

This CVPR 2020 workshop paper is the Open Access version, provided by the Computer Vision Foundation. Except for this watermark, it is identical to the accepted version;

the final published version of the proceedings is available on IEEE Xplore.

## Learning Ordered Top-k Adversarial Attacks via Adversarial Distillation

Zekun Zhang and Tianfu Wu\*

Department of ECE and the Visual Narrative Initiative, NC State University

{zzhang56, tianfu\_wu}@ncsu.edu

## Abstract

Deep Neural Networks (DNNs) are vulnerable to adversarial attacks, especially white-box targeted attacks. This paper studies the problem of how aggressive white-box targeted attacks can be to go beyond widely used Top-1 attacks. We propose to learn ordered Top-k attacks  $(k \ge 1)$ , which enforce the Top-k predicted labels of an adversarial example to be the k (randomly) selected and ordered labels (the ground-truth label is exclusive). Two methods are presented. First, we extend the vanilla Carlini-Wagner (C&W) method and use it as a strong baseline. Second, we present an Adversarial Distillation (AD) framework consisting of two components: (i) Computing an adversarial probability distribution for a given ordered Top-k targeted labels. (ii) Learning adversarial examples by minimizing the Kullback-Leibler (KL) divergence between the adversarial distribution and the predicted distribution, together with the perturbation energy penalty. In computing adversarial distributions, we explore how to leverage label semantic similarities, leading to knowledge-oriented attacks. In experiments, we test Top-k (k = 1, 2, 5, 10) attacks in the ImageNet-1000 val dataset using three representative DNNs trained with the clean ImageNet-1000 train dataset, ResNet-50 [11], DenseNet-121 [14] and AOGNet-12M [21]. Overall, the proposed AD approach obtains the best results, especially by a large margin when computation budget is limited. It reduces the perturbation energy consistently with the same attack success rate on all the four k's, and improve the attack success rate by a large margin against the modified C&W method for k = 10.

## 1. Introduction

Despite the recent dramatic progress, deep neural networks (DNNs) [20, 18, 11, 36] trained for visual recognition tasks (*e.g.*, image classification) can be easily fooled by socalled **adversarial attacks** which utilize visually imperceptible, carefully-crafted perturbations to cause networks to misclassify inputs in arbitrarily chosen ways in the close set

of labels used in training [25, 38, 2, 4], even with onepixel attacks [35]. The existence of adversarial attacks hinders the deployment of DNNs-based visual recognition systems in a wide range of applications such as autonomous driving and smart medical diagnosis in the long-run.

In this paper, we are interested in learning visuallyimperceptible targeted attacks under the whitebox and limited



Figure 1. Comparisons of Top-*k* attack success rates (ASR, higher is better) between the modified C&W method (CW<sup>k</sup>) and the proposed Adversarial Distillation (AD) method (k = 1, 2, 5, 10). The thickness of plotted lines represents the reverse  $\ell_2$  energy of the learned perturbation (thicker is better).  $9 \times 30$  and  $9 \times 1000$ represent the computing budgets in learning (see Section 4 for experimental settings).

computing budget setting in image classification tasks. In the literature, most methods address targeted attacks in the Top-1 manner, in which an adversarial attack is said to be successful if a randomly selected label (not the ground-truth label) is predicted as the Top-1 label with the added perturbation satisfying to be visually-imperceptible. One question arises,

• *The "robustness" of an attack method itself*: How far is the attack method able to push the underlying ground-truth label in the prediction of the learned adversarial examples?

Table 1 shows the evaluation results of the "robustness" of different attack methods. The widely used C&W method [4] does not push the GT labels very far, especially when smaller perturbation energy is aimed using larger search range (*e.g.*, the average rank of the GT label is 2.6

<sup>\*</sup>T. Wu is the corresponding author.



Figure 2. Learned adversarial examples of ordered Top-10 adversarial attacks for ResNet-50 [11] pretrained with clean images. The proposed AD method has smaller perturbation energies and "cleaner" (lower-entropy) prediction distributions than the proposed modified C&W method ( $CW^k$ ). Note that for Top-10 attacks, the 9 × 30 search scheme does not work (see Table. 2). See text for detail.

Method	ASR	Proport	ion of GT	Average Rank of GT			
		Top-3	Top-5	Top-10	Top-50	Top-100	Labels (larger is better)
C&W <sub>9×30</sub> [4]	99.9	36.9	50.5	66.3	90.0	95.1	20.4
$C\&W_{9 \times 1000}$ [4]	100	71.9	87.0	96.1	99.9	100	2.6
FGSM [10]	80.7	25.5	37.8	52.8	81.2	89.2	44.2
PGD10 [24]	100	3.3	6.7	12	34.7	43.9	306.5
MIFGSM <sub>10</sub> [7]	99.9	0.7	1.9	6.0	22.5	32.3	404.4

Table 1. Results showing where the ground-truth (GT) labels are in the prediction of learned adversarial examples for different attack methods. The test is performed in ImageNet-1000 val dataset using ResNet-50 [11] pretrained with clean images. The subscripts of methods indicate the computing budgets used. See Section 4 for experimental settings.

for C&W<sub>9×1000</sub>). Consider Top-5, if the ground-truth labels of adversarial examples still largely appear in the Top-5 of the prediction, we may be over-confident about the 100% ASR, especially when some downstream modules may rely on Top-5 predictions in their decision making. On the contrary, the three untargeted attack approaches are much better in terms of pushing the GT labels since they usually move against the GT label explicitly in the optimization, but the energies of learned perturbation are usually much larger since no explicit constraints are posed. As we shall show, more "robust" attack methods can be developed by harnessing the advantages of the two types of attack methods. In addition, the targeted Top-1 attack setting could limit the flexibility of attacks, and may lead to less rich perturbations.

To facilitate explicit control of targeted attacks and enable more "robust" attack methods, one natural solution, which is *the focus of this paper*, is to develop **ordered Top-** k targeted attacks which enforce the Top-k predicted labels of an adversarial example to be the k (randomly) selected and ordered labels (k > 1, the GT label is exclusive). In this paper, we *present two methods* of learning ordered Top-k attacks. The basic idea is to design proper adversarial objective functions that result in imperceptible perturbations for a testing image (whose original prediction by a model is correct) through iterative gradient-based backpropagation. First, we extend the vanilla Carlini-Wagner (C&W) method [4], denoted by  $CW^k$ , and use it as a strong baseline. Second, we present an Adversarial Distillation (AD) framework consisting of two components: (i) Computing an adversarial probability distribution for any given ordered Top-k targeted labels. (ii) Learning adversarial examples by minimizing the Kullback-Leibler (KL) divergence between the adversarial distribution and the predicted distribution, together with the perturbation energy penalty.

