster the attacks' robustness in dynamic physical environments. Simulation engines like Unreal Engine[3] and Unity[4] can provide a variety of environmental conditions for attack methods, including lighting, backgrounds, camera-to-object distances, view angles, *etc*. Most existing methods employ these simulators to evaluate the effectiveness of their attacks [23], [65]. However, due to non-differentiability concerns, they cannot be employed within an end-to-end optimization process for adversarial perturbations. In addition to basic operations

---

3. https://www.unrealengine.com/
4. https://unity.com/

like rotation, adding noise, affine transformations, and occlusions [13], [22], [94], more advanced physical scene simulation methods should be integrated into the attack pipeline, enabling the consideration of these dynamic settings during adversarial perturbation design.

### 6.2.4 Transition from Physical Space to Digital Space

Adversarial mediums carry elaborate adversarial perturbations that are captured by cameras in the real world, resulting in digital images, which are then used to attack the DNNs-based model. Throughout this process, there exists a gap between the transformations from the real world to digital images. For instance, different camera settings can result in variations in digital images. Phan *et al.* [125] have studied physical adversarial attacks under specific ISP conditions, but they did not explore the performance of attacks across different imaging devices. This domain transformation gap needs further consideration.

### 6.2.5 Physical Adversarial Attacks on New Tasks

As described in this survey, the current mainstream physical adversarial attack methods are oriented to tasks such as person detection [22], [91], traffic sign detection [99], [111], face recognition [14], [159], *etc*. Although many fields have been covered, there are some tasks that have not yet been explored. For example, [80] recently proposed an adversarial patch attack against monocular depth estimation (MDE), which is a critical vision task in real-world driving scenarios. It is the first time to propose an attack on MDE. Besides, we consider that domains with the following two characteristics can be explored for physical adversarial attacks: **1)** using the DNNs techniques, and **2)** applying in the physical world. Such as trajectory prediction [176], pose estimation [177], action recognition [178], *etc*.

## 7 CONCLUSION

Physical adversarial attacks have cast a shadow over the reliability of deep neural networks, raising security concerns. Consequently, extensive research has proposed various methods for real-world attacks across multiple tasks. We have provided an overview of the field of physical adversarial attacks on computer vision tasks, covering classification, detection, re-identification, and some niche tasks, with a focus on the adversarial mediums and a comprehensive evaluation. We first propose a general workflow for launching a physical adversarial attack, underlining the important role of the adversarial medium. Additionally, we have devised a new metric termed *hiPPA*, systematically quantifying and assessing attack methods from six distinct perspectives. Correspondingly, we present comparative results for existing methods, offering valuable insights for future improvements. Many challenges remain ahead, and we hope that this paper can motivate further discussion in this field and provides important guidance for future research, ultimately advancing the safety and reliability of machine vision systems.

## REFERENCES

[1] A. Voulodimos, N. Doulamis, A. Doulamis, and E. Protopapadakis, "Deep learning for computer vision: A brief review," *Computational intelligence and neuroscience*, vol. 2018, 2018. 1

[2] D. W. Otter, J. R. Medina, and J. K. Kalita, "A survey of the usages of deep learning for natural language processing," *IEEE transactions on neural networks and learning systems*, vol. 32, no. 2, pp. 604–624, 2020. 1

[3] A. B. Nassif, I. Shahin, I. Attili, M. Azzeh, and K. Shaalan, "Speech recognition using deep neural networks: A systematic review," *IEEE access*, vol. 7, pp. 19 143–19 165, 2019. 1

[4] C. Szegedy, W. Zaremba, I. Sutskever, J. Bruna, D. Erhan, I. J. Goodfellow, and R. Fergus, "Intriguing properties of neural networks," in *2nd International Conference on Learning Representations, ICLR 2014, Banff, AB, Canada, April 14-16, 2014, Conference Track Proceedings*, Y. Bengio and Y. LeCun, Eds., 2014. 1

[5] N. Carlini and D. Wagner, "Towards evaluating the robustness of neural networks," in *2017 ieee symposium on security and privacy (sp)*. Ieee, 2017, pp. 39–57. 1

[6] C. Xiao, B. Li, J.-Y. Zhu, W. He, M. Liu, and D. Song, "Generating adversarial examples with adversarial networks," *arXiv preprint arXiv:1801.02610*, 2018. 1

[7] N. Inkawhich, W. Wen, H. H. Li, and Y. Chen, "Feature space perturbations yield more transferable adversarial examples," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2019, pp. 7066–7074. 1

[8] X. Yuan, P. He, Q. Zhu, and X. Li, "Adversarial examples: Attacks and defenses for deep learning," *IEEE transactions on neural networks and learning systems*, vol. 30, no. 9, pp. 2805–2824, 2019. 1

[9] Z. Zhao, Z. Liu, and M. Larson, "Towards large yet imperceptible adversarial image perturbations with perceptual color distance," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2020, pp. 1039–1048. 1

[10] Y. Diao, T. Shao, Y.-L. Yang, K. Zhou, and H. Wang, "Basar: Blackbox attack on skeletal action recognition," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2021, pp. 7597–7607. 1

[11] Z. Cai, S. Rane, A. E. Brito, C. Song, S. V. Krishnamurthy, A. K. Roy-Chowdhury, and M. S. Asif, "Zero-query transfer attacks on context-aware object detectors," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2022, pp. 15 024–15 034. 1

[12] Z. Wei, J. Chen, M. Goldblum, Z. Wu, T. Goldstein, and Y.-G. Jiang, "Towards transferable adversarial attacks on vision transformers," in *Proceedings of the AAAI Conference on Artificial Intelligence*, vol. 36, no. 3, 2022, pp. 2668–2676. 1

[13] S. Thys, W. Van Ranst, and T. Goedemé, "Fooling automated surveillance cameras: adversarial patches to attack person detection," in *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition workshops*, 2019, pp. 0–0. 1, 5, 6, 10, 11, 15

