

Advances in adversarial attacks and defenses in computer vision: A survey

Naveed Akhtar, Ajmal Mian, Navid Kardan and Mubarak Shah

Abstract—Deep Learning (DL) is the most widely used tool in the contemporary field of computer vision. Its ability to accurately solve complex problems is employed in vision research to learn deep neural models for a variety of tasks, including security critical applications. However, it is now known that DL is vulnerable to adversarial attacks that can manipulate its predictions by introducing visually imperceptible perturbations in images and videos. Since the discovery of this phenomenon in 2013 [1], it has attracted significant attention of researchers from multiple sub-fields of machine intelligence. In [2], we reviewed the contributions made by the computer vision community in adversarial attacks on deep learning (and their defenses) until the advent of year 2018. Many of those contributions have inspired new directions in this area, which has matured significantly since witnessing the first generation methods. Hence, as a legacy sequel of [2], this literature review focuses on the advances in this area since 2018. To ensure authenticity, we mainly consider peer-reviewed contributions published in the prestigious sources of computer vision and machine learning research. Besides a comprehensive literature review, the article also provides concise definitions of technical terminologies for non-experts in this domain. Finally, this article discusses challenges and future outlook of this direction based on the literature reviewed herein and [2].

Index Terms—Deep learning, adversarial examples, adversarial machine learning, perturbation, black-box attack, white-box attack, adversarial defense.

I. INTRODUCTION

DEEP LEARNING (DL) [3] is a data driven technology that can precisely model complex mathematical functions over large data sets. It has recently provided scientists with numerous breakthroughs in machine intelligence applications. From analysing mutations in DNA [4] to reconstruction of brain circuits [5] and exploring cell data [6]; deep learning methods are currently advancing our knowledge for many cutting-edge scientific problems. Thus, it is not surprising that multiple contemporary sub-fields of machine intelligence are fast adopting this technology as ‘the tool’ to solve their long-standing problems. Along speech recognition [7] and natural language processing [8], computer vision is one of the sub-fields that currently relies heavily on deep learning.

The rise of deep learning in computer vision was triggered by the seminal work of Krizhevsky et al. [9] in 2012, reporting a record performance improvement on a hard image recognition task [10] using a Convolutional Neural Network

Naveed Akhtar and Ajmal Mian are with the Department of Computer Science and Software Engineering, University of Western Australia, 35 Stirling Highway, Crawley 6009, WA, Australia.

Mubarak Shah and Navid Kardan are with the Center for Research in Computer Vision, University of Central Florida, Orlando, FL 32816, United States.



Fig. 1. Attacking a deep visual model (GoogLeNet [16] here) by imperceptible image manipulation results in incorrect prediction with high confidence. FGSM attack [17] is used here to manipulate the image.

(CNN) [11]. Since [9], the computer vision community has contributed significantly to deep learning research, which has led to increasingly powerful neural networks [12], [13], [14] that can handle a large number of layers in their architectures - establishing the essence of ‘deep’ learning. The advances made in the context of computer vision have also enabled deep learning to solve complex problems of Artificial Intelligence (AI). For instance, one of the crowning achievements of the modern AI, i.e. tabula-rasa learning [15] owes a fair share to Residual Learning [12], which originated in the field of computer vision.

Owing to the (apparent) super-human abilities of deep learning [15], computer vision-based AI is believed to have reached the maturity required for deployment in safety and security critical systems. Auto-pilots of vehicles [18], facial recognition in ATMs [19] and Face ID technology of mobile devices [20] are a few fore-running real-world examples that portray the developing faith of modern societies in computer vision solutions. With highly active deep learning-based vision research for autonomous vehicles [21], face recognition [22], [23], robotics [24] and surveillance systems [25] etc., we can anticipate the *omnipresence* of deep learning in security critical computer vision applications. However, serious concerns are now emerging for this prospect due to an unsought discovery of adversarial vulnerability of deep learning [1].

It was discovered by Szegedy et al. [1] that deep neural network predictions can be manipulated with extremely low magnitude input perturbations. For images, these perturbations can be restricted to the imperceptible regime of human vision system, yet they can completely alter the output predictions of a deep visual model (see Fig. 1). Originally, these manipulative signals were discovered for the image classification task [1]. However, their existence is now well-established for a variety of mainstream computer vision problems, e.g. semantic segmentation [27], [28]; object detection [29], [30]; and object tracking [31], [32]. The literature highlights numerous

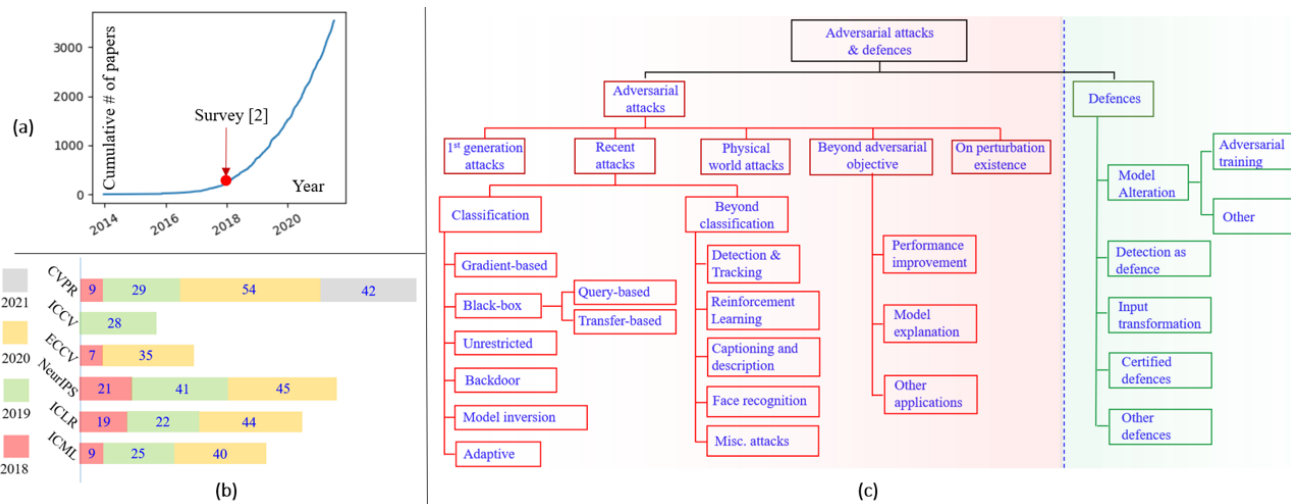


Fig. 2. (a) Cumulative number of adversarial attacks and defense papers appearing on arXiv in recent years (data from [26]). Over 3,000 papers have appeared since the first survey article [2]. (b) An increasing number of publications in this direction is experienced by the leading research sources of computer vision and machine learning. The bar chart indicates the total number of papers appearing per year which include ‘adversarial’, ‘attack’ or ‘defense’ keyword in their title, while the paper-content directly focuses on adversarial attack or defense problem. (c) Structuring of the literature reviewed in the article. The survey covers both aspects of attacks and defenses with emphasis on the attack methods.

characteristics of adversarial perturbations, that make them a real threat to deep learning as a pragmatic technology. For instance, it is repeatedly observed that the attacked models generally show high confidence on the wrong predictions of the manipulated images [2], [17]. It is also established that the same perturbation can often fool multiple models [33], [34]. The literature has also witnessed pre-computed perturbations, known as universal perturbations, that can be added to ‘any’ image to fool a given model with high probability [35], [36]. These facts have profound implications for security critical applications, especially when it is widely believed that deep learning solutions have predictive prowess that can surpass human abilities [15], [37].

Due to its critical nature, the topic of adversarial attacks (and their defenses) has received considerable attention of the research community in the last five years. In [2], we surveyed the contributions surfaced in this direction until the advent of 2018. Most of those works can be seen as the *first-generation* techniques that explore the core algorithms and techniques to fool deep learning or defend it against the adversarial attacks. Some of those algorithms have inspired streams of followup methods that further refine and adapt the core attack and defense techniques. These *second-generation* methods are also found to focus more on other vision tasks instead of just the classification problem, which is the main topic of interest in early contributions in this direction.

Since 2018, there has been an ever increasing number of publications in this research direction (see Fig. 2-a,b). Naturally, these publications also include instances of literature reviews, e.g. [38], [39], [40], [41], [42]. The literature survey we provide here differs from the existing reviews in many ways. This article is unique in that it is a legacy sequel of [2] - the first-ever peer-reviewed literature survey on this topic. Subsequent reviews, e.g. [41] are often found to be closely following [2]; or building on [2] for specific problems [42]. In recent years, this direction has matured significantly within

the field of computer vision. By building on the insights of [2] and subsequent literature, we are able to provide more precise definitions of the technical terminologies for this fast developing research direction. This also resulted in a more coherent structure of literature reviewed in this article, for which we provide concise discussions based on the current understanding of the terminologies by the research community. Moreover, we focus on peer-reviewed publications appearing in the prestigious research publication venues of computer vision and machine learning. Focusing on the leading contributions allows us to provide a more clear outlook of this direction for computer vision and machine learning researchers. Not to mention, this article reviews the most recent contributions of this fast evolving area to provide the most comprehensive review in this direction to date.

The rest of the article is organized as follows. In Section II, we provide definitions of technical terminologies used in the rest of the article. In Section III, we formulate the broader problem of adversarial attacks. The first generation of the attacks are discussed in Section IV, followed by the recent attacks focusing on the classification problem in Section V. We focus on recent attacks beyond classification problem in Section VI, and on the attacks tailored to the Physical world in Section VII. Contributions focusing more on the theoretical aspect of the existence of adversarial examples are discussed in Section IX. Recent defense methods are the topic of Section X. The article reflects on the literature trends in Section XI, where it also provides a discussion on the outlook of this research direction and future venues. Finally, we conclude in Section XII.

II. DEFINITION OF TERMS

To provide a clear discussion on the literature, it is imperative to first specify precise definitions of the technical terminologies commonly appearing in publications. Currently, the domain of adversarial machine learning is evolving rapidly.

Hence, understanding of the related technical terms is also evolving in the research community. Arranged alphabetically below, we provide definitions of the frequently encountered terminologies in the related literature, as widely understood by the computer vision (and machine learning) community. The same definitions of the concepts are followed in the rest of this article.

- *Adversarial example/image* is an image that is intentionally manipulated to cause incorrect model prediction. It is generally computed by adding *adversarial perturbation* to a natural image. *Clean, natural* or *benign* image are the commonly used terms to describe the opposite of an *adversarial* image.
- *Adversarial perturbation* is the component of an adversarial image that causes the incorrect prediction. Commonly, it is a low magnitude additive noise-like signal. However, exceptions are possible.
- *Adversarial training* is a process that injects adversarial examples in the training data of a model to make it adversarially robust.
- *Adversary* is the agent (i.e. the attacker) creating an adversarial example. Alternatively, the adversarial signal/perturbation is also referred to as the adversary, albeit much less often.
- *Attack detector* is an external mechanism for a model to (only) identify an input as adversarial or clean.
- *Black-box attack* assumes no knowledge of the target model. More strictly, the adversary is unaware of its training process and parameters. One category of black-box attacks allows *probing* the deployed target models with queries. This setup is more commonly known as *query-based* attack. To distinguish from the query-based attacks, other black-box attacks are sometimes also referred to as *Zero-knowledge attacks*. Opposite of black-box attack is *white-box* attack - see the definition below.
- *Data membership attack* aims to identify if a sample was used in the training of a model or not.
- *Defense/adversarial defense* is a broader term used for any mechanism of inducing inherent robustness in a model, or external/internal mechanisms to detect adversarial signals, or image processing to negate adversarial effects of input manipulations. *Adversarial robustness* is the preferred alternate term for the techniques focusing on inducing inherent resilience in the models, e.g. with adversarial training.
- *Digital attack* assumes that the adversary has full access to the actual digital input to the model. Most of the existing adversarial attacks are digital attacks. The opposite of digital attack is *Physical (world) attack* - see definition below.
- *Evasion attack* is a broader term for the adversarial attacks that fool pre-trained models into misclassifying input images at ‘test time’. Poisoning attack (see below) is its close antonym that poisons a model during ‘training’.
- *Fooling rate/ratio* is the commonly used evaluation metric, defined as the percentage of adversarial images on which the target model prediction is incorrect.
- *Gradient-based attacks* involve gradient computation of a model’s cost surface (or intermediate internal representation) with respect to the input. White-box attacks are predominately gradient-based.
- *Gradient-free attacks* do not involve gradient computation of any model.
- *Gray-box attack* assumes partial knowledge of the target model. However, since partial knowledge may actually lead to more knowledge, we prefer the term *restricted knowledge white-box* over the *gray-box* in this review. Under this nomenclature, gray-box attacks form a subcategory of the white-box attacks.
- *Image-specific attack* is computed to fool a target model on a specific image. Close antonyms for this term are *universal attack* and *label universal attack*.
- *Insertion attacks* insert an adversarial object (or a well-localised visible pattern) in an image, e.g. adversarial patch to alter the model prediction.
- *Label universal (adversarial) attack* aims at class-specific fooling. It computes an additive perturbation that has a pronounced effect on all samples of a selected class.
- *Model extraction attack* aims at recovering information about a target model (e.g. its classification boundaries) to subsequently use the information for fooling it.
- *Model inversion attack* aims at reconstructing individual training samples of the target model.
- *Norm-bounded perturbations* restrict the ℓ_p -norm of additive adversarial perturbations to control their perceptibility in adversarial examples. An overwhelming majority of the additive adversarial perturbations is norm-bounded.
- *One-step methods* compute perturbations in a single step, as opposed to *iterative methods* that use multiple iterations in their algorithm. These terms are generally more relevant to white-box attacks.
- *Physical (world) attacks* do not assume any access to the digital representation of the target model’s input. Adversarial examples are ‘clean’ images of e.g. physically modified or adversarially illuminated objects.
- *Poisoning attack* causes a model to misbehave when exposed to a trigger in the input. This (mis-)behavior is programmed into the model by manipulating the training process with tampered training data or algorithm. Generally, *trojan* or *backdoor* attack are used as synonyms for poisoning attack. This article largely focuses on the attacks (and their defenses) launched on clean pre-trained models. Hence, poisoning attack is not a direct topic for this survey.
- *Quasi-imperceptible perturbations* introduce slight visual impairment to images. This is in contrast to the imperceptible changes induced by imperceptible perturbations.
- *Query-based attack* is a form of black-box attack where the attacker is able to query the target model and exploit its output to optimize adversarial image(s). It either treats the target model as an *oracle* or learns a substitute model (see below) to be used as an oracle to subsequently generate adversarial images. A *decision/boundary-based* attack is a specific form of query-based attacks that assumes knowledge of only the predicted labels (not

confidence scores) of the target model. The query-based attacks that also exploit confidence scores of the target model are termed *score-based* attacks.

- *Real-world attacks* are evaluated in practical conditions by attacking real-world systems, as opposed to the bare models in laboratory setup. These attacks may still be digital or physical.
- *Targeted attack* forces the output of a model to pre-specified prediction of adversary’s choice, as opposed to random incorrect prediction in the case of *non-targeted attack*.
- *Target image* is the clean image being manipulated by the adversary.
- *Target model* is the model under attack.
- *Target label* is the (desired) incorrect label of the adversarial example. The term is more relevant for targeted attacks.
- *Threat model* refers to the assumed collective adversarial conditions against which a defense mechanism is designed and tested to verify its effectiveness.
- *Transferability* is the ability of an adversarial example/perturbation to generalise beyond the model for which it was originally computed.
- *Substitute model* is a model trained by an adversary to replicate the prediction behavior of the target model. *Surrogate model* and *auxiliary model* are the commonly used synonyms for the term substitute model.
- *Universal (adversarial) perturbations* are image-agnostic manipulative signals that can alter the model prediction on any input with high probability.
- *Unrestricted adversarial attacks* replace a natural image with a (synthetically) generated adversarial image, such that the latter has the same semantic meaning as the former for humans but not for the target model¹.
- *White-box attack* assumes complete knowledge of the target model. We refer to the attacks that assume partial knowledge of the target model or its training process, as *restricted knowledge white-box attacks*. Such attacks differ from the black-box attacks in that the latter only assume the knowledge of ‘prediction’ made by the model. The prediction may include a single/set of labels or a single/set of confidence scores. Any further, but incomplete knowledge makes the attack restricted knowledge white-box attack.

III. ADVERSARIAL ATTACKS: THE FORMAL PROBLEM

Let $\mathcal{M}(\cdot)$ be the target deep visual model such that $\mathcal{M}(\mathbf{I}) : \mathbf{I} \rightarrow \ell$, where $\mathbf{I} \in \mathbb{R}^m$ is a natural image and $\ell \in \mathbb{Z}^+$ is the output of the model. In the most common form of adversarial attacks, the adversary seeks a signal $\boldsymbol{\rho} \in \mathbb{R}^m$ to achieve $\mathcal{M}(\mathbf{I} + \boldsymbol{\rho}) \rightarrow \tilde{\ell}$, where $\tilde{\ell} \neq \ell$. To ensure that the manipulation to a clean image is humanly imperceptible, the perturbation $\boldsymbol{\rho}$ is often norm-bounded, e.g. by enforcing $\|\boldsymbol{\rho}\|_p < \eta$, where $\|\cdot\|_p$

denotes the ℓ_p -norm of a vector and ‘ η ’ is a pre-defined scalar. More concisely, the adversary seeks $\boldsymbol{\rho}$ that satisfies:

$$\mathcal{M}(\mathbf{I} + \boldsymbol{\rho}) \rightarrow \tilde{\ell} \text{ s.t. } \tilde{\ell} \neq \ell, \|\boldsymbol{\rho}\|_p < \eta. \quad (1)$$

The formulation above underpins the most prevailing contemporary understanding of the adversarial attacks. Yet, it does not encompass all attacks. For instance, unrestricted adversarial examples [43], [44], where the adversary is neither restricted to manipulate the original image (i.e. the image itself can be replaced) nor concerned with limiting the perturbation norm, can not be described by the constraint in (1). Similarly, the addition of a localized, but perceivable adversarial pattern in an image (e.g. adversarial patch [45]) is not accounted for by (1). Hence, for comprehensiveness, we also consider a more broader constraint, given as

$$\mathcal{M}(\tilde{\mathbf{I}}) \rightarrow \tilde{\ell} \text{ s.t. } \tilde{\ell} \neq \ell, \tilde{\mathbf{I}} \in \mathcal{S}_{\mathbf{I}}, \mathcal{M}(\mathbf{I} \sim \{\mathcal{S}_{\mathbf{I}} - \tilde{\mathbf{I}}\}) = \ell, \quad (2)$$

where $\mathcal{S}_{\mathbf{I}}$ is the set of images *perceived* as clean or allowed by humans to produce the desired output ℓ . For the sake of brevity, we are assuming a single adversarial sample in $\mathcal{S}_{\mathbf{I}}$ in (2). The conventional view of additive perturbations (in Eq. 1) becomes a special case of this constraint where $\tilde{\mathbf{I}} = \mathbf{I} + \boldsymbol{\rho}$ and $\tilde{\mathbf{I}} \in \mathcal{S}_{\mathbf{I}}$ is ensured by restricting the perturbation norm. Since (2) does not deal with $\boldsymbol{\rho}$ explicitly, one must articulate any additional constraint over $\boldsymbol{\rho}$ to specify an attack under (2) - as we have done above for the imperceptible perturbation.

Adversarial examples for deep visual models were originally discovered for the image classification task [1], where *additive* perturbations were used to launch the attack. Consequently, a vast majority of the existing attacks leverage some form of the additive perturbations to manipulate the model output. Moreover, image classifiers still remain the most popular target models for attacks. This trend partially owes to the fact that classification is one of the fundamental tasks in pattern recognition. Thus, it is important to explicitly, though briefly, discuss the broad concept of adversarial attacks on deep image classifiers under the above formulation.

For the image classifiers, an output is a class label $\ell \in \mathbb{Z}^+$. The nature of the task makes it more interesting to change this label to a pre-specified incorrect label $\tilde{\ell} \in \mathbb{Z}^+$ by the attack, which motivates the *targeted* adversarial attacks on classifiers. A *non-targeted* attack on a classifier can also be considered as a special case of the targeted attacks, where $\tilde{\ell}$ is chosen at random. Whereas *image-specific* attacks lead to misclassification of individual images, it is also possible to compute additive perturbations $\boldsymbol{\rho}$ that cause incorrect label predictions on a large number of images. Such *universal* perturbations were first reported by Moosavi-Dezfooli [35]. Here, we discuss the notions of image-specific vs universal, and targeted vs non-targeted attacks in the context of classifiers for a clear understanding of the text to follow immediately. Nevertheless, these concepts are more general and can also be applied to other computer vision tasks.

IV. FIRST-GENERATION ATTACKS

In the context of this survey, as a legacy sequel of [2], the first generation of adversarial attacks include the most

¹This understanding of the term is slightly different from [43] and relates more to [44] that allows a clearer delineation between the unrestricted and conventional adversarial examples.

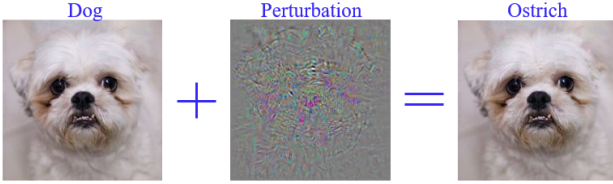


Fig. 3. Szegedy et al. [1] were the first to demonstrate imperceptible perturbations to images to fool deep learning. Here, the image of a ‘dog’ is confused as ‘ostrich’ by AlexNet [9] when the shown perturbation is added to it. The perturbation is exaggerated for visualisation.

influential contributions surfacing before 2018, which inspired series of followup methods. These attacks focus more on the fundamental algorithms to compute adversarial images, using image classification task as the test bed. We discuss these methods upfront as a separate section for two main reasons. First, by organizing the discussion on these methods in a (roughly) chronological order, we also provide the readers with a historical account of this research direction. Second, describing these seminal works early provides a more clear understanding of the inspiration of the more recent techniques.

The L-BFGS Attack: Szegedy et al. [1] first discovered the vulnerability of deep visual models to adversarial perturbations by solving for the following optimisation problem:

$$\min_{\rho} \|\rho\|_2 \quad \text{s.t. } \mathcal{M}(\mathbf{I} + \rho) = \tilde{\ell}; \quad \mathbf{I} + \rho \in [0, 1]^m. \quad (3)$$

The above is a hard problem, for which an approximate solution is computed by Szegedy et al. with the Limited Memory Broyden–Fletcher–Goldfarb–Shanno (L-BFGS) algorithm, which is a quasi-Newton algorithm involving computation of inverse Hessian [46] - inspiring the name of the attack adopted in the subsequent literature. To solve (3), the constraint $\min_{\rho} \|\rho\|_2$ is combined using a Lagrangian multiplier ‘ c ’ and the solution is computed by estimating the smallest $c > 0$ for which the minimizer ρ of the problem (4) satisfies $\mathcal{M}(\mathbf{I} + \rho) = \tilde{\ell}$.

$$\min_{\rho} c|\rho| + \mathcal{L}(\mathbf{I} + \rho, \tilde{\ell}) \quad \text{s.t. } \mathbf{I} + \rho \in [0, 1]^m, \quad (4)$$

where $\mathcal{L}(\cdot, \cdot)$ is the classifier loss. Manipulation of a clean image with the additive perturbation resulting from (4) remains imperceptible to the human visual system, see Fig. 3. This observation had a profound impact on the vision research community, which was fast developing the impression that deep visual features well approximate the perceptual differences in images with Euclidean distances. Discovery of adversarial perturbations that could completely alter the decisions of deep visual models with minuscule Euclidean norm revised this impression. Szegedy et al. also demonstrated that their adversarial attack transfers well between different deep visual classifiers. This intriguing vulnerability of deep learning to adversarial attacks attracted a wide interest of researchers in the subsequent years.

The FGSM Attack: It was originally observed by Szegedy et al. [1] that including adversarial images in the training data of a classifier improves its robustness to adversarial examples. Reinforced by multiple followup works, this observation is the

main motivation behind the idea of *adversarial training* in the literature. However, solving (4) for a large number of images is computationally prohibitive. This inspired the Fast Gradient Sign Method (FGSM) [17] to efficiently compute adversarial perturbations as:

$$\rho = \epsilon \text{sign}(\nabla_{\mathbf{I}} \mathcal{J}_{\theta}(\mathbf{I}, \ell)), \quad (5)$$

where $\mathcal{J}_{\theta}(\cdot, \cdot)$ is the cost for the model with parameters θ , $\nabla_{\mathbf{I}}$ computes its gradient *w.r.t.* \mathbf{I} , $\text{sign}(\cdot)$ denotes the sign function which is applied to each element in a vector, and ϵ is a prefixed scalar value to control perturbation perceptibility. The adversarial image is finally computed as $\tilde{\mathbf{I}} = \mathbf{I} + \rho$.

The FGSM is a one-step gradient-based method that computes norm-bounded perturbations, focusing on the ‘efficiency’ of perturbation computation rather than achieving high fooling rates. Goodfellow et al. [17] also used this attack to corroborate their linearity hypothesis, which considers the linear behavior of the modern neural networks in high dimension spaces (induced by e.g. ReLUs) as a sufficient reason for their vulnerability to adversarial perturbations. They also advocated this behavior as a major cause of transferability of the attacks between different modern networks, as their architectures pervasively allow such linearity for training efficiency. At the time, the linearity hypothesis was in sharp contrast to the developing idea that adversarial vulnerability was a result of high ‘non-linearity’ of the complex modern networks.

The FGSM [17] is among the most influential attacks in the existing literature, especially in the white-box setup. Its core concept of performing gradient ascend over the model’s loss surface to fool it, is the basis for a plethora of adversarial attacks. Many follow-up attacks can be strongly related to the original idea of FGSM. For instance, the Fast Gradient Value Method (FGVM) of Rozsa et al. [47] mainly removes the sign function from (5) to launch the attack. Similarly, ignoring the sign function, Miyato et al. [48] normalised the gradient with its ℓ_2 -norm to launch the attack. Kurakin et al. [49] also analysed the normalisation with ℓ_{∞} -norm. They also extended the FGSM to I-FGSM - its iterative variant, which is subsequently enhanced to incorporate momentum during the iterative optimisation by Dong et al. [50]. Their technique is known as Momentum Iterative (MI)-FGSM. Diverse Input I-FGSM, i.e. DI²-FGSM [51] is another example of the attacks that directly builds on FGSM. The main idea of [51] is to diversify the input used in each iteration of the iterative FGSM by applying image transformations, such as random resizing and padding, with a fixed probability. This diversification is claimed to facilitate better transferability of the resulting attack in a black-box setup. The authors also extend DI²-FGSM to M-DI²-FGSM by incorporating the momentum following [50].

In the above discussion, we consider FGSM as the first generation attack that inspired the followup works. It is emphasized that the discussed follow-up contributions do not form an exhaustive list of the methods that largely build on FGSM by far. Other such methods will keep appearing in the remaining article.