The proposed AD framework can be viewed as applying the network distillation frameworks [12, 3, 28] for "the bad" induced by target adversarial distributions. To compute a proper adversarial distribution for a given ordered Top-k targeted labels, the AD framework is motivated by two aspects: (i) The difference between the objective functions used by the C&W method and the three untargeted attack methods (Table 1) respectively. The former maximizes the margin of the logits between the target and the runner-up (either GT or not), while the latter maximizes the cross-entropy between the prediction probabilities (softmax of logits) and the one-hot distribution of the ground-truth. (ii) The label smoothing methods [37, 30], which are often used to improve the performance of DNNs by addressing the over-confidence issue in the one-hot vector encoding of labels. We explore how to leverage label semantic similarities in computing "smoothed" adversarial distributions, leading to **knowledge-oriented attacks**. We measure label semantic similarities using the cosine distance between some off-the-shelf word2vec embedding of labels such as the pretrained Glove embedding [29]. Along this direction, another question of interest is further investigated: *Are all Top-k targets equally challenging for an attack approach?* The answer is no (as intuitively perceived), and we observe some meaningful "blind spots" in experiments.

In experiments, we test Top-k (k = 1, 2, 5, 10) in the ImageNet-1000 [33] val dataset using three representative DNNs trained with clean ImageNet-1000 train dataset, ResNet-50 [11], DenseNet-121 [14] and AOGNet-12M [21] respectively. Overall, the proposed AD approach obtains the best results. It reduces the perturbation energy consistently with the same attack success rate on all the four k's, and improve the attack success rate by a large margin against the modified  $CW^k$  method for k = 10 (Fig. 1 and learned adversarial examples in Fig. 2). We observe that Top-k targets that are distant from the GT label in terms of either label semantic distance or prediction scores of clean images are actually more difficult to attack. In summary, not only can ordered Top-k attacks improve the "robustness" of attacks, but also they provide insights on how aggressive adversarial attacks can be (under limited computing budgets).

## 2. Related Work and Our Contributions

The growing ubiquity of DNNs in advanced machine learning and AI systems dramatically increases their capabilities, but also increases the potential for new vulnerabilities to attacks [39, 17, 34, 9, 31, 22]. This situation has become critical as many powerful approaches have been developed where imperceptible perturbations to DNN inputs could deceive a well-trained DNN, significantly altering its prediction. Such results have initiated a rapidly proliferating field of research characterized by ever more complex attacks [27, 23, 41, 8] that prove increasingly strong against defensive countermeasures [10, 16, 24, 40]. Please refer to [1] for a comprehensive survey of attack methods in computer vision. We review some related work that motivate our work and show the difference. Assuming full access to DNNs pretrained with clean images, white-box targeted attacks are powerful ways of investigating the brittleness of DNNs and their sensitivity to non-robust yet wellgeneralizing features in the data, and of exploiting adversarial examples as useful features [15].

**Distillation.** The central idea of our proposed AD method is built on distillation. Network distillation [3, 12]

is a powerful training scheme proposed to train a new, usually lightweight model (a.k.a., the student) to mimic another already trained model (a.k.a. the teacher). It takes a functional viewpoint of the knowledge learned by the teacher as the conditional distribution it produces over outputs given an input. It teaches the student to keep up or emulate by adding some regularization terms to the loss in order to encourage the two models to be similar directly based on the distilled knowledge, replacing the training labels. Label smoothing [37] can be treated as a simple hand-crafted knowledge to help improve model performance. Distillation has been exploited to develop defense models [28, 26] to improve model robustness. Our proposed adversarial distillation method utilizes the distillation idea in an opposite direction, leveraging label semantic knowledge for learning ordered Top-k attacks and improving attack robustness.

Adversarial Attack. For image classification tasks using DNNs, the discovery of the existence of visuallyimperceptible adversarial attacks [38] was a big shock in developing DNNs. White-box attacks provide a powerful way of evaluating model brittleness. In a plain and loose explanation, DNNs are universal function approximator [13] and capable of even fitting random labels [43] in large scale classification tasks as ImageNet-1000 [33]. Thus, adversarial attacks are generally learnable provided proper objective functions are given, especially when DNNs are trained with fully differentible back-propagation. Many white-box attack methods focus on norm-ball constrained objective functions [38, 19, 4, 7, 32]. The C&W method investigates 7 different loss functions. The best performing loss function found by the C&W method has been applied in many attack methods and achieved strong results [6, 24, 5]. By introducing momentum in the MIFGSM method [7] and the  $\ell_p$  gradient projection in the PGD method [24], they usually achieve better performance in generating adversarial examples. In the meanwhile, some other attack methods such as the StrAttack [42] also investigate different loss functions for better interpretability of attacks. Our proposed method leverages label semantic knowledge in the loss function design for the first time.

**Our Contributions.** This paper makes three main contributions to the field of learning adversarial attacks: (i) *The problem in study is novel*. Learning ordered Top-k adversarial attacks is an important problem that reflects the robustness of attacks themselves, but has not been addressed in the literature. (ii) *The proposed adversarial distillation framework is effective,* especially when k is large (such as k = 5, 10). (iii) The proposed knowledge-oriented adversarial distillation framework for a novel problem (ordered Top-k adversarial attacks) with some novel modifications (knowledge-oriented target distributions as "teachers").

## **3. Problem Formulation**

In this section, we first define the white-box attack setting and the widely used C&W method [4] under the Top-1 protocol, to be self-contained. Then we present the ordered Top-k attack formulation. To learn ordered Top-k attacks, we present details of the modified C&W method, CW<sup>k</sup>, as a strong baseline and the proposed AD framework.

# **3.1.** Background on White-Box Targeted Attack and the Top-1 Setting

We focus on image classification tasks using DNNs. Denote by (x, y) a clean input,  $x \in \mathcal{X}$  and its ground-truth label  $y \in \mathcal{Y}$ . For example, in the ImageNet-1000 classification task, x represents a RGB image defined in the lattice of  $224 \times 224$  and we have  $\mathcal{X} \triangleq R^{3 \times 224 \times 224}$ . y is the category label and we have  $\mathcal{Y} \triangleq \{1, \dots, 1000\}$ . Let  $f(\cdot; \Theta)$  be a DNN pretrained with clean training images where  $\Theta$  collects all estimated parameters and is **frozen** in learning adversarial examples. For notation simplicity, we denote by  $f(\cdot)$  a pretrained DNN. The prediction for an input x from  $f(\cdot)$  is usually defined using softmax function by,

$$P = f(x) = softmax(z(x)), \tag{1}$$

where  $P \in R^{|\mathcal{Y}|}$  represents the estimated confidence/probability vector ( $P_c \geq 0$  and  $\sum_c P_c = 1$ ) and z(x) is the logit vector. The predicted label is then inferred by  $\hat{y} = \arg \max_{c \in [1, |\mathcal{Y}|]} P_c$ .