[14] M. Sharif, S. Bhagavatula, L. Bauer, and M. K. Reiter, "Accessorize to a crime: Real and stealthy attacks on state-of-the-art face recognition," in *Proceedings of the 2016 acm sigsac conference on computer and communications security*, 2016, pp. 1528–1540. 1, 6, 12, 14, 15

[15] L. Sun, M. Tan, and Z. Zhou, "A survey of practical adversarial example attacks," *Cybersecurity*, vol. 1, pp. 1–9, 2018. 1, 2

[16] X. Wei, B. Pu, J. Lu, and B. Wu, "Physically adversarial attacks and defenses in computer vision: A survey," *arXiv preprint arXiv:2211.01671*, 2022. 1, 2

[17] D. Wang, W. Yao, T. Jiang, G. Tang, and X. Chen, "A survey on physical adversarial attack in computer vision," *arXiv preprint arXiv:2209.14262*, 2022. 1, 2

[18] K. Nguyen, T. Fernando, C. Fookes, and S. Sridharan, "Physical adversarial attacks for surveillance: A survey," *arXiv preprint arXiv:2305.01074*, 2023. 1, 2

[19] Q. Xu, G. Tao, S. Cheng, and X. Zhang, "Towards feature space adversarial attack by style perturbation," in *Proceedings of the AAAI Conference on Artificial Intelligence*, vol. 35, no. 12, 2021, pp. 10 523–10 531. 1

[20] Z. Cai, X. Xie, S. Li, M. Yin, C. Song, S. V. Krishnamurthy, A. K. Roy-Chowdhury, and M. S. Asif, "Context-aware transfer attacks for object detection," in *Proceedings of the AAAI Conference on Artificial Intelligence*, vol. 36, no. 1, 2022, pp. 149–157. 1

[21] T. B. Brown, D. Mané, A. Roy, M. Abadi, and J. Gilmer, "Adversarial patch," in *Proceedings of the Advances in Neural Information Processing Systems Workshop*, 2017. 1, 5, 6, 7, 8, 9, 14

[22] J. Tan, N. Ji, H. Xie, and X. Xiang, "Legitimate adversarial patches: Evading human eyes and detection models in the physical world," in *Proceedings of the ACM Multimedia*, 2021, pp. 5307–5315. 1, 5, 7, 10, 11, 15

[23] Y. Zhang, H. Foroosh, P. David, and B. Gong, "Camou: Learning physical vehicle camouflages to adversarially attack detectors in the wild," in *International Conference on Learning Representations*, 2018. 1, 6, 9, 10, 14

[24] S. Komkov and A. Petiushko, "Advhat: Real-world adversarial attack on arcface face id system," in *2020 25th International Conference on Pattern Recognition (ICPR)*. IEEE, 2021, pp. 819–826. 1, 6, 12

[25] O. Russakovsky, J. Deng, H. Su, J. Krause, S. Satheesh, S. Ma, Z. Huang, A. Karpathy, A. Khosla, M. Bernstein *et al.*, "Imagenet large scale visual recognition challenge," *International Journal of Computer Vision*, vol. 115, no. 3, pp. 211–252, 2015. 2

[26] A. Dosovitskiy, L. Beyer, A. Kolesnikov, D. Weissenborn, X. Zhai, T. Unterthiner, M. Dehghani, M. Minderer, G. Heigold, S. Gelly *et al.*, "An image is worth 16x16 words: Transformers for image recognition at scale," in *International Conference on Learning Representations*, 2021. 2

[27] Z. Liu, Y. Lin, Y. Cao, H. Hu, Y. Wei, Z. Zhang, S. Lin, and B. Guo, "Swin transformer: Hierarchical vision transformer using shifted windows," in *Proceedings of the IEEE International Conference on Computer Vision*, 2021, pp. 10 012–10 022. 2

[28] J. Yu, Z. Wang, V. Vasudevan, L. Yeung, M. Seyedhosseini, and Y. Wu, "Coca: Contrastive captioners are image-text foundation models," *Transactions on Machine Learning Research*, 2022. [Online]. Available: https://openreview.net/forum?id=Ee277P3AYC 2

[29] O. Ronneberger, P. Fischer, and T. Brox, "U-net: Convolutional networks for biomedical image segmentation," in *International Conference on Medical image computing and computer-assisted intervention*. Springer, 2015, pp. 234–241. 2

[30] K. He, G. Gkioxari, P. Dollár, and R. Girshick, "Mask r-cnn," in *Proceedings of the IEEE International Conference on Computer Vision*, 2017, pp. 2961–2969. 2, 12, 13

[31] V. Badrinarayanan, A. Kendall, and R. Cipolla, "Segnet: A deep convolutional encoder-decoder architecture for image segmentation," *IEEE Transactions on Pattern Analysis and Machine Intelligence*, vol. 39, no. 12, pp. 2481–2495, 2017. 2

[32] J. Redmon and A. Farhadi, "Yolov3: An incremental improvement," *arXiv preprint arXiv:1804.02767*, 2018. 2, 12

[33] Z. Tian, C. Shen, H. Chen, and T. He, "Fcos: Fully convolutional one-stage object detection," in *Proceedings of the IEEE International Conference on Computer Vision*, 2019, pp. 9627–9636. 2

[34] H. Zhang, F. Li, S. Liu, L. Zhang, H. Su, J. Zhu, L. M. Ni, and H.-Y. Shum, "Dino: Detr with improved denoising anchor boxes for end-to-end object detection," *arXiv preprint arXiv:2203.03605*, 2022. 2

[35] T. Karras, S. Laine, and T. Aila, "A style-based generator architecture for generative adversarial networks," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2019, pp. 4401–4410. 2

[36] J. Ho, A. Jain, and P. Abbeel, "Denoising diffusion probabilistic models," *Proceedings of the Advances in Neural Information Processing Systems*, vol. 33, pp. 6840–6851, 2020. 2