The BIM & ILCM Attacks: Though closely building on the FGSM [17] as the original concept of iterative FGSM,

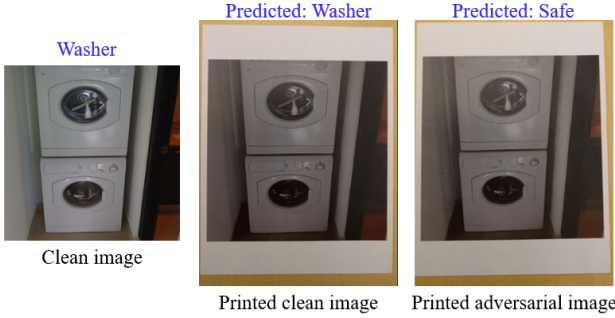


Fig. 4. Kurakin et al. [49] first demonstrated adversarial attack in the physical world by fooling a classifier on a printed adversarial image. Printed clean image of ‘Washer’ is predicted correctly, but printed adversarial image is predicted as ‘Safe’ by the TensorFlow Camera app used by [49].

the Basic Iterative Method (BIM) [49] is also an influential contribution that introduced the Physical World attacks. The attack, which is essentially the iterative FGSM algorithm, computes an adversarial image by repeating

$$\tilde{\mathbf{I}}_{i+1} = \text{Clip}_\epsilon \left\{ \tilde{\mathbf{I}}_i + \alpha \text{sign}(\nabla_{\mathbf{I}} \mathcal{J}_\theta(\tilde{\mathbf{I}}_i, \ell)) \right\}, \quad (6)$$

where ‘ i ’ indicates the i^{th} iteration, $\text{Clip}_\epsilon\{\cdot\}$ performs clipping at ϵ , and α is a pre-selected fixed scalar. Kurakin et al. [49] fooled the ImageNet inception model [52] on a mobile device by imaging *printed* adversarial images in the physical world, see Fig. 4. This idea also played its role in inspiring physical world attacks. The notion of targeted adversarial attacks can also be traced back to [49] and [53], where it is shown that the log-probability of prediction for a target class of adversarial image can be maximised by modifying (6) by changing addition to subtraction and replacing ℓ by $\tilde{\ell}$ as:

$$\tilde{\mathbf{I}}_{i+1} = \text{Clip}_\epsilon \left\{ \tilde{\mathbf{I}}_i - \alpha \text{sign}(\nabla_{\mathbf{I}} \mathcal{J}_\theta(\tilde{\mathbf{I}}_i, \tilde{\ell})) \right\}. \quad (7)$$

For a classifier with cross-entropy loss, solving (7) maximizes the confidence of the model on $\tilde{\ell}$ for the image $\tilde{\mathbf{I}}$. Originally, the authors proposed to use the label of the least-likely class of the clean image (as predicted by the model) as $\tilde{\ell}$ to compute interesting fooling outcomes. Hence, the technique is also referred to as Iterative Least-likely Class Method (ILCM).

The PGD Attack: The Projected Gradient Descent (PGD) attack is widely considered as one of the most powerful attacks in the literature, while referring to the seminal work of Madry et al. [54] as its origin. However, Madry et al. also refer to the iterative FGSM ([49], [53]) as a PGD method because Projected Gradient Descent is a standard optimization technique that projects gradients to a ball. Specifically, the authors see the iterative FGSM as the ℓ_∞ -bounded PGD, in which the ℓ_∞ -norm of the perturbation is bounded by the clipping operation - the projection. The main contribution of [54] comes in the form of looking at the adversarial robustness of deep models through the lens of robust optimisation, thereby defining adversarial training of deep models as a formal min-max optimisation problem below:

$$\min_{\theta} \rho(\theta), \quad \text{s.t.} \quad \rho(\theta) = \mathbb{E}_{(\mathbf{I}, \ell) \sim \mathcal{I}} \left[\max_{\rho} \mathcal{L}(\theta, \tilde{\mathbf{I}}, \ell) \right], \quad (8)$$

where $\mathbb{E}[\cdot]$ is the Expectation operator and \mathcal{I} is a distribution defined over the input images. This view allowed the authors to identify PGD as possibly the strongest first-order attack.

From the above view, we can also look at the variants of I-FGSM discussed in the previous sections as variants of PGD. In turn, PGD can be related to FGSM. However, a crucial finding by Madry et al. [54] makes PGD more appealing than FGSM for adversarial training. That is, the phenomenon of ‘label leaking’, observed in FGSM-based adversarial training [53], does not occur for PGD-based adversarial training. In plain words, label leaking occurs when adversarially trained model ends up with higher prediction accuracy for adversarial images, as compared to the clean images. FGSM results in a restricted set of adversarial examples, which can lead to overfitting in adversarial training, thereby causing label leaking. Considering that a major objective of FGSM is to compute samples for better adversarial training, avoiding label leaking is a significant advantage of PGD. Madry et al. [54] also showed that adversarial training with PGD - the strongest first-order attack - automatically makes the model robust against the weaker first-order attacks, e.g. FGSM. Nevertheless, being an iterative technique, PGD is computationally expensive.

JSMA & One-pixel Attack: Whereas most of the early attacks focused on perturbing a clean image holistically while enforcing perturbation imperceptibility by restricting the ℓ_2 or ℓ_∞ norms of the perturbations, the Jacobian-based Saliency Map Attack (JSMA) [55] and One-pixel attack [56] deviate from this practice by restricting the perturbations to smaller regions of the image. Contrary to the convention of computing the backward-gradient of the network for perturbation estimation, as done by e.g. FGSM and its variants, JSMA computes the forward-gradient of a network $\mathcal{M}(\cdot)$ as:

$$\nabla \mathcal{M}(\mathbf{I}) = \frac{\partial \mathcal{M}(\mathbf{I})}{\partial \mathbf{I}} = \left[\frac{\partial \mathcal{M}_j(\mathbf{I})}{\partial x_i} \right], \quad (9)$$

where $j \in 1, \dots, M$ for the M -dimensional function represented by $\mathcal{M}(\cdot)$, $i \in 1, \dots, N$ for the N -dimensional vectorized form of \mathbf{I} , whose i^{th} element is denoted as x_i . Essentially, (9) computes the Jacobian of the function learned by the network. Later, an adversarial extension of the saliency map [57] is used by [55] to modify only a few selected pixels that are most influential in terms of altering the model prediction.

Su et al. [56] demonstrated that a deep visual model can even be fooled by restricting the perturbation to a single pixel. However, this is generally more effective for the smaller image sizes, e.g. 64×64 or smaller. They used Differential Evolution (DE) [58] to estimate the location and RGB value of the pixel to be modified in the image to create an adversarial image, where the fitness criterion of the evolution is defined by accounting for the model prediction. Interestingly, the use of DE in contrast to model gradients, inherently makes their attack a query-based black-box attack. The authors also analysed the cases of a few pixel modifications, e.g. altering 5 instead of a single pixel for fooling. Although not originally emphasized as such, both JSMA and One-pixel attacks can be casted as optimisation problems with external constraints over the ℓ_0 -pseudo norm of the perturbations.

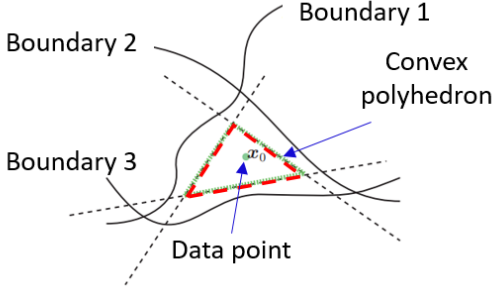


Fig. 5. The DeepFool algorithm [59] linearizes the decision boundaries around a data point to form a convex polyhedron to gradually push the point over the closest boundary for minimal perturbation.

The DeepFool Attack: Instead of restricting the perturbation norms to pre-fixed values, Moosavi-Dezfooli et al. [59] specifically aimed at minimising the norm of the adversarial perturbation by solving:

$$\Delta(\mathbf{I}; \ell) := \min_{\rho} \|\rho\|_2 \quad \text{s.t.} \quad \tilde{\ell} \neq \ell. \quad (10)$$

The main motivation behind computing the perturbations with minimal norm was to effectively quantify the adversarial robustness of the target models, where the robustness measure was defined as:

$$\rho_{\text{adv}} = \mathbb{E}_{\mathbf{I}} \frac{\Delta(\mathbf{I}; \ell)}{\|\mathbf{I}\|_2}, \quad (11)$$

where $\mathbb{E}_{\mathbf{I}}$ is the expectation over the data distribution.

DeepFool is the algorithm that computes ρ in (10) to compute the robustness defined by (11). The iterative algorithm linearizes the class boundaries around the current image to form a convex polyhedron and pushes the image towards the closest hyperplane to change the class label, see Fig. 5. The image gets updated in each iteration with the additive perturbation. Though originally proposed to quantify model robustness, DeepFool is now generally seen as an effective image-specific adversarial attack, while overlooking the quantification aspect.

The C&W Attack: The discovery of adversarial vulnerability of deep learning [1] also started a parallel research direction of defenses against adversarial attacks on deep learning in 2015-16. Defensive distillation [60] was a prominent technique that promised an effective solution to the problem, by building on the insights of knowledge distillation in deep networks [61]. However Carlini & Wagner [62] developed a set of attacks that computes norm-restricted additive perturbations that completely break defensive distillation. It is also shown that their attack is successful in fooling a defensively distilled network under black-box settings, where the perturbation is computed using an unsecured white-box model. Transferability of their attack in this setting significantly undermines the efficacy of defensive distillation.

To compute the adversarial perturbation, Carlini and Wagner solve the following optimisation problem:

$$\min \|\rho\|_p + c \cdot f(\mathbf{I} + \rho), \quad \text{s.t.} \quad \mathbf{I} + \rho \in [0, 1]^m, \quad (12)$$

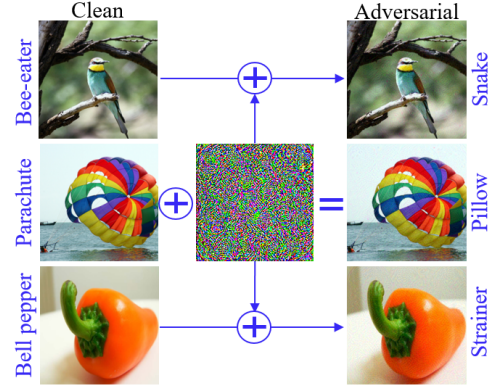


Fig. 6. A single Universal Adversarial Perturbation [35] can fool a model on multiple images. Fooling of GoogLeNet is shown here. These perturbations often transfer well across different models.

where $f(\cdot)$ is a function satisfying $\mathcal{M}(\tilde{\mathbf{I}}) \rightarrow \tilde{\ell}$, $\iff f(\mathbf{I} + \rho) \leq 0$. A range of analytical forms of $f(\cdot)$ are discussed by the authors to compute the desired perturbations. Carlini and Wagner [62] bounded the perturbations in their ℓ_2 , ℓ_∞ and ℓ_0 -pseudo norms, which gave rise to a set of attacks. The authors later showed that their attacks are also effective against other defense techniques [63]. The Carlini & Wagner (C&W) attack is generally considered a very strong attack, however, it does have a higher computational cost.

Universal Adversarial Perturbations: The above-mentioned methods compute adversarial perturbations that fool a target model on a specific image. Moosavi-Dezfooli et al. [35] focused on computing image-agnostic perturbations that could fool the model on *any* image with a high probability, see Fig. 6. Dubbed ‘universal’ for their transferability across different images (as opposed to models), these perturbations aim at satisfying the following constraint:

$$\mathbb{P}_{\mathbf{I} \sim \mathcal{I}} (\mathcal{M}(\mathbf{I}) \neq \mathcal{M}(\mathbf{I} + \rho)) \geq \delta \quad \text{s.t.} \quad \|\rho\|_p \leq \eta, \quad (13)$$

where $\mathbb{P}(\cdot)$ is the probability, \mathcal{I} denotes the distribution of clean images and $\delta \in (0, 1]$ is a predefined scalar, deciding the acceptable fooling ratio for the perturbations. The resulting universal adversarial perturbations are shown to be effective with both ℓ_2 and ℓ_∞ bounds over their respective norms. It can be observed from the experiments of [35] that perturbations bounded to around 4% of the respective image norms are able to achieve a significant fooling ratio (of $\sim 80\%$) for popular ImageNet models, e.g. ResNet [12], Inception [14]. However, a 4% distortion in an image is often slightly perceivable to the human visual system, hence the authors termed the perturbations to be quasi-imperceptible.

The universal adversarial perturbations are also able to transfer well across different models. In a sense, this property makes them ‘doubly universal’, as suggested by the authors [35]. However, since the estimated perturbations depend on parameters δ and η , (δ, η) -universal perturbations is a more qualified term for these signals. Moosavi-Dezfooli et al. compute these perturbations by building on the concept of Deepfool [59], where a single image is gradually pushed out of the decision boundary of its class. In the case of universal

perturbations, the iterative algorithm sequentially pushes all the data points out of their respective class regions, while accumulating the (label changing) perturbations by back-projection them onto an ℓ_p -ball of radius η . It is shown in the original paper that computing universal perturbations with as little as 2000 training images can still achieve $\sim 50\%$ fooling ratio for ImageNet models.

V. RECENT ATTACKS ON CLASSIFIERS

Mainly building on the core concepts of the first-generation attacks, there have been a multitude of more recent attacks on image classifiers. We cover those attacks in this section as per the structure illustrated in Fig. 2(c).

A. Advanced gradient based attacks

There is still a variety of contributions that are intended to improve the core strategy of gradient ascend for adversarial attacks. Naturally, these methods can be seen as downstream fine-tuning of first generation attacks like FGSM or PGD. For instance, Dong et al. [64] proposed to focus the gradient-based perturbation computed in an FGSM-like manner on the salient regions of images with the help of super-pixel guided attention. Such perturbations are claimed to be more robust against image processing based defenses. Similarly, Guo et al. [65] focused on improving the transferability of gradient-based attacks by backpropagating the computed gradients *linearly* through the model. Their gradient backpropagation mimics the scenario in which nonlinear activations are not encountered in the forward pass. Their modification is claimed to achieve better transferability of gradient-based attacks on large scale models.

Dong et al. [66] proposed a so-called GreedyFool algorithm that performs a sparse distortion in the input image based on gradients of its pixels. With improved sparsity, the perceptibility of their gradient-based perturbations becomes lower. Sriramanan et al. [67] proposed a Guided Adversarial Margin Attack (GAMA) that introduces a relaxation term in the standard losses (e.g. cross-entropy) of gradient-based attacks, e.g. PGD. It is claimed that this modification allows the attack to find better gradient directions, thereby increasing its efficacy. Similarly, Tohsiro et al. [68] devised a gradient-based strategy called Output Diversity Sampling (ODS) that is claimed to improve attacks in both white and black-box setups. Many adversarial attacks use random sampling of distributions, e.g. for initializing optimization process or updating query (in black-box setup). The ODS is mainly directed to provide a better sampling scheme for such attacks.

In [69], decoupling of the direction and norm of ℓ_2 -norm bounded gradient-based perturbations was proposed to make the attack more lethal. The resulting attack is commonly referred to as Decoupled Direction and Norm (DDN) attack. In [70], Yao et al. recommended to upgrade the first generation gradient-based attacks with Trust Regions [71]. During optimization, trust regions around the current point in the loss landscape finds descent/ascent directions that reduce errors due to the local nature of decisions. It is shown that multiple first-generation attacks can be improved for norm reduction and

computational efficiency using trust regions. In [72], Phan et al. argue to also consider the influence of image processing pipeline of cameras in attacks. They develop a gradient-based attack by differential approximation of this pipeline such that their perturbations are able to fool classifiers by images from one camera pipeline and not for another.

We emphasize that although we categorize only a few methods under advanced gradient attacks, nearly all white box (and transfer-based) attacks can be placed under this title, because those attacks inadvertently deal with model gradients rather directly. However, we introduce those attacks under subcategories more suited to their objective or threat models. Our intention to include a separate subsection for ‘advanced’ gradient attacks is to emphasize on the fact that improving the core gradient ascend scheme for attacks is still an active direction in this domain. The gradient based attacks, which are inherently white box, are generally the easiest to compute. Hence, they are the hardest to defend against. This makes them a useful tool to analyze model robustness.

B. Black-box attacks

From a pure adversarial perspective, black-box attacks form the most pragmatic category, because they assume no (or minimal) knowledge of the target model. Their practicality is making them highly popular in the recent literature. We review the recent black-box attacks along the directions of query-based and transfer-based attacks.

1) *Query-based attacks*: These attacks query the target model and use their outputs to construct adversarial images. Generally, their objective is to achieve minimal distortion in adversarial samples while maintaining model fooling. Queries are normally utilized for refining stronger perturbations for imperceptibility, see Fig. 7. Due to their practicality, decision/boundary-based attacks in this category are overwhelmingly popular as compared to their score-based counterpart.

Recently, Rahmati et al. [73] introduced a framework exploiting the decision boundary geometry to launch a black-box attack with a small number of queries to the target model that returns only the top-1 label. The attack exploits the smaller ‘mean’ curvature of the decision boundaries near the data point to estimate the normal vector, along which the data point can be efficiently nudged to the other side of the decision boundary by adding perturbations with small ℓ_p -norm for $p \geq 1$. The authors also show that the computed perturbation converges to the minimal norm for $p = 2$ for curvature-bounded decision boundaries. Better performance in terms of the number of queries and perturbation norm are reported as compared to the Boundary attack [74], HopSkipJump attack [75] and the qFool attack [76].

The Customized Adversarial Boundary (CAB) attack [77] reduces the number of queries by customising adversarial noise distribution with the queries in query-history, and initializing with perturbations already aimed for transferable attacks. Similarly, to improve query efficiency, a technique to extract generalizable prior using the earlier queries with meta learning is proposed in [78]. Another effort to improve query efficiency

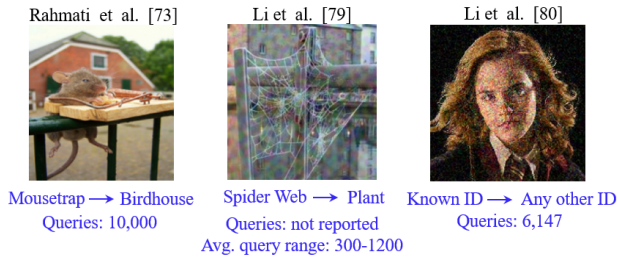


Fig. 7. Representative examples of query-based adversarial examples (attacks selected randomly): Generally, a larger number of queries is required for smaller perturbation perceptibility. For [79], we provide average query range, as queries for the shown image are not reported. Perturbed images are taken directly from [73], [79], [80].

includes Projection & Probability-Driven Black-box Attack (PPBA) [79] that restricts the solution space of the problem with low-frequency constrained sensing matrix - a concept inspired by compressive sensing theory. Li et al. [80] proposed a Query Efficient Boundary-based Black-box Attack (QEBA), that iteratively adds perturbation to a source image to retain its original label, but alters the image to form a perceptibly clean target image of a different object.

In [81], a Bayesian optimisation based attack is proposed. One method to reduce the number of queries it to search for adversarial images in a lower dimensional latent space as compared to the original image space. In that case, estimating the correct dimensionality of the latent space becomes a problem of its own. Ru et al. [81] employ non-parametric Bayesian strategy to resolve that by exploiting Gaussian Processes [82] based surrogate models to generate queries. Cheng et al. [83] also claimed a query efficient attack, altering the optimization objective of their previous work [84] that performed a binary search to estimate the gradient of the target model using query results. Later, they improved the attack by drastically reducing the number of queries by focusing on estimating gradient signs instead of gradients [83]. In another attempt to decrease the number of queries, Cheng et al. [85] also introduced a prior-guided random gradient-free method.

A TRansferable EMbedding based Black-box Attack (TREMBA) [86] trains an encoder-decoder model to learn a low-dimensional embedding space, where an adversarial example is searched for a given target model in a query-based setup. This process is claimed to reduce the number of queries significantly due to reduction in the search space of queries. Another method looking at the problem from search space perspective performs the attack as a progressive binary search using the gradient signs (instead of magnitude) [87]. The attack shows fooling of MNIST models with as little as 12 queries. Ilyas et al. [88] revisited the zeroth-order optimization (zoo) and proposed a query-based attack using bandit optimization that exploits prior information about the target model gradient. From zoo perspective, Zhao et al. [89] also proposed to augment the optimization with an ADMM-based framework.

Query-based black-box attacks are attracting significant interest of the research community in the recent literature. There are multiple other recent works that deal with these kinds of attacks, e.g., [90], [91], [92], [84], [93]. Mostly, the



Fig. 8. Typical examples of transfer-based perturbations (chosen randomly). Due to the harder objective of targeted transfer-based attacks, success rates are generally low (e.g. $< 50\%$) while perturbations are often perceptible. The reported average success rates across ImageNet models are taken from the original papers, which fall in the typical range of transfer-based attack success rates in the literature. Images are taken from [107] and [108].

current literature is dealing with decision-based attacks [94], [95], [96], [97], [98]. However, score-based attack schemes are also frequently encountered in the recent literature [99], [100].

2) *Transfer-based attacks*: Among the black-box attacks, transfer-based attacks are even more popular than the query-based attacks. This is because transfer based attacks do not require to query the black-box model and hence avoid suspicion altogether. The core idea behind transfer-based attacks is to compute perturbation on local surrogate models such that the perturbations will also effectively fool the remote target model. Popularity of these attacks also owes to the fact that the insights from white-box setup can often be readily leveraged for these attacks. The main objective of the methods appearing in this direction then is to amplify the intrinsic transferability of perturbations, for which different strategies are adopted.

Recently, Wu et al. [101] proposed to boost the transferability of perturbations by focusing them more on the salient image regions, where the regions are computed with Grad-CAM [102]. Improving perturbation transferability by manipulating the internal representation of the models is studied in [103]. Similarly, Huang et al. [104] fine-tuned adversarial examples using representations of pre-specified layers of the source model to improve attack transferability. A concept of ‘Adversarial Example Game’ is introduced in [105] that trains a generator for a transfer-based attack by training it against a discriminator for a hypothesis class of the target classifier. Since the underlying attack generation method does not assume details of the target (remote model), this setup is termed No-box attack in [106].

From the perspective of enhancing perturbation transferability, Lin et al. [107] exploit Nesterov gradient acceleration [109] with iterative FGSM for computing more generalizable, and hence transferable perturbations. The authors also introduced a scale-invariant attack method that induces an ensemble of models from an original model using data transformations that preserve the original loss of the model. Adversarial examples computed with these models are shown to exhibit better transferability. Lu et al. [110] demonstrated the possibility of fooling deep learning across different computer vision tasks. Analysing image classification, object detection, semantic segmentation and content detection as the tasks, they showed transferability of adversarial examples across them with rather modest perturbations. This is mainly achieved by reducing the

dispersion in the feature maps of the internal layers of the surrogate model with the help of a specialized loss. Inkawich et al. [111] also claimed that feature space perturbations are particularly helpful in computing adversarial examples that are more transferable across models.

There are also examples of targeted transferable attacks. For instance, Li et al. [108] proposed to make gradient-based targeted attacks more transferable by identifying two characteristics of white-box targeted attacks that restrict their transferability. First, reduction in gradient magnitude across iterations - leading to noise curing. Second, proximity of the adversarial examples to the true class region. The first issue is handled in [108] by allowing adaptive gradient magnitude in optimisation. Whereas the second is mitigated by metric-learning based regularization. Inkawich et al. [112] claimed state-of-the-art results for transferable targeted attacks on pristine ImageNet models. Instead of the classification layer, their method focuses on modeling layer-wise and class-wise feature distributions of a white box model and uses this information to alter the label of an adversarial image.

The direction of transferable attacks is also expanding in terms of the target tasks and underlying objective. For example, Wang et al. [113] demonstrated successful transferable attacks for the task of person re-identification. We also find an example of improving transferability of universal adversarial perturbations [114]. Moreover, training surrogate models in a data-free manner for transferable attacks was proposed in [115]. The core idea is to use a generator to construct synthetic images and label those with the target model (similar to query-based attacks), and train the substitute model with those images to better replicate the decision boundaries of the target model. Other recent examples focusing directly or indirectly on improving transferability of perturbations include [116], [117], [118], [119], [120].

C. Unrestricted adversarial attacks

Whereas the majority of mainstream attacks induce perturbation imperceptibility in adversarial images by restricting the ℓ_p -norm of perturbations, it is sometimes argued that the perturbation norm is not a good indicator of the perceptual difference between the two images [121]. The works related to achieving perturbation imperceptibility based on preserving semantics of the target image [122], [123], [124], [125] and preserving the structural information [126], [127] are motivated by this argument, see Fig. 9. In [128], unrestricted perturbations are introduced by manipulating the image color and texture to make them adversarial. It is claimed that such unrestricted perturbations are generally robust to defenses like feature squeezing, JPEG compression and adversarial training. On the other hand, compression and adversarial training are sometimes found effective against norm-bounded attacks like FGSM [17]. Shamsabadi et al. [129] demonstrated that it is possible to selectively manipulate image colors imperceptibly by operating on the decorrelated a , b channels of the Lab color space [130]. By changing image colors only to natural colors, and restricting manipulation to perceptually less sensitive regions in images, they computed transferable unrestricted adversarial examples that appear natural to humans.

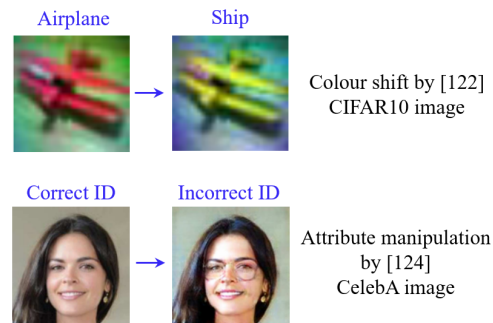


Fig. 9. Examples of unrestricted attacks. Images taken from [122], [124].

Zhao et al. [131] recently proposed to use the perceptual color distance CIEDE2000 [132] to control the imperceptibility of perturbations. The CIEDE2000 distance is known to align better with human perception. Zhao et al. demonstrated that accounting for the perceptual color distance while perturbing images can allow larger perturbations (having higher ℓ_p -norm) to remain imperceptible. The authors extended the C&W attack [62] to its variant that accounts for the perceptual color distance, Per-C&W. Their results show that higher confidence on incorrect labels and better transferability of attacks is possible by considering the perceptual color distance, without sacrificing perturbation imperceptibility. The proposed Per-C&W attack still computes a norm-bounded perturbation though, and the resulting image is not an unrestricted adversarial example. Another example of unrestricted perturbation attack is the semantic adversarial attack that manipulates image attributes with parametric conditional generative models [124], [133]. Incidentally, we can also categorise the emerging deep-fakes [134], [135] as unrestricted attacks. In a recent work, Hendrycks et al. [136] also reported two sets of natural images for which ImageNet models have extremely low accuracies ($< 5\%$). Named ImageNet-A (for adversarial) and ImageNet-O (for out-of-distribution), these images are termed natural adversarial examples by the author.

D. Backdoor attacks

Whereas adversarial attacks manipulate images during test time, backdoor attacks embed a backdoor or Trojan in the model during training. The targeted model normally shows high accuracy for clean input, however, its output is easily manipulated by embedding an attacker-defined trigger in the input. Although this article does not directly deal with backdoor or Trojan attacks, we still include recent papers in the surveyed venues due to the proximity of this research direction to adversarial attacks and for the sake of comprehensiveness of our survey. For a more detailed review of backdoor attacks appearing in other venues, we refer to [137], [138].