The traditional **Top-1** protocol of learning targeted attacks. For an input (x, y), given a target label  $t \neq y$ , we seek to compute some visually-imperceptible perturbation  $\delta(x, t, f)$  using the pretrained and fixed DNN  $f(\cdot)$  under the white-box setting. White-box attacks assume the complete knowledge of the pretrained DNN f, including its parameter values, architecture, training method, etc. The perturbed example is defined by,

$$x' = x + \delta(x, t, f), \tag{2}$$

which is called **an adversarial example** of x if  $t = \hat{y'} = \arg \max_c f(x')_c$  and the perturbation  $\delta(x, t, f)$  is sufficiently small according to some energy metric.

The C&W Method [4]. Learning  $\delta(x, t, f)$  under the Top-1 protocol is posed as a constrained optimization problem [2, 4],

minimize 
$$\mathcal{E}(\delta) = ||\delta||_p$$
, (3)  
subject to  $t = \arg\max_c f(x+\delta)_c$ ,  
 $x+\delta \in [0,1]^n$ ,

where  $\mathcal{E}(\cdot)$  is defined by a  $\ell_p$  norm (e.g., the  $\ell_2$  norm) and n the size of the input domain (e.g., the number of pixels). To overcome the difficulty (non-linear and nonconvex constraints) of directly solving Eqn. 3, the C&W method expresses it in a different form by designing some loss functions  $L(x') = L(x+\delta)$  such that the first constraint  $t = \arg \max_c f(x')_c$  is satisfied if and only if  $L(x') \leq 0$ . The best loss function proposed by the C&W method is defined by the hinge loss,

$$L_{CW}(x') = \max(0, \max_{c \neq t} z(x')_c - z(x')_t).$$
(4)

which induces penalties when the logit of the target label is not the maximum among all labels.

Then, the learning problem is formulated by,

minimize 
$$||\delta||_p + \lambda \cdot L(x+\delta),$$
 (5)  
subject to  $x+\delta \in [0,1]^n,$ 

which can be solved via back-propagation with the constraint satisfied via introducing a tanh layer.

**Computing Budget:** For the trade-off parameter  $\lambda$  (Eqn. 5), a binary search will be performed during the learning. For example, typically,  $9 \times 30$  and  $9 \times 1000$  are used by the C&W method, which represents 9 tries of different values for  $\lambda$  and 30 or 1000 back-propagation iterations used for each  $\lambda$ .

## 3.2. The Proposed Ordered Top-k Attack Setting

We extend Eqn. 3 for learning ordered Top-k attacks  $(k \ge 1)$ . Denote by  $(t_1, \dots, t_k)$  the ordered Top-k targets  $(t_i \ne y)$ . We have,

minimize 
$$\mathcal{E}(\delta) = ||\delta||_p$$
, (6)  
subject to  $t_i = \arg \max_{c \in [1, |\mathcal{V}|], c \notin \{t_1, \cdots, t_{i-1}\}} f(x+\delta)_c$ ,  
 $i \in \{1, \cdots, k\},$   
 $x + \delta \in [0, 1]^n$ .

Directly solving Eqn. 6 is a challenging task and proper loss functions are entailed, similar in spirit to the approximation approaches widely adopted in the Top-1 protocol, to ensure the first constraint can be satisfied once the optimization is converged (note that the optimization may fail with a given computing budget).

## 3.3. Learning Ordered Top-k Attacks

#### **3.3.1** The Modified C&W Method, $CW^k$

We modify the loss function (Eqn. 4) of the C&W method accordingly to solve Eqn. 6. We have,

$$L_{CW}^{(k)}(x') = \sum_{i=1}^{k} \max\left(0, \max_{j \notin \{t_1, \cdots, t_i\}} z(x')_j - \min_{j \in \{t_1, \cdots, t_i\}} z(x')_j\right)$$
(7)

which covers the vanilla C&W loss (Eqn. 4), *i.e.*, when k = 1,  $L_{CW}(x') = L_{CW}^{(1)}(x')$ . The C&W loss function does not care where the underlying GT label will be as long as it is not in the Top-k. On the one hand, it is powerful in terms of attack success rate. On the other hand, the GT label may be very close to the Top-k, leading to over-confident attacks (Tabel. 1). In addition, it is generic for any given Top-k targets. As we shall show, they are less effective if we select the Top-k targets from the sub-set of labels which are least like the ground-truth label in terms of label semantics.

#### 3.3.2 Knowledge-Oriented Adversarial Distillation

To overcome the shortcomings of the C&W loss function and In our adversarial distillation framework, we adopt the view of point proposed in the network distillation method [12] that the full confidence/probability distribution summarizes the knowledge of a trained DNN. We hypothesize that we can leverage the network distillation framework to learn the ordered Top-k attacks by designing a proper adversarial probability distribution across the entire set of labels that satisfies the specification of the given ordered Top-k targets, and facilitates explicit control of placing the GT label, as well as top-down integration of label semantics.

Consider a given set of Top-k targets,  $\{t_1, \dots, t_k\}$ , denoted by  $P^{AD}$  the adversarial probability distribution in which  $P_{t_i}^{AD} > P_{t_j}^{AD}$  ( $\forall i < j$ ) and  $P_{t_i}^{AD} > P_l^{AD}$  ( $\forall l \notin \{t_1, \dots, t_k\}$ ). The space of candidate distributions are huge. We present a simple knowledge-oriented approach to define the adversarial distribution. We first specify the logit distribution and then compute the probability distribution using softmax. Denote by Z the maximum logit (*e.g.*, Z = 10 in our experiments).

We define the adversarial logits for the ordered Top-k targets by,

 $z_{t_i}^{AD} = Z - (i - 1) \times \gamma, \quad i \in [1, \dots, k],$  (8) where  $\gamma$  is an empirically chosen decreasing factor (*e.g.*,  $\gamma = 0.3$  in our experiments).

For the remaining categories  $l \notin \{t_1, \cdots, t_k\}$ , we define the adversarial logit by,

$$z_l^{AD} = \alpha \times \frac{1}{k} \sum_{i=1}^k s(t_i, l) + \epsilon, \qquad (9)$$

where  $0 \le \alpha < z_{t_k}^{AD}$  is the maximum logit that can be assigned to any j, s(a, b) is the semantic similarity between the label a and label b, and  $\epsilon$  is a small position for numerical consideration (e.g.,  $\epsilon = 1e$ -5). We compute s(a, b) using the cosine distance between the Glove [29] embedding vectors of category names and  $-1 \le s(a, b) \le 1$ . Here, when  $\alpha = 0$ , we discard the semantic knowledge and treat all the remaining categories equally. Note that our design of  $P^{AD}$  is similar in spirit to the label smoothing technique and its variants [37, 30] except that we target attack labels and exploit label semantic knowledge. The design choice is still preliminary, although we observe its effectiveness in experiments. We hope this can simulate more sophisticated work to be explored.