[37] P. Dhariwal and A. Nichol, "Diffusion models beat gans on image synthesis," *Proceedings of the Advances in Neural Information Processing Systems*, vol. 34, pp. 8780–8794, 2021. 2

[38] M. Ye, Z. Wang, X. Lan, and P. C. Yuen, "Visible thermal person re-identification via dual-constrained top-ranking." in *Proceedings of the International Joint Conferences on Artificial Intelligence*, vol. 1, 2018, p. 2. 2

[39] Z. Wang, M. Ye, F. Yang, X. Bai, and S. S. 0001, "Cascaded sr-gan for scale-adaptive low resolution person re-identification." in *Proceedings of the International Joint Conferences on Artificial Intelligence*, vol. 1, no. 2, 2018, p. 4. 2

[40] Q. Meng, S. Zhao, Z. Huang, and F. Zhou, "Magface: A universal representation for face recognition and quality assessment," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2021, pp. 14 225–14 234. 2

[41] T. Gu, B. Dolan-Gavitt, and S. Garg, "Badnets: Identifying vulnerabilities in the machine learning model supply chain," *arXiv preprint arXiv:1708.06733*, 2017. 3

[42] Y. Liu, S. Ma, Y. Aafer, W.-C. Lee, J. Zhai, W. Wang, and X. Zhang, "Trojaning attack on neural networks," 2017. 3

[43] E. Wenger, J. Passananti, A. N. Bhagoji, Y. Yao, H. Zheng, and B. Y. Zhao, "Backdoor attacks against deep learning systems in the physical world," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2021, pp. 6206–6215. 3

[44] X. Qi, T. Xie, R. Pan, J. Zhu, Y. Yang, and K. Bu, "Towards practical deployment-stage backdoor attack on deep neural networks," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2022, pp. 13 347–13 357. 3, 4

[45] M. Barreno, B. Nelson, R. Sears, A. D. Joseph, and J. D. Tygar, "Can machine learning be secure?" in *Proceedings of the 2006 ACM Symposium on Information, computer and communications security*, 2006, pp. 16–25. 3

[46] A. Shafahi, W. R. Huang, M. Najibi, O. Suciu, C. Studer, T. Dumitras, and T. Goldstein, "Poison frogs! targeted clean-label poisoning attacks on neural networks," *Proceedings of the Advances in Neural Information Processing Systems*, vol. 31, 2018. 3

[47] X. Zhang, X. Zhu, and L. Lessard, "Online data poisoning attacks," in *Learning for Dynamics and Control*. PMLR, 2020, pp. 201–210. 3

[48] A. Oprea, A. Singhal, and A. Vassilev, "Poisoning attacks against machine learning: Can machine learning be trustworthy?" *Computer*, vol. 55, no. 11, pp. 94–99, 2022. 3

[49] Y. Gao, B. G. Doan, Z. Zhang, S. Ma, J. Zhang, A. Fu, S. Nepal, and H. Kim, "Backdoor attacks and countermeasures on deep learning: A comprehensive review," *arXiv preprint arXiv:2007.10760*, 2020. 3

[50] E. Bagdasaryan and V. Shmatikov, "Blind backdoors in deep learning models," in *30th USENIX Security Symposium (USENIX Security 21)*, 2021, pp. 1505–1521. 3

[51] K. Doan, Y. Lao, W. Zhao, and P. Li, "Lira: Learnable, imperceptible and robust backdoor attacks," in *Proceedings of the IEEE/CVF international conference on computer vision*, 2021, pp. 11 966–11 976. 3

[52] J. Su, D. V. Vargas, and K. Sakurai, "One pixel attack for fooling deep neural networks," *IEEE Transactions on Evolutionary Computation*, vol. 23, no. 5, pp. 828–841, 2019. 4, 8

[53] H. Ma, Y. Li, Y. Gao, Z. Zhang, A. Abuadbba, A. Fu, S. F. Al-Sarawi, N. Surya, and D. Abbott, "Macab: Model-agnostic clean-annotation backdoor to object detection with natural trigger in real-world," 2022. 4

[54] M. P. Van Albada and A. Lagendijk, "Observation of weak localization of light in a random medium," *Physical review letters*, vol. 55, no. 24, p. 2692, 1985. 4

[55] D. Y. Yang, J. Xiong, X. Li, X. Yan, J. Raiti, Y. Wang, H. Wu, and Z. Zhong, "Building towards" invisible cloak": Robust physical adversarial attack on yolo object detector," in *2018 9th IEEE Annual Ubiquitous Computing, Electronics & Mobile Communication Conference (UEMCON)*. IEEE, 2018, pp. 368–374. 6, 10, 11

[56] J. Li, F. Schmidt, and Z. Kolter, "Adversarial camera stickers: A physical camera-based attack on deep learning systems," in *International Conference on Machine Learning*. PMLR, 2019, pp. 3896–3904. 6, 8, 9

[57] M. Sharif, S. Bhagavatula, L. Bauer, and M. K. Reiter, "A general framework for adversarial examples with objectives," *ACM Transactions on Privacy and Security (TOPS)*, vol. 22, no. 3, pp. 1–30, 2019. 6, 12

[58] I. Singh, S. Momiyama, K. Kakizaki, and T. Araki, "On brightness agnostic adversarial examples against face recognition systems," in *2021 International Conference of the Biometrics Special Interest Group (BIOSIG)*. IEEE, 2021, pp. 1–5. 6, 12

[59] J. Wang, A. Liu, Z. Yin, S. Liu, S. Tang, and X. Liu, "Dual attention suppression attack: Generate adversarial camouflage in physical world," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2021, pp. 8565–8574. 5, 6, 10