Generally, backdoors are embedded in the victim model by including trigger patterns in the training data so that the model learns a false association of a label with the trigger pattern. An issue with such triggers in training data is that the trigger patterns are often conspicuous, leading to easy detection of triggers with visual inspection. Liu et al. [139] recently proposed to use reflection patterns as triggers. Casting the triggers as natural looking shadows makes them harder

to detect. Often, triggers in the backdoor attacks are uniform across input images. However, Nguyen and Tran [140] proposed a generator-based backdoor attack that allows using different trigger patterns based on the context in the image. This makes detection of the trigger pattern even harder. Xie et al. [141] introduced a distributed backdoor attack on Federated Learning [142] in which the trigger is distributed among different parties providing the training data. This is in contrast to the centralized poisoning of data that appears in conventional supervised learning [143], [144].

A method claiming effective targeted poisoning was proposed in [145] for the practical setups where minimal assumptions can be made about the target network. That technique uses a pre-trained network to learn an attack model that can be directly used to generate images that would fool the victim model. In another example of backdoor attacks, Rakin et al. [146] generated a Trojan trigger to locate and flip the vulnerable bits of a DNN in DRAM to make it misbehave. It is noted in [147] that static backdoor attack on images do not work well for videos. Hence, a specialized backdoor attack for the task of video recognition was proposed. Similar to the concept of universal adversarial perturbations [35], their method uses a universal trigger to perform Trojan attack on video models.

The current literature is also witnessing multiple methods to secure deep learning models against the backdoor attacks. For instance, Kolouri et al. [148] introduced a ‘universal litmus patterns’ (ULP) for detecting a backdoor in pre-trained models. The detection is done by binary classification of the response of logit layers of the model in question for multiple geometric ULPs. The geometric ULPs are pre-defined, which are computed by an optimization problem inspired by universal adversarial perturbations [35]. Along the line of defense against Trojan attacks, Wang et al. [149] analyzed the possibility of detecting backdoors in the context of Federated Learning. They claimed that the detection is “unlikely” - assuming first-order oracles or polynomial time. Building on this theoretical insight they introduced a new family of backdoor attacks, called edge-case backdoors, which forces model fooling on the inputs living on the tail of the input distribution. By doing so, they make the detection of their attack very hard.

E. Model inversion

Model inversion aims to reconstruct training data or its markers from a trained model [150]. These attacks raise serious privacy concerns. Although model inversion problem is currently not as popular in the computer vision literature as additive perturbations, the attack is highly relevant for visual models in practical adversarial setups.

Since the discovery of the model inversion phenomenon [150], there have been multiple attempts to formalize the underlying problem for systematic investigation of the issue. For example, [151] uses the notion of influence from Boolean analysis to characterize inversion of Boolean functions. Similarly, [152] formulates the risk faced by the model to reveal individuals in training data. It is shown that the risk increases with over-fitting. To an extent, the model

inversion problem can be related to feature visualization [153] or the recently introduced attack to explain [154]. However, the overall adversarial objective of model inversion remains different from these frameworks which are more focused on model explanation.

Recently, Zhang et al. [155] proposed a generative model-inversion attack that trains a GAN to estimate the distributional prior of the target model’s training data. Combining the prediction loss of the target model with the loss of the discriminator, the trained generator is shown to produce high quality training samples of the target model, especially for the face models. Interestingly, the authors showed that highly predictive models establish stronger correlation between learned features and the sample labels. This is exactly what can be leveraged to do better in launching an inversion attack. An implication of this fact is that more accurate models might be more vulnerable to inversion attacks.

F. Adaptive attacks

It is now known that defenses against adversarial attacks can also be broken with counter-counter measures. For instance, in [63] and [156], we see multiple defenses broken with subsequent stronger attacks. This has prompted the research community to evaluate defenses against adaptive attacks [157]. Adaptive attacks are designed to specifically fool a defense mechanism. Although, research community is fast adapting the convention of evaluating defenses against adaptive attacks, Tramer et al. [157] showed that these evaluations are still far from providing guaranteed robustness against such attacks. The authors demonstrate this by circumventing thirteen recent defenses published in the proceedings of ICML, ICLR and NeurIPS. These defenses include [158], [159], [160], [161], [162], [163], [164], [165], [166], [167], [168], [169] and [170]. A key takeaway from [157] is that adaptive attacks should be hand-designed to specific defenses to be more effective, instead of automated attack adaption.

Although certifiable defenses are sometimes assumed robust to adaptive attacks, we also witness instances in the literature for adaptive attacks against certified defenses (see Section X-D). For example, Ghiasi et al [171] proposed a “Shadow Attack” to break certifiable defenses. Through generating the perturbation outside the certified ℓ_p bounds, their method produces a “spoofed” certificate, which results in visually imperceptible adversarial perturbations to break the defense. Whereas the underlying tools to develop adaptive attacks (i.e. counter-counter measures) are generally similar to conventional adversarial attacks, it is normally the objective of circumventing a specific (type of) defense that characterizes adaptive attacks. In recent years, these attacks are studied mainly in the context of developing robust defense techniques.

G. Miscellaneous attacks

The above sections reviewed literature related to the attacks on classifiers along popular directions. There are also other multiple interesting attacks related to the classification problem that do not fall under the above-mentioned subcategories. We provide a summary of those attacks in this section.

In [172], the authors devise an adversarial ranking attack, where the attacker can raise or reduce the rank of the potential label for the image. The unique objective of this attack differentiates it from the conventional fooling attacks. The literature has also seen attempts to fool deep neural networks by exploiting their storage on Dynamic Random Access Memory (DRAM) [173], [174], [175]. These attacks are particularly interesting for deep learning in practice. In another interesting work, Rezaei and Liu [176] demonstrated the possibility of adversarial manipulation of predictions for transfer learning, without the knowledge of the target domain. Similarly, Mor et al. [177] study optimal strategies against generative adversarial attacks. A Feature Disruptive Attack was proposed in [178] that is targeted at disrupting the internal representation of models for the adversarial samples, instead of simply focusing on altering the prediction.

The literature has also seen attacks to disrupt classifiers for point clouds. For instance, Zhou et al. [179] proposed a label guided GAN-based method for targeted attack on 3D point clouds in real-time. The proposed Label-Guided employs a multi-branch adversarial network for input feature extraction and then embeds the target label information in the features with an encoder. Vulnerability of deep 3D point cloud models to isometry transformations has also been exposed [180]. Other recent examples of attacks on point clouds and 3D data (and their defenses) include [181], [182], [183], [184], [185].

We also witness adversarial examples for video classifiers. Due to the additional time dimension, attacks on images can often not be readily translated to video. Hence, specialized attacks for videos are devised. Zhang et al. [186] used the movement patterns in the video frames to compute a noise prior that can help in gradient estimation for fooling video classifiers in the context of query-based attacks. A spatio-temporal attack is also introduced for embodied agents in [187]. Liu et al. [188] proposed an FGSM-like attack to fool skeleton-based human action recognition models. Their attack also accounts for multiple task-specific constraints in the optimization problem, e.g. anthropomorphic plausibility of adversarial inputs. Another example of skeleton action recognition attack is [189]. Wang et al. [190] have also provided an analysis of adversarial robustness of skeleton-based action recognition. We also see exploitation of task-specific constraints in developing attacks and defenses against such attacks. For example, Pony et al. [191] introduce flickering across the temporal dimension to fool video recognition systems. Xiao et al. [192] proposed a defense against attacks on videos that detects adversarial inputs by analysing temporal consistency property of the videos.

We also see examples in the literature that devise attacks for specific types of network architectures. For instance, attacks specifically devised for GCNs are studied in [188] and [193]. Jin et al. [194] also analyze certified robustness of GCNs under topological perturbations. Similarly, an attack specialized to binarized neural networks can be found in [195]. Another example of attacks on quantized networks can be found in [196]. We also see other unique ways of rendering inputs adversarial for deep learning models. Alaifari et al. [197] deformed image planes to construct adversarial examples. The

techniques in [198] and [126] aim at perceptibility reduction of adversarial perturbations by directly focusing on ℓ_0 -norm reduction of perturbation vector. In [199], we witness an unsupervised universal attack method to compute perturbations that exploit model uncertainty. The method uses a Monte Carlo scheme to activate more neurons to increase model uncertainty during perturbation computation with a stochastic gradient descent technique. It also exploits a textural bias prior. A steganography based universal adversarial perturbation method is proposed in [200] that embeds a secret natural image in another image to render the latter adversarial. In another example of universal attacks, Rampini et al. [201] extended the notion to deformable geometric shapes. They compute the attack in the spectral domain by perturbing eigenvalue sequence of the representation. Recovering shape from spectrum then leads to adversarial samples.

VI. ATTACKS BEYOND CLASSIFICATION

In this section, we focus on the contributions that fool deep visual models for tasks other than classification. Whereas the fundamental tools to generate perturbations for these tasks are the same as those used to fool classifiers, the unique objectives of these tasks results in more specialized attack algorithms.

A. Object detection and tracking

Object detection and tracking are longstanding computer vision problems. Their wide application in practical deep learning has led to numerous specialized techniques for these tasks. From the adversarial perspective, it has also resulted in specialized attacks. Interestingly, many of those attacks foray into the realm of physical world attacks (Section VII) due to the practical nature of these tasks. For instance, Eykholt et al. [123] and Zhong et al. [202] have analyzed adversarial stickers on stop signs in the context of autonomous driving to fool YOLO [203] - a popular object detector. Jia et al. [31] have recently developed a ‘tracker hijacking’ technique to fool multiple object trackers with adversarial examples computed for object detectors in the perceptual pipeline of autonomous driving. We note that the original tracker used by [31] follows tracking-by-detection paradigm [204]. Adversarial training of detectors for their robustification is discussed in [205]. The authors also proposed a class-aware adversarial training that uses universal perturbations to eventually compute class-weighted loss for improved robustness.

Yan et al. [206] developed an attacking technique to deceive single object trackers, in specific SiamRPN++ [207]. Their method trains a generator model to construct adversarial frames under a ‘cooling-shrinking’ loss. The loss is designed to cool down the hot target regions and forcing the bounding boxes to shrink during online tracking. A Fast-Attack-Network is also developed in [208] to attack the trackers based on Siamese network. In [209], the authors introduced an adversarial pattern that can be printed on a poster in the physical world. When a target moves in front of that poster, the tracker locks itself to the poster pattern instead of tracking the target.

Huang et al. [210] studied physical attacks on object detectors in-the-wild by developing a universal camouflage



Fig. 10. Representative attacks on detectors - randomly selected. **(Left)** Universal adversarial camouflage [210] incorrectly detects the target class at the cost of conspicuous patterns. **(Right)** Invisibility cloak [211] makes target object invisible with high probability. Images are taken from [210] and [211].

for object categories. The hard objective of the problem resulted in conspicuous patterns for their attack though, see Fig. 10. Robustness of object detectors are also explicitly studied in [30]. Whereas it appears that object detectors are relatively hard to fool, [30] shows that their robustness can also be improved with adversarial training. Another example of deceiving object detector in the real-world can be found in [211]. Zolfi et al. [212] develop a physical translucent patch that can be placed on camera lens to deceive detectors operating down the stream. A one-shot adversarial attack is proposed in [32] for single object tracking where inserting a patch in the first frame of the video results in losing the target in the subsequent frames. A spatial-aware attack (SPARK) is proposed in [213] to fool online trackers. This method imposes an $L_{1,2}$ regularization constraint over perturbations while computing them incrementally based on previous frames. It is shown that their perturbations are able to fool multiple state-of-the-art trackers. An example of black-box attack (decision-based) on trackers can be found in [214] that focuses on IoU reduction by accounting for current and previous frames.

B. Reinforcement learning

Reinforcement Learning (RL) is a major research direction in AI. Although it is not a mainstream topic in computer vision research, adversarial attacks on RL systems are often inspired by attacks on visual models. Hence, we find it imperative to briefly touch upon the advances made in adversarial attacks on RL in our literature survey.

Huang et al. [215] were among the first to demonstrate that FGSM-like perturbations can also be used to degrade the performance of trained policies in RL. They considered adversaries that can manipulate raw input of policies. Their experiments prove success of adversaries even in black-box scenarios. Xiang et al. [216] developed a PCA-based model for predicting adversarial examples in the context of Q-learning based path-finding. In another related work, Bai et al. [217] also attacked the Deep Q Network (DQN) [218] for robotic path-finding in a white-box setup. Similarly, Chen et al. [219] also explored adversarial attacks for the same problem, and devised a so-called Common Dominant Adversarial Examples Generation Method for computing adversarial examples for a given map. In light of their threat models, we can categorize [216], [217] and [219] as white-box attacks within the RL context. We can also find early instances of black-box attacks for RL. For example, [220] showed successful transferability of attacks across different DQN models. Additional early examples of black-box attacks on RL include [221] and [222].

For the interested readers, we refer to [223] for a more thorough review of the literature in adversarial attacks and defenses on RL up until 2019.

More recently, Gleave et al. [224] showed the existence of adversarial policies in zero-sum games between robots, especially in high dimensional environments. The victim of their policies are robust opponents trained with self-play. It is claimed that their adversarial policies defeat the victims reliably, while generating apparently random behavior. Rakhsha et al. [225] analysed a training time attack on RL where the adversary can poison the environment of an agent to enforce execution of a target policy in a stealthy manner. Zhang et al. [226] also developed ‘adaptive’ reward-poisoning attack that allows perturbation to reward at every step to cause learning of adversarial policies. From the perspective of endowing robustness to deep reinforcement learning from adversarial observations in an agent’s environment, Zhang et al. [227] showed that directly applying robustification methods, e.g. adversarial training is insufficient. They proposed a State Adversarial Markov Decision Process (SA-MDP) method for regularizing the policies. It is claimed that this method is applicable to a large family of popular deep RL techniques, including DQN.

C. Image Captioning/Description

Image and video captioning/description [228] is a multi-model task that normally involves a visual model (e.g. CNN) to extract visual information from the input, followed by a language model (e.g. RNN). Due to the temporal dependency in captions, attacking such a captioning/description framework is more challenging than attacking a visual model alone. Nevertheless, we do find recent examples that successfully fool these frameworks. For instance, Xu et al. [229] fool an image captioning framework by treating the generated sentences as individual labels. Their focus is on fooling the language model (i.e. RNN) while keeping the CNN embeddings of input image intact. In a related work, Chen et al. proposed ‘Show and Fool’ method [230] that fools the ‘Show and Tell’ model [231]. Their technique can generate a pre-specified target caption for any image, or embed adversarial keywords in the caption, see Fig. 11. Recently, Xu et al. [232] also proposed a targeted partial caption attack that formulates the underlying task of generating adversarial partial captions as a structured output learning task with latent variables. The problem is solved under a generalized expectation maximization method and structural SVMs with latent variables. In [233], an adversarial optimization-based attack is developed for scene text recog-

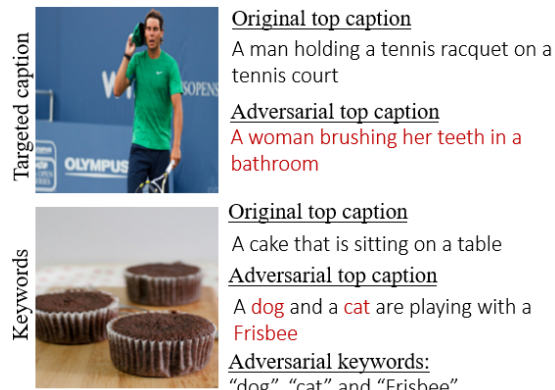


Fig. 11. Example of ‘Show and Fool’ [230] generating a targeted incorrect caption (top) and embedding adversarial keywords in a caption (bottom).

nition that employs sequential models. However, the method focuses on the related problem of text recognition, not directly on caption generation. Another remotely related work to captioning is FreeLB [234] that promotes adversarial robustness in language models with adversarial training. Whereas attacks on captioning are currently not as popular as attacks on other mainstream computer vision tasks, following the trends of other tasks, we can anticipate a gradual rise in the popularity of captioning attack in the future.

D. Face recognition

Face recognition is also a long-standing problem in computer vision. Although the task is closely related to classification, due to specific data properties, it is often treated separately from classification. From adversarial perspective, treating face recognition separately is even more meaningful due to the serious implications of adversarial attacks on these systems, which are generally not relevant to general purpose visual classifiers.

Although deep learning era has witnessed highly accurate face recognition models [22], [235] these systems are also vulnerable to adversarial attacks. Goswami et al. [236] provided an analyses of face recognition systems’ robustness against adversarial attacks. They ascertained the susceptibility of popular model OpenFace [237] and VGG-Face [238]. Dong et al. [239] also reported adversarial vulnerability of of face recognition in black-box setups, specifically to decision-based attacks. They adapted a popular evolutionary strategy [240] to perform search over the perturbation in the black-box setup, where the search is guided by the local geometry of the searched directions for efficiency. Zhong et al. [241] used *transferability* to fool face recognition systems in another black-box scenario. They devised a so-called Drop-out Face Attacking Network (DFANet) that focuses on matching the internal representation of an identify (image) with another identity to confuse the target model between the two. An FGSM-like method, Penalized Fast Gradient Value Method was introduced in [242] to demonstrate fooling of face recognition models. A friend-safe attack on face recognition systems was also introduced in [243], which computes images that are adversarial for ‘enemy’ models, but benign for ‘friend’ models.

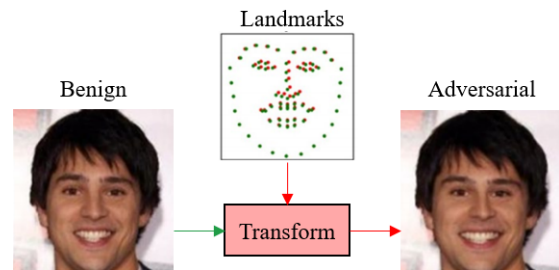


Fig. 12. An example of fooling face recognition system by landmark manipulation [244].

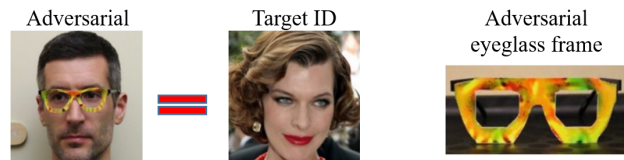


Fig. 13. An example of fooling a face recognition model by wearing an adversarial eyeglass fame [250].

The above methods mainly compute additive perturbations without explicitly accounting for geometric information of faces. In contrast, Dabouei et al. [244] devised a facial landmark manipulation method to mislead recognition systems. Their technique computes adversarial landmarks to perform spatial distortions in images that result in incorrect recognition, see Fig. 12. Adversarial patches for faces are also studied for their transferability in [245]. In another related example, Yang et al. [246] proposed an Attentional Adversarial Attack Generative Network (A^3GN) for targeted fooling of face recognition models. It is claimed that their network is able to exploit geometric and context information of the target with the help of a conditional VAE and attention modules to achieve this feat. Another example of using GAN for deceiving face recognition systems is AdvFaces [247] that manipulate geometric features of the face in the image. Similarly, Li et al. [248] generate a fake face image by matching the latent representation of the image with its adversarial counterpart that can fool fake image detectors. Along the line of utilizing GANs for manipulating faces in images and videos, an interesting research direction of DeepFakes is emerging. Interested readers are referred to [249] for a recent survey of that direction.

Due to the practical nature of face recognition task, the literature has also witnessed fooling attempts through manipulating faces in the physical world. For instance, Sharif et al. [250] demonstrated the possibility of physically realizable attacks to impersonate an identity or evade the face recognition system. They devised an eyeglass frame for fooling the target network, see Fig. 13. Their technique was further improved in [125] for attack robustness. On a similar line, Zhou et al. [251] developed a cap that illuminates face of the person wearing it to fool the recognition system. They compute the adversarial illumination pattern on a image of the identity and use the cap to project that pattern on the face in physical world while presenting the face to the vision system. A related concept of ‘adversarial light projection’ is studied in [252]

that projects a rather conspicuous pattern on faces to evade FaceNet model [253] in white-box settings. Other examples of physical world attack on face recognition systems include AdvHat [254] and adversarial patches for faces [255]. Face presentation attacks are also studied in [256]. We refer to our first survey [2] for the more classic attacks on face recognition.

E. Miscellaneous attacks

Even for the tasks beyond classification, multiple other attacks exist that do not fall under the categories described above. For example, Nakka et al. [257] devised an attack to demonstrate the vulnerability of semantic segmentation networks against holistic perturbations and localized ones. Similarly, for the problem of segmentation, a data membership attack is devised in [258]. Choi et al. [259] also observed that deep learning models for super-resolution are also vulnerable to adversarial attacks. This is demonstrated by introducing unnoticeable distortions in the low-resolution images, which adversely affect the super resolution results. Mehra et al. [260] proposed a poisoning attacks for reducing the average certified radius of a given class for certified defenses.

Deep neural networks are often successfully applied to predict depth in monocular scenes. Recently, [261] showed that adversarial attacks can be used to manipulate the predicted distance from the camera. The method in [261] can match the predicted distance to a different target scene or directly fabricate the depth of specific instances in the scene. Targeted attacks on hashing based retrieval are proposed in [262], [119], whereas a universal perturbation for image retrieval systems is computed in [263]. An example of adversarial attack on Graph Matching can be found in [264]. There has also been enhancements and variants of patch attacks for multiple vision tasks. For example, Yang et al. [265] improved the patch attack in a blackbox setup by reducing the required number of queries with reinforcement learning. Similarly, a universal patch is proposed for face recognition in [266] and the patch attack is extended to optical flow in [267]. It is shown that a patch as small as 1% of the image size can disrupt optical flow networks.

VII. PHYSICAL WORLD ATTACKS

We already reviewed some of the literature performing physical world attacks in § VI-A for the tasks of object detection and tracking, and for face recognition in § VI-D. Below, we further expand on the literature in this direction by focusing on the practical physical world application of autonomous driving and general purpose object detection and classification attacks.

In the context of autonomous driving, Tu et al. [29] proposed a technique to compute physically realizable adversarial examples using LiDAR data to fool object detectors in simulated autonomous vehicle scenarios. It is claimed that placing an adversarial object (with underlying adversarial mesh computed from their technique) on the rooftop of a target vehicle can make the vehicle undetectable. The mesh surface of the computed object generally remains unnatural though. In another study, Cao et al. [268] showed that CNN-based



Fig. 14. Representative examples of successful physical world attacks to fool recognition systems with AdvCam [273], RP2 [123] and adversarial patch [45].

object detectors can be fooled in vehicle detection scenarios. An adversarial pattern computed by their technique serves as a camouflage to evade detectors in their work. The notion of camouflage is also explored in the physical world settings in [269], [270].

Kong et al. [271] used a GAN-based setup to generate norm-bounded adversarial images, which when printed, demonstrate resilience to changes in the physical world conditions, e.g. lighting condition, viewing angle. Their method, PhysGAN is specifically designed to fool steering models of autonomous vehicles, under a regression-based formulation of the angle prediction problem. PhysGAN computes perturbations for a stream of visual features of driving video while ignoring the scene background. It is claimed that this strategy allows effective perturbations for dynamic scene conditions, nullifying the need of static scene assumption appearing in earlier literature [123]. In [272], it is shown that with camera shake and pose variation while imaging physical world objects, one can acquire images that can easily fool deep learning models. Here, the imperceptibility of the perturbation comes in the form of semantic-imperceptibility i.e. contextually, the pose or shake appears natural.

The use of ‘adversarial patch’ [45] is another effective method to launch a physical world attack. An adversarial patch is normally a clearly visible, but well localized pattern - i.e. a patch that can be placed beside an object to cause model fooling, see Fig. 14. More recently, Duan et al. [273] proposed a neural style transfer [274] based technique to compute unrestricted perturbations that can take effect as a physical world attack to camouflage a target object. Their proposed AdvCam is able to compute patterns that are claimed to be more stealthy than earlier related techniques (e.g. adversarial patch [45], RP2 [123]) in that the adversarial pattern appears more natural to humans (Fig. 14). To compute the physical world pattern, their technique captures image of the scene with a given camera and estimates the perturbation digitally in a restricted knowledge white box setup. The substitute model used to compute the adversarial pattern has the same architecture as the target model. The pattern is then placed in the physical world and the same camera is again used to capture the adversarial example. The technique requires manual specification of the region to place the adversarial pattern and the target style. In another example related to adversarial patch, Liu et al. [275] constructed a universal patch and used it to deceive automatic checkout models.

As also noted in § VI-A, deceiving object detectors is a particularly interesting problem for the physical world attacks. We are already witnessing interesting methods in this direction.

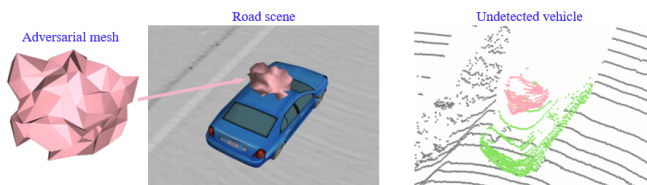


Fig. 15. Representative example of fooling object detector with LiDAR data [29]. The adversarial mesh placed on top of vehicle makes it undetectable for visual detector. Image taken from [29].

Recently, Xu et al. [276] proposed a technique to fabricate adversarial T-shirt for evading detectors. Another example of fooling object detectors on similar lines is [211]. The authors compute a so-mentioned ‘invisibility cloak’ that contains the patterns causing misdetections for state-of-the-art detectors, see Fig. 10. Considering the implications of this direction of research, a dataset for adversarial attacks on object detectors in the physical world is also introduced in [277]. Whereas currently attacks on object detectors are not as popular as attacks on classifiers, we can anticipate much larger interest of the research community for this problem due to many interesting, and sometimes security-critical applications.

There are also works in the literature that fool classifiers by distorted illumination in the scene. For instance, Sayles et al. [278] distort leverage radiometric rolling shutter effect for distortions that cause misclassification. Similarly, Duan et al. [279] proposed an adversarial laser beam attack, which computes adversarial parameters for a laser that can be used to distort illumination such that the captured image is adversarial. We have already seen multiple examples for fooling face recognition systems with adversarial illumination patterns in Section VI-D.

VIII. BEYOND ADVERSARIAL OBJECTIVE

Although the primary objective of adversarial attacks in the literature is to fool deep learning models, there are also instances where adversarial perturbations are exploited under more constructive objectives of improving model performance, interpreting it, or estimating the performance. Note that, for the former, we are not alluding to the works leveraging adversarial training to robustify models - explained shortly.

A. Improving model performance

Xie et al. [280] recently demonstrated that adversarial examples can actually help in performance gain in fully supervised setups for large-scale models, e.g. ImageNet [10]. To demonstrate that, the authors propose Adversarial Propagation (AdvProp) technique that is applied to EfficientNet-B7 [281] to achieve performance gains of 0.7%, 6.5%, 7.0% and 4.8% for ImageNet, ImageNet-C, ImageNet-A and Stylized ImageNet datasets. Moreover, after enhancing the network to EfficientNet-B8, their method sets the new state-of-the-art of 85.5% on ImageNet top-1 accuracy without extra training data. The key insight used by AdvProp is that the underlying distribution of adversarial images is different from natural images. This calls for disentangling the normalisation statistics for the networks in the Batch Normalization (BN) layers.