With the adversarial probability distribution  $P^{AD}$  defined above as the target, we use the KL divergence as the loss function in our adversarial distillation framework as done in network distillation [12] and we have,

$$L_{AD}^{(k)}(x') = KL(f(x')||P^{AD}),$$
(10)

and then we follow the same optimization scheme as done in the C&W method (Eqn. 5).

#### 4. Experiments

In this section, we evaluate ordered Top-k attacks with k = 1, 2, 5, 10 in the ImageNet-1000 benchmark [33] using three representative pretrained DNNs: (i) ResNet-50 [11]. ResNets are the most widely used DNNs. (ii) DenseNet-121 [14]. DenseNets are also popular in practice. (iii) AOGNet-12M [21]. AOGNets are grammar-guided networks with the vanilla Bottleneck building bock, which represent an interesting direction of network architecture engineering. So, the attacking results by the proposed methods will be broadly useful for ResNets and DenseNets based deployment in practice and potentially insightful for on-going and future development of more powerful and robust DNNs. Due to computational cost, we choose the small versions of the three different DNNs. Pretrained models of ResNet-50 and DenseNet-121 are from the PyTorch model zoo<sup>1</sup>. Pretrained AOGNet-12M is from their Github repo<sup>2</sup>. We implement the proposed two methods using the AdverTorch toolkit<sup>3</sup>. Our reproducible source code will be released.

**Data.** In ImageNet-1000 [33], there are 50,000 images for validation. To study attacks, we utilize the subset of images for which the predictions of all the three networks, ResNet-50, DenseNet-121 and AOGNet-12M, are correct. To reduce the computational demand, we randomly sample a smaller subset, as commonly done in the literature. We first randomly select 500 categories to enlarge the coverage of categories, and then randomly chose 2 images per selected categories, resulting in 1000 test images in total.

**Settings.** We follow the protocol used in the C&W method. We only test  $\ell_2$  norm as the energy penalty for perturbations in learning (Eqn. 5). But, we evaluate learned adversarial examples in terms of three norms ( $\ell_1$ ,  $\ell_2$  and  $\ell_{\infty}$ ). We test two search schema for the trade-off parameter  $\lambda$  in optimization: both use 9 steps of binary search, and 30 and 1000 iterations of optimization are performed for each trial of  $\lambda$ . In practice, computation budget is an important factor and less computationally expensive ones are usually preferred. Only  $\alpha = 1$  is used in Eqn. 9 in experiments for simplicity due to computational demand. We compare the results under three scenarios proposed in the C&W method [4]: The Best Case settings test the attack against all incorrect classes, and report the target class(es) that was least difficult to attack. The Worst Case settings test the attack against all incorrect classes, and report the target class(es) that was most difficult to attack. The Av*erage Case* settings select the target class(es) uniformly at random among the labels that are not the GT.

<sup>&</sup>lt;sup>1</sup>https://github.com/pytorch/vision/tree/master/torchvision/models

<sup>&</sup>lt;sup>2</sup>https://github.com/iVMCL/AOGNets

<sup>&</sup>lt;sup>3</sup>https://github.com/BorealisAI/advertorch



Figure 3. Learned adversarial examples of ordered Top-5 (top) and vanilla Top-1 (bottom) attacks for ResNet-50 [11] pretrained with clean images. The proposed AD method has smaller perturbation energies and "cleaner" (lower-entropy) prediction distributions than both the modified  $CW^k$  method in Top-5 attacks and the vanilla C&W method in the Top-1 attacks.

#### 4.1. Results for ResNet-50

We first test ordered Top-k attacks using ResNet-50 for the four selected k's. Fig. 3 shows some learned adversarial examples of ordered Top-5 and Top-1 attacks (Top-10 attack examples in Fig. 2). The visualizations are similar across different networks. Table. 2 summarizes the quantitative results and comparisons. For Top-10 attacks, the proposed AD method obtains significantly better results in terms of both ASR and the  $\ell_2$  energy of the added perturbation. For example, the proposed AD method has *relative* 362.3% ASR *improvement* over the strong C&W baseline for the worst case setting. For Top-5 attacks, the AD method obtains significantly better results when the search budget is relatively low (i.e.,  $9 \times 30$ ). For Top-k (k = 1, 2) attacks, both the C&W method and the AD method can achieve 100% ASR, but the AD method has consistently lower energies of the added perturbation, i.e., finding more effective attacks and richer perturbations.

#### 4.2. Are all Top-k targets equally difficult to attack?

Intuitively, we understand that they should not be equally difficult. We conduct some experiments to test this hypothesis. In particular, we test whether the label semantic knowledge can help identify the weak spots of different attack

Protocol	Attack Method		st Case		Averag	ge Case		Worst Case					
11010001	- maint method	ASR	$\ell_1$	$\ell_2$	$\ell_{\infty}$	ASR	$\ell_1$	$\ell_2$	$\ell_{\infty}$	ASR	$\ell_1$	$\ell_2$	$\ell_{\infty}$
	$CW^{k}_{9 \times 30}$	0	N.A.	N.A.	N.A.	0	N.A.	N.A.	N.A.	0	N.A.	N.A.	N.A.
	$AD_{9 \times 30}$	0.8	2579	8.18	0.096	0.16	2579	8.18	0.096	0	N.A.	N.A.	N.A.
Top-10	$CW^{k}_{9\times 100}$	43.4	2336	7.83	0.109	11.8	2330	7.82	0.109	0.1	2479	8.26	0.119
-	$AD_{9 \times 100}$	91.8	1677	5.56	0.088	51.2	1867	6.14	0.098	5.6	2021	6.62	0.110
	$CW^{k}_{9 \times 1000}$	97.7	1525	5.26	0.092	64.5	1742	5.99	0.103	20.4	1898	6.61	0.120
	$AD_{9 \times 1000}$	99.8	678	2.45	0.060	98.4	974	3.45	0.081	94.3	1278	4.48	0.103
	Improvement	2.1 (3.0%)		2.81 (53.4%)		33.9 (52.6%)		2.54 (42.4%)		73.9 (362.3%)		1.13 (17.1%)	
Top-5	$CW_{9\times 30}^{k}$	75.8	2370	7.76	0.083	29.34	2425	7.94	0.086	0.7	2553	8.37	0.094
	$AD_{9 \times 30}$	96.1	1060	3.58	0.056	80.68	1568	5.13	0.070	49.8	2215	7.07	0.087
	$CW^{k}_{9 \times 1000}$	100	437	1.59	0.044	100	600	2.16	0.058	100	779	2.77	0.074
	$AD_{9 \times 1000}$	100	285	1.09	0.034	100	359	1.35	0.043	100	456	1.68	0.055
	$CW^{k}_{9 \times 30}$	99.9	1002	3.40	0.037	99.36	1504	4.95	0.050	97.9	2007	6.52	0.065
Ton-2	$AD_{9 \times 30}$	99.9	308	1.12	0.028	99.5	561	1.94	0.037	98.4	873	2.92	0.049
100 2	$CW^{k}_{9 \times 1000}$	100	185	0.72	0.025	100	241	0.91	0.033	100	303	1.12	0.042
	$AD_{9 \times 1000}$	100	137	0.56	0.022	100	174	0.70	0.028	100	220	0.85	0.035
	$C\&W_{9\times 30}$	100	209.7	0.777	0.022	99.92	354.1	1.273	0.031	99.9	560.9	1.987	0.042
	$AD_{9 \times 30}$	100	140.9	0.542	0.018	99.9	184.6	0.696	0.025	99.9	238.6	0.880	0.032
<b>T</b> 1	$C\&W_{9\times 1000}$	100	95.6	0.408	0.017	100	127.2	0.516	0.023	100	164.1	0.635	0.030
Top-1	$AD_{9 \times 1000}$	100	81.3	0.380	0.016	100	109.6	0.472	0.023	100	143.9	0.579	0.029
	FGSM	2.3	9299	24.1	0.063	0.46	9299	24.1	0.063	0	N.A.	N.A.	N.A.
	$PGD_{10}$	99.6	4691	14.1	0.063	88.1	4714	14.2	0.063	57.1	4748	14.3	0.063
	$MIFGSM_{10}$	100	5961	17.4	0.063	99.98	6082	17.6	0.063	99.9	6211	17.9	0.063