[60] P. Jing, Q. Tang, Y. Du, L. Xue, X. Luo, T. Wang, S. Nie, and S. Wu, "Too good to be safe: Tricking lane detection in autonomous driving with crafted perturbations," in *30th USENIX Security Symposium (USENIX Security 21)*. USENIX Association, Aug. 2021, pp. 3237–3254. [Online]. Available: https://www.usenix.org/conference/usenixsecurity21/presentation/jing 6, 11

[61] D. Wang, T. Jiang, J. Sun, W. Zhou, Z. Gong, X. Zhang, W. Yao, and X. Chen, "Fca: Learning a 3d full-coverage vehicle camouflage for multi-view physical adversarial attack," in *Proceedings*

*of the AAAI Conference on Artificial Intelligence*, vol. 36, no. 2, 2022, pp. 2414–2422. 6, 9, 10

[62] N. Suryanto, Y. Kim, H. Kang, H. T. Larasati, Y. Yun, T.-T.-H. Le, H. Yang, S.-Y. Oh, and H. Kim, "Dta: Physical camouflage attacks using differentiable transformation network," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2022, pp. 15 305–15 314. 6, 10

[63] A. Zolfi, S. Avidan, Y. Elovici, and A. Shabtai, "Adversarial mask: Real-world adversarial attack against face recognition models," *arXiv preprint arXiv:2111.10759*, 2021. 6, 12

[64] X. Wei, Y. Guo, and J. Yu, "Adversarial sticker: A stealthy attack method in the physical world," *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 2022. 6, 12

[65] Y. Duan, J. Chen, X. Zhou, J. Zou, Z. He, J. Zhang, W. Zhang, and Z. Pan, "Learning coated adversarial camouflages for object detectors," in *Proceedings of the International Joint Conferences on Artificial Intelligence*, 2022, pp. 891–897. 6, 9, 10, 14

[66] K. Eykholt, I. Evtimov, E. Fernandes, B. Li, A. Rahmati, C. Xiao, A. Prakash, T. Kohno, and D. Song, "Robust physical-world attacks on deep learning visual classification," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2018, pp. 1625–1634. 6, 9, 11

[67] Y. Zhao, H. Zhu, R. Liang, Q. Shen, S. Zhang, and K. Chen, "Seeing isn't believing: Towards more robust adversarial attack against real world object detectors," in *Proceedings of the 2019 ACM SIGSAC Conference on Computer and Communications Security*, 2019, pp. 1989–2004. 6, 11

[68] A. Liu, X. Liu, J. Fan, Y. Ma, A. Zhang, H. Xie, and D. Tao, "Perceptual-sensitive gan for generating adversarial patches," in *Proceedings of the AAAI Conference on Artificial Intelligence*, vol. 33, no. 01, 2019, pp. 1028–1035. 5, 6, 9

[69] R. R. Wiyatno and A. Xu, "Physical adversarial textures that fool visual object tracking," in *Proceedings of the IEEE/CVF International Conference on Computer Vision*, 2019, pp. 4822–4831. 6, 13, 14

[70] Z. Wang, S. Zheng, M. Song, Q. Wang, A. Rahimpour, and H. Qi, "advpattern: physical-world attacks on deep person re-identification via adversarially transformable patterns," in *Proceedings of the IEEE International Conference on Computer Vision*, 2019, pp. 8341–8350. 6, 12, 13

[71] A. Ranjan, J. Janai, A. Geiger, and M. J. Black, "Attacking optical flow," in *Proceedings of the IEEE International Conference on Computer Vision*, 2019, pp. 2404–2413. 6, 13, 14

[72] A. Liu, J. Wang, X. Liu, B. Cao, C. Zhang, and H. Yu, "Bias-based universal adversarial patch attack for automatic checkout," in *Proceedings of the European Conference on Computer Vision*. Springer, 2020, pp. 395–410. 6, 7, 9

[73] A. Zolfi, M. Kravchik, Y. Elovici, and A. Shabtai, "The translucent patch: A physical and universal attack on object detectors," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2021, pp. 15 232–15 241. 6, 11, 12

[74] L. Ding, Y. Wang, K. Yuan, M. Jiang, P. Wang, H. Huang, and Z. J. Wang, "Towards universal physical attacks on single object tracking," in *Proceedings of the AAAI Conference on Artificial Intelligence*, vol. 35, no. 2, 2021, pp. 1236–1245. 6, 13, 14

[75] Z. Xiao, X. Gao, C. Fu, Y. Dong, W. Gao, X. Zhang, J. Zhou, and J. Zhu, "Improving transferability of adversarial patches on face recognition with generative models," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2021, pp. 11 845–11 854. 6, 12

[76] J. Wang, A. Liu, X. Bai, and X. Liu, "Universal adversarial patch attack for automatic checkout using perceptual and attentional bias," *IEEE Transactions on Image Processing*, vol. 31, pp. 598–611, 2021. 6, 7, 9

[77] P. A. Sava, J.-P. Schulze, P. Sperl, and K. Böttinger, "Assessing the impact of transformations on physical adversarial attacks," in *Proceedings of the 15th ACM Workshop on Artificial Intelligence and Security*, 2022, pp. 79–90. 6, 11

[78] S. Casper, M. Nadeau, D. Hadfield-Menell, and G. Kreiman, "Robust feature-level adversaries are interpretability tools," *Advances in Neural Information Processing Systems*, vol. 35, pp. 33 093–33 106, 2022. 6, 7, 8, 9

[79] B. G. Doan, M. Xue, S. Ma, E. Abbasnejad, and D. C. Ranasinghe, "Tnt attacks! universal naturalistic adversarial patches against deep neural network systems," *IEEE Transactions on Information Forensics and Security*, vol. 17, pp. 3816–3830, 2022. 6, 7, 9

[80] Z. Cheng, J. Liang, H. Choi, G. Tao, Z. Cao, D. Liu, and X. Zhang, "Physical attack on monocular depth estimation with optimal