Hence, the authors proposed an auxiliary BN layer that is explicitly used for adversarial examples during training, and dropped during testing. During training, the loss is computed by propagating the clean and adversarial images separately through their respective BN layers.

The AdvProp is unique in that successfully aims at performance gain for large-scale models on clean images with adversarial examples. This is different from adversarial training, which generally results in sacrificing model accuracy on the clean images to gain robustness to adversarial examples [2], [282]. There are also other instances that report performance gain on clean data by accounting for adversarial image in training. For instance, [1], [283] report improved model accuracy for a small dataset (MNIST) under a fully supervised setup. Similarly, [48] and [284] improve model performance with adversarial examples for large models in a semi-supervised setup. Ho and Nuno et al. [285] also found use of adversarial example in Contrastive Learning for self-supervised learning. They used adversarial examples to augment data for pretext learning of embeddings.

It is also claimed by Salman et al. [286] that adversarially trained models, while less accurate than the standard models, often perform better for transfer learning. In another study, Gan et al. [287] propose VILLA, a representation learning approach based on large-scale adversarial training on vision-and-language data. They perform a task-agnostic adversarial training followed by a task-specific adversarial fine-tuning in the embedding space. This method is claimed to achieve state-of-the-art performance on a variety of tasks, including Visual Question Answering, Visual Commonsense Reasoning, and Image-Text Retrieval. Using the anti-adversarial directions for weakly supervised models, Lee et al. [288] claimed improvement in the semantic segmentation performance.

B. The link between attacks and model interpretation

Jalwana et al. [289] developed a technique to visually reveal the understanding of human-defined semantic concepts by deep learning perceptual models, see Fig. 16. By expanding the domain of the adversarial perturbation and iteratively refining it, the authors demonstrate the presence of human-understandable patterns in the perturbations. A more clear relation between the adversarial and explanation character of their perturbations is later established in [154]. The authors also utilize their ‘attack to explain’ to perform low-level vision tasks by attacking robust classifiers. This concept builds on [290]. In a related approach, Augustine et al. [291] associate model explainability to its adversarial robustness, demonstrating generative properties of their adversarially robust model similar to [290]. Elliott et al. [292] also attempts to bridge the gap between adversarial perturbations and counter-factual explanation of deep models. They localized their perturbations to salient regions of inputs to demonstrate that perceptually regularized counterfactuals provide useful model explanation.

There is also a line of research that considers interpretability of the induced perturbation patterns themselves. For instance, Xu et al. [293] applied group sparsity over the perturbation vector and showed that the resulting perturbations are more

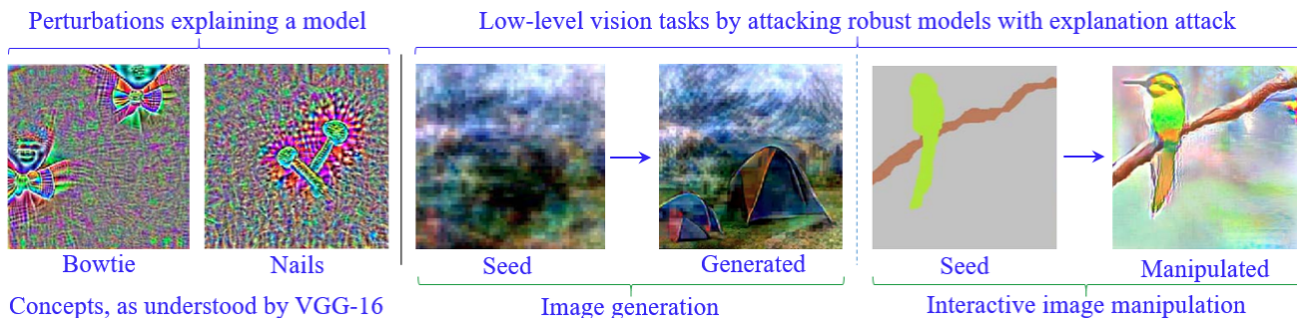


Fig. 16. Jalwana et al. [289] proposed an ‘attack to explain’ which uses perturbations to visualize human-defined semantic concepts as understood by a model. They also utilized their attack to perform low-level vision tasks by attacking robust classifiers. **(Left)** Perturbations computed with attack to explain VGG-16. The signals visualize understanding of VGG-16 for human-defined concepts of Bowtie and Nails. **(Right)** The attack is used to generate an image of ‘mountain tent’ with a random seed image. For interactive image manipulation, the seed image is refined into a bird with the attack.

interpretable. Nevertheless, this method does not offer interpretability of the model itself like [289], [154].

C. Other applications

Among other interesting related applications, Elsayed et al. [294] showed that adversarial perturbations can be used to reprogram a target model. For instance, with embedded perturbations, they successfully converted a classifier into a box-counting machine. Finally, Sakaguchi et al. [295] recently proposed an algorithm, called AFLITE, to adversarially reduce task- or dataset-specific biases in head distribution, while preserving complexity of the tail. This bias reduction mitigates overestimation of model performance, which is evident by their performance on out-of-distribution and adversarial examples. In a followup work, Le et al. [296] theoretically studied AFLITE and provided extensive evidence that AFLITE reduces measurable dataset biases. They showed that the models trained on the filtered dataset generalize better to out-of-distribution data.

IX. ON THE EXISTENCE OF ADVERSARIAL EXAMPLES

The existence of adversarial examples for otherwise highly accurate deep visual models has confounded the research community since the discovery of this phenomenon. The literature has witnessed numerous hypotheses to explain the adversarial vulnerability of deep learning. However, many of those fail to generalise, and the remaining often conflict with each other. It can be argued that there is still no consensus on the reasons of the existence of adversarial examples. Whereas it was common among the earlier contributions to also hypothesize about generic causes of the adversarial susceptibility of neural networks, the recent attack methods are more concerned with achieving higher fooling rates and better transferability etc. Nevertheless, the contributions that analyze the causes of adversarial vulnerability as the core topic, are still attractive because the wider impression is that this phenomenon is still not fully understood. Below, we review contributions and major hypotheses in this direction along the lines of input-specific perturbations, input-agnostic perturbations and other prevailing topics.

A. On input-specific perturbations

One of the first popular hypotheses on the existence of adversarial examples for modern deep network was the linearity hypothesis [17] - see the FGSM attack in § IV for details. However, it was later shown by Tanay and Griffin [297] that this hypothesis does not generalize as there are classes that do not suffer from adversarial examples for linear classifiers. Nevertheless, Kortov and Hopfiled [298] later provided another evidence based on Dense Associative Memory [299] that supports the role of linearity in neural model susceptibility to adversarial examples. Similarly, linearity is also blamed for adversarial vulnerability in [300]. In [301], ‘inherent prediction uncertainty’ of neural networks is blamed for their adversarial vulnerability. The claim is corroborated by computing a functional form of the prediction uncertainty which remains independent of the architecture and training of the model. It is also argued that clean image accuracy of models correlates with their adversarial robustness, which resonates with the findings of [302] and other earlier observations [17], [54]. Evolutionary stalling [303] is another interesting hypothesis, according to which, the inability of training samples to contribute beyond a certain capacity leaves their representation very close to the model decision boundaries. This allows adversaries to easily nudge those and similar representations out of the correct classification regions.

Exploring the space of adversarial examples, Tabacof and Valle [304] showed that adversarial examples reside in large regions in the pixel space of images. Their findings suggest that weak shallow networks are as susceptible to adversarial examples as the complex deep networks. On a similar note, Tramer et al. [305] claimed that adversarial examples span a contiguous high dimensional space. The high dimensionality of this space and subspaces of different classifiers results in their intersections which causes transferability of the attacks across different models. More recently, [306] claimed that model’s gradient leakage along the perpendicular to a tangent space of training data manifold contributes to adversarial vulnerability of the models.

In another work, Jacobsen et al. [307] claimed that deep neural networks are highly invariant to a variety of task-relevant changes in the input that causes vast input space regions to be vulnerable to adversarial perturbations. This is

in addition to the high sensitivity of the models to the task-irrelevant changes to the input. Along the lines of analysing the existence of adversarial examples from the robustness perspective, Reddy et al. [308] studied the biological visual system of primates. They showed that non-uniform sampling done by primate retina and the existence of multiple receptive fields (having a range of field sizes) improves the robustness of neural networks to adversarial perturbations.

Pal and Vidal [309] proposed a game-theoretic framework for analysing attacks and defenses that exist in equilibrium. They proved that under a locally linear decision boundary model, FGSM and the randomized smoothing [310] exhibit a Nash Equilibrium [311]. Daniely et al. [312] provided a theoretical analysis that studies the vulnerability of ReLU networks against adversarial perturbations, concluding that most ReLU networks suffer from ℓ_2 perturbations. We also find a similar but broader claim that adversarial examples are inevitable for certain types of problems in [313].

Similar to the linearity hypothesis of Goodfellow et al. [17], another popular concept related to the existence of adversarial examples is ‘manifold assumption’ [297], which argues that adversarial examples tend to leave the clean data manifold. Nevertheless, there is also evidence of on-manifold adversarial examples [314], [315], [43]. There is also a debate in the literature of associating robustness of neural models to their generalization [316]. For instance, Trsipras et al. [317] provide systematic evidence of clash between adversarial robustness and generalization of a model. This is also partially supported by the empirical study in [318]. However, we also find works in the literature that suggest the opposite [314], [302], i.e. improved generalization results in better robustness.

B. On Input-agnostic perturbations

Analysing the existence of universal perturbation, Moosavi-Dezfooli et al. [35] claimed that these signals leverage the geometric correlations between the decision boundaries of classifiers. The authors theoretically demonstrate the existence of common directions for multiple data points along which a classifier’s decision boundaries can be highly curved [319], [320]. It is argued that such directions allow the universal perturbations to effectively fool the classifier across multiple samples. On a similar note, Jetley et al. [321] demonstrated that the directions (in image space) used by neural networks to achieve higher performance are the same that make them vulnerable to adversarial attacks. Thus, the high accuracy and adversarial vulnerability of neural networks are related phenomena, which allow the existence of universal perturbations.

Analysing the Pearson correlation between the coefficients of logit vectors of a classifier for clean and adversarial images, Zhang et al. [322] showed that for universal perturbations, adversarial examples are strongly correlated with the perturbations. On the other hand, a low correlation is observed between the adversarial and clean images. This leads to the conclusion that universal perturbations hold more dominant features as compared to clean images despite their low power and visual (quasi-)imperceptibility. The authors also leverage this insight to introduce a method to compute universal perturbations using random clean images.

C. Adversarial examples as features & other sources

Inline with the findings of Jetley et al. [321] (discussed above), Ilyas et al. [323] claimed adversarial examples to be essential data features for neural networks, as opposed to unwarranted bugs. They demonstrated that the existence of adversarial examples can be attributed to non-robust features that are pervasive in datasets and are an effective source of achieving higher accuracy for the neural perception models. The authors also demonstrate the possibility of disentangling robust and non-robust features and showed that robust features align more to human perception than their non-robust counterparts. This insight is exploited by Santurkar et al. [290] to use adversarially robust models to perform visually appealing image synthesis. Tsipras et al. [317] also noted the tension between adversarial robustness and classifier accuracy under the idea that adversarial examples are non-robust features used by models to achieve better performance.

Bubeck et al. [324] argued that adversarial vulnerability of classifiers in high dimension is “likely not due to information theoretic limitations, but rather it could be due to computational constraints”. They provided evidence to support the hypothesis that “identifying a robust classifier from limited training data is information theoretically possible but computationally impossible”. Interestingly, their evidence weakens the notion that identifying a robust classifier requires huge amount of training data. In a study more focused on ResNet inspired architectures, Wu et al. [325] identified adversarial vulnerabilities of skip-connections. The authors claimed that the use of skip connection results in more transferable adversarial examples for models. They introduced a Skip Gradient Method (SGM) that relies on the gradient flow from skip connections to compute more transferable examples for the models that employ skip connections.

X. DEFENSE AGAINST ADVERSARIAL ATTACKS

Akhtar and Mian [2] organized defenses against adversarial attacks into three broad categories. They incorporated defenses resulting from (1) modifying the target models for robustification, (2) modifying input for perturbation removal and (3) adding external modules (mainly detectors) to the model. Since 2018, the research direction of adversarial defenses has evolved mainly along the same three lines. Hence, we first review the recent literature along the same directions. However, there is a subclass of defenses that is gaining rapid popularity in the recent literature, known as ‘certified defenses’. Although most of the works in this subclass follow (1), we review these methods separately in § X-D due to their unique common objective of providing certificates/guarantees for the developed defenses. In § X-E, we also provide a bird’s-eye view of other recent defense techniques that either combine more than one strategies noted above, or develop specialized defenses, e.g. for specific tasks or network types.

A. Model alteration for defense

The most common framework that modifies the (potentially) targeted model itself for robustification against adversarial attacks is ‘adversarial training’. Hence, we review techniques

focusing on this framework separately, before discussing other recent methods in the category.

1) *Adversarial training*: The adversarial training framework is considered among the strongest principled defenses against adversarial attacks. It exposes the model to adversarial examples during training to obtain some level of immunity against them. Adversarial training was originally employed in [1], [17]. However, Madry et al. [54] are the first to theoretically study and justify it through the lens of robust optimization for deep learning. Since [54], adversarial training has attracted significant interest of the research community. This also resulted in multiple contributions highlighting weaknesses of this framework. For instance, Zhang et al. [326] showed that adversarially trained models are still vulnerable to ‘blind-spot’ attacks. Arguments against the robustness induced by adversarial training can also be found in [327]. It is also claimed that adversarial training is sensitive to the training data distribution in [328]. Moreover, poor generalization of adversarial training is also often highlighted in the literature [329], [330], [331], [332].

Despite its shortcomings, adversarial training is still favored by the research community due to its principled nature. Over the last few years, multiple variants and enhancements of adversarial training have surfaced. For example, a Misclassification Aware adversarial Training (MART) is proposed in [333] to incorporate distinctive influence of clean misclassified examples in the training process. Gowal et al. [334] improved adversarial training by varying the Style-GAN-based [335] disentangled representations of original images. This can be considered a defense against unrestricted perturbations. A margin-maximization variant of adversarial training was proposed in [336] that creates adversarial examples using sample-specific η (see Eq. 1) instead of a fixed η value across the training samples. The η value corresponds to the “shortest successful perturbation” that fools the model.

Among other variants of adversarial training frameworks, we have [337] where in each training iteration, the model is verified for robustness using convex relaxation and adversarial examples are computed under that relaxation for training purpose. Vivek et al. [338] also proposed a dropout scheduling method to improve the efficacy of adversarial training with single-step methods. To improve generalization of adversarially trained models, Song et al. [339] proposed Robust Local Features for Adversarial Training (RLFAT) that employs random block shuffle of the input during training. Farnia et al. [340] also proposed a spectral normalization based regularization for adversarial training to address the generalization issue. In [341], enhancement is suggested by using adversarial examples generated by attacking a model other than the model to be defended. To make adversarial training more efficient, Zheng et al. [33] proposed to use the same adversarial perturbations across multiple epochs during the training. This reduces the number of computations in the overall training process while achieving acceptable performance.

Naseer et al. [342] proposed self-supervised adversarial training, whereas adversarial training is independently analyzed for self-supervision by incorporating it in pretraining

in [343]. Similarly, [344] used a generator in adversarial training to generate more diverse adversarial examples. A Dual Manifold Adversarial Training (DMAT) is proposed in [345], which uses perturbations in the image space as well as latent space of StyleGAN to make training more effective. In another related work, Wang et al. [346] proposed a bilateral adversarial training that not only perturbs input images during training, but also their labels. The authors claimed improvements in state-of-the-art adversarial training results with this modification. It is often argued that adversarial training leads to the requirement of larger models. Ye et al. [347] proposed a concurrent adversarial training and weight pruning strategy to address this specific issue.

Considering further variants of adversarial training, Don et al. [348] proposed an adversarial distributional training. Their method also formulates adversarial training as a minimax problem, however, the inner maximization is aimed at learning an adversarial distribution under an entropic regularizer. The outer minimization problem minimizes the loss over the worst-case adversarial distributions. Madaan et al. [349] proposed a vulnerability suppression loss that minimizes the expected difference between latent features of the network on clean and their corresponding adversarial images. They further learned a pruning mask that explicitly minimizes adversarial loss by pruning features with high distortion. In an attempt to address multiple perturbation models for adversarial training, Maini et al. [350] proposed to incorporate several perturbations into a single attack by taking the worst-case over the entire steepest directions as an extension to the standard PGD. It is claimed that their approach produces state-of-the-art robust classifiers against ℓ_1 , ℓ_2 , and ℓ_∞ norm bounded perturbations simultaneously.

In the literature, we also find methods that conceptually relate to adversarial training closely without presenting themselves as such. For example, in [351], a mixup training of neural networks was introduced. The main concept of the method is to augment training data with additional samples that are created as convex linear combination of the already available samples. The same is done to the labels of the combined samples to provide the label of the resulting samples. It is shown that besides improving accuracy of the original model, this practice also helps in robustness against adversarial samples [352]. Pang et al. [167] takes this notion further by also applying mixup of samples in the inference phase. A related adversarial vertex mixup method is adopted in [353] to achieve better adversarial generalization of the models.

We also find multiple contributions in the literature that focus on *analyzing* adversarial training instead of devising its variants. For example, Xie et al. [354] have reported some interesting properties of adversarial training. The most intriguing ones include an improvement in the adversarial robustness for the process with separate Batch Normalization for clean and adversarial images, and a consistent improvement in the adversarial robustness with even deeper models as compared to the popular depth limits among the visual models. Li et al. [355] also analyzed the implicit bias of gradient descent on adversarial training on separable data. Their findings theoretically back the efficacy of adversarial

training for robustness. In [356], it is also demonstrated that transfer learning on adversarially robust models retains (to an extent) the robustness effect for the target domain. Sehwag et al. [357] also devised a method for an adversarial training-aware model pruning in resource constrained environment.

Wong et al. [358] showed that adversarial training with FGSM combined with random initialization is as effective as adversarial training with the first order PGD attack. On their computational setup, they trained a robust CIFAR10 classifier with 45% robust accuracy in 6 minutes as compared to the 10 hours training of PGD-based counterpart that achieves similar results. Their improvement of adding randomization with FGSM-based adversarial training is, however, contradicted to an extent by [359]. Zhao et al. [360] study the mode connectivity of loss landscape of adversarially robust and regular models, demonstrating the existence of robustness loss barrier for the former. Wu et al. [361] showed that many adversarial training improvements appearing in the literature, e.g. early stopping, new objective functions, or exploiting unlabeled data, implicitly flatten the weight loss landscape (i.e. loss change w.r.t. weights). Hence, they proposed an Adversarial Weight Perturbation (AWP) that directly regularizes the flatness of weight loss landscape, and can be used to improve adversarial training.

Even though Madry et al. [54] have justified adversarial training by robust optimization theory, it is still unclear how adversarial training results in low robust training loss. Gau et al. [362] provide a theoretical analysis of adversarial training to explain its success using Neural Tangent Kernel and tools from online learning. They also prove that more model capacity is required for robust interpolation. However, their approach is limited to networks with exponential width and run time. Zhang et al. [363] extend their work for situations where the width of the network and its run time is polynomial in input dimension. They also extend the results to ReLU activation function. Another related method [364] proposes to boost adversarial training by embedding a hypersphere method in the training process by regularizing features onto a compact manifolds.

The literature also contains instances in which adversarial training is moulded to specific task requirements. For example, Wu et al. [365] proposed an adversarial training method in which the adversarial samples are generated specifically keeping in view the physical world attacks. It is noted in [365] that commonly used adversarial training and randomized smoothing for the digital attacks do not perform well for the physical world attacks. Hence, the modification was proposed. Instead of focusing on robustness against adversarial attacks, Zhu et al. [234] employ adversarial training in natural language understanding for achieving higher embedding space invariance by perturbing the word embeddings. This is reported to result in better generalization of language models. This result also resonates with the observations of [366].

2) *Other model modifications*: Besides adversarial training that focuses on modifying model weights through alternate training samples, there are multiple approaches that alter the basic building blocks of the model to incorporate adversarial robustness through regular training data. For instance, Pang

et al. [162] suggested to replace the softmax cross-entropy loss with a new loss, called Max-Mahalanobis center loss to induce adversarial robustness in the model. Xiao et al. [158] proposed to alter the ReLU activations with a k -winner-takes-all C^0 discontinuous to secure models against the gradient-based attacks. There are also works that advocate on modifying the networks in a holistic manner. For instance, in [165], the authors suggest using quantized models for robustness against gradient-based attacks. Guo et al. [367] also proposed RobNets, designed with neural architecture search, which are claimed to provide up to 5% gain in robust accuracy on large datasets, e.g. ImageNet. Bui et al. [368] propose an Adversary Divergence Reduction Network that can be used in conjunction with adversarial training for improved robustness. Similarly, a Bayesian neural network is proposed in [369] for adversarial robustness.

From the viewpoint of altering internal components of models, Jeddi et al. [370] proposed perturbation-injection modules in the internal layers of model during training and testing and used an alternating back-propagation scheme to train the network. A 4-7% improvement in robustness over adversarial training with FGSM and PGD (ℓ_∞) is claimed by the authors. Li et al [371] introduced image restoration and denoising modules in the network and constrained its classification layer's Lipschitz constant for adversarial robustness. In the context of defense against the universal perturbations [35], [372] devised a method to identify adversarially vulnerable convolutional filters in a model and introduces 'regeneration units' to generate resilient features for those filters to avoid fooling. Wang et al. [373] proposed to model adversarial noise with a generator that is trained jointly with a discriminator classifier and showed its effectiveness against black-box attacks. Xie et al. [374] suggested that adversarial perturbations result in noisy features of the networks. Hence, they proposed networks containing denoising blocks with non-local means or other filters. Building on the idea of injecting noise in the network while training [375], [376], He et al. [377] proposed a trainable Gaussian model for injecting the noise. A family of CNNs that alternate between the Euclidean convolutions and graph convolutions to leverage the information from the graph of peer samples is proposed in [378].

Another emerging model alteration approach to defend against adversarial attacks is through search for robust architectures. Following this paradigm, Hosseini et al. [379] propose DSRNA to search for robust architectures via two differentiable metrics for robustness. Moreover, Cazenavette et al. [380] proposed a deep pursuit algorithm that formulate the architecture search as a global sparse coding problem that jointly computes all network activations.

We also witness techniques that approach at adversarial robustness from the model regularization perspective. For example, based on the observation that Jacobian of adversarially robust models are more salient and interpretable as compared to their non-robust counterparts [317], Chan et al. [381] proposed a Jacobian-based GAN-like regularization scheme to show improved robustness. A joint gradient phase and magnitude regularization was proposed in [382] to improve robustness of ensemble models. A concept of biologically

inspired post-learning sleep phase of neural networks was introduced in [383]. The proposed technique allows a trained network to reflect on its statistics in an unsupervised manner and alter the weights to avoid over-fitting to the training data. Addepalli et al. [384] proposed to regularize models with Bit-plane consistency as an efficient alternate for adversarial training.

We note that whereas we discuss the above methods separately from adversarial training to provide a better structure to the literature, the boundary separating these two lines of research is often abstract. One can understand adversarial training as a more fundamental framework that can generally be combined with other defenses, including those discussed in the subsequent sections, for improved robustness. Other methods discussed above often demand model modifications that are less generic.

B. Detection for defense

Instead of proactively inducing a robust model during the training phase, there are also techniques that provide add-on mechanisms and modules for pre-trained models to defend them against adversarial attacks. Mostly, these methods are limited to detecting the presence of adversarial perturbations in the input during inference. Based on our earlier survey [2] and recent literature, we can say that this line of research is getting slightly less popular (as compared to its earlier years) in the leading research sources of computer vision and machine learning. A possible reason for that is their ad-hoc nature as compared to defenses like adversarial training. Nevertheless, we still witness interesting techniques of adversarial detection using add-on mechanisms in the recent literature.

Qin et al. [385] proposed a mechanism of class-conditional reconstruction of images to detect adversarial examples during test time. The authors also introduced an attack to overcome this defense, demonstrating better robustness of CapsNet [386] over CNNs for their attack. More importantly, their attack shows more visual similarity between the adversarial examples and target object category for CapsNet. In essence, this demonstrates a larger perceptual alignment between CapsNet representation and human visual system as compared to CNNs. For reference, perceptual alignment between deep visual models and human vision is also discussed at length in [154], [289]. In [387], the authors proposed to leverage Lightweight Bayesian neural networks for task agnostic detection of adversarial perturbations in inputs using Bayes principle. The technique replaces last few layers of the attacked model with Bayesian module and performs detection-oriented fine-tuning that allows to maintain original performance while enabling detection.

Li et al. [388] proposed to use context inconsistency of adversarial patterns in images for their detection using an external mechanism. For face recognition, Tao et al. [389] proposed a method to identify internal neurons corresponding to critical facial attributes. By amplifying activation of these neurons, they construct an attribute-steered model. Later, they detect adversarial examples by identifying inconsistencies between the original and the attribute-steered models. In [390], the authors proposed a mechanism to trace the activation

paths of clean and adversarial images and detect adversarial perturbations based on the different characteristics of these paths. Liu et al. [391] proposed to detect adversarial examples by analysing inputs from steganography point of view. Their method estimates the probability of modification to images keeping in view adversarial perturbations. Yin et al. [392] introduced a so-called generative adversarial training method that learns an adversarial example detector. To robustify the detector against adaptive attacks, the authors employed asymmetric adversarial training.

C. Input transformations for defense

Instead of focusing on model robustness to ‘adversarial’ inputs, transformation based methods aim at cleaning inputs to make them benign for the target model. For instance, JPEG-based compression of input has been studied for removing adversarial perturbations from images [393], [394], [395]. Compressed adversarial images have been found to significantly lose their fooling abilities. Generally, input transformation provides the benefit that it can be easily used in conjunction with other defense mechanisms, e.g. with adversarially trained models. In some cases, different input transformations are also combined to improve their collective strength. For example, in [396], Raff et al. proposed to stochastically combine multiple input transformations to also secure their defense against adaptive attacks. However, it is also observed in [396] that more transformations undesirably lead to significant reduction in model performance on clean images. Similarly, [397] also proposed to utilize a set of random input transformations as an adversarial defense. The main idea behind this method is that the ‘key’ controlling the randomization of transformations is assumed to be kept secret during test time. This mitigates the risk of potential adaptive attacks on their defense.

Instead of directly using standard image compression, learnable compression methods that use neural models are also proposed in the literature for adversarial defense [398], [399]. In [399], an external defender module for a deployed model is learned that projects inputs to a so-called adversarial-free data zone for the target model. We can also categorize learning-based compression techniques as defense mechanisms altering the models by appending add-on modules to them. In another related work, Sun et al. [400] transformed an input image using convolutional sparse coding. Their method use a ‘Sparse Transformation Layer’ to project input to a quasi-natural space that is claimed to be less sensitive to adversarial manipulation.

