



DeepSearch: A Simple and Effective Blackbox Attack for Deep Neural Networks

Fuyuan Zhang
MPI-SWS, Germany
fuyuan@mpi-sws.org

Sankalan Pal Chowdhury
MPI-SWS, Germany
sankalan@mpi-sws.org

Maria Christakis
MPI-SWS, Germany
maria@mpi-sws.org

ABSTRACT

Although deep neural networks have been very successful in image-classification tasks, they are prone to adversarial attacks. To generate adversarial inputs, there has emerged a wide variety of techniques, such as black- and whitebox attacks for neural networks. In this paper, we present DeepSearch, a novel fuzzing-based, query-efficient, blackbox attack for image classifiers. Despite its simplicity, DeepSearch is shown to be more effective in finding adversarial inputs than state-of-the-art blackbox approaches. DeepSearch is additionally able to generate the most subtle adversarial inputs in comparison to these approaches.

CCS CONCEPTS

• **Computing methodologies** → **Neural networks**; • **Software and its engineering** → **Software testing and debugging**.

KEYWORDS

Adversarial Attack, Deep Neural Networks, Blackbox Fuzzing

ACM Reference Format:

Fuyuan Zhang, Sankalan Pal Chowdhury, and Maria Christakis. 2020. DeepSearch: A Simple and Effective Blackbox Attack for Deep Neural Networks. In *Proceedings of the 28th ACM Joint European Software Engineering Conference and Symposium on the Foundations of Software Engineering (ESEC/FSE '20)*, November 8–13, 2020, Virtual Event, USA. ACM, New York, NY, USA, 13 pages. <https://doi.org/10.1145/3368089.3409750>

1 INTRODUCTION

Deep neural networks have been impressively successful in pattern recognition and image classification [30, 39, 42]. However, it is intriguing that deep neural networks are extremely vulnerable to adversarial attacks [72]. In fact, even very subtle perturbations of a correctly classified image, imperceptible to the human eye, may cause a deep neural network to change its prediction. This poses serious security risks to deploying deep neural networks in safety critical applications.

Various adversarial attacks have been developed to evaluate the vulnerability of neural networks against adversarial perturbations. Early work on generating adversarial examples has focused on whitebox attacks [12, 24, 40, 49, 51, 57]. In the whitebox setting, attackers have full access to the network under evaluation, which

enables them to calculate gradients of the network. Many gradient-based attacks have been shown to be highly effective. However, in several real-world scenarios, having complete access to network parameters is not realistic. This has motivated the development of blackbox adversarial attacks.

In the blackbox setting, attackers assume no knowledge about the network structure or its parameters and may only query the target network for its prediction when given particular inputs. One important metric to measure the efficiency of blackbox attacks is the number of queries needed, because queries are essentially time and monetary costs for attackers, e.g., each query to an online, commercial machine-learning service costs money. Evaluating the robustness of deep neural networks in a query-limited blackbox setting is already standard. Gradient-estimation-based blackbox attacks [5, 13], although effective, require a huge number of queries, which makes generating an attack too costly. Various state-of-the-art blackbox attacks (e.g., [27, 32, 33, 50]) can already achieve successful attacks with low number of queries. However, constructing query-efficient blackbox attacks is still open and challenging.

In this paper, we develop a blackbox fuzzing-based technique for evaluating adversarial robustness of neural networks. The two key challenges of applying fuzzing here are (1) to maintain a high attack success rate, and (2) to require a low number of queries. In many cases, without careful guidance while searching, a naive fuzzing approach, e.g., random fuzzing, is not able to find adversarial examples even after a huge number of queries. To improve attack success rate, we introduce carefully designed feedback to guide our search so that images are efficiently fuzzed toward the decision boundaries. To reduce the number of queries, we adapt hierarchical grouping [50] to our setting so that multiple dimensions can be fuzzed simultaneously, which dramatically reduces query numbers. Furthermore, a refinement step, which can be viewed as a backward search step for fuzzing, can effectively reduce distortion of adversarial examples. Therefore, we extend fuzz testing and show how to apply it on neural networks in a blackbox setting.

Our approach. Inspired by the linear explanation of adversarial examples [24], we develop DeepSearch, a simple, yet effective, query-efficient, blackbox attack, which is based on feedback-directed fuzzing. DeepSearch targets deep neural networks for image classification. Our attack is constrained by the L_∞ distance and only queries the attacked network for its prediction scores on perturbed inputs. The design of our approach is based on the following three aspects:

- (1) **Feedback-directed fuzzing:** Starting from a correctly classified image, DeepSearch strategically mutates its pixels to values that are more likely to lead to an adversarial input. The



This work is licensed under a Creative Commons Attribution International 4.0 License.

ESEC/FSE '20, November 8–13, 2020, Virtual Event, USA

© 2020 Copyright held by the owner/author(s).

ACM ISBN 978-1-4503-7043-1/20/11.

<https://doi.org/10.1145/3368089.3409750>

fuzzing process continues until it either finds an adversarial input or it reaches the query limit.

- (2) *Iterative refinement*: Once an adversarial input is found, our approach starts a refinement step to reduce the L_∞ distance of this input. The iteration of refinement continues until either the query limit is reached or some termination criterion is met. Our evaluation shows that iterative refinement is able to find subtle adversarial inputs generated by only slightly perturbing pixels in the original image.
- (3) *Query reduction*: By utilizing the spatial regularities in input images, DeepSearch adapts an existing hierarchical-grouping strategy [50] to our setting and dramatically reduces the number of queries for constructing successful attacks. The query-reduction step significantly improves the efficiency of our fuzzing and refinement process.

We evaluate DeepSearch against four state-of-the-art blackbox attacks in a query-limited setting, where attackers have only a limited query budget to construct attacks. For our evaluation, we use three popular datasets, namely SVHN [53], CIFAR-10 [38], and ImageNet [61]. For SVHN and CIFAR-10, we further attack neural networks with state-of-the-art defenses based on adversarial training [49]. Our experimental results show that DeepSearch is the most effective in attacking both defended and undefended neural networks. Moreover, it outperforms the other four attacks. Although it is important to develop defense techniques against blackbox adversarial attacks, it is not the focus of this paper and we leave it for future work.

Contributions. We make the following contributions:

- (1) We present a simple, yet very effective, fuzzing-based blackbox attack for deep neural networks.
- (2) We perform an extensive evaluation demonstrating that DeepSearch is more effective in finding adversarial examples than state-of-the-art blackbox approaches.
- (3) We show that the refinement step in our approach gives DeepSearch the advantage of finding the most subtle adversarial examples in comparison to related approaches.
- (4) We show that the hierarchical-grouping strategy is effective for query reduction in our setting.

Outline. The next section briefly introduces background. In Sect. 3, we present DeepSearch for binary classifiers, that is, networks that classify inputs into two classes. Sect. 4 generalizes the technique to multiclass classifiers, which classify inputs into multiple classes. In Sect. 5, we extend our technique with iterative refinement such that the generated adversarial examples are more subtle. We adapt hierarchical grouping for query reduction in Sect. 6. We present our experimental evaluation in Sect. 7, discuss related work in Sect. 8, and conclude in Sect. 9.

2 BACKGROUND

In this section, we introduce some notation and terminology. Let \mathbb{R}^n be the n -dimensional vector space for input images. We represent images as column vectors $\mathbf{x} = (x_1, \dots, x_n)^T$, where $x_i \in \mathbb{R}$ ($1 \leq i \leq n$) is the i th coordinate of \mathbf{x} . We also write $\mathbf{x}(i)$ to denote the i th coordinate x_i , i.