[81] F. Nesti, G. Rossolini, S. Nair, A. Biondi, and G. Buttazzo, "Evaluating the robustness of semantic segmentation for autonomous driving against real-world adversarial patch attacks," in *Proceedings of the IEEE/CVF Winter Conference on Applications of Computer Vision*, 2022, pp. 2280–2289. 6, 13, 14

[82] S. Liu, J. Wang, A. Liu, Y. Li, Y. Gao, X. Liu, and D. Tao, "Harnessing perceptual adversarial patches for crowd counting," in *Proceedings of the 2022 ACM SIGSAC Conference on Computer and Communications Security*, ser. CCS '22. New York, NY, USA: Association for Computing Machinery, 2022, p. 2055–2069. [Online]. Available: https://doi.org/10.1145/3548606.3560566 6, 13, 14

[83] Z. Chen, B. Li, S. Wu, J. Xu, S. Ding, and W. Zhang, "Shape matters: deformable patch attack," in *Proceedings of the European Conference on Computer Vision*. Springer, 2022, pp. 529–548. 6, 7, 9

[84] X. Wei, Y. Guo, J. Yu, and B. Zhang, "Simultaneously optimizing perturbations and positions for black-box adversarial patch attacks," *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 2022. 6, 12

[85] A. Du, B. Chen, T.-J. Chin, Y. W. Law, M. Sasdelli, R. Rajasegaran, and D. Campbell, "Physical adversarial attacks on an aerial imagery object detector," in *Proceedings of the IEEE/CVF Winter Conference on Applications of Computer Vision*, 2022, pp. 1796–1806. 6, 10

[86] X. Wei, J. Yu, and Y. Huang, "Physically adversarial infrared patches with learnable shapes and locations," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, June 2023, pp. 12 334–12 342. 6, 11

[87] H. Huang, Z. Chen, H. Chen, Y. Wang, and K. Zhang, "T-sea: Transfer-based self-ensemble attack on object detection," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2023. 6, 10, 14

[88] X. Wei, Y. Huang, Y. Sun, and J. Yu, "Unified adversarial patch for visible-infrared cross-modal attacks in the physical world," 2023. 5, 6, 11

[89] S. Zhang, Y. Cheng, W. Zhu, X. Ji, and W. Xu, "{CAPatch}: Physical adversarial patch against image captioning systems," in *32nd USENIX Security Symposium (USENIX Security 23)*, 2023, pp. 679–696. 6, 13, 14

[90] K. Xu, G. Zhang, S. Liu, Q. Fan, M. Sun, H. Chen, P.-Y. Chen, Y. Wang, and X. Lin, "Adversarial t-shirt! evading person detectors in a physical world," in *Proceedings of the European Conference on Computer Vision*. Springer, 2020, pp. 665–681. 6, 7, 10, 11, 14

[91] L. Huang, C. Gao, Y. Zhou, C. Xie, A. L. Yuille, C. Zou, and N. Liu, "Universal physical camouflage attacks on object detectors," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2020, pp. 720–729. 5, 7, 10, 11, 15

[92] Z. Wu, S.-N. Lim, L. S. Davis, and T. Goldstein, "Making an invisibility cloak: Real world adversarial attacks on object detectors," in *Proceedings of the European Conference on Computer Vision*. Springer, 2020, pp. 1–17. 7, 10, 11

[93] Y.-C.-T. Hu, B.-H. Kung, D. S. Tan, J.-C. Chen, K.-L. Hua, and W.-H. Cheng, "Naturalistic physical adversarial patch for object detectors," in *Proceedings of the IEEE International Conference on Computer Vision*, 2021, pp. 7848–7857. 6, 7, 10, 11

[94] Z. Hu, S. Huang, X. Zhu, F. Sun, B. Zhang, and X. Hu, "Adversarial texture for fooling person detectors in the physical world," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2022, pp. 13 307–13 316. 5, 6, 7, 10, 11, 15

[95] X. Zhu, Z. Hu, S. Huang, J. Li, and X. Hu, "Infrared invisible clothing: Hiding from infrared detectors at multiple angles in real world," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2022, pp. 13 317–13 326. 7, 11

[96] Z. Hu, W. Chu, X. Zhu, H. Zhang, B. Zhang, and X. Hu, "Physically realizable natural-looking clothing textures evade person detectors via 3d modeling," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, June 2023, pp. 16 975–16 984. 7, 11

[97] A. Kurakin, I. J. Goodfellow, and S. Bengio, "Adversarial examples in the physical world," 2017. [Online]. Available: https://openreview.net/forum?id=HJGU3Rodl 7, 9

[98] S.-T. Chen, C. Cornelius, J. Martin, and D. H. P. Chau, "Shapeshifter: Robust physical adversarial attack on faster r-cnn object detector," in *Joint European Conference on Machine Learning*

*and Knowledge Discovery in Databases.* Springer, 2018, pp. 52–68. 7, 10, 11

[99] D. Song, K. Eykholt, I. Evtimov, E. Fernandes, B. Li, A. Rahmati, F. Tramèr, A. Prakash, and T. Kohno, "Physical adversarial examples for object detectors," in *12th USENIX Workshop on Offensive Technologies (WOOT 18).* Baltimore, MD: USENIX Association, Aug. 2018. [Online]. Available: https://www.usenix.org/conference/woot18/presentation/eykholt 7, 11, 15

[100] S. T. Jan, J. Messou, Y.-C. Lin, J.-B. Huang, and G. Wang, "Connecting the digital and physical world: Improving the robustness of adversarial attacks," in *Proceedings of the AAAI Conference on Artificial Intelligence*, vol. 33, no. 01, 2019, pp. 962–969. 7, 8, 9, 14

[101] Q. Guo, F. Juefei-Xu, X. Xie, L. Ma, J. Wang, B. Yu, W. Feng, and Y. Liu, "Watch out! motion is blurring the vision of your deep neural networks," *Advances in Neural Information Processing Systems*, vol. 33, pp. 975–985, 2020. 7, 8, 9