In [401], Samangouei et al. presented one of the first examples of input transformation using GANs. Their method, Defense-GAN learns the distribution of clean images. For inference, it computes an output close to the input image, which does not contain the potential adversarial perturbation. A denoising based defense is proposed by [402] that selectively denoises high attention regions of an image to recover the correct label. Kuo et al. [403] noted that when input transformation is employed as a defense technique [395], the softmax distribution characteristics can be used to improve the clean image accuracy of the classifier with the help of an external lightweight classifier trained on the softmax distribution of clean images.

In [404], an ensemble generative cleaning with feedback loop is proposed to clean the image from adversarial patterns. Their method also relies on external generative modules to denoise adversarial images. Cohen et al. [405] developed an external detector mechanism for adversarial samples using a so-called influence function. This function measures the impact of all training samples on the validation data to provide sample influence scores. Supportive training instances for validation samples are identified with their scores. A k-nearest neighbor (k-NN) model is also fitted on the models activations to compute a ranking of the supportive training samples. Supportive samples are claimed to be highly correlated with the nearest neighbors of clean test sample, while the correlation is found to be weaker for adversarial inputs.

D. Certified defenses

Although the literature is witnessing multiple defense techniques, it is shown that stronger attacks can be formed to defeat the existing defense methods [156], [63], see § V-F for more examples. Even adversarial training has its problems despite being widely considered a reliable defense strategy. For instance, adversarially trained models with ℓ_∞ -norm bounded perturbations are still found vulnerable to ℓ_p -norm perturbations, where $p \neq \infty$ [327], [406]. Certified defenses attempt to provide guarantee that the target model can not be fooled within an ℓ_p -ball of the clean image. This guarantee is either achieved by computing the minimal ℓ_p -norm of the perturbation to break the provided defense [407], [408]; or by providing a lower bound on the norm [409], [410], [411]. There are also other methods that aim to both enhance network robustness and produce models that are more amenable to robustness verification techniques [412], [413]. Nevertheless, most of the certified defenses are able to prove their robustness against only one kind of bound on the perturbation, e.g. ℓ_2 , ℓ_∞ , struggling to provide generic bounds for multiple ℓ_p -norms simultaneously [414], with a few exceptions [327], [415].

Corce and Hein [414] recently proposed a regularization scheme for ReLU networks to enforce robustness against ℓ_1 and ℓ_∞ attacks and showed that it results in provable robust models for any ℓ_p norm, where $p \geq 1$. As opposed to providing certified robustness for top-1 predictions, Jia et al. [416] derived tight robustness in ℓ_2 -norm using Gaussian randomized smoothing for top-k predictions. Their method builds on the notion of randomized smoothing introduced in [417] and [310]. Zhai et al. [418] also built on the insights of [310] to develop a method for MAXimizing the CERtified Radius (MACER) of the models that is claimed to be scalable to large models.

Fischer et al. [419] also extended the notion of randomized smoothing to incorporate parameterized transformations (e.g., translations, rotations) and certified the robustness of models in parameter space (e.g., rotation angle). Another example of using randomised smoothing for a certifiable defense can be found in [420]. This defense is aimed at patch attacks, and it provides certificate against given image and patch size. For patch attacks, more certified defenses are studied in [45], [421], [422]. Zhang et al. [423] extended the Gaussian

smoothing noise in randomized classifiers to non-Gaussian noise. They designed a family of non-Gaussian smoothing distributions that works more efficiently against ℓ_1 , ℓ_2 , and ℓ_∞ attacks.

As noted earlier, the direction of certified defenses is gradually becoming quite popular in adversarial machine learning literature. Incidentally, the problem is attracting more interest of machine learning community as compared to the compute vision community. Nevertheless, inspirations for the techniques developed along this line of research are coming from different directions. For instance, Rahnama et al. [424] treat the networks from a control theory perspective to provide tight bounds on any layer’s response to adversarial examples. Generally, randomised smoothing is the most commonly utilised tool be certified defense tools, which relates to adversarial in essence. Further recent examples of certified defenses can be found in [425], [426], [427], [428], [429].

E. Miscellaneous methods

Among defenses, there are numerous works that either propose methods for specialized tasks, networks or attack types. There are also techniques that mainly focus on improving the defense strength by combining multiple defense strategies discussed above. This section provides a summary of such works in the recent literature.

Cemgil et al. [430] analyzed the susceptibility of Variational Auto-Encoders (VAEs) to adversarial examples. They identified ‘evidence lower bound’ as one of its major causes, which is addressed by a data augmentation strategy during training in their work. He et al. [431] specifically proposed a binarization aware training method to defend against the Bit Flip Attack [175]. Robustness of Bayesian networks to gradient-based attacks is studied in [432]. Similarly, inherent robustness of spiking neural networks is the main topic of discussion in [433]. Defending Graph Neural Networks (GNNs) is studied in [434]. Among other specialized defenses, differential privacy is used in [435] to detect poisoning samples for backdoor attacks. There are also methods that focus entirely on defending neural models against the universal adversarial perturbations [436], [437] Cost sensitive adversarial robustness is studied in [438], whereas a so-called ‘guided complement entropy’ loss is proposed in [439], claiming better robustness over the standard cross entropy loss.

There are also examples providing specialized defenses for computer vision tasks other than standard classification, e.g. tracking [440], open-set recognition [441], face recognition [442]. Goldblum et al. [443] proposed a method to infer robust models for few-shot classification tasks based on adversarially robust meta-learners. A prediction poisoning attack is adopted as a defense against the model stealing attacks in [444]. The technique systematically alters the prediction of a target model to maintain the original performance but poison any model trained to steal the target model. A similar approach is taken in [445] by selectively making incorrect prediction for out-of-distribution queries to avoid model stealing.

Methods analyzing defense mechanisms and robustness instead of proposing new specialized defenses are also found

for this category of defenses. For instance, [446] analyzes the robustness of sparse coding to adversarial examples. It is observed in [447] that multitask learning generally results in improving adversarial robustness of the models. Kim et al. [448] claim that by leveraging sparsity and other perceptual biological mechanisms, adversarial robustness of models can be improved. Wang et al. [449] studied how to calibrate a trained model in-situ, in order to analyze the achievable trade-offs between the standard and robust accuracy of the model. Trade-off between the backdoor and adversarial robustness of models is studied in [450]. Chen et al. [451] proposed to use Neural Architecture Search to find adversarially robust architectures. Adversarial robustness in the more practical scenario of long-tailed data distribution is analyzed in [452]. The authors combine adversarial training with the existing recognition methods for imbalanced and long-tailed data to highlight interesting properties of models. For instance, it is shown that unreliable evaluation can easily give fake robustness gain impression for these models.

XI. DISCUSSION

Since its advent in 2013, the problem of adversarial attacks and lack of defenses for deep learning has intrigued the computer vision community considerably. Currently, this research direction is more active than ever. We found an ever-increasing number of papers appearing in the leading research sources of computer vision. The mainstream venues of machine learning research are also publishing papers with almost the same frequency as the computer vision venues. Interestingly, we found that most works in the direction of adversarial attacks and defenses appearing in machine learning sources still use ‘visual models’ as their test-bed. Nevertheless, we find a particular interest of the machine learning community in robustification of the models instead of devising new methods of fooling them. Of 400+ papers identified among the top six computer vision and machine learning venues in the last three years, we find around 74% papers dealing with defense techniques in machine learning venues. In contrast, only 40% of the papers in computer vision venues make adversarial defense as their central topic.

Among many interesting sub-problems in this area, the problem of ‘black-box’ attacks under better transferability and query-based setup is gaining significant popularity in computer vision research sources. In parallel, the topics of ‘adversarial training’ and ‘provable/certified’ defenses currently stand out in the literature appearing in machine learning sources. Below we summarize a few general trends that we observed in the literature. We intentionally keep the discussion at a higher-level of abstraction while covering the broader direction. The reader is encouraged to visit the related sections of the article to observe these trends with specific instances.

A. General Trends and Challenges

a) Adversarial Attacks: Whereas the first generation of attacks explored new core tools to fool deep visual classifiers, the more recent attacks are concerned with utilizing those tools for more specific fooling objectives. Gradient ascend over

the loss surface of the model is arguably the most common (and effective) tool for adversarial attacks in the literature. An overwhelming majority of the existing white-box attacks and transfer-based black-box attacks use this tool in some form to compute the additive adversarial perturbations. Model gradients are sometimes also utilized to satisfy linearization assumptions used by the attacks that aim at exploiting the geometry of the classification regions of the model. The observation that model gradients are the central tool for adversarial attacks resounds with the fact that deep models are, in the end, differentiable programs. Nevertheless, there are also other tools and heuristics, e.g. evolutionary algorithms, color-space search, that have been shown to find effective adversarial examples. As compared to model-gradients, such techniques are found to be more ad-hoc though.

The more recent core attack methods often aim at making the attacks more threatening by further reducing the norm of the perturbations and amplifying the transferability of the adversarial examples in black-box setups. Although universal perturbations can be considered a more serious threat from a practical viewpoint, the vast majority of the existing literature (> 95%) is concerned with image-specific attacks. A major reason for that is, from the defense perspective, securing models against (stronger) image-specific attacks already provides some robustness against the universal attacks, because the adversarial objective of the latter is already more challenging than that of the former. Nevertheless, we still find active investigations related to specifically securing models against the universal perturbations.

Since black-box attacks are gaining considerable popularity in the recent literature, it is worth summarizing some of the trends specifically in this direction. For the transfer-based black-box attacks, currently an accuracy reduction (of the target model) in the range 40 – 50% with ℓ_∞ perturbation norm of 15/255 is generally considered a good achievement in the recent literature for ImageNet models. This is true only for untargeted fooling though. The norm-bound is often considerably relaxed for the targeted black-box fooling (e.g. up to 32/255) without achieving fooling ratios at par with untargeted fooling with half the perturbation norm. We also observe that black-box attacks are reported to transfer better between the models with architectural similarity. For instance, one can expect to see an accuracy reduction of $\sim 50\%$ when transferring perturbations computed on inception-V3 to inception-V4. This number is expected to be $\sim 25\%$ when those perturbations are transferred to a ResNet-50 model. These numbers are not hypothetical. We provide them by observing multiple contributions. However, since each attack method has its own specific algorithm, the exact transfer rates may vary. We intentionally do not associate these numbers to specific methods, and only provide a rough estimate as a general guide to the readers.

One surprising trend we observed in the literature is about the evaluation of transfer-based black-box attacks. The term ‘black-box’ is understood by the community as a setup where the attacker does not have ‘any’ information about the target model (except its output in query-based setup). However, the existing methods generally report the attack transfer rates on

‘ImageNet’ models while also computing the perturbations on the ‘ImageNet’ models. In essence, this setup entails complete knowledge of the training data of the target model, which violates the definition of ‘black-box’ setup. Strictly speaking, the target models must be trained on unseen data, and should have, e.g. unknown number of output labels. We suggest that the research community considers this aspect in evaluating the transfer-based black-box attacks.

Among the query-based black-box attacks, the boundary attacks are more popular - outweighing their score-based counterpart by ~ 5 to 1. Generally, the query-based attacks optimize for two contradictory objectives of (a) achieving high fooling rates with stronger perturbations that use minimum number of queries, (b) keeping the perturbations imperceptible by restricting their norm. The most widely used strategy is to first query the black-box model with large perturbations, and then reduce the perturbation norm with a refinement mechanism while maintaining the incorrect prediction. We witness a large variation in the achieved fooling ratios and the number of queries utilized by different methods in this direction. It is clear from the reported results that these values depend rather strongly on image size. For ImageNet sized images, current literature considers 20K to 100K+ queries per image to still be reasonable to achieve imperceptible perturbations. This number drastically reduces to ~ 1 K for image sizes of 32×32 .

From the perspective of threat of adversarial attacks in the physical world, we do not find research to be as active as in the digital domain. One reason for that is the processes involved in physically realizing the computed adversarial patterns are often cumbersome and time taking. This does not mix well with the extremely fast pace of this research direction in the digital domain. Hence, even the attacks devised for the physical world applications are often just evaluated in simulated environments, e.g. camouflage cars for autonomous driving. Whereas we did not find any convincing argument that could suggest that physical world attacks are not a real concern to vision systems, we do find the adversarial samples for the physical world to be more conspicuous. Generally, such samples can be marked by obvious unnatural/irrelevant geometry or texture. Such a compromise over the stealthiness of the attack directly comes from the fact that visual sensors digitize only the ‘visible’ information. Hence, the adversarial patterns have to be visible to the sensor. Hiding the adversarial patterns from humans by semantically blending them in the scene environment is then the obvious choice for imperceptibility of the attacks in the physical world.

One interesting emerging utility of adversarial attacks is in explaining deep visual models. Considering that model gradients are utilized by both attack methods and popular model explanation methods, e.g., Grad-CAM [102], CAM-ERAS [453], it is not surprising that this overlap is emerging. Since deep learning models are differentiable programs, one can expect adversarial perturbations (which are a processed form of gradient information) can carry a signature of that program. From another perspective, perturbations can also be expected to focus (in some form) more on the salient regions of object to take model’s attention away from those regions. This

notion also resounds with image saliency. Hence, researchers are also getting interested in explainability of the perturbation itself.

Beyond the above-mentioned trends and challenges, the recent attacks are gradually getting more and more specialized to specific vision tasks and data modalities. Nevertheless, since the discovery of adversarial perturbations, their theoretical understanding has always been a topic of debate. Though a number of hypotheses exist in the literature on the susceptibility of deep learning to adversarial attacks, there is no single theory to fully explain all the observed phenomena in this direction. The adversarial vulnerability of deep models seem to emerge from a number of processes, and the debate on its existence and theory to explain all its aspect can be expected to stay as a long-standing problem for this research direction.

b) Adversarial defenses: Whereas a large number of defenses against adversarial attacks are appearing in the literature, arguably the most promising stream of works still concerns itself with ‘adversarial training’. Interestingly, the concept of adversarial training was presented simultaneously with adversarial perturbations in the original work of Szegedy et al. [1]. Most of the later literature significantly digressed from this original idea of robustifying the models. However, the later defense strategies mostly rely on ad-hoc rules and heuristics. Many of those are also shown to be broken with stronger attacks or different attack conditions [454]. In fact, recently, Tramer et al. [157] also show that thirteen different defenses that actually account for adaptive attack strategies can also be broken. From the defense perspective, the research community (especially machine learning community) is focusing more on adversarial training and certified defenses due to their principled nature. Nevertheless, reduction in the accuracy of the robust models on clean images is a major challenge at this front. It is easy to observe in the literature that methods withstanding stronger attacks have proportionally low accuracy on clean images.

c) Future outlook: Considering an ever-increasing influx of research papers in adversarial attacks (and defenses) since the advent of this direction, we can easily predict high research activity in this direction in the near-future. From the attacks perspective, whereas white-box attacks are likely to keep building around the tools used by the first-generation attack methods, a variety of new techniques for black-box setup can be anticipated. This is especially true for query-based attacks that is gaining increasing interest of the research community. Naturally, we can also anticipate the attacks to soon circumscribe Transformer models in vision, which are gaining popularity in computer vision community [455].

Based on the existing literature, we can argue that topics like understanding the existence of adversarial examples, intrinsically robust models, robustness-accuracy trade-off, adversarial training, certified defenses; are gradually adapting into long-standing problems of this direction. Hence, we also expect a multitude of works directed to address these problems in the future. We are likely to see mergence of adversarial perturbation techniques with other related directions, e.g. deep-Fakes [249], backdoor attacks [456]. Specifically, a potential interesting scenario is adopting the adversarial objective of

perturbations to independently fool the detectors of deepFakes and backdoor attacks. We are also likely to witness more activity in terms of expansion of adversarial attacks through visual models to multi-model tasks, e.g. image/video captioning [228] which combines visual models with language models, providing the opportunity to control the latter by attacking the former.

Since adversarial examples question the core utility of deep learning of making ‘reliable’ automated decisions, the research direction of adversarial attacks (and their defenses) seems to be here to stay with deep learning research. Just like deep learning is finding utilities in all applications, adversarial attacks are gradually adapting to those applications as its nemesis. From the viewpoint of this research direction, this arm race is promising, but not so much for deep learning in practice.

XII. CONCLUSION

In this article, we reviewed the research direction of adversarial attacks and defenses for deep learning models, focusing on the visual models. Since its advent in 2013, this direction has particularly intrigued the computer vision community, which has led to a large influx of papers in the recent years. To ensure the authenticity and quality of the discussed contributions, the survey mainly focused on the papers published in the top-ranked sources of computer vision and machine learning research. For standardising technical terminologies in this relatively new research direction, the survey also provided a list of definitions of the frequently used terms in the related literature. It also presented a detailed discussion on the early contributions in adversarial attacks to provide a historical account of the overall direction. The presented review builds on the first-ever peer-reviewed survey in this direction [2] - co-authored by the authors of this survey - as a legacy sequel. In [2], literature until 2018 is covered thoroughly. Hence, this article focused on the more recent literature, published after 2018. The covered literature is divided into attacks and defenses methods, which are further broken down into sub-topics by clustering the papers. This provided a clear indication of the current and emerging trends in the literature, that we discussed and reflected upon explicitly after reviewing the literature.

ACKNOWLEDGMENT

This material is based upon work supported by the Defense Advanced Research Projects Agency (DARPA) under Agreement No. HR00112090095, and the Australian Research Council Discovery Grant DP190102443. Dr. Naveed Akhtar is the recipient of an Office of National Intelligence National Intelligence Postdoctoral Grant funded by the Australian Government.

REFERENCES

- [1] C. Szegedy, W. Zaremba, I. Sutskever, J. Bruna, D. Erhan, I. Goodfellow, and R. Fergus, “Intriguing properties of neural networks,” *arXiv preprint arXiv:1312.6199*, 2013.
- [2] N. Akhtar and A. Mian, “Threat of adversarial attacks on deep learning in computer vision: A survey,” *IEEE Access*, vol. 6, pp. 14410–14430, 2018.
- [3] Y. LeCun, Y. Bengio, and G. Hinton, “Deep learning,” *nature*, vol. 521, no. 7553, pp. 436–444, 2015.
- [4] H. Y. Xiong, B. Alipanahi, L. J. Lee, H. Bretschneider, D. Merico, R. K. Yuen, Y. Hua, S. Gueroussov, H. S. Najafabadi, T. R. Hughes *et al.*, “The human splicing code reveals new insights into the genetic determinants of disease,” *Science*, vol. 347, no. 6218, 2015.
- [5] M. Helmstaedter, K. L. Briggman, S. C. Turaga, V. Jain, H. S. Seung, and W. Denk, “Connectomic reconstruction of the inner plexiform layer in the mouse retina,” *Nature*, vol. 500, no. 7461, pp. 168–174, 2013.
- [6] M. Amodio, D. Van Dijk, K. Srinivasan, W. S. Chen, H. Mohsen, K. R. Moon, A. Campbell, Y. Zhao, X. Wang, M. Venkataswamy *et al.*, “Exploring single-cell data with deep multitasking neural networks,” *Nature methods*, pp. 1–7, 2019.
- [7] G. Hickok and D. Poeppel, “The cortical organization of speech processing,” *Nature reviews neuroscience*, vol. 8, no. 5, pp. 393–402, 2007.
- [8] C. Manning and H. Schütze, *Foundations of statistical natural language processing*. MIT press, 1999.
- [9] A. Krizhevsky, I. Sutskever, and G. E. Hinton, “Imagenet classification with deep convolutional neural networks,” in *Advances in neural information processing systems*, 2012, pp. 1097–1105.
- [10] 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.
- [11] Y. LeCun, B. Boser, J. S. Denker, D. Henderson, R. E. Howard, W. Hubbard, and L. D. Jackel, “Backpropagation applied to handwritten zip code recognition,” *Neural computation*, vol. 1, no. 4, pp. 541–551, 1989.
- [12] K. He, X. Zhang, S. Ren, and J. Sun, “Deep residual learning for image recognition,” in *Proceedings of the IEEE conference on computer vision and pattern recognition*, 2016, pp. 770–778.
- [13] G. Huang, Z. Liu, L. Van Der Maaten, and K. Q. Weinberger, “Densely connected convolutional networks,” in *Proceedings of the IEEE conference on computer vision and pattern recognition*, 2017, pp. 4700–4708.
- [14] C. Szegedy, S. Ioffe, V. Vanhoucke, and A. Alemi, “Inception-v4, inception-resnet and the impact of residual connections on learning,” *arXiv preprint arXiv:1602.07261*, 2016.
- [15] D. Silver, J. Schrittwieser, K. Simonyan, I. Antonoglou, A. Huang, A. Guez, T. Hubert, L. Baker, M. Lai, A. Bolton *et al.*, “Mastering the game of go without human knowledge,” *nature*, vol. 550, no. 7676, pp. 354–359, 2017.
- [16] C. Szegedy, W. Liu, Y. Jia, P. Sermanet, S. Reed, D. Anguelov, D. Erhan, V. Vanhoucke, and A. Rabinovich, “Going deeper with convolutions,” in *Proceedings of the IEEE conference on computer vision and pattern recognition*, 2015, pp. 1–9.
- [17] I. J. Goodfellow, J. Shlens, and C. Szegedy, “Explaining and harnessing adversarial examples,” *arXiv preprint arXiv:1412.6572*, 2014.
- [18] Tesla, “Future of driving,” 2020 (accessed August 25, 2020), https://www.tesla.com/en_AU/autopilot.
- [19] C. Middlehurst, “China unveils world’s first facial recognition atm,” 2020 (accessed August 25, 2020), <https://www.telegraph.co.uk/news/worldnews/asia/china/11643314/China-unveils-worlds-first-facial-recognition-ATM.html>.
- [20] Apple, “About face id advanced technology,” 2020 (accessed August 25, 2020), <https://support.apple.com/en-au/HT208108>.
- [21] S. Grigorescu, B. Trasnea, T. Cocias, and G. Macesanu, “A survey of deep learning techniques for autonomous driving,” *Journal of Field Robotics*, vol. 37, no. 3, pp. 362–386, 2020.
- [22] S. Zulfarnain Gilani and A. Mian, “Learning from millions of 3d scans for large-scale 3d face recognition,” in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2018, pp. 1896–1905.
- [23] I. Masi, Y. Wu, T. Hassner, and P. Natarajan, “Deep face recognition: A survey,” in *2018 31st SIBGRAPI conference on graphics, patterns and images (SIBGRAPI)*. IEEE, 2018, pp. 471–478.
- [24] N. Sünderhauf, O. Brock, W. Scheirer, R. Hadsell, D. Fox, J. Leitner, B. Upcroft, P. Abbeel, W. Burgard, M. Milford *et al.*, “The limits and potentials of deep learning for robotics,” *The International Journal of Robotics Research*, vol. 37, no. 4-5, pp. 405–420, 2018.
- [25] M. M. Najafabadi, F. Villanustre, T. M. Khoshgoftaar, N. Seliya, R. Wald, and E. Muharemagic, “Deep learning applications and challenges in big data analytics,” *Journal of Big Data*, vol. 2, no. 1, p. 1, 2015.
- [26] N. Carlini, *A Complete List of All (arXiv) Adversarial Example Papers*, 2020 (accessed October 1, 2020). [Online]. Available: <https://nicholas.carlini.com/writing/2019/all-adversarial-example-papers.html>