Table 2. Results and comparisons under the ordered Top-k targeted attack protocol using randomly selected and ordered k targets (GT exclusive) in ImageNet using **ResNet-50**. For Top-1 attacks, we also compare with three state-of-the-art untargeted attack methods, FGSM [10], PGD [24] and MIFGSM [7]. 10 iterations are used for both PGD and MIFGSM.

methods, and whether the proposed AD method can gain more in those weak spots. We test Top-5 using ResNet-50. Table. 3 summarizes the results. Similar results are observed for DenseNets and AOGNets. We observe that for the  $9 \times 30$  search budget, attacks are more challenging if the Top-5 targets are selected from the least-like set in terms of the label semantic similarity (Eqn. 9), or from the lowestscore set in terms of prediction scores on clean images.

## 4.3. Results for DenseNet-121 and AOGNet-12M

Overall, we obtain similar results for DenseNet-121 summarized in Table. 4. For AOGNet-12M (Table. 5), the proposed AD does not show improvement as significant as for the other two networks, especially for Top-10.

## 5. Conclusions and Discussions

This paper proposes to extend the traditional Top-1 targeted attack setting to the ordered Top-k setting  $(k \ge 1)$ under the white-box attack protocol. The ordered Top-ktargeted attacks can improve the robustness of attacks themselves. To our knowledge, it is the first work studying this ordered Top-k attacks. To learn the ordered Top-k attacks, we present a conceptually simple yet effective adversarial distillation framework motivated by network distillation. We also develop a modified C&W method as the strong baseline for the ordered Top-k targeted attacks. In experiments, the proposed method is tested in ImageNet-1000 using two popular DNNs, ResNet-50 and DenseNet-121, with consistently better results obtained. We investigate the effectiveness of label semantic knowledge in designing the adversarial distribution for distilling the ordered Top-k targeted attacks.

Protocol	Similarity	Method	ASR	$\ell_1$	$\ell_2$	$\ell_{\infty}$
	Most like	${ m CW}^k{}_{9 imes 30} \ { m AD}_{9 imes 30}$	80 <b>96.5</b>	1922 1286	6.30 4.20	0.066 0.054
Label similarity		$\begin{array}{c} \mathrm{CW}^{k}_{9\times1000} \\ \mathrm{AD}_{9\times1000} \end{array}$	100 100	392 <b>277</b>	1.43 <b>1.05</b>	0.042 <b>0.035</b>
	Least like	$\mathrm{CW}^{k}_{9 imes 30}$ AD <sub>9 imes 30</sub>	27.1 <b>77.1</b>	2418 1635	7.90 5.35	0.085 0.072
	Louse into	$\begin{array}{c} \mathrm{CW}^{k}_{9\times1000} \\ \mathrm{AD}_{9\times1000} \end{array}$	100 100	596 <b>370</b>	2.15 <b>1.39</b>	0.060 <b>0.045</b>
	Highest	$\mathrm{CW}^{k}_{9 imes 30}$ $\mathrm{AD}_{9 imes 30}$	93 <b>99.9</b>	1546 1182	4.98 3.78	0.042 0.039
Prediction Score	mgneor	$\mathrm{CW}^{k}_{9 imes 1000}$ $\mathrm{AD}_{9 imes 1000}$	100 100	205 <b>170</b>	0.75 <b>0.65</b>	0.025 <b>0.023</b>
	Lowest	${{ m CW}^k}_{9 imes 30} \ { m AD}_{9 imes 30}$	13.4 <b>68.6</b>	2231 1791	7.30 5.86	0.082 0.077
		$\mathrm{CW}^{k}_{9 imes 1000}$ AD <sub>9 imes 1000</sub>	100 100	621 <b>392</b>	2.25 1.47	0.064 <b>0.047</b>

Table 3. Results of ordered Top-5 targeted attacks for ResNet-50 with targets being selected based on (Top) label similarity, which uses 5 most-like labels and 5 least-like labels as targets respectively, and (Bottom) prediction score of clean image, which uses 5 highest-score labels and 5-lowest score labels. In both cases, GT labels are exclusive.

**Discussions.** We have shown that the proposed AD method is generally applicable to learn ordered Top-k attacks. But, we note that the two components of the AD framework are in their simplest forms in this paper, and need to be more thoroughly studied: designing more informative adversarial distributions to guide the optimization to learn adversarial examples better and faster, and investigating loss functions other than KL divergence such as the Jensen-Shannon (JS) divergence or the Earth-Mover distance. On the other hand, we observed that the proposed