[102] Z. Kong, J. Guo, A. Li, and C. Liu, "Physgan: Generating physical-world-resilient adversarial examples for autonomous driving," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2020, pp. 14 254–14 263. 5, 7, 13, 14

[103] R. Duan, X. Ma, Y. Wang, J. Bailey, A. K. Qin, and Y. Yang, "Adversarial camouflage: Hiding physical-world attacks with natural styles," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2020, pp. 1000–1008. 7, 8, 9

[104] K. Yang, T. Tsai, H. Yu, T.-Y. Ho, and Y. Jin, "Beyond digital domain: Fooling deep learning based recognition system in physical world," in *Proceedings of the AAAI Conference on Artificial Intelligence*, vol. 34, no. 01, 2020, pp. 1088–1095. 7, 11

[105] W. Feng, B. Wu, T. Zhang, Y. Zhang, and Y. Zhang, "Meta-attack: Class-agnostic and model-agnostic physical adversarial attack," in *Proceedings of the IEEE International Conference on Computer Vision*, 2021, pp. 7787–7796. 7, 8, 9

[106] Y. Dong, S. Ruan, H. Su, C. Kang, X. Wei, and J. Zhu, "Viewfool: Evaluating the robustness of visual recognition to adversarial viewpoints," in *Advances in Neural Information Processing Systems*, A. H. Oh, A. Agarwal, D. Belgrave, and K. Cho, Eds., 2022. [Online]. Available: https://openreview.net/forum?id=X0m9q0IcsmX 7, 8, 9

[107] W. Feng, N. Xu, T. Zhang, B. Wu, and Y. Zhang, "Robust and generalized physical adversarial attacks via meta-gan," *IEEE Transactions on Information Forensics and Security*, 2023. 7, 8, 9

[108] N. Nichols and R. Jasper, "Projecting trouble: Light based adversarial attacks on deep learning classifiers," 2018. 7, 8, 9, 14

[109] Y. Man, M. Li, and R. Gerdes, "Poster: Perceived adversarial examples," in *IEEE Symposium on Security and Privacy*, no. 2019, 2019. 7, 8, 9, 14

[110] D.-L. Nguyen, S. S. Arora, Y. Wu, and H. Yang, "Adversarial light projection attacks on face recognition systems: A feasibility study," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition Workshops*, 2020, pp. 814–815. 7, 12, 14

[111] G. Lovisotto, H. Turner, I. Sluganovic, M. Strohmeier, and I. Martinovic, "SLAP: Improving physical adversarial examples with Short-Lived adversarial perturbations," in *30th USENIX Security Symposium (USENIX Security 21).* USENIX Association, Aug. 2021, pp. 1865–1882. [Online]. Available: https://www.usenix.org/conference/usenixsecurity21/presentation/lovisotto 7, 11, 12, 15

[112] A. Gnanasambandam, A. M. Sherman, and S. H. Chan, "Optical adversarial attack," in *Proceedings of the IEEE International Conference on Computer Vision*, 2021, pp. 92–101. 7, 8, 9, 12

[113] R. Duan, X. Mao, A. K. Qin, Y. Chen, S. Ye, Y. He, and Y. Yang, "Adversarial laser beam: Effective physical-world attack to dnns in a blink," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2021, pp. 16 062–16 071. 7, 8, 9, 14

[114] Y. Zhong, X. Liu, D. Zhai, J. Jiang, and X. Ji, "Shadows can be dangerous: Stealthy and effective physical-world adversarial attack by natural phenomenon," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2022, pp. 15 345–15 354. 7, 9

[115] B. Huang and H. Ling, "Spaa: Stealthy projector-based adversarial attacks on deep image classifiers," in *2022 IEEE Conference on Virtual Reality and 3D User Interfaces (VR).* IEEE, 2022, pp. 534–542. 7, 8, 9

[116] H. Wen, S. Chang, and L. Zhou, "Light projection-based physical-world vanishing attack against car detection," in *ICASSP 2023 -*

*2023 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, 2023, pp. 1–5. 7, 10

[117] C. Hu, Y. Wang, K. Tiliwalidi, and W. Li, "Adversarial laser spot: Robust and covert physical-world attack to dnns," in *Asian Conference on Machine Learning.* PMLR, 2023, pp. 483–498. 7, 11

[118] Y. Li, Y. Li, X. Dai, S. Guo, and B. Xiao, "Physical-world optical adversarial attacks on 3d face recognition," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, June 2023, pp. 24 699–24 708. 7, 12

[119] N. Hingun, C. Sitawarin, J. Li, and D. Wagner, "Reap: A large-scale realistic adversarial patch benchmark," *Proceedings of the IEEE International Conference on Computer Vision*, 2023. 5

[120] N. Wang, Y. Luo, T. Sato, K. Xu, and Q. A. Chen, "Does physical adversarial example really matter to autonomous driving? towards system-level effect of adversarial object evasion attack," *Proceedings of the IEEE International Conference on Computer Vision*, 2023. 5

[121] S. Li, S. Zhang, G. Chen, D. Wang, P. Feng, J. Wang, A. Liu, X. Yi, and X. Liu, "Towards benchmarking and assessing visual naturalness of physical world adversarial attacks," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2023, pp. 12 324–12 333. 5, 6

[122] X. Zhu, X. Li, J. Li, Z. Wang, and X. Hu, "Fooling thermal infrared pedestrian detectors in real world using small bulbs," in *Proceedings of the AAAI Conference on Artificial Intelligence*, vol. 35, no. 4, 2021, pp. 3616–3624. 7, 8, 11

[123] H. Wei, Z. Wang, X. Jia, Y. Zheng, H. Tang, S. Satoh, and Z. Wang, "Hotcold block: Fooling thermal infrared detectors with a novel wearable design," in *Proceedings of the AAAI Conference on Artificial Intelligence*, vol. 37, no. 12, 2023, pp. 15 233–15 241. 7, 8, 11