- [27] A. Arnab, O. Miksik, and P. H. Torr, "On the robustness of semantic segmentation models to adversarial attacks," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2018, pp. 888–897.
- [28] Y. He, S. Rahimian, B. Schiele, and M. Fritz, "Segmentations-leak: Membership inference attacks and defenses in semantic image segmentation," *arXiv preprint arXiv:1912.09685*, 2019.
- [29] J. Tu, M. Ren, S. Manivasagam, M. Liang, B. Yang, R. Du, F. Cheng, and R. Urtasun, "Physically realizable adversarial examples for lidar object detection," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2020, pp. 13 716–13 725.
- [30] H. Zhang and J. Wang, "Towards adversarially robust object detection," in *Proceedings of the IEEE International Conference on Computer Vision*, 2019, pp. 421–430.
- [31] Y. Jia, Y. Lu, J. Shen, Q. A. Chen, H. Chen, Z. Zhong, and T. Wei, "Fooling detection alone is not enough: Adversarial attack against multiple object tracking," in *International Conference on Learning Representations*, 2019.
- [32] X. Chen, X. Yan, F. Zheng, Y. Jiang, S.-T. Xia, Y. Zhao, and R. Ji, "One-shot adversarial attacks on visual tracking with dual attention," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2020, pp. 10 176–10 185.
- [33] H. Zheng, Z. Zhang, J. Gu, H. Lee, and A. Prakash, "Efficient adversarial training with transferable adversarial examples," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2020, pp. 1181–1190.
- [34] W. Zhou, X. Hou, Y. Chen, M. Tang, X. Huang, X. Gan, and Y. Yang, "Transferable adversarial perturbations," in *Proceedings of the European Conference on Computer Vision (ECCV)*, 2018, pp. 452–467.
- [35] S.-M. Moosavi-Dezfooli, A. Fawzi, O. Fawzi, and P. Frossard, "Universal adversarial perturbations," in *Proceedings of the IEEE conference on computer vision and pattern recognition*, 2017, pp. 1765–1773.
- [36] N. Akhtar, M. A. Jalwana, M. Bennamoun, and A. Mian, "Label universal targeted attack," *arXiv preprint arXiv:1905.11544*, 2019.
- [37] O. Vinyals, I. Babuschkin, W. M. Czarnecki, M. Mathieu, A. Dudzik, J. Chung, D. H. Choi, R. Powell, T. Ewalds, P. Georgiev *et al.*, "Grandmaster level in starcraft ii using multi-agent reinforcement learning," *Nature*, vol. 575, no. 7782, pp. 350–354, 2019.
- [38] 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.
- [39] H. X. Y. M. Hao-Chen, L. D. Deb, H. L. J.-L. T. Anil, and K. Jain, "Adversarial attacks and defenses in images, graphs and text: A review," *International Journal of Automation and Computing*, vol. 17, no. 2, pp. 151–178, 2020.
- [40] M. Ozdag, "Adversarial attacks and defenses against deep neural networks: a survey," *Procedia Computer Science*, vol. 140, pp. 152–161, 2018.
- [41] Y. Zhou, M. Han, L. Liu, J. He, and X. Gao, "The adversarial attacks threats on computer vision: A survey," in *2019 IEEE 16th International Conference on Mobile Ad Hoc and Sensor Systems Workshops (MASSW)*. IEEE, 2019, pp. 25–30.
- [42] F. Vakhshiteh, R. Ramachandra, and A. Nickabadi, "Threat of adversarial attacks on face recognition: A comprehensive survey," *arXiv preprint arXiv:2007.11709*, 2020.
- [43] T. B. Brown, N. Carlini, C. Zhang, C. Olsson, P. Christiano, and I. Goodfellow, "Unrestricted adversarial examples," *arXiv preprint arXiv:1809.08352*, 2018.
- [44] Y. Song, R. Shu, N. Kushman, and S. Ermon, "Constructing unrestricted adversarial examples with generative models," in *Advances in Neural Information Processing Systems*, 2018, pp. 8312–8323.
- [45] T. B. Brown, D. Mané, A. Roy, M. Abadi, and J. Gilmer, "Adversarial patch," *arXiv preprint arXiv:1712.09665*, 2017.
- [46] R. Fletcher, *Practical methods of optimization*. John Wiley & Sons, 2013.
- [47] A. Rozsa, E. M. Rudd, and T. E. Boulton, "Adversarial diversity and hard positive generation," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition Workshops*, 2016, pp. 25–32.
- [48] T. Miyato, S.-i. Maeda, M. Koyama, and S. Ishii, "Virtual adversarial training: a regularization method for supervised and semi-supervised learning," *IEEE transactions on pattern analysis and machine intelligence*, vol. 41, no. 8, pp. 1979–1993, 2018.
- [49] A. Kurakin, I. Goodfellow, and S. Bengio, "Adversarial examples in the physical world," *arXiv preprint arXiv:1607.02533*, 2016.
- [50] Y. Dong, F. Liao, T. Pang, H. Su, J. Zhu, X. Hu, and J. Li, "Boosting adversarial attacks with momentum," in *Proceedings of the IEEE conference on computer vision and pattern recognition*, 2018, pp. 9185–9193.
- [51] C. Xie, Z. Zhang, Y. Zhou, S. Bai, J. Wang, Z. Ren, and A. L. Yuille, "Improving transferability of adversarial examples with input diversity," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2019, pp. 2730–2739.
- [52] 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.
- [53] A. Kurakin, I. Goodfellow, and S. Bengio, "Adversarial machine learning at scale," *arXiv preprint arXiv:1611.01236*, 2016.
- [54] A. Madry, A. Makelov, L. Schmidt, D. Tsipras, and A. Vladu, "Towards deep learning models resistant to adversarial attacks," *arXiv preprint arXiv:1706.06083*, 2017.
- [55] N. Papernot, P. McDaniel, S. Jha, M. Fredrikson, Z. B. Celik, and A. Swami, "The limitations of deep learning in adversarial settings," in *2016 IEEE European symposium on security and privacy (EuroS&P)*. IEEE, 2016, pp. 372–387.
- [56] 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.
- [57] K. Simonyan, A. Vedaldi, and A. Zisserman, "Deep inside convolutional networks: Visualising image classification models and saliency maps," *arXiv preprint arXiv:1312.6034*, 2013.
- [58] S. Das and P. N. Suganthan, "Differential evolution: A survey of the state-of-the-art," *IEEE transactions on evolutionary computation*, vol. 15, no. 1, pp. 4–31, 2010.
- [59] S.-M. Moosavi-Dezfooli, A. Fawzi, and P. Frossard, "Deepfool: a simple and accurate method to fool deep neural networks," in *Proceedings of the IEEE conference on computer vision and pattern recognition*, 2016, pp. 2574–2582.
- [60] N. Papernot, P. McDaniel, X. Wu, S. Jha, and A. Swami, "Distillation as a defense to adversarial perturbations against deep neural networks," in *2016 IEEE Symposium on Security and Privacy (SP)*. IEEE, 2016, pp. 582–597.
- [61] G. Hinton, O. Vinyals, and J. Dean, "Distilling the knowledge in a neural network," *arXiv preprint arXiv:1503.02531*, 2015.
- [62] 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.
- [63] —, "Adversarial examples are not easily detected: Bypassing ten detection methods," in *Proceedings of the 10th ACM workshop on artificial intelligence and security*, 2017, pp. 3–14.
- [64] X. Dong, J. Han, D. Chen, J. Liu, H. Bian, Z. Ma, H. Li, X. Wang, W. Zhang, and N. Yu, "Robust superpixel-guided attentional adversarial attack," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2020, pp. 12 895–12 904.
- [65] Y. Guo, Q. Li, and H. Chen, "Backpropagating linearly improves transferability of adversarial examples," *NeurIPS*, 2020.
- [66] X. Dong, D. Chen, J. Bao, C. Qin, L. Yuan, W. Zhang, N. Yu, and D. Chen, "Greedyfool: Distortion-aware sparse adversarial attack," *arXiv preprint arXiv:2010.13773*, 2020.
- [67] G. Sriraman, S. Addepalli, A. Baburaj, and R. V. Babu, "Guided adversarial attack for evaluating and enhancing adversarial defenses," *NeurIPS*, 2020.
- [68] Y. Tashiro, Y. Song, and S. Ermon, "Diversity can be transferred: Output diversification for white-and black-box attacks," *Advances in Neural Information Processing Systems*, vol. 33, 2020.
- [69] J. Rony, L. G. Hafemann, L. S. Oliveira, I. B. Ayed, R. Sabourin, and E. Granger, "Decoupling direction and norm for efficient gradient-based l2 adversarial attacks and defenses," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2019, pp. 4322–4330.
- [70] Z. Yao, A. Gholami, P. Xu, K. Keutzer, and M. W. Mahoney, "Trust region based adversarial attack on neural networks," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2019, pp. 11 350–11 359.
- [71] A. R. Conn, N. I. Gould, and P. L. Toint, *Trust region methods*. SIAM, 2000.
- [72] B. Phan, F. Mannan, and F. Heide, "Adversarial imaging pipelines," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, June 2021, pp. 16 051–16 061.
- [73] A. Rahmati, S.-M. Moosavi-Dezfooli, P. Frossard, and H. Dai, "Geoda: a geometric framework for black-box adversarial attacks," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2020, pp. 8446–8455.

- [74] W. Brendel, J. Rauber, and M. Bethge, "Decision-based adversarial attacks: Reliable attacks against black-box machine learning models," *ICLR*, 2018.
- [75] J. Chen, M. I. Jordan, and M. J. Wainwright, "Hopskipjumpattack: A query-efficient decision-based attack," in *2020 IEEE Symposium on Security and Privacy (SP)*. IEEE, 2020, pp. 1277–1294.
- [76] Y. Liu, S.-M. Moosavi-Dezfooli, and P. Frossard, "A geometry-inspired decision-based attack," in *Proceedings of the IEEE International Conference on Computer Vision*, 2019, pp. 4890–4898.
- [77] Y. Shi, Y. Han, and Q. Tian, "Polishing decision-based adversarial noise with a customized sampling," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2020, pp. 1030–1038.
- [78] J. Du, H. Zhang, J. T. Zhou, Y. Yang, and J. Feng, "Query-efficient meta attack to deep neural networks," *ICLR*, 2020.
- [79] J. Li, R. Ji, H. Liu, J. Liu, B. Zhong, C. Deng, and Q. Tian, "Projection & probability-driven black-box attack," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2020, pp. 362–371.
- [80] H. Li, X. Xu, X. Zhang, S. Yang, and B. Li, "Qeba: Query-efficient boundary-based blackbox attack," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2020, pp. 1221–1230.
- [81] B. Ru, A. Cobb, A. Blaas, and Y. Gal, "Bayesopt adversarial attack," in *International Conference on Learning Representations*, 2019.
- [82] N. Akhtar and A. Mian, "Hyperspectral recovery from rgb images using gaussian processes," *IEEE transactions on pattern analysis and machine intelligence*, vol. 42, no. 1, pp. 100–113, 2018.
- [83] M. Cheng, S. Singh, P. Chen, P.-Y. Chen, S. Liu, and C.-J. Hsieh, "Sign-opt: A query-efficient hard-label adversarial attack," *ICLR*, 2020.
- [84] M. Cheng, T. Le, P.-Y. Chen, J. Yi, H. Zhang, and C.-J. Hsieh, "Query-efficient hard-label black-box attack: An optimization-based approach," *ICLR*, 2019.
- [85] S. Cheng, Y. Dong, T. Pang, H. Su, and J. Zhu, "Improving black-box adversarial attacks with a transfer-based prior," *arXiv preprint arXiv:1906.06919*, 2019.
- [86] Z. Huang and T. Zhang, "Black-box adversarial attack with transferable model-based embedding," *ICLR*, 2020.
- [87] A. Al-Dujaili and U.-M. O'Reilly, "Sign bits are all you need for black-box attacks," in *International Conference on Learning Representations*, 2019.
- [88] A. Ilyas, L. Engstrom, and A. Madry, "Prior convictions: Black-box adversarial attacks with bandits and priors," *ICLR*, 2019.
- [89] P. Zhao, S. Liu, P.-Y. Chen, N. Hoang, K. Xu, B. Kailkhura, and X. Lin, "On the design of black-box adversarial examples by leveraging gradient-free optimization and operator splitting method," in *Proceedings of the IEEE/CVF International Conference on Computer Vision*, 2019, pp. 121–130.
- [90] M. Andriushchenko, F. Croce, N. Flammarion, and M. Hein, "Square attack: a query-efficient black-box adversarial attack via random search," in *European Conference on Computer Vision*. Springer, 2020, pp. 484–501.
- [91] W. Chen, Z. Zhang, X. Hu, and B. Wu, "Boosting decision-based black-box adversarial attacks with random sign flip," in *European Conference on Computer Vision*. Springer, 2020, pp. 276–293.
- [92] S.-T. Xia and W. Guo, "Improving query efficiency of black-box adversarial attack," in *European Conference on Computer Vision*. Springer, 2020.
- [93] W. Wang, B. Yin, T. Yao, L. Zhang, Y. Fu, S. Ding, J. Li, F. Huang, and X. Xue, "Delving into data: Effectively substitute training for black-box attack," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, June 2021, pp. 4761–4770.
- [94] T. Brunner, F. Diehl, M. T. Le, and A. Knoll, "Guessing smart: Biased sampling for efficient black-box adversarial attacks," in *Proceedings of the IEEE/CVF International Conference on Computer Vision*, 2019, pp. 4958–4966.
- [95] G. Tolias, F. Radenovic, and O. Chum, "Targeted mismatch adversarial attack: Query with a flower to retrieve the tower," in *Proceedings of the IEEE/CVF International Conference on Computer Vision*, 2019, pp. 5037–5046.
- [96] H. M. Dolatabadi, S. Erfani, and C. Leckie, "Advflow: Inconspicuous black-box adversarial attacks using normalizing flows," *NeruIPS*, 2020.
- [97] T. Maho, T. Furon, and E. Le Merrer, "Surf-free: A fast surrogate-free black-box attack," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, June 2021, pp. 10430–10439.
- [98] X. Li, J. Li, Y. Chen, S. Ye, Y. He, S. Wang, H. Su, and H. Xue, "Qair: Practical query-efficient black-box attacks for image retrieval," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, June 2021, pp. 3330–3339.
- [99] J. Yang, Y. Jiang, X. Huang, B. Ni, and C. Zhao, "Learning black-box attackers with transferable priors and query feedback," *Advances in Neural Information Processing Systems*, vol. 33, 2020.
- [100] C. Ma, L. Chen, and J.-H. Yong, "Simulating unknown target models for query-efficient black-box attacks," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, June 2021, pp. 11835–11844.
- [101] W. Wu, Y. Su, X. Chen, S. Zhao, I. King, M. R. Lyu, and Y.-W. Tai, "Boosting the transferability of adversarial samples via attention," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2020, pp. 1161–1170.
- [102] R. R. Selvaraju, M. Cogswell, A. Das, R. Vedantam, D. Parikh, and D. Batra, "Grad-cam: Visual explanations from deep networks via gradient-based localization," in *Proceedings of the IEEE international conference on computer vision*, 2017, pp. 618–626.
- [103] Q. Li, Y. Guo, and H. Chen, "Yet another intermediate-level attack," in *European Conference on Computer Vision*. Springer, 2020, pp. 241–257.
- [104] Q. Huang, I. Katsman, H. He, Z. Gu, S. Belongie, and S.-N. Lim, "Enhancing adversarial example transferability with an intermediate level attack," in *Proceedings of the IEEE/CVF International Conference on Computer Vision*, 2019, pp. 4733–4742.
- [105] A. J. Bose, G. Gidel, H. Berrard, A. Cianflone, P. Vincent, S. Lacoste-Julien, and W. L. Hamilton, "Adversarial example games," *arXiv preprint arXiv:2007.00720*, 2020.
- [106] Q. Li, Y. Guo, and H. Chen, "Practical no-box adversarial attacks against dnns," *NeurIPS*, 2020.
- [107] J. Lin, C. Song, K. He, L. Wang, and J. E. Hopcroft, "Nesterov accelerated gradient and scale invariance for adversarial attacks," in *International Conference on Learning Representations*, 2019.
- [108] M. Li, C. Deng, T. Li, J. Yan, X. Gao, and H. Huang, "Towards transferable targeted attack," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2020, pp. 641–649.
- [109] Y. Nesterov, "A method for unconstrained convex minimization problem with the rate of convergence $O(1/k^2)$," in *Doklady an ussr*, vol. 269, 1983, pp. 543–547.
- [110] Y. Lu, Y. Jia, J. Wang, B. Li, W. Chai, L. Carin, and S. Velipasalar, "Enhancing cross-task black-box transferability of adversarial examples with dispersion reduction," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2020, pp. 940–949.
- [111] N. Inkawhich, W. Wen, H. H. Li, and Y. Chen, "Feature space perturbations yield more transferable adversarial examples," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2019, pp. 7066–7074.
- [112] N. Inkawhich, K. J. Liang, L. Carin, and Y. Chen, "Transferable perturbations of deep feature distributions," *ICLR*, 2020.
- [113] H. Wang, G. Wang, Y. Li, D. Zhang, and L. Lin, "Transferable, controllable, and inconspicuous adversarial attacks on person re-identification with deep mis-ranking," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2020, pp. 342–351.
- [114] Y. Li, S. Bai, C. Xie, Z. Liao, X. Shen, and A. L. Yuille, "Regional homogeneity: Towards learning transferable universal adversarial perturbations against defenses," *ECCV*, 2020.
- [115] M. Zhou, J. Wu, Y. Liu, S. Liu, and C. Zhu, "Dast: Data-free substitute training for adversarial attacks," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2020, pp. 234–243.
- [116] J. Zou, Z. Pan, J. Qiu, X. Liu, T. Rui, and W. Li, "Improving the transferability of adversarial examples with resized-diverse-inputs, diversity-ensemble and region fitting," in *European Conference on Computer Vision*. Springer, 2020, pp. 563–579.
- [117] Y. Dong, T. Pang, H. Su, and J. Zhu, "Evading defenses to transferable adversarial examples by translation-invariant attacks," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2019, pp. 4312–4321.
- [118] Y. Shi, S. Wang, and Y. Han, "Curls & whey: Boosting black-box adversarial attacks," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2019, pp. 6519–6527.
- [119] X. Wang and K. He, "Enhancing the transferability of adversarial attacks through variance tuning," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, June 2021, pp. 1924–1933.

- [120] W. Wu, Y. Su, M. R. Lyu, and I. King, "Improving the transferability of adversarial samples with adversarial transformations," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, June 2021, pp. 9024–9033.
- [121] M. Sharif, L. Bauer, and M. K. Reiter, "On the suitability of l_p -norms for creating and preventing adversarial examples," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition Workshops*, 2018, pp. 1605–1613.
- [122] H. Hosseini and R. Poovendran, "Semantic adversarial examples," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition Workshops*, 2018, pp. 1614–1619.
- [123] 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.
- [124] A. Joshi, A. Mukherjee, S. Sarkar, and C. Hegde, "Semantic adversarial attacks: Parametric transformations that fool deep classifiers," in *Proceedings of the IEEE International Conference on Computer Vision*, 2019, pp. 4773–4783.
- [125] 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.
- [126] F. Croce and M. Hein, "Sparse and imperceptible adversarial attacks," in *Proceedings of the IEEE International Conference on Computer Vision*, 2019, pp. 4724–4732.
- [127] E. Wong, F. R. Schmidt, and J. Z. Kolter, "Wasserstein adversarial examples via projected sinkhorn iterations," *arXiv preprint arXiv:1902.07906*, 2019.
- [128] A. Bhattad, M. J. Chong, K. Liang, B. Li, and D. A. Forsyth, "Unrestricted adversarial examples via semantic manipulation," *ICLR*, 2020.
- [129] A. S. Shamsabadi, R. Sanchez-Matilla, and A. Cavallaro, "Colorfool: Semantic adversarial colorization," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2020, pp. 1151–1160.
- [130] D. L. Ruderman, T. W. Cronin, and C.-C. Chiao, "Statistics of cone responses to natural images: implications for visual coding," *JOSA A*, vol. 15, no. 8, pp. 2036–2045, 1998.
- [131] Z. Zhao, Z. Liu, and M. Larson, "Towards large yet imperceptible adversarial image perturbations with perceptual color distance," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2020, pp. 1039–1048.
- [132] M. R. Luo, G. Cui, and B. Rigg, "The development of the cie 2000 colour-difference formula: Ciede2000," *Color Research & Application: Endorsed by Inter-Society Color Council, The Colour Group (Great Britain), Canadian Society for Color, Color Science Association of Japan, Dutch Society for the Study of Color, The Swedish Colour Centre Foundation, Colour Society of Australia, Centre Français de la Couleur*, vol. 26, no. 5, pp. 340–350, 2001.
- [133] H. Qiu, C. Xiao, L. Yang, X. Yan, H. Lee, and B. Li, "Semanticadv: Generating adversarial examples via attribute-conditioned image editing," in *European Conference on Computer Vision*. Springer, 2020, pp. 19–37.
- [134] P. Korshunov and S. Marcel, "Deepfakes: a new threat to face recognition? assessment and detection," *arXiv preprint arXiv:1812.08685*, 2018.
- [135] Z. Chen, L. Xie, S. Pang, Y. He, and B. Zhang, "Magdr: Mask-guided detection and reconstruction for defending deepfakes," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, June 2021, pp. 9014–9023.
- [136] D. Hendrycks, K. Zhao, S. Basart, J. Steinhardt, and D. Song, "Natural adversarial examples," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, June 2021, pp. 15 262–15 271.
- [137] Y. Li, B. Wu, Y. Jiang, Z. Li, and S.-T. Xia, "Backdoor learning: A survey," *arXiv preprint arXiv:2007.08745*, 2020.
- [138] Y. Liu, A. Mondal, A. Chakraborty, M. Zuzak, N. Jacobsen, D. Xing, and A. Srivastava, "A survey on neural trojans," in *2020 21st International Symposium on Quality Electronic Design (ISQED)*. IEEE, 2020, pp. 33–39.
- [139] Y. Liu, X. Ma, J. Bailey, and F. Lu, "Reflection backdoor: A natural backdoor attack on deep neural networks," in *European Conference on Computer Vision*. Springer, 2020, pp. 182–199.
- [140] A. Nguyen and A. Tran, "Input-aware dynamic backdoor attack," *arXiv preprint arXiv:2010.08138*, 2020.
- [141] C. Xie, K. Huang, P.-Y. Chen, and B. Li, "Dba: Distributed backdoor attacks against federated learning," in *International Conference on Learning Representations*, 2019.
- [142] V. Smith, C.-K. Chiang, M. Sanjabi, and A. S. Talwalkar, "Federated multi-task learning," in *Advances in neural information processing systems*, 2017, pp. 4424–4434.
- [143] A. N. Bhagoji, S. Chakraborty, P. Mittal, and S. Calo, "Analyzing federated learning through an adversarial lens," in *International Conference on Machine Learning*. PMLR, 2019, pp. 634–643.
- [144] E. Bagdasaryan, A. Veit, Y. Hua, D. Estrin, and V. Shmatikov, "How to backdoor federated learning," in *International Conference on Artificial Intelligence and Statistics*. PMLR, 2020, pp. 2938–2948.
- [145] J. Guo and C. Liu, "Practical poisoning attacks on neural networks," in *European Conference on Computer Vision*. Springer, 2020, pp. 142–158.
- [146] A. S. Rakin, Z. He, and D. Fan, "Tbt: Targeted neural network attack with bit trojan," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2020, pp. 13 198–13 207.
- [147] S. Zhao, X. Ma, X. Zheng, J. Bailey, J. Chen, and Y.-G. Jiang, "Clean-label backdoor attacks on video recognition models," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2020, pp. 14 443–14 452.
- [148] S. Kolouri, A. Saha, H. Pirsiavash, and H. Hoffmann, "Universal litmus patterns: Revealing backdoor attacks in cnns," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2020, pp. 301–310.
- [149] H. Wang, K. Sreenivasan, S. Rajput, H. Vishwakarma, S. Agarwal, J.-y. Sohn, K. Lee, and D. Papailiopoulos, "Attack of the tails: Yes, you really can backdoor federated learning," *arXiv preprint arXiv:2007.05084*, 2020.
- [150] M. Fredrikson, E. Lantz, S. Jha, S. Lin, D. Page, and T. Ristenpart, "Privacy in pharmacogenetics: An end-to-end case study of personalized warfarin dosing," in *23rd {USENIX} Security Symposium ({USENIX} Security 14)*, 2014, pp. 17–32.
- [151] X. Wu, M. Fredrikson, S. Jha, and J. F. Naughton, "A methodology for formalizing model-inversion attacks," in *2016 IEEE 29th Computer Security Foundations Symposium (CSF)*. IEEE, 2016, pp. 355–370.
- [152] S. Yeom, I. Giacomelli, M. Fredrikson, and S. Jha, "Privacy risk in machine learning: Analyzing the connection to overfitting," in *2018 IEEE 31st Computer Security Foundations Symposium (CSF)*. IEEE, 2018, pp. 268–282.
- [153] A. Nguyen, A. Dosovitskiy, J. Yosinski, T. Brox, and J. Clune, "Synthesizing the preferred inputs for neurons in neural networks via deep generator networks," *Advances in neural information processing systems*, vol. 29, pp. 3387–3395, 2016.
- [154] N. Akhtar, M. Jalwana, M. Bennamoun, and A. S. Mian, "Attack to fool and explain deep networks," *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 2021.
- [155] Y. Zhang, R. Jia, H. Pei, W. Wang, B. Li, and D. Song, "The secret revealer: generative model-inversion attacks against deep neural networks," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2020, pp. 253–261.
- [156] A. Athalye, N. Carlini, and D. Wagner, "Obfuscated gradients give a false sense of security: Circumventing defenses to adversarial examples," *ICML*, 2018.
- [157] F. Tramer, N. Carlini, W. Brendel, and A. Madry, "On adaptive attacks to adversarial example defenses," *NeurIPS*, 2020.
- [158] C. Xiao, P. Zhong, and C. Zheng, "Enhancing adversarial defense by k-winners-take-all," *arXiv preprint arXiv:1905.10510*, 2020.
- [159] K. Roth, Y. Kilcher, and T. Hofmann, "The odds are odd: A statistical test for detecting adversarial examples," in *International Conference on Machine Learning*. PMLR, 2019, pp. 5498–5507.
- [160] Y. Li, J. Bradshaw, and Y. Sharma, "Are generative classifiers more robust to adversarial attacks?" in *International Conference on Machine Learning*. PMLR, 2019, pp. 3804–3814.
- [161] M. Bafna, J. Murtagh, and N. Vyas, "Thwarting adversarial examples: An l_0 -robustsparse fourier transform," *NeurIPS*, 2018.
- [162] T. Pang, K. Xu, Y. Dong, C. Du, N. Chen, and J. Zhu, "Rethinking softmax cross-entropy loss for adversarial robustness," *ICLR*, 2020.
- [163] G. Verma and A. Swami, "Error correcting output codes improve probability estimation and adversarial robustness of deep neural networks," *Advances in Neural Information Processing Systems*, vol. 32, pp. 8646–8656, 2019.
- [164] T. Pang, K. Xu, C. Du, N. Chen, and J. Zhu, "Improving adversarial robustness via promoting ensemble diversity," in *International Conference on Machine Learning*. PMLR, 2019, pp. 4970–4979.