Protocol	Method	Best Case					Av	erage Case		Worst Case			
11000001	moulou	ASR	$\ell_1$	$\ell_2$	$\ell_{\infty}$	ASR	$\ell_1$	$\ell_2$	$\ell_{\infty}$	ASR	$\ell_1$	$\ell_2$	$\ell_{\infty}$
Top-10	$CW^{k}_{9 \times 30}$ AD <sub>9×30</sub>	0	N.A. 2674	N.A. 8.48	N.A. 0.102	0 0.26	N.A. 2674	N.A. 8.48	N.A. 0.102	0	N.A. N.A.	N.A. N.A.	N.A. N.A.
	${\mathop{\rm CW^{k}}_{9 imes100}} {\mathop{\rm AD}_{9 imes100}}$	45.3 96.9	2320 1864	7.77 6.21	0.108 0.101	12.42 61.4	2338 2331	7.83 7.71	0.109 0.123	0 13.2	N.A. 2878	N.A. 9.58	N.A. 0.156
	$CW^{k}_{9 \times 1000}$ AD <sub>9 × 1000</sub> Improvement	100 100	1163 627	4.09 <b>2.26</b> 1.83 (44.7%)	0.085 <b>0.058</b>	97.16 100	1640 <b>900</b>	5.67 <b>3.19</b> 2.48 (43.7%)	0.107 <b>0.076</b>	85.8 100	2157 1250	7.39 <b>4.37</b> 3.02 (40.9%)	0.130 0.100
Top-5	$\begin{array}{c} \mathrm{CW}^{k}{}_{9\times 30} \\ \mathrm{AD}_{9\times 30} \end{array}$	96.6 97.7	2161 6413	7.09 2.14	0.071 0.043	73.68	2329 1063	7.65 3.57	0.080 0.057	35.6 83.3	2530 1636	8.28 5.35	0.088 0.072
	$\begin{array}{c} \mathrm{CW}^{k}{}_{9\times 1000} \\ \mathrm{AD}_{9\times 1000} \end{array}$	100 100	392 273	1.42 1.05	0.040 <b>0.033</b>	100 100	527 <b>344</b>	1.89 1.29	0.052 0.042	100 100	669 <b>425</b>	2.37 1.57	0.065 0.052
Top-2	$\begin{array}{c} \mathrm{CW}^k{}_{9 imes 30} \\ \mathrm{AD}_{9 imes 30} \end{array}$	99.9 99.9	549 199	1.92 0.74	0.033 0.023	99.72 99.8	1058 249	3.54 0.92	0.042 0.029	99.4 99.7	1640 308	5.35 1.12	0.051 0.037
	$\begin{array}{c} \mathrm{CW}^{k}_{9\times1000} \\ \mathrm{AD}_{9\times1000} \end{array}$	100 100	146 121	0.58 <b>0.52</b>	0.022 <b>0.021</b>	100 100	187 153	0.72 <b>0.62</b>	0.029 <b>0.027</b>	100 100	230 187	0.86 <b>0.74</b>	0.037 <b>0.034</b>
	$\begin{array}{c} C\&W_{9\times 30}\\ AD_{9\times 30} \end{array}$	99.9 99.9	188.6 136.4	0.694 0.523	0.019 0.017	99.9 99.9	279.4 181.8	1.008 0.678	0.028 0.024	99.9 99.9	396.5 240.0	1.404 0.870	0.037 0.031
Top-1	$\begin{array}{c} C\&W_{9\times 1000}\\ AD_{9\times 1000} \end{array}$	100 100	98.5 <b>83.8</b>	0.415 <b>0.384</b>	0.016 0.016	100 100	132.3 115.9	0.528 <b>0.485</b>	0.023 0.023	100 100	174.8 <b>158.69</b>	0.657 <b>0.610</b>	0.030 0.030
	FGSM PGD <sub>10</sub> MIFGSM <sub>10</sub>	6.4 100 100	9263 4617 5979	24.0 14.2 17.6	0.063 0.063 0.063	1.44 97.2 100	9270 4716 6095	24.0 14.2 17.6	0.063 0.063 0.063	0 87.6 100	N.A. 4716 6218	N.A. 14.2 17.9	N.A. 0.063 0.063

Table 4. Results and comparisons using **DenseNet-121** [14] under the ordered Top-*k* targeted attack protocol using randomly selected and ordered 10 targets (GT exclusive). For Top-1 attacks, we also compare with three state-of-the-art untargeted attack methods, FGSM [10], PGD [24] and MIFGSM [7]. 10 iterations are used for both PGD and MIFGSM.

Protocol	Method		Bes	st Case			Avera	ge Case		Worst Case				
11010001	Method	ASR	$\ell_1$	$\ell_2$	$\ell_{\infty}$	ASR	$\ell_1$	$\ell_2$	$\ell_{\infty}$	ASR	$\ell_1$	$\ell_2$	$\ell_{\infty}$	
	$\mathrm{CW}^k{}_{9 imes 30}$	0	N.A.	N.A.	N.A.	0	N.A.	N.A.	N.A.	0	N.A.	N.A.	N.A.	
	$AD_{9 \times 30}$	0	N.A.	N.A.	N.A.	0	N.A.	N.A.	N.A.	0	N.A.	N.A.	N.A.	
Top-10	$CW^k_{9 \times 100}$	43.5	2221	7.50	0.11	12.56	2246	7.57	0.11	0.1	2657	8.56	0.13	
	$AD_{9 \times 100}$	42.8	1920	6.47	0.11	11.64	1947	6.56	0.11	0	N.A.	N.A.	N.A.	
	$CW^{k}_{9 \times 1000}$	100	1018	3.68	0.084	99.94	1393	4.93	0.10	99.9	1816	6.31	0.12	
	$AD_{9 \times 1000}$	100	651	2.40	0.065	99.84	947	3.43	0.088	99.3	1317	4.72	0.12	
Top-5	$CW^{k}_{9 \times 30}$	37.1	1809	6.08	0.078	10.84	1918	6.41	0.079	0.1	2290	7.80	0.090	
	$AD_{9 \times 30}$	82.7	1126	3.84	0.063	46.26	1599	5.28	0.074	10.2	2396	7.71	0.093	
	$CW^k_{9 \times 1000}$	100	355	1.33	0.041	100	468	1.73	0.053	100	588	2.15	0.068	
	$AD_{9 \times 1000}$	100	199	0.81	0.029	100	258	1.01	0.037	100	327	1.26	0.047	
	$CW^{k}_{9 \times 30}$	99.6	644	2.33	0.046	97.22	1176	4.04	0.059	91	1918	6.38	0.074	
Ton-2	$AD_{9 \times 30}$	99.9	205	0.78	0.024	98.72	283	1.06	0.032	96.2	391	1.43	0.041	
100 -	$CW^k_{9 \times 1000}$	100	131	0.53	0.021	100	178	0.70	0.029	100	232	0.90	0.038	
	$AD_{9\times 1000}$	100	96	0.44	0.019	100	124	0.54	0.025	100	154	0.65	0.032	
	$C\&W_{9\times 30}$	99.9	256	0.946	0.025	99.64	532	1.90	0.039	99.1	932	3.25	0.054	
	$AD_{9 \times 30}$	99.9	141	0.541	0.019	99.8	190	0.720	0.026	99.5	246	0.922	0.034	
Top-1	$C\&W_{9\times 1000}$	100	80	0.36	0.016	100	114	0.48	0.023	100	153	0.61	0.030	
	$AD_{9\times 1000}$	100	62	0.312	0.016	100	86	0.398	0.021	100	115	0.496	0.028	
	FGSM	1.7	9254	24.0	0.0625	0.34	9254	24.0	0.0625	0	N.A.	N.A.	N.A.	
	$PGD_{10}$	100	4685	14.1	0.0625	98.16	4698	14.2	0.0625	91.2	4714	14.2	0.0625	
	$MIFGSM_{10}$	100	5940	17.3	0.0625	99.92	6046	17.5	0.0625	99.6	6165	17.8	0.0625	

Table 5. Results and comparisons under the ordered Top-k targeted attack protocol using randomly selected and ordered 10 targets (GT exclusive) in ImageNet using **AOGNet-12M**.