[124] A. Sayles, A. Hooda, M. Gupta, R. Chatterjee, and E. Fernandes, "Invisible perturbations: Physical adversarial examples exploiting the rolling shutter effect," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2021, pp. 14 666–14 675. 8, 9

[125] B. Phan, F. Mannan, and F. Heide, "Adversarial imaging pipelines," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2021, pp. 16 051–16 061. 8, 9, 14, 15

[126] B. Yin, W. Wang, T. Yao, J. Guo, Z. Kong, S. Ding, J. Li, and C. Liu, "Adv-makeup: A new imperceptible and transferable attack on face recognition," 2021. 8, 12, 14

[127] A. Athalye, L. Engstrom, A. Ilyas, and K. Kwok, "Synthesizing robust adversarial examples," in *International conference on machine learning.* PMLR, 2018, pp. 284–293. 8, 9, 13, 14

[128] X. Zeng, C. Liu, Y.-S. Wang, W. Qiu, L. Xie, Y.-W. Tai, C.-K. Tang, and A. L. Yuille, "Adversarial attacks beyond the image space," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2019, pp. 4302–4311. 8, 9

[129] X. Yang, C. Liu, L. Xu, Y. Wang, Y. Dong, N. Chen, H. Su, and J. Zhu, "Towards effective adversarial textured 3d meshes on physical face recognition," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, June 2023, pp. 4119–4128. 8, 12

[130] I. Goodfellow, J. Pouget-Abadie, M. Mirza, B. Xu, D. Warde-Farley, S. Ozair, A. Courville, and Y. Bengio, "Generative adversarial networks," *Communications of the ACM*, vol. 63, no. 11, pp. 139–144, 2020. 7

[131] J. Deng, W. Dong, R. Socher, L.-J. Li, K. Li, and L. Fei-Fei, "Imagenet: A large-scale hierarchical image database," in *2009 IEEE conference on computer vision and pattern recognition.* Ieee, 2009, pp. 248–255. 8

[132] C. Szegedy, V. Vanhoucke, S. Ioffe, J. Shlens, and Z. Wojna, "Rethinking the inception architecture for computer vision," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2016, pp. 2818–2826. 8

[133] J.-Y. Zhu, T. Park, P. Isola, and A. A. Efros, "Unpaired image-to-image translation using cycle-consistent adversarial networks," in *Proceedings of the IEEE international conference on computer vision*, 2017, pp. 2223–2232. 8

[134] T. Wu, X. Ning, W. Li, R. Huang, H. Yang, and Y. Wang, "Physical adversarial attack on vehicle detector in the carla simulator," *arXiv preprint arXiv:2007.16118*, 2020. 9

[135] A. Dosovitskiy, G. Ros, F. Codevilla, A. Lopez, and V. Koltun, "Carla: An open urban driving simulator," in *Conference on robot learning.* PMLR, 2017, pp. 1–16. 9, 13

[136] H. Kato, Y. Ushiku, and T. Harada, "Neural 3d mesh renderer," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2018, pp. 3907–3916. 9

[137] J. Thies, M. Zollhöfer, and M. Nießner, "Deferred neural rendering: Image synthesis using neural textures," *ACM Transactions on Graphics*, vol. 38, no. 4, pp. 1–12, 2019. 9

[138] K. Rematas and V. Ferrari, "Neural voxel renderer: Learning an accurate and controllable rendering tool," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2020, pp. 5417–5427. 9

[139] J. Redmon, S. Divvala, R. Girshick, and A. Farhadi, "You only look once: Unified, real-time object detection," in *Proceedings of the IEEE conference on computer vision and pattern recognition*, 2016, pp. 779–788. 10

[140] J. Redmon and A. Farhadi, "Yolo9000: better, faster, stronger," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2017, pp. 7263–7271. 10, 11

[141] N. Dalal and B. Triggs, "Histograms of oriented gradients for human detection," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, vol. 1. Ieee, 2005, pp. 886–893. 10

[142] Z. Wang, Q. She, and T. E. Ward, "Generative adversarial networks in computer vision: A survey and taxonomy," *ACM Computing Surveys (CSUR)*, vol. 54, no. 2, pp. 1–38, 2021. 10

[143] A. Brock, J. Donahue, and K. Simonyan, "Large scale gan training for high fidelity natural image synthesis," *arXiv preprint arXiv:1809.11096*, 2018. 10

[144] T. Karras, M. Aittala, S. Laine, E. Härkönen, J. Hellsten, J. Lehtinen, and T. Aila, "Alias-free generative adversarial networks," *Proceedings of the Advances in Neural Information Processing Systems*, vol. 34, pp. 852–863, 2021. 10

[145] M. Andriluka, L. Pishchulin, P. Gehler, and B. Schiele, "2d human pose estimation: New benchmark and state of the art analysis," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2014, pp. 3686–3693. 10

[146] J. Long, E. Shelhamer, and T. Darrell, "Fully convolutional networks for semantic segmentation," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2015, pp. 3431–3440. 11

[147] J. T. Springenberg, A. Dosovitskiy, T. Brox, and M. Riedmiller, "Striving for simplicity: The all convolutional net," *arXiv preprint arXiv:1412.6806*, 2014. 11

[148] A. Hatcher, *Algebraic topology*. Cambridge University Press, 2002. 11

[149] M. Martin, A. Roitberg, M. Haurilet, M. Horne, S. Reiß, M. Voit, and R. Stiefelhagen, "Drive&act: A multi-modal dataset for fine-grained driver behavior recognition in autonomous vehicles," in *Proceedings of the IEEE International Conference on Computer Vision*, 2019, pp. 2801–2810. 11

[150] M. Ye, C. Chen, J. Shen, and L. Shao, "Dynamic tri-level relation mining with attentive graph for visible infrared re-identification," *IEEE Transactions on Information Forensics and Security*, vol. 17, pp. 386–398, 2021. 11