- [165] S. Sen, B. Ravindran, and A. Raghunathan, "Empir: Ensembles of mixed precision deep networks for increased robustness against adversarial attacks," *ICLR*, 2020.
- [166] Z. Yang, B. Li, P.-Y. Chen, and D. Song, "Characterizing audio adversarial examples using temporal dependency," *ICLR*, 2019.
- [167] T. Pang, K. Xu, and J. Zhu, "Mixup inference: Better exploiting mixup to defend adversarial attacks," *ICLR*, 2020.
- [168] Y. Yang, G. Zhang, D. Katabi, and Z. Xu, "Me-net: Towards effective adversarial robustness with matrix estimation," *ICLR*, 2019.
- [169] X. Yin, S. Kolouri, and G. K. Rohde, "Adversarial example detection and classification with asymmetrical adversarial training," *ICLR*, 2020.
- [170] T. Yu, S. Hu, C. Guo, W.-L. Chao, and K. Q. Weinberger, "A new defense against adversarial images: Turning a weakness into a strength," *NeurIPS*, 2019.
- [171] A. Ghiasi, A. Shafahi, and T. Goldstein, "Breaking certified defenses: Semantic adversarial examples with spoofed robustness certificates," *ICLR*, 2020.
- [172] M. Zhou, Z. Niu, L. Wang, Q. Zhang, and G. Hua, "Adversarial ranking attack and defense," *ECCV*, 2020.
- [173] K. Razavi, B. Gras, E. Bosman, B. Preneel, C. Giuffrida, and H. Bos, "Flip feng shui: Hammering a needle in the software stack," in *25th {USENIX} Security Symposium ({USENIX} Security 16)*, 2016, pp. 1–18.
- [174] S. Hong, P. Frigo, Y. Kaya, C. Giuffrida, and T. Dumitras, "Terminal brain damage: Exposing the graceless degradation in deep neural networks under hardware fault attacks," in *28th {USENIX} Security Symposium ({USENIX} Security 19)*, 2019, pp. 497–514.
- [175] A. S. Rakin, Z. He, and D. Fan, "Bit-flip attack: Crushing neural network with progressive bit search," in *Proceedings of the IEEE International Conference on Computer Vision*, 2019, pp. 1211–1220.
- [176] S. Rezaei and X. Liu, "A target-agnostic attack on deep models: Exploiting security vulnerabilities of transfer learning," *ICLR*, 2020.
- [177] R. Mor, E. Peterfreund, M. Gavish, and A. Globerson, "Optimal strategies against generative attacks," in *International Conference on Learning Representations*, 2019.
- [178] A. Ganeshan and R. V. Babu, "Fda: Feature disruptive attack," in *Proceedings of the IEEE/CVF International Conference on Computer Vision*, 2019, pp. 8069–8079.
- [179] H. Zhou, D. Chen, J. Liao, K. Chen, X. Dong, K. Liu, W. Zhang, G. Hua, and N. Yu, "Lg-gan: Label guided adversarial network for flexible targeted attack of point cloud based deep networks," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2020, pp. 10356–10365.
- [180] Y. Zhao, Y. Wu, C. Chen, and A. Lim, "On isometry robustness of deep 3d point cloud models under adversarial attacks," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2020, pp. 1201–1210.
- [181] A. Hamdi, S. Rojas, A. Thabet, and B. Ghanem, "Advpc: Transferable adversarial perturbations on 3d point clouds," in *European Conference on Computer Vision*. Springer, 2020, pp. 241–257.
- [182] H. Zhou, K. Chen, W. Zhang, H. Fang, W. Zhou, and N. Yu, "Dup-net: Denoiser and upsampler network for 3d adversarial point clouds defense," in *Proceedings of the IEEE/CVF International Conference on Computer Vision*, 2019, pp. 1961–1970.
- [183] M. Wicker and M. Kwiatkowska, "Robustness of 3d deep learning in an adversarial setting," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2019, pp. 11767–11775.
- [184] C. Xiang, C. R. Qi, and B. Li, "Generating 3d adversarial point clouds," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2019, pp. 9136–9144.
- [185] 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 Conference on Computer Vision and Pattern Recognition*, 2019, pp. 4302–4311.
- [186] H. Zhang, L. Zhu, Y. Zhu, and Y. Yang, "Motion-excited sampler: Video adversarial attack with sparked prior," *ECCV*, 2020.
- [187] A. Liu, T. Huang, X. Liu, Y. Xu, Y. Ma, X. Chen, S. J. Maybank, and D. Tao, "Spatiotemporal attacks for embodied agents," in *European Conference on Computer Vision*. Springer, 2020, pp. 122–138.
- [188] J. Liu, N. Akhtar, and A. Mian, "Adversarial attack on skeleton-based human action recognition," *IEEE Transactions on Neural Networks and Learning Systems*, 2020.
- [189] Y. Diao, T. Shao, Y.-L. Yang, K. Zhou, and H. Wang, "Basar:black-box attack on skeletal action recognition," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, June 2021, pp. 7597–7607.
- [190] H. Wang, F. He, Z. Peng, T. Shao, Y.-L. Yang, K. Zhou, and D. Hogg, "Understanding the robustness of skeleton-based action recognition under adversarial attack," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, June 2021, pp. 14656–14665.
- [191] R. Pony, I. Naei, and S. Mannor, "Over-the-air adversarial flickering attacks against video recognition networks," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, June 2021, pp. 515–524.
- [192] C. Xiao, R. Deng, B. Li, T. Lee, B. Edwards, J. Yi, D. Song, M. Liu, and I. Molloy, "Advit: Adversarial frames identifier based on temporal consistency in videos," in *Proceedings of the IEEE/CVF International Conference on Computer Vision*, 2019, pp. 3968–3977.
- [193] D. Zügner and S. Günnemann, "Adversarial attacks on graph neural networks via meta learning," *ICLR*, 2019.
- [194] H. Jin, Z. Shi, V. J. S. A. Peruri, and X. Zhang, "Certified robustness of graph convolution networks for graph classification under topological attacks," *Advances in Neural Information Processing Systems*, vol. 33, 2020.
- [195] E. B. Khalil, A. Gupta, and B. Dilkina, "Combinatorial attacks on binarized neural networks," *arXiv preprint arXiv:1810.03538*, 2018.
- [196] J. Lin, C. Gan, and S. Han, "Defensive quantization: When efficiency meets robustness," *arXiv preprint arXiv:1904.08444*, 2019.
- [197] R. Alaifari, G. S. Alberti, and T. Gauksson, "Adef: an iterative algorithm to construct adversarial deformations," *arXiv preprint arXiv:1804.07729*, 2018.
- [198] Y. Fan, B. Wu, T. Li, Y. Zhang, M. Li, Z. Li, and Y. Yang, "Sparse adversarial attack via perturbation factorization," in *Proceedings of European Conference on Computer Vision*, 2020.
- [199] H. Liu, R. Ji, J. Li, B. Zhang, Y. Gao, Y. Wu, and F. Huang, "Universal adversarial perturbation via prior driven uncertainty approximation," in *Proceedings of the IEEE/CVF International Conference on Computer Vision*, 2019, pp. 2941–2949.
- [200] S. U. Din, N. Akhtar, S. Younis, F. Shafait, A. Mansoor, and M. Shafique, "Steganographic universal adversarial perturbations," *Pattern Recognition Letters*, vol. 135, pp. 146–152, 2020.
- [201] A. Rampini, F. Pestarini, L. Cosmo, S. Melzi, and E. Rodola, "Universal spectral adversarial attacks for deformable shapes," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, June 2021, pp. 3216–3226.
- [202] Z. Zhong, W. Xu, Y. Jia, and T. Wei, "Perception deception: Physical adversarial attack challenges and tactics for dnn-based object detection," *Black Hat Europe*, 2018.
- [203] 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.
- [204] S. Sun, N. Akhtar, H. Song, A. Mian, and M. Shah, "Deep affinity network for multiple object tracking," *IEEE transactions on pattern analysis and machine intelligence*, vol. 43, no. 1, pp. 104–119, 2019.
- [205] P.-C. Chen, B.-H. Kung, and J.-C. Chen, "Class-aware robust adversarial training for object detection," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, June 2021, pp. 10420–10429.
- [206] B. Yan, D. Wang, H. Lu, and X. Yang, "Cooling-shrinking attack: Blinding the tracker with imperceptible noises," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2020, pp. 990–999.
- [207] B. Li, W. Wu, Q. Wang, F. Zhang, J. Xing, and J. Yan, "Siamrpn++: Evolution of siamese visual tracking with very deep networks," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2019, pp. 4282–4291.
- [208] S. Liang, X. Wei, S. Yao, and X. Cao, "Efficient adversarial attacks for visual object tracking," in *European Conference on Computer Vision*. Springer, 2020, pp. 34–50.
- [209] 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.
- [210] 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/CVF Conference on Computer Vision and Pattern Recognition*, 2020, pp. 720–729.
- [211] Z. Wu, S.-N. Lim, L. S. Davis, and T. Goldstein, "Making an invisibility cloak: Real world adversarial attacks on object detectors," in *European Conference on Computer Vision*. Springer, 2020, pp. 1–17.
- [212] A. Zolfi, M. Kravchik, Y. Elovici, and A. Shabtai, "The translucent patch: A physical and universal attack on object detectors," in *Pro-*

- ceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, June 2021, pp. 15 232–15 241.
- [213] Q. Guo, X. Xie, F. Juefei-Xu, L. Ma, Z. Li, W. Xue, W. Feng, and Y. Liu, “Spark: Spatial-aware online incremental attack against visual tracking,” in *Proceedings of the European Conference on Computer Vision (ECCV)*, vol. 2. Springer, 2020.
- [214] S. Jia, Y. Song, C. Ma, and X. Yang, “You attack: Towards temporally coherent black-box adversarial attack for visual object tracking,” in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, June 2021, pp. 6709–6718.
- [215] S. Huang, N. Papernot, I. Goodfellow, Y. Duan, and P. Abbeel, “Adversarial attacks on neural network policies,” *arXiv preprint arXiv:1702.02284*, 2017.
- [216] Y. Xiang, W. Niu, J. Liu, T. Chen, and Z. Han, “A pca-based model to predict adversarial examples on q-learning of path finding,” in *2018 IEEE Third International Conference on Data Science in Cyberspace (DSC)*. IEEE, 2018, pp. 773–780.
- [217] X. Bai, W. Niu, J. Liu, X. Gao, Y. Xiang, and J. Liu, “Adversarial examples construction towards white-box q table variation in dqn pathfinding training,” in *2018 IEEE Third International Conference on Data Science in Cyberspace (DSC)*. IEEE, 2018, pp. 781–787.
- [218] V. Mnih, K. Kavukcuoglu, D. Silver, A. Graves, I. Antonoglou, D. Wierstra, and M. Riedmiller, “Playing atari with deep reinforcement learning,” *arXiv preprint arXiv:1312.5602*, 2013.
- [219] T. Chen, W. Niu, Y. Xiang, X. Bai, J. Liu, Z. Han, and G. Li, “Gradient band-based adversarial training for generalized attack immunity of a3c path finding,” *arXiv preprint arXiv:1807.06752*, 2018.
- [220] V. Behzadan and A. Munir, “Vulnerability of deep reinforcement learning to policy induction attacks,” in *International Conference on Machine Learning and Data Mining in Pattern Recognition*. Springer, 2017, pp. 262–275.
- [221] Y.-C. Lin, Z.-W. Hong, Y.-H. Liao, M.-L. Shih, M.-Y. Liu, and M. Sun, “Tactics of adversarial attack on deep reinforcement learning agents,” *arXiv preprint arXiv:1703.06748*, 2017.
- [222] J. Liu, W. Niu, J. Liu, J. Zhao, T. Chen, Y. Yang, Y. Xiang, and L. Han, “A method to effectively detect vulnerabilities on path planning of vin,” in *International Conference on Information and Communications Security*. Springer, 2017, pp. 374–384.
- [223] T. Chen, J. Liu, Y. Xiang, W. Niu, E. Tong, and Z. Han, “Adversarial attack and defense in reinforcement learning-from ai security view,” *Cybersecurity*, vol. 2, no. 1, pp. 1–22, 2019.
- [224] A. Gleave, M. Dennis, C. Wild, N. Kant, S. Levine, and S. Russell, “Adversarial policies: Attacking deep reinforcement learning,” *ICLR*, 2020.
- [225] A. Rakhsha, G. Radanovic, R. Devidze, X. Zhu, and A. Singla, “Policy teaching via environment poisoning: Training-time adversarial attacks against reinforcement learning,” in *International Conference on Machine Learning*. PMLR, 2020, pp. 7974–7984.
- [226] X. Zhang, Y. Ma, A. Singla, and X. Zhu, “Adaptive reward-poisoning attacks against reinforcement learning,” in *International Conference on Machine Learning*. PMLR, 2020, pp. 11 225–11 234.
- [227] H. Zhang, H. Chen, C. Xiao, B. Li, M. Liu, D. Boning, and C.-J. Hsieh, “Robust deep reinforcement learning against adversarial perturbations on state observations,” *NeurIPS*, 2020.
- [228] N. Aafaq, A. Mian, W. Liu, S. Z. Gilani, and M. Shah, “Video description: A survey of methods, datasets, and evaluation metrics,” *ACM Computing Surveys (CSUR)*, vol. 52, no. 6, pp. 1–37, 2019.
- [229] X. Xu, X. Chen, C. Liu, A. Rohrbach, T. Darrell, and D. Song, “Fooing vision and language models despite localization and attention mechanism,” in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2018, pp. 4951–4961.
- [230] H. Chen, H. Zhang, P.-Y. Chen, J. Yi, and C.-J. Hsieh, “Attacking visual language grounding with adversarial examples: A case study on neural image captioning,” *arXiv preprint arXiv:1712.02051*, 2017.
- [231] O. Vinyals, A. Toshev, S. Bengio, and D. Erhan, “Show and tell: A neural image caption generator,” in *Proceedings of the IEEE conference on computer vision and pattern recognition*, 2015, pp. 3156–3164.
- [232] Y. Xu, B. Wu, F. Shen, Y. Fan, Y. Zhang, H. T. Shen, and W. Liu, “Exact adversarial attack to image captioning via structured output learning with latent variables,” in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2019, pp. 4135–4144.
- [233] X. Xu, J. Chen, J. Xiao, L. Gao, F. Shen, and H. T. Shen, “What machines see is not what they get: Fooing scene text recognition models with adversarial text images,” in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2020, pp. 12 304–12 314.
- [234] C. Zhu, Y. Cheng, Z. Gan, S. Sun, T. Goldstein, and J. Liu, “FreeLb: Enhanced adversarial training for natural language understanding,” in *International Conference on Learning Representations*, 2019.
- [235] J. Deng, J. Guo, N. Xue, and S. Zafeiriou, “Arcface: Additive angular margin loss for deep face recognition,” in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2019, pp. 4690–4699.
- [236] G. Goswami, N. Ratha, A. Agarwal, R. Singh, and M. Vatsa, “Unravelling robustness of deep learning based face recognition against adversarial attacks,” in *Proceedings of the AAAI Conference on Artificial Intelligence*, vol. 32, no. 1, 2018.
- [237] B. Amos, B. Ludwiczuk, J. Harkes, P. Pillai, K. Elgazzar, and M. Satyanarayanan, “Openface: Face recognition with deep neural networks,” in *IEEE Winter Conference on Applications of Computer Vision*, vol. 1, no. 2, 2016, p. 6.
- [238] O. M. Parkhi, A. Vedaldi, and A. Zisserman, “Deep face recognition,” in *Proceedings of the British Machine Vision Conference (BMVC)*. BMVA Press, September 2015, pp. 41.1–41.12.
- [239] Y. Dong, H. Su, B. Wu, Z. Li, W. Liu, T. Zhang, and J. Zhu, “Efficient decision-based black-box adversarial attacks on face recognition,” in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2019, pp. 7714–7722.
- [240] N. Hansen and A. Ostermeier, “Completely derandomized self-adaptation in evolution strategies,” *Evolutionary computation*, vol. 9, no. 2, pp. 159–195, 2001.
- [241] Y. Zhong and W. Deng, “Towards transferable adversarial attack against deep face recognition,” *IEEE Transactions on Information Forensics and Security*, vol. 16, pp. 1452–1466, 2020.
- [242] E. Chatzikyriakidis, C. Papaioannidis, and I. Pitas, “Adversarial face de-identification,” in *2019 IEEE International Conference on Image Processing (ICIP)*. IEEE, 2019, pp. 684–688.
- [243] H. Kwon, O. Kwon, H. Yoon, and K.-W. Park, “Face friend-safe adversarial example on face recognition system,” in *2019 Eleventh International Conference on Ubiquitous and Future Networks (ICUFN)*. IEEE, 2019, pp. 547–551.
- [244] A. Dabouei, S. Soleymani, J. Dawson, and N. Nasrabadi, “Fast geometrically-perturbed adversarial faces,” in *2019 IEEE Winter Conference on Applications of Computer Vision (WACV)*. IEEE, 2019, pp. 1979–1988.
- [245] 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 (CVPR)*, June 2021, pp. 11 845–11 854.
- [246] L. Yang, Q. Song, and Y. Wu, “Attacks on state-of-the-art face recognition using attentional adversarial attack generative network,” *Multimedia Tools and Applications*, pp. 1–21, 2020.
- [247] D. Deb, J. Zhang, and A. K. Jain, “Advfaces: Adversarial face synthesis,” in *2020 IEEE International Joint Conference on Biometrics (IJCB)*. IEEE, 2019, pp. 1–10.
- [248] D. Li, W. Wang, H. Fan, and J. Dong, “Exploring adversarial fake images on face manifold,” in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, June 2021, pp. 5789–5798.
- [249] R. Tolosana, R. Vera-Rodriguez, J. Fierrez, A. Morales, and J. Ortega-Garcia, “Deepfakes and beyond: A survey of face manipulation and fake detection,” *Information Fusion*, vol. 64, pp. 131–148, 2020.
- [250] 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.
- [251] Z. Zhou, D. Tang, X. Wang, W. Han, X. Liu, and K. Zhang, “Invisible mask: Practical attacks on face recognition with infrared,” *arXiv preprint arXiv:1803.04683*, 2018.
- [252] 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.
- [253] 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*, 2015, pp. 815–823.
- [254] S. Komkov and A. Petiushko, “Advhat: Real-world adversarial attack on arcfac face id system,” *arXiv preprint arXiv:1908.08705*, 2019.
- [255] M. Pautov, G. Melnikov, E. Kaziakhmedov, K. Kireev, and A. Petiushko, “On adversarial patches: real-world attack on arcfac-100 face recognition system,” in *2019 International Multi-Conference on*

- Engineering, Computer and Information Sciences (SIBIRCON)*. IEEE, 2019, pp. 0391–0396.
- [256] R. Shao, X. Lan, J. Li, and P. C. Yuen, “Multi-adversarial discriminative deep domain generalization for face presentation attack detection,” in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2019, pp. 10 023–10 031.
- [257] K. K. Nakka and M. Salzmann, “Indirect local attacks for context-aware semantic segmentation networks,” in *European Conference on Computer Vision*. Springer, 2020, pp. 611–628.
- [258] Y. He, S. Rahimian, B. Schiele, and M. Fritz, “Segmentations-leak: Membership inference attacks and defenses in semantic image segmentation,” in *European Conference on Computer Vision*. Springer, 2020, pp. 519–535.
- [259] J.-H. Choi, H. Zhang, J.-H. Kim, C.-J. Hsieh, and J.-S. Lee, “Evaluating robustness of deep image super-resolution against adversarial attacks,” in *Proceedings of the IEEE/CVF International Conference on Computer Vision*, 2019, pp. 303–311.
- [260] A. Mehra, B. Kailkhura, P.-Y. Chen, and J. Hamm, “How robust are randomized smoothing based defenses to data poisoning?” in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, June 2021, pp. 13 244–13 253.
- [261] A. Wong, S. Cicek, and S. Soatto, “Targeted adversarial perturbations for monocular depth prediction,” *arXiv preprint arXiv:2006.08602*, 2020.
- [262] J. Bai, B. Chen, Y. Li, D. Wu, W. Guo, S.-t. Xia, and E.-h. Yang, “Targeted attack for deep hashing based retrieval,” in *European Conference on Computer Vision*. Springer, 2020, pp. 618–634.
- [263] J. Li, R. Ji, H. Liu, X. Hong, Y. Gao, and Q. Tian, “Universal perturbation attack against image retrieval,” in *Proceedings of the IEEE/CVF International Conference on Computer Vision*, 2019, pp. 4899–4908.
- [264] Z. Zhang, Z. Zhang, Y. Zhou, Y. Shen, R. Jin, and D. Dou, “Adversarial attacks on deep graph matching,” *Advances in Neural Information Processing Systems*, vol. 33, 2020.
- [265] C. Yang, A. Kortylewski, C. Xie, Y. Cao, and A. Yuille, “Patchattack: A black-box texture-based attack with reinforcement learning,” in *European Conference on Computer Vision*. Springer, 2020, pp. 681–698.
- [266] X. Yang, F. Wei, H. Zhang, X. Ming, and J. Zhu, “Design and interpretation of universal adversarial patches in face detection,” *ECCV*, 2020.
- [267] A. Ranjan, J. Janai, A. Geiger, and M. J. Black, “Attacking optical flow,” in *Proceedings of the IEEE/CVF International Conference on Computer Vision*, 2019, pp. 2404–2413.
- [268] Y. Cao, C. Xiao, B. Cyr, Y. Zhou, W. Park, S. Rampazzi, Q. A. Chen, K. Fu, and Z. M. Mao, “Adversarial sensor attack on lidar-based perception in autonomous driving,” in *Proceedings of the 2019 ACM SIGSAC Conference on Computer and Communications Security*, 2019, pp. 2267–2281.
- [269] 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/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, June 2021, pp. 8565–8574.
- [270] 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*, 2019.
- [271] Z. Kong, J. Guo, A. Li, and C. Liu, “Physgan: Generating physical-world-resilient adversarial examples for autonomous driving,” in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2020, pp. 14 254–14 263.
- [272] C.-H. Ho, B. Leung, E. Sandstrom, Y. Chang, and N. Vasconcelos, “Catastrophic child’s play: easy to perform, hard to defend adversarial attacks,” in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2019, pp. 9229–9237.
- [273] 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/CVF Conference on Computer Vision and Pattern Recognition*, 2020, pp. 1000–1008.
- [274] Y. Jing, Y. Yang, Z. Feng, J. Ye, Y. Yu, and M. Song, “Neural style transfer: A review,” *IEEE transactions on visualization and computer graphics*, 2019.
- [275] A. Liu, J. Wang, X. Liu, B. Cao, C. Zhang, and H. Yu, “Bias-based universal adversarial patch attack for automatic check-out,” in *Proc. Eur. Conf. Comput. Vis.* Springer, 2020, pp. 395–410.
- [276] 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 *European Conference on Computer Vision*. Springer, 2020, pp. 665–681.
- [277] A. Brauneegg, A. Chakraborty, M. Krumdick, N. Lape, S. Leary, K. Manville, E. Merkhofer, L. Strickhart, and M. Walmer, “Apricot: A dataset of physical adversarial attacks on object detection,” in *European Conference on Computer Vision*. Springer, 2020, pp. 35–50.
- [278] 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/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, June 2021, pp. 14 666–14 675.
- [279] R. Duan, X. Mao, A. K. Qin, Y. Chen, S. Ye, Y. He, and Y. Yang, “Adversarial laser beam: Effective physical-world attack to dnms in a blink,” in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, June 2021, pp. 16 062–16 071.
- [280] C. Xie, M. Tan, B. Gong, J. Wang, A. L. Yuille, and Q. V. Le, “Adversarial examples improve image recognition,” in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2020, pp. 819–828.
- [281] M. Tan and Q. V. Le, “Efficientnet: Rethinking model scaling for convolutional neural networks,” *arXiv preprint arXiv:1905.11946*, 2019.
- [282] A. Raghunathan, S. M. Xie, F. Yang, J. C. Duchi, and P. Liang, “Adversarial training can hurt generalization,” *arXiv preprint arXiv:1906.06032*, 2019.
- [283] Y. Li, E. X. Fang, H. Xu, and T. Zhao, “Inductive bias of gradient descent based adversarial training on separable data,” *arXiv preprint arXiv:1906.02931*, 2019.
- [284] S. Qiao, W. Shen, Z. Zhang, B. Wang, and A. Yuille, “Deep co-training for semi-supervised image recognition,” in *Proceedings of the european conference on computer vision (eccv)*, 2018, pp. 135–152.
- [285] C.-H. Ho and N. Vasconcelos, “Contrastive learning with adversarial examples,” *NeurIPS*.
- [286] H. Salman, A. Ilyas, L. Engstrom, A. Kapoor, and A. Madry, “Do adversarially robust imagenet models transfer better?” *arXiv preprint arXiv:2007.08489*, 2020.
- [287] Z. Gan, Y.-C. Chen, L. Li, C. Zhu, Y. Cheng, and J. Liu, “Large-scale adversarial training for vision-and-language representation learning,” *arXiv preprint arXiv:2006.06195*, 2020.
- [288] J. Lee, E. Kim, and S. Yoon, “Anti-adversarially manipulated attributions for weakly and semi-supervised semantic segmentation,” in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, June 2021, pp. 4071–4080.
- [289] M. A. Jalwana, N. Akhtar, M. Bennamoun, and A. Mian, “Attack to explain deep representation,” in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2020, pp. 9543–9552.
- [290] S. Santurkar, A. Ilyas, D. Tsipras, L. Engstrom, B. Tran, and A. Madry, “Image synthesis with a single (robust) classifier,” in *Advances in Neural Information Processing Systems*, 2019, pp. 1262–1273.
- [291] M. Augustin, A. Meinke, and M. Hein, “Adversarial robustness on in- and out-distribution improves explainability,” in *European Conference on Computer Vision*. Springer, 2020, pp. 228–245.
- [292] A. Elliott, S. Law, and C. Russell, “Explaining classifiers using adversarial perturbations on the perceptual ball,” in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, June 2021, pp. 10 693–10 702.
- [293] K. Xu, S. Liu, P. Zhao, P.-Y. Chen, H. Zhang, Q. Fan, D. Erdogmus, Y. Wang, and X. Lin, “Structured adversarial attack: Towards general implementation and better interpretability,” *ICLR*, 2019.
- [294] G. F. Elsayed, I. Goodfellow, and J. Sohl-Dickstein, “Adversarial reprogramming of neural networks,” *arXiv preprint arXiv:1806.11146*, 2018.
- [295] K. Sakaguchi, R. Le Bras, C. Bhagavatula, and Y. Choi, “Winogrande: An adversarial winograd schema challenge at scale,” in *Proceedings of the AAAI Conference on Artificial Intelligence*, vol. 34, no. 05, 2020, pp. 8732–8740.
- [296] R. Le Bras, S. Swayamdipta, C. Bhagavatula, R. Zellers, M. Peters, A. Sabharwal, and Y. Choi, “Adversarial filters of dataset biases,” in *International Conference on Machine Learning*. PMLR, 2020, pp. 1078–1088.
- [297] T. Tanay and L. Griffin, “A boundary tilting perspective on the phenomenon of adversarial examples,” *arXiv preprint arXiv:1608.07690*, 2016.
- [298] D. Krotov and J. Hopfield, “Dense associative memory is robust to adversarial inputs,” *Neural computation*, vol. 30, no. 12, pp. 3151–3167, 2018.