AD method is more effective when computation budget is limited (*e.g.*, using the  $9 \times 30$  search scheme). This leads to the theoretically and computationally interesting question whether different attack methods all will work comparably well if the computation budget is not limited. Of course, in practice, we prefer more powerful ones when only limited computation budget is allowed. Furthermore, we observed that both the modified C&W method and the AD method largely do not work in learning Top-k ( $k \geq 20$ ) attacks with the two search schema ( $9 \times 30$  and  $9 \times 1000$ ). We are working on addressing the aforementioned issues to test the Top-k ( $k \ge 20$ ) cases, thus providing a thorough empirical answer to the question: how aggressive can adversarial attacks be?

### Acknowledgement

This work was supported in part by ARO Grant W911NF1810295 and NSF IIS-1909644. The views presented in this paper are those of the authors and should not be interpreted as representing any funding agencies.

## References

- Naveed Akhtar and Ajmal S. Mian. Threat of adversarial attacks on deep learning in computer vision: A survey. *IEEE Access*, 6:14410–14430, 2018. 3
- [2] Anish Athalye and Ilya Sutskever. Synthesizing robust adversarial examples. *CoRR*, abs/1707.07397, 2017. 1, 4
- [3] Cristian Bucila, Rich Caruana, and Alexandru Niculescu-Mizil. Model compression. In Proceedings of the Twelfth ACM SIGKDD International Conference on Knowledge Discovery and Data Mining, Philadelphia, PA, USA, August 20-23, 2006, pages 535–541, 2006. 2, 3
- [4] Nicholas Carlini and David A. Wagner. Towards evaluating the robustness of neural networks. In 2017 IEEE Symposium on Security and Privacy, SP 2017, San Jose, CA, USA, May 22-26, 2017, pages 39–57, 2017. 1, 2, 3, 4, 5
- [5] Pin-Yu Chen, Yash Sharma, Huan Zhang, Jinfeng Yi, and Cho-Jui Hsieh. EAD: elastic-net attacks to deep neural networks via adversarial examples. In *AAAI*, pages 10–17. AAAI Press, 2018. 3
- [6] Pin-Yu Chen, Huan Zhang, Yash Sharma, Jinfeng Yi, and Cho-Jui Hsieh. ZOO: zeroth order optimization based blackbox attacks to deep neural networks without training substitute models. In *AISec* @*CCS*, pages 15–26. ACM, 2017. 3
- [7] Yinpeng Dong, Fangzhou Liao, Tianyu Pang, Hang Su, Jun Zhu, Xiaolin Hu, and Jianguo Li. Boosting adversarial attacks with momentum. In *CVPR*, pages 9185–9193. IEEE Computer Society, 2018. 2, 3, 7, 8
- [8] Yinpeng Dong, Tianyu Pang, Hang Su, and Jun Zhu. Evading defenses to transferable adversarial examples by translation-invariant attacks. In *IEEE Conference on Computer Vision and Pattern Recognition, CVPR 2019, Long Beach, CA, USA, June 16-20, 2019*, pages 4312–4321, 2019.
   3
- [9] Javid Ebrahimi, Anyi Rao, Daniel Lowd, and Dejing Dou. Hotflip: White-box adversarial examples for text classification. In Proceedings of the 56th Annual Meeting of the Association for Computational Linguistics, ACL 2018, Melbourne, Australia, July 15-20, 2018, Volume 2: Short Papers, pages 31–36, 2018. 3
- [10] Ian J. Goodfellow, Jonathon Shlens, and Christian Szegedy. Explaining and harnessing adversarial examples. In 3rd International Conference on Learning Representations, ICLR 2015, San Diego, CA, USA, May 7-9, 2015, Conference Track Proceedings, 2015. 2, 3, 7, 8
- [11] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image recognition. In *IEEE Conference on Computer Vision and Pattern Recognition* (CVPR), 2016. 1, 2, 3, 5, 6
- [12] Geoffrey E. Hinton, Oriol Vinyals, and Jeffrey Dean. Distilling the knowledge in a neural network. *CoRR*, abs/1503.02531, 2015. 2, 3, 5
- [13] Kurt Hornik, Maxwell B. Stinchcombe, and Halbert White. Multilayer feedforward networks are universal approximators. *Neural Networks*, 2(5):359–366, 1989. 3
- [14] Gao Huang, Zhuang Liu, Laurens van der Maaten, and Kilian Q Weinberger. Densely connected convolutional net-

works. In Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition, 2017. 1, 3, 5, 8