[151] S. Adams, T. Bucknall, and A. Kouzani, "An initial study on the agreement of body temperatures measured by infrared cameras and oral thermometry," *Scientific reports*, vol. 11, no. 1, pp. 1–10, 2021. 11

[152] FLIR, "Teledyne flir free adas thermal datasets v2," [EB/OL], 2022, https://adas-dataset-v2.flirconservator.com/. 11

[153] S. Ren, K. He, R. Girshick, and J. Sun, "Faster r-cnn: Towards real-time object detection with region proposal networks," *Proceedings of the Advances in Neural Information Processing Systems*, vol. 28, 2015. 11

[154] G. Jocher, "Yolov5 detector," https://github.com/ultralytics/yolov5, 2020, accessed: 2023-08-15. 11

[155] MobileSec, "Mobilesec android authentication framework," https://github.com/mobilesec/authentication-framework-module-face, 2022. 12

[156] N. Technology, "Sentiveillance sdk," http://www.neurotechnology.com/sentiveillance.html, 2022. 12

[157] M. Inc, "Face++," http://www.faceplusplus.com/, 2022. 12

[158] P. Viola and M. Jones, "Rapid object detection using a boosted cascade of simple features," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, vol. 1. Ieee, 2001, pp. I–I. 12

[159] M. Pautov, G. Melnikov, E. Kaziakhmedov, K. Kireev, and A. Petiushko, "On adversarial patches: Real-world attack on arcface-100 face recognition system," in *2019 International Multi-Conference on Engineering, Computer and Information Sciences (SIBIRCON)*, 2019, pp. 0391–0396. 12, 15

[160] J. Deng, J. Guo, N. Xue, and S. Zafeiriou, "Arcface: Additive angular margin loss for deep face recognition," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, June 2019. 12

[161] F. Schroff, D. Kalenichenko, and J. Philbin, "Facenet: A unified embedding for face recognition and clustering," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, June 2015. 12

[162] W. Liu, Y. Wen, Z. Yu, M. Li, B. Raj, and L. Song, "Sphereface: Deep hypersphere embedding for face recognition," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, July 2017. 12

[163] I. Kemelmacher-Shlizerman, S. M. Seitz, D. Miller, and E. Brossard, "The megaface benchmark: 1 million faces for recognition at scale," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2016, pp. 4873–4882. 12

[164] P. Grother and M. Ngan, "Face recognition vendor test ( frvt ) performance of face identification algorithms," *NIST Interagency/Internal Report (NISTIR) - 8009*, 2014. 12

[165] M. Jaderberg, K. Simonyan, A. Zisserman, and k. kavukcuoglu, "Spatial transformer networks," in *Proceedings of the Advances in Neural Information Processing Systems*, C. Cortes, N. Lawrence, D. Lee, M. Sugiyama, and R. Garnett, Eds., vol. 28. Curran Associates, Inc., 2015. [Online]. Available: https://proceedings.neurips.cc/paper/2015/file/33ceb07bf4eeb3da587e268d663aba1a-Paper.pdf 12

[166] Z. Wang, Z. Wang, Y. Zheng, Y.-Y. Chuang, and S. Satoh, "Learning to reduce dual-level discrepancy for infrared-visible person re-identification," in *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, 2019, pp. 618–626. 13

[167] E. Ilg, N. Mayer, T. Saikia, M. Keuper, A. Dosovitskiy, and T. Brox, "Flownet 2.0: Evolution of optical flow estimation with deep networks," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2017, pp. 2462–2470. 13

[168] X. Jiang, L. Zhang, M. Xu, T. Zhang, P. Lv, B. Zhou, X. Yang, and Y. Pang, "Attention scaling for crowd counting," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2020, pp. 4706–4715. 13

[169] Y. Wang, W.-L. Chao, D. Garg, B. Hariharan, M. Campbell, and K. Q. Weinberger, "Pseudo-lidar from visual depth estimation: Bridging the gap in 3d object detection for autonomous driving," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, June 2019. 13

[170] H. Zhao, X. Qi, X. Shen, J. Shi, and J. Jia, "Icnet for real-time semantic segmentation on high-resolution images," in *Proceedings of the European Conference on Computer Vision*, 2018, pp. 405–420. 13

[171] C. Chen, A. Seff, A. Kornhauser, and J. Xiao, "Deepdriving: Learning affordance for direct perception in autonomous driving," in *Proceedings of the IEEE International Conference on Computer Vision*, 2015, pp. 2722–2730. 13

[172] A. Prakash, K. Chitta, and A. Geiger, "Multi-modal fusion transformer for end-to-end autonomous driving," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2021, pp. 7077–7087. 13

[173] K. Chitta, A. Prakash, and A. Geiger, "Neat: Neural attention fields for end-to-end autonomous driving," in *Proceedings of the IEEE International Conference on Computer Vision*, 2021, pp. 15 793–15 803. 13

[174] A. Yilmaz, O. Javed, and M. Shah, "Object tracking: A survey," *Acm computing surveys (CSUR)*, vol. 38, no. 4, pp. 13–es, 2006. 13

[175] M. Z. Hossain, F. Sohel, M. F. Shiratuddin, and H. Laga, "A comprehensive survey of deep learning for image captioning," *ACM Computing Surveys (CsUR)*, vol. 51, no. 6, pp. 1–36, 2019. 14

[176] J. Gu, C. Sun, and H. Zhao, "Densetnt: End-to-end trajectory prediction from dense goal sets," in *Proceedings of the IEEE International Conference on Computer Vision*, 2021, pp. 15 303–15 312. 15

[177] W. Liu and T. Mei, "Recent advances of monocular 2d and 3d human pose estimation: A deep learning perspective," *ACM Computing Surveys (CSUR)*, 2022. 15

[178] Z. Wang, Q. She, and A. Smolic, "Action-net: Multipath excitation for action recognition," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2021, pp. 13 214–13 223. 15