- [299] D. Krotov and J. J. Hopfield, "Dense associative memory for pattern recognition," *Advances in neural information processing systems*, vol. 29, pp. 1172–1180, 2016.
- [300] S. A. Taghanaki, K. Abhishek, S. Azizi, and G. Hamarneh, "A kernelized manifold mapping to diminish the effect of adversarial perturbations," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2019, pp. 11 340–11 349.
- [301] E. D. Cubuk, B. Zoph, S. S. Schoenholz, and Q. V. Le, "Intriguing properties of adversarial examples," *ICLR*, 2018.
- [302] A. Rozsa, M. Günther, and T. E. Boult, "Are accuracy and robustness correlated," in *2016 15th IEEE international conference on machine learning and applications (ICMLA)*. IEEE, 2016, pp. 227–232.
- [303] A. Rozsa, M. Gunther, and T. E. Boult, "Towards robust deep neural networks with bang," *WACV*, 2018.
- [304] P. Tabacof and E. Valle, "Exploring the space of adversarial images," in *2016 International Joint Conference on Neural Networks (IJCNN)*. IEEE, 2016, pp. 426–433.
- [305] F. Tramèr, N. Papernot, I. Goodfellow, D. Boneh, and P. McDaniel, "The space of transferable adversarial examples," *arXiv preprint arXiv:1704.03453*, 2017.
- [306] Y. Li, S. Cheng, H. Su, and J. Zhu, "Defense against adversarial attacks via controlling gradient leaking on embedded manifolds," in *European Conference on Computer Vision*. Springer, 2020, pp. 753–769.
- [307] J.-H. Jacobsen, J. Behrmann, R. Zemel, and M. Bethge, "Excessive invariance causes adversarial vulnerability," *ICLR*, 2019.
- [308] M. V. Reddy, A. Banburski, N. Pant, and T. Poggio, "Biologically inspired mechanisms for adversarial robustness," *arXiv preprint arXiv:2006.16427*, 2020.
- [309] A. Pal and R. Vidal, "A game theoretic analysis of additive adversarial attacks and defenses," *NeurIPS*, 2020.
- [310] J. M. Cohen, E. Rosenfeld, and J. Z. Kolter, "Certified adversarial robustness via randomized smoothing," *ICML*, 2019.
- [311] J. F. Nash *et al.*, "Equilibrium points in n-person games," *Proceedings of the national academy of sciences*, vol. 36, no. 1, pp. 48–49, 1950.
- [312] A. Daniely and H. Schacham, "Most relu networks suffer from l2 adversarial perturbations," *NeurIPS*, 2020.
- [313] A. Shafahi, W. R. Huang, C. Studer, S. Feizi, and T. Goldstein, "Are adversarial examples inevitable?" *arXiv preprint arXiv:1809.02104*, 2018.
- [314] J. Gilmer, L. Metz, F. Faghri, S. S. Schoenholz, M. Raghu, M. Wattenberg, and I. Goodfellow, "Adversarial spheres," *arXiv preprint arXiv:1801.02774*, 2018.
- [315] Y. Song, R. Shu, N. Kushman, and S. Ermon, "Generative adversarial examples," *arXiv preprint arXiv:1805.07894*, 2018.
- [316] D. Stutz, M. Hein, and B. Schiele, "Disentangling adversarial robustness and generalization," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2019, pp. 6976–6987.
- [317] D. Tsipras, S. Santurkar, L. Engstrom, A. Turner, and A. Madry, "Robustness may be at odds with accuracy," *ICLR*, 2019.
- [318] D. Su, H. Zhang, H. Chen, J. Yi, P.-Y. Chen, and Y. Gao, "Is robustness the cost of accuracy?—a comprehensive study on the robustness of 18 deep image classification models," in *Proceedings of the European Conference on Computer Vision (ECCV)*, 2018, pp. 631–648.
- [319] S.-M. Moosavi-Dezfooli, A. Fawzi, O. Fawzi, P. Frossard, and S. Soatto, "Analysis of universal adversarial perturbations," *arXiv preprint arXiv:1705.09554*, 2017.
- [320] —, "Robustness of classifiers to universal perturbations: A geometric perspective," in *International Conference on Learning Representations*, 2018.
- [321] S. Jetley, N. Lord, and P. Torr, "With friends like these, who needs adversaries?" in *Advances in neural information processing systems*, 2018, pp. 10 749–10 759.
- [322] C. Zhang, P. Benz, T. Imtiaz, and I. S. Kweon, "Understanding adversarial examples from the mutual influence of images and perturbations," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2020, pp. 14 521–14 530.
- [323] A. Ilyas, S. Santurkar, D. Tsipras, L. Engstrom, B. Tran, and A. Madry, "Adversarial examples are not bugs, they are features," in *Advances in Neural Information Processing Systems*, 2019, pp. 125–136.
- [324] S. Bubeck, Y. T. Lee, E. Price, and I. Razenshteyn, "Adversarial examples from computational constraints," in *International Conference on Machine Learning*. PMLR, 2019, pp. 831–840.
- [325] D. Wu, Y. Wang, S.-T. Xia, J. Bailey, and X. Ma, "Skip connections matter: On the transferability of adversarial examples generated with resnets," *arXiv preprint arXiv:2002.05990*, 2020.
- [326] H. Zhang, H. Chen, Z. Song, D. Boning, I. S. Dhillon, and C.-J. Hsieh, "The limitations of adversarial training and the blind-spot attack," *ICLR*, 2020.
- [327] L. Schott, J. Rauber, M. Bethge, and W. Brendel, "Towards the first adversarially robust neural network model on mnist," *ICLR*, 2019.
- [328] G. W. Ding, K. Y. C. Lui, X. Jin, L. Wang, and R. Huang, "On the sensitivity of adversarial robustness to input data distributions," in *ICLR (Poster)*, 2019.
- [329] C. Song, K. He, L. Wang, and J. E. Hopcroft, "Improving the generalization of adversarial training with domain adaptation," *arXiv preprint arXiv:1810.00740*, 2018.
- [330] R. Geirhos, P. Rubisch, C. Michaelis, M. Bethge, F. A. Wichmann, and W. Brendel, "Imagenet-trained cnns are biased towards texture; increasing shape bias improves accuracy and robustness," *ICLR*, 2019.
- [331] T. Zhang and Z. Zhu, "Interpreting adversarially trained convolutional neural networks," *ICML*, 2019.
- [332] C. Gong, T. Ren, M. Ye, and Q. Liu, "Maxup: Lightweight adversarial training with data augmentation improves neural network training," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, June 2021, pp. 2474–2483.
- [333] Y. Wang, D. Zou, J. Yi, J. Bailey, X. Ma, and Q. Gu, "Improving adversarial robustness requires revisiting misclassified examples," in *International Conference on Learning Representations*, 2019.
- [334] S. Gowal, C. Qin, P.-S. Huang, T. Cemgil, K. Dvijotham, T. Mann, and P. Kohli, "Achieving robustness in the wild via adversarial mixing with disentangled representations," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2020, pp. 1211–1220.
- [335] 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.
- [336] G. W. Ding, Y. Sharma, K. Y. C. Lui, and R. Huang, "Mma training: Direct input space margin maximization through adversarial training," in *International Conference on Learning Representations*, 2019.
- [337] M. Balunovic and M. Vechev, "Adversarial training and provable defenses: Bridging the gap," in *International Conference on Learning Representations*, 2019.
- [338] B. Vivek and R. V. Babu, "Single-step adversarial training with dropout scheduling," in *2020 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*. IEEE, 2020, pp. 947–956.
- [339] C. Song, K. He, J. Lin, L. Wang, and J. E. Hopcroft, "Robust local features for improving the generalization of adversarial training," *ICLR*, 2020.
- [340] F. Farnia, J. M. Zhang, and D. Tse, "Generalizable adversarial training via spectral normalization," *ICLR*, 2019.
- [341] C. Xiao and C. Zheng, "One man's trash is another man's treasure: Resisting adversarial examples by adversarial examples," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2020, pp. 412–421.
- [342] M. Naseer, S. Khan, M. Hayat, F. S. Khan, and F. Porikli, "A self-supervised approach for adversarial robustness," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2020, pp. 262–271.
- [343] T. Chen, S. Liu, S. Chang, Y. Cheng, L. Amini, and Z. Wang, "Adversarial robustness: From self-supervised pre-training to fine-tuning," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2020, pp. 699–708.
- [344] Y. Jang, T. Zhao, S. Hong, and H. Lee, "Adversarial defense via learning to generate diverse attacks," in *Proceedings of the IEEE/CVF International Conference on Computer Vision*, 2019, pp. 2740–2749.
- [345] W.-A. Lin, C. P. Lau, A. Levine, R. Chellappa, and S. Feizi, "Dual manifold adversarial robustness: Defense against lp and non-lp adversarial attacks," *arXiv preprint arXiv:2009.02470*, 2020.
- [346] J. Wang and H. Zhang, "Bilateral adversarial training: Towards fast training of more robust models against adversarial attacks," in *Proceedings of the IEEE/CVF International Conference on Computer Vision*, 2019, pp. 6629–6638.
- [347] S. Ye, K. Xu, S. Liu, H. Cheng, J.-H. Lambrechts, H. Zhang, A. Zhou, K. Ma, Y. Wang, and X. Lin, "Adversarial robustness vs. model compression, or both?" in *Proceedings of the IEEE/CVF International Conference on Computer Vision*, 2019, pp. 111–120.
- [348] Y. Dong, Z. Deng, T. Pang, H. Su, and J. Zhu, "Adversarial distributional training for robust deep learning," *NeurIPS*, 2020.
- [349] D. Madaan, J. Shin, and S. J. Hwang, "Adversarial neural pruning with latent vulnerability suppression," in *International Conference on Machine Learning*. PMLR, 2020, pp. 6575–6585.

- [350] P. Maini, E. Wong, and Z. Kolter, "Adversarial robustness against the union of multiple perturbation models," in *International Conference on Machine Learning*. PMLR, 2020, pp. 6640–6650.
- [351] H. Zhang, M. Cisse, Y. N. Dauphin, and D. Lopez-Paz, "mixup: Beyond empirical risk minimization," *arXiv preprint arXiv:1710.09412*, 2017.
- [352] V. Verma, A. Lamb, C. Beckham, A. Courville, I. Mitliagkis, and Y. Bengio, "Manifold mixup: Encouraging meaningful on-manifold interpolation as a regularizer," *ICML*, 2019.
- [353] S. Lee, H. Lee, and S. Yoon, "Adversarial vertex mixup: Toward better adversarially robust generalization," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2020, pp. 272–281.
- [354] C. Xie and A. Yuille, "Intriguing properties of adversarial training at scale," *ICLR*, 2020.
- [355] Y. Li, E. X. Fang, H. Xu, and T. Zhao, "Implicit bias of gradient descent based adversarial training on separable data," in *International Conference on Learning Representations*, 2019.
- [356] A. Shafahi, P. Saadatpanah, C. Zhu, A. Ghiassi, C. Studer, D. Jacobs, and T. Goldstein, "Adversarially robust transfer learning," *arXiv preprint arXiv:1905.08232*, 2019.
- [357] V. Sehwag, S. Wang, P. Mittal, and S. Jana, "Hydra: Pruning adversarially robust neural networks," *Advances in Neural Information Processing Systems (NeurIPS)*, vol. 7, 2020.
- [358] E. Wong, L. Rice, and J. Z. Kolter, "Fast is better than free: Revisiting adversarial training," *ICLR*, 2020.
- [359] M. Andriushchenko and N. Flammarion, "Understanding and improving fast adversarial training," *NeurIPS*, 2020.
- [360] P. Zhao, P.-Y. Chen, P. Das, K. N. Ramamurthy, and X. Lin, "Bridging mode connectivity in loss landscapes and adversarial robustness," *ICLR*, 2020.
- [361] D. Wu, S.-T. Xia, and Y. Wang, "Adversarial weight perturbation helps robust generalization," *Advances in Neural Information Processing Systems*, vol. 33, 2020.
- [362] R. Gao, T. Cai, H. Li, C.-J. Hsieh, L. Wang, and J. D. Lee, "Convergence of adversarial training in overparametrized neural networks," *Advances in Neural Information Processing Systems*, vol. 32, pp. 13 029–13 040, 2019.
- [363] Y. Zhang, O. Plevrakis, S. S. Du, X. Li, Z. Song, and S. Arora, "Overparameterized adversarial training: An analysis overcoming the curse of dimensionality," *arXiv preprint arXiv:2002.06668*, 2020.
- [364] T. Pang, X. Yang, Y. Dong, K. Xu, J. Zhu, and H. Su, "Boosting adversarial training with hypersphere embedding," *arXiv preprint arXiv:2002.08619*, 2020.
- [365] T. Wu, L. Tong, and Y. Vorobeychik, "Defending against physically realizable attacks on image classification," *ICLR*, 2020.
- [366] Y. Cheng, L. Jiang, and W. Macherey, "Robust neural machine translation with doubly adversarial inputs," *arXiv preprint arXiv:1906.02443*, 2019.
- [367] M. Guo, Y. Yang, R. Xu, Z. Liu, and D. Lin, "When nas meets robustness: In search of robust architectures against adversarial attacks," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2020, pp. 631–640.
- [368] A. Bui, T. Le, H. Zhao, P. Montague, O. deVel, T. Abraham, and D. Phung, "Improving adversarial robustness by enforcing local and global compactness," *ECCV*, 2020.
- [369] X. Liu, Y. Li, C. Wu, and C.-J. Hsieh, "Adv-bnn: Improved adversarial defense through robust bayesian neural network," *ICLR*, 2019.
- [370] A. Jeddi, M. J. Shafiee, M. Karg, C. Scharfenberger, and A. Wong, "Learn2perturb: an end-to-end feature perturbation learning to improve adversarial robustness," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2020, pp. 1241–1250.
- [371] G. Li, S. Ding, J. Luo, and C. Liu, "Enhancing intrinsic adversarial robustness via feature pyramid decoder," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2020, pp. 800–808.
- [372] T. Borkar, F. Heide, and L. Karam, "Defending against universal attacks through selective feature regeneration," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2020, pp. 709–719.
- [373] H. Wang and C.-N. Yu, "A direct approach to robust deep learning using adversarial networks," *arXiv preprint arXiv:1905.09591*, 2019.
- [374] C. Xie, Y. Wu, L. v. d. Maaten, A. L. Yuille, and K. He, "Feature denoising for improving adversarial robustness," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2019, pp. 501–509.
- [375] M. Lecuyer, V. Atlidakis, R. Geambasu, D. Hsu, and S. Jana, "Certified robustness to adversarial examples with differential privacy," in *2019 IEEE Symposium on Security and Privacy (SP)*. IEEE, 2019, pp. 656–672.
- [376] X. Liu, M. Cheng, H. Zhang, and C.-J. Hsieh, "Towards robust neural networks via random self-ensemble," in *Proceedings of the European Conference on Computer Vision (ECCV)*, 2018, pp. 369–385.
- [377] Z. He, A. S. Rakin, and D. Fan, "Parametric noise injection: Trainable randomness to improve deep neural network robustness against adversarial attack," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2019, pp. 588–597.
- [378] J. Svoboda, J. Masci, F. Monti, M. M. Bronstein, and L. Guibas, "Peernets: Exploiting peer wisdom against adversarial attacks," *ICLR*, 2019.
- [379] R. Hosseini, X. Yang, and P. Xie, "Dsma: Differentiable search of robust neural architectures," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2021, pp. 6196–6205.
- [380] G. Cazenavette, C. Murdock, and S. Lucey, "Architectural adversarial robustness: The case for deep pursuit," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2021, pp. 7150–7158.
- [381] A. Chan, Y. Tay, Y. S. Ong, and J. Fu, "Jacobian adversarially regularized networks for robustness," *ICLR*, 2020.
- [382] A. Dabouei, S. Soleymani, F. Taherkhani, J. Dawson, and N. M. Nasrabadi, "Exploiting joint robustness to adversarial perturbations," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2020, pp. 1122–1131.
- [383] T. Tadros, G. Krishnan, R. Ramyaa, and M. Bazhenov, "Biologically inspired sleep algorithm for increased generalization and adversarial robustness in deep neural networks," in *International Conference on Learning Representations*, 2019.
- [384] S. Addepalli, A. Baburaj, G. Sriramanan, and R. V. Babu, "Towards achieving adversarial robustness by enforcing feature consistency across bit planes," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2020, pp. 1020–1029.
- [385] Y. Qin, N. Frosst, S. Sabour, C. Raffel, G. Cottrell, and G. Hinton, "Detecting and diagnosing adversarial images with class-conditional capsule reconstructions," *ICLR*, 2020.
- [386] S. Sabour, N. Frosst, and G. E. Hinton, "Dynamic routing between capsules," in *Advances in neural information processing systems*, 2017, pp. 3856–3866.
- [387] Z. Deng, X. Yang, S. Xu, H. Su, and J. Zhu, "Libre: A practical bayesian approach to adversarial detection," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, June 2021, pp. 972–982.
- [388] S. Li, S. Zhu, S. Paul, A. Roy-Chowdhury, C. Song, S. Krishnamurthy, A. Swami, and K. S. Chan, "Connecting the dots: Detecting adversarial perturbations using context inconsistency," in *European Conference on Computer Vision*. Springer, 2020, pp. 396–413.
- [389] G. Tao, S. Ma, Y. Liu, and X. Zhang, "Attacks meet interpretability: Attribute-steered detection of adversarial samples," *arXiv preprint arXiv:1810.11580*, 2018.
- [390] Y. Qiu, J. Leng, C. Guo, Q. Chen, C. Li, M. Guo, and Y. Zhu, "Adversarial defense through network profiling based path extraction," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2019, pp. 4777–4786.
- [391] J. Liu, W. Zhang, Y. Zhang, D. Hou, Y. Liu, H. Zha, and N. Yu, "Detection based defense against adversarial examples from the steganalysis point of view," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2019, pp. 4825–4834.
- [392] X. Yin, S. Kolouri, and G. K. Rohde, "Gat: Generative adversarial training for adversarial example detection and robust classification," in *International Conference on Learning Representations*, 2019.
- [393] Z. Liu, Q. Liu, T. Liu, N. Xu, X. Lin, Y. Wang, and W. Wen, "Feature distillation: Dnn-oriented jpeg compression against adversarial examples," in *2019 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*. IEEE, 2019, pp. 860–868.
- [394] N. Das, M. Shanbhogue, S.-T. Chen, F. Hohman, L. Chen, M. E. Kounavis, and D. H. Chau, "Keeping the bad guys out: Protecting and vaccinating deep learning with jpeg compression," *arXiv preprint arXiv:1705.02900*, 2017.
- [395] C. Guo, M. Rana, M. Cisse, and L. Van Der Maaten, "Countering adversarial images using input transformations," *ICLR*, 2018.
- [396] E. Raff, J. Sylvester, S. Forsyth, and M. McLean, "Barrage of random transforms for adversarially robust defense," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2019, pp. 6528–6537.

- [397] O. Taran, S. Rezaeifar, T. Holotyak, and S. Voloshynovskiy, “Defending against adversarial attacks by randomized diversification,” in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2019, pp. 11 226–11 233.
- [398] X. Jia, X. Wei, X. Cao, and H. Foroosh, “Comdefend: An efficient image compression model to defend adversarial examples,” in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2019, pp. 6084–6092.
- [399] R. Theagarajan, M. Chen, B. Bhanu, and J. Zhang, “Shieldnets: Defending against adversarial attacks using probabilistic adversarial robustness,” in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2019, pp. 6988–6996.
- [400] B. Sun, N.-h. Tsai, F. Liu, R. Yu, and H. Su, “Adversarial defense by stratified convolutional sparse coding,” in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2019, pp. 11 447–11 456.
- [401] P. Samangouei, M. Kabkab, and R. Chellappa, “Defense-gan: Protecting classifiers against adversarial attacks using generative models,” *ICLR*, 2018.
- [402] P. Gupta and E. Rahtu, “Ciidefence: Defeating adversarial attacks by fusing class-specific image inpainting and image denoising,” in *Proceedings of the IEEE/CVF International Conference on Computer Vision*, 2019, pp. 6708–6717.
- [403] C. Kou, H. K. Lee, E.-C. Chang, and T. K. Ng, “Enhancing transformation-based defenses against adversarial attacks with a distribution classifier,” in *International Conference on Learning Representations*, 2019.
- [404] J. Yuan and Z. He, “Ensemble generative cleaning with feedback loops for defending adversarial attacks,” in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2020, pp. 581–590.
- [405] G. Cohen, G. Sapiro, and R. Giryes, “Detecting adversarial samples using influence functions and nearest neighbors,” in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2020, pp. 14 453–14 462.
- [406] F. Croce, J. Rauber, and M. Hein, “Scaling up the randomized gradient-free adversarial attack reveals overestimation of robustness using established attacks,” *International Journal of Computer Vision*, pp. 1–19, 2019.
- [407] G. Katz, C. Barrett, D. L. Dill, K. Julian, and M. J. Kochenderfer, “Reluplex: An efficient smt solver for verifying deep neural networks,” in *International Conference on Computer Aided Verification*. Springer, 2017, pp. 97–117.
- [408] V. Tjeng, K. Xiao, and R. Tedrake, “Evaluating robustness of neural networks with mixed integer programming,” *ICLR*, 2019.
- [409] M. Hein and M. Andriushchenko, “Formal guarantees on the robustness of a classifier against adversarial manipulation,” in *Advances in Neural Information Processing Systems*, 2017, pp. 2266–2276.
- [410] A. Raghunathan, J. Steinhardt, and P. Liang, “Certified defenses against adversarial examples,” *ICLR*, 2018.
- [411] E. Wong and Z. Kolter, “Provable defenses against adversarial examples via the convex outer adversarial polytope,” in *International Conference on Machine Learning*. PMLR, 2018, pp. 5286–5295.
- [412] M. Mirman, T. Gehr, and M. Vechev, “Differentiable abstract interpretation for provably robust neural networks,” in *International Conference on Machine Learning*, 2018, pp. 3578–3586.
- [413] K. Y. Xiao, V. Tjeng, N. M. Shafiqullah, and A. Madry, “Training for faster adversarial robustness verification via inducing relu stability,” *ICLR*, 2019.
- [414] F. Croce and M. Hein, “Provable robustness against all adversarial l_p -perturbations for $p \geq 1$,” in *ICLR*, 2020.
- [415] F. Tramèr and D. Boneh, “Adversarial training and robustness for multiple perturbations,” in *Advances in Neural Information Processing Systems*, 2019, pp. 5866–5876.
- [416] J. Jia, X. Cao, B. Wang, and N. Z. Gong, “Certified robustness for top-k predictions against adversarial perturbations via randomized smoothing,” *ICLR*, 2020.
- [417] X. Cao and N. Z. Gong, “Mitigating evasion attacks to deep neural networks via region-based classification,” in *Proceedings of the 33rd Annual Computer Security Applications Conference*, 2017, pp. 278–287.
- [418] R. Zhai, C. Dan, D. He, H. Zhang, B. Gong, P. Ravikumar, C.-J. Hsieh, and L. Wang, “Macer: Attack-free and scalable robust training via maximizing certified radius,” *ICLR*, 2020.
- [419] M. Fischer, M. Baader, and M. Vechev, “Certified defense to image transformations via randomized smoothing,” *NeurIPS*, 2020.
- [420] A. Levine and S. Feizi, “(de) randomized smoothing for certifiable defense against patch attacks,” *arXiv preprint arXiv:2002.10733*, 2020.
- [421] P.-y. Chiang, R. Ni, A. Abdelkader, C. Zhu, C. Studor, and T. Goldstein, “Certified defenses for adversarial patches,” *ICLR*, 2020.
- [422] P. Awasthi, H. Jain, A. S. Rawat, and A. Vijayaraghavan, “Adversarial robustness via robust low rank representations,” *NeurIPS*, 2020.
- [423] D. Zhang, M. Ye, C. Gong, Z. Zhu, and Q. Liu, “Black-box certification with randomized smoothing: A functional optimization based framework,” *arXiv preprint arXiv:2002.09169*, 2020.
- [424] A. Rahnama, A. T. Nguyen, and E. Raff, “Robust design of deep neural networks against adversarial attacks based on Lyapunov theory,” in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2020, pp. 8178–8187.
- [425] S. Saralajew, L. Holdijk, T. Villmann, and U. Mittweida, “Fast adversarial robustness certification of nearest prototype classifiers for arbitrary seminorms,” *Advances in Neural Information Processing Systems*, vol. 33, 2020.
- [426] R. Y. Zhang, “On the tightness of semidefinite relaxations for certifying robustness to adversarial examples,” *arXiv preprint arXiv:2006.06759*, 2020.
- [427] S. Goldwasser, A. T. Kalai, Y. T. Kalai, and O. Montasser, “Beyond perturbations: Learning guarantees with arbitrary adversarial test examples,” 2020.
- [428] H. Salman, M. Sun, G. Yang, A. Kapoor, and J. Z. Kolter, “Denoised smoothing: A provable defense for pretrained classifiers,” *Advances in Neural Information Processing Systems*, vol. 33, 2020.
- [429] P. Awasthi, G. Yu, C.-S. Ferng, A. Tomkins, and D.-C. Juan, “Adversarial robustness across representation spaces,” in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, June 2021, pp. 7608–7616.
- [430] T. Cemgil, S. Ghaisas, K. D. Dvijotham, and P. Kohli, “Adversarially robust representations with smooth encoders,” in *International Conference on Learning Representations*, 2019.
- [431] Z. He, A. S. Rakin, J. Li, C. Chakrabarti, and D. Fan, “Defending and harnessing the bit-flip based adversarial weight attack,” in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2020, pp. 14 095–14 103.
- [432] G. Carbone, M. Wicker, L. Laurenti, A. Patane, L. Bortolussi, and G. Sanguinetti, “Robustness of bayesian neural networks to gradient-based attacks,” *arXiv preprint arXiv:2002.04359*, 2020.
- [433] S. Sharmir, N. Rathi, P. Panda, and K. Roy, “Inherent adversarial robustness of deep spiking neural networks: Effects of discrete input encoding and non-linear activations,” in *European Conference on Computer Vision*. Springer, 2020, pp. 399–414.
- [434] X. Zhang and M. Zitnik, “Gnnguard: Defending graph neural networks against adversarial attacks,” *NeurIPS*, 2020.
- [435] M. Du, R. Jia, and D. Song, “Robust anomaly detection and backdoor attack detection via differential privacy,” *ICLR*, 2020.
- [436] C. K. Mummadi, T. Brox, and J. H. Metzen, “Defending against universal perturbations with shared adversarial training,” in *Proceedings of the IEEE/CVF International Conference on Computer Vision*, 2019, pp. 4928–4937.
- [437] N. Akhtar, J. Liu, and A. Mian, “Defense against universal adversarial perturbations,” in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2018, pp. 3389–3398.
- [438] X. Zhang and D. Evans, “Cost-sensitive robustness against adversarial examples,” *arXiv preprint arXiv:1810.09225*, 2018.
- [439] H.-Y. Chen, J.-H. Liang, S.-C. Chang, J.-Y. Pan, Y.-T. Chen, W. Wei, and D.-C. Juan, “Improving adversarial robustness via guided complement entropy,” in *Proceedings of the IEEE/CVF International Conference on Computer Vision*, 2019, pp. 4881–4889.
- [440] S. Jia, C. Ma, Y. Song, and X. Yang, “Robust tracking against adversarial attacks,” in *European Conference on Computer Vision*. Springer, 2020, pp. 69–84.
- [441] R. Shao, P. Perera, P. C. Yuen, and V. M. Patel, “Open-set adversarial defense,” *ECCV*, 2020.
- [442] J. Zhou, C. Liang, and J. Chen, “Manifold projection for adversarial defense on face recognition,” in *European Conference on Computer Vision*. Springer, 2020, pp. 288–305.
- [443] M. Goldblum, L. Fowl, and T. Goldstein, “Adversarially robust few-shot learning: A meta-learning approach,” *Advances in Neural Information Processing Systems*, vol. 33, 2020.
- [444] T. Orekondy, B. Schiele, and M. Fritz, “Prediction poisoning: Towards defenses against dnn model stealing attacks,” in *International Conference on Learning Representations*, 2019.

- [445] S. Kariyappa and M. K. Qureshi, “Defending against model stealing attacks with adaptive misinformation,” in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2020, pp. 770–778.
- [446] J. Sulam, R. Muthumukar, and R. Arora, “Adversarial robustness of supervised sparse coding,” *arXiv preprint arXiv:2010.12088*, 2020.
- [447] J. Yang and C. Vondrick, “Multitask learning strengthens adversarial robustness,” in *European Conference on Computer Vision*. Springer, 2020.
- [448] E. Kim, J. Rego, Y. Watkins, and G. T. Kenyon, “Modeling biological immunity to adversarial examples,” in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2020, pp. 4666–4675.
- [449] H. Wang, T. Chen, S. Gui, T.-K. Hu, J. Liu, and Z. Wang, “Calibratable adversarial training: In-situ tradeoff between robustness and accuracy for free,” *NeurIPS*, 2020.
- [450] C.-H. Weng, Y.-T. Lee, and S.-H. B. Wu, “On the trade-off between adversarial and backdoor robustness,” *Advances in Neural Information Processing Systems*, vol. 33, 2020.
- [451] H. Chen, B. Zhang, S. Xue, X. Gong, H. Liu, R. Ji, and D. Doermann, “Anti-bandit neural architecture search for model defense,” in *European Conference on Computer Vision*. Springer, 2020, pp. 70–85.
- [452] T. Wu, Z. Liu, Q. Huang, Y. Wang, and D. Lin, “Adversarial robustness under long-tailed distribution,” in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, June 2021, pp. 8659–8668.
- [453] M. A. Jalwana, N. Akhtar, M. Bennamoun, and A. Mian, “Cameras: Enhanced resolution and sanity preserving class activation mapping for image saliency,” in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2021.
- [454] Y. Yu, X. Gao, and C.-Z. Xu, “Lafeat: Piercing through adversarial defenses with latent features,” in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, June 2021, pp. 5735–5745.
- [455] S. Khan, M. Naseer, M. Hayat, S. W. Zamir, F. S. Khan, and M. Shah, “Transformers in vision: A survey,” *arXiv preprint arXiv:2101.01169*, 2021.
- [456] Y. Liu, S. Ma, Y. Aafer, W.-C. Lee, J. Zhai, W. Wang, and X. Zhang, “Trojaning attack on neural networks,” 2017.