- [15] Andrew Ilyas, Shibani Santurkar, Dimitris Tsipras, Logan Engstrom, Brandon Tran, and Aleksander Madry. Adversarial examples are not bugs, they are features. *CoRR*, abs/1905.02175, 2019. 3
- [16] Harini Kannan, Alexey Kurakin, and Ian J. Goodfellow. Adversarial logit pairing. *CoRR*, abs/1803.06373, 2018. 3
- [17] Jernej Kos, Ian Fischer, and Dawn Song. Adversarial examples for generative models. In 2018 IEEE Security and Privacy Workshops, SP Workshops 2018, San Francisco, CA, USA, May 24, 2018, pages 36–42, 2018. 3
- [18] Alex Krizhevsky, Ilya Sutskever, and Geoffrey E. Hinton. Imagenet classification with deep convolutional neural networks. In *Neural Information Processing Systems (NIPS)*, pages 1106–1114, 2012. 1
- [19] Alexey Kurakin, Ian J. Goodfellow, and Samy Bengio. Adversarial examples in the physical world. In *ICLR (Workshop)*. OpenReview.net, 2017. 3
- [20] Yann LeCun, Leon Bottou, Yoshua Bengio, and Patrick Haffner. Gradient-based learning applied to document recognition. *Proceedings of the IEEE*, 86(11):2278–2324, 1998.
- [21] Xilai Li, Xi Song, and Tianfu Wu. Aognets: Compositional grammatical architectures for deep learning. In *IEEE Conference on Computer Vision and Pattern Recognition, CVPR* 2019, Long Beach, CA, USA, June 16-20, 2019, pages 6220– 6230, 2019. 1, 3, 5
- [22] Yen-Chen Lin, Zhang-Wei Hong, Yuan-Hong Liao, Meng-Li Shih, Ming-Yu Liu, and Min Sun. Tactics of adversarial attack on deep reinforcement learning agents. In Proceedings of the Twenty-Sixth International Joint Conference on Artificial Intelligence, IJCAI 2017, Melbourne, Australia, August 19-25, 2017, pages 3756–3762, 2017. 3
- [23] Yanpei Liu, Xinyun Chen, Chang Liu, and Dawn Song. Delving into transferable adversarial examples and blackbox attacks. In 5th International Conference on Learning Representations, ICLR 2017, Toulon, France, April 24-26, 2017, Conference Track Proceedings, 2017. 3
- [24] Aleksander Madry, Aleksandar Makelov, Ludwig Schmidt, Dimitris Tsipras, and Adrian Vladu. Towards deep learning models resistant to adversarial attacks. In 6th International Conference on Learning Representations, ICLR 2018, Vancouver, BC, Canada, April 30 - May 3, 2018, Conference Track Proceedings, 2018. 2, 3, 7, 8
- [25] Anh Mai Nguyen, Jason Yosinski, and Jeff Clune. Deep neural networks are easily fooled: High confidence predictions for unrecognizable images. In *IEEE Conference on Computer Vision and Pattern Recognition, CVPR 2015, Boston, MA, USA, June 7-12, 2015*, pages 427–436, 2015. 1
- [26] Nicolas Papernot and Patrick D. McDaniel. Extending defensive distillation. *CoRR*, abs/1705.05264, 2017. 3
- [27] Nicolas Papernot, Patrick D. McDaniel, Ian J. Goodfellow, Somesh Jha, Z. Berkay Celik, and Ananthram Swami. Practical black-box attacks against machine learning. In Proceedings of the 2017 ACM on Asia Conference on Computer and Communications Security, AsiaCCS 2017, Abu

Dhabi, United Arab Emirates, April 2-6, 2017, pages 506–519, 2017. 3

- [28] Nicolas Papernot, Patrick D. McDaniel, Xi Wu, Somesh Jha, and Ananthram Swami. Distillation as a defense to adversarial perturbations against deep neural networks. In *IEEE Symposium on Security and Privacy*, pages 582–597. IEEE Computer Society, 2016. 2, 3
- [29] Jeffrey Pennington, Richard Socher, and Christopher D. Manning. Glove: Global vectors for word representation. In *In EMNLP*, 2014. 3, 5
- [30] Gabriel Pereyra, George Tucker, Jan Chorowski, Lukasz Kaiser, and Geoffrey E. Hinton. Regularizing neural networks by penalizing confident output distributions. *CoRR*, abs/1701.06548, 2017. 3, 5
- [31] Yao Qin, Nicholas Carlini, Garrison W. Cottrell, Ian J. Good-fellow, and Colin Raffel. Imperceptible, robust, and targeted adversarial examples for automatic speech recognition. In *Proceedings of the 36th International Conference on Machine Learning, ICML 2019, 9-15 June 2019, Long Beach, California, USA*, pages 5231–5240, 2019. 3
- [32] Jérôme Rony, Luiz G. Hafemann, Luiz S. Oliveira, Ismail Ben Ayed, Robert Sabourin, and Eric Granger. Decoupling direction and norm for efficient gradient-based L2 adversarial attacks and defenses. In *IEEE Conference on Computer Vision and Pattern Recognition, CVPR 2019, Long Beach, CA, USA, June 16-20, 2019*, pages 4322–4330, 2019. 3
- [33] Olga Russakovsky, Jia Deng, Hao Su, Jonathan Krause, Sanjeev Satheesh, Sean Ma, Zhiheng Huang, Andrej Karpathy, Aditya Khosla, Michael Bernstein, Alexander C. Berg, and Li Fei-Fei. ImageNet Large Scale Visual Recognition Challenge. *International Journal of Computer Vision (IJCV)*, 115(3):211–252, 2015. 3, 5
- [34] Mahmood Sharif, Sruti Bhagavatula, Lujo Bauer, and Michael 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, Vienna, Austria, October 24-28, 2016, pages 1528–1540, 2016. 3
- [35] Jiawei Su, Danilo Vasconcellos Vargas, and Kouichi Sakurai. One pixel attack for fooling deep neural networks. *IEEE Trans. Evolutionary Computation*, 23(5):828–841, 2019.
- [36] Christian Szegedy, Sergey Ioffe, and Vincent Vanhoucke. Inception-v4, inception-resnet and the impact of residual connections on learning. *CoRR*, abs/1602.07261, 2016.
- [37] Christian Szegedy, Vincent Vanhoucke, Sergey Ioffe, Jonathon Shlens, and Zbigniew Wojna. Rethinking the inception architecture for computer vision. *CoRR*, abs/1512.00567, 2015. 3, 5
- [38] Christian Szegedy, Wojciech Zaremba, Ilya Sutskever, Joan Bruna, Dumitru Erhan, Ian J. Goodfellow, and Rob Fergus. Intriguing properties of neural networks. In *ICLR*, 2014. 1, 3
- [39] Cihang Xie, Jianyu Wang, Zhishuai Zhang, Yuyin Zhou, Lingxi Xie, and Alan L. Yuille. Adversarial examples for semantic segmentation and object detection. In *IEEE International Conference on Computer Vision, ICCV 2017, Venice, Italy, October 22-29, 2017*, pages 1378–1387, 2017. 3

- [40] Cihang Xie, Yuxin Wu, Laurens van der Maaten, Alan L. Yuille, and Kaiming He. Feature denoising for improving adversarial robustness. In *IEEE Conference on Computer Vision and Pattern Recognition, CVPR 2019, Long Beach, CA, USA, June 16-20, 2019*, pages 501–509, 2019. 3
- [41] Cihang Xie, Zhishuai Zhang, Yuyin Zhou, Song Bai, Jianyu Wang, Zhou Ren, and Alan L. Yuille. Improving transferability of adversarial examples with input diversity. In *IEEE Conference on Computer Vision and Pattern Recognition*, *CVPR 2019, Long Beach, CA, USA, June 16-20, 2019*, pages 2730–2739, 2019. 3
- [42] Kaidi Xu, Sijia Liu, Pu Zhao, Pin-Yu Chen, Huan Zhang, Deniz Erdogmus, Yanzhi Wang, and Xue Lin. Structured adversarial attack: Towards general implementation and better interpretability. *CoRR*, abs/1808.01664, 2018. 3
- [43] Chiyuan Zhang, Samy Bengio, Moritz Hardt, Benjamin Recht, and Oriol Vinyals. Understanding deep learning requires rethinking generalization. *CoRR*, abs/1611.03530, 2016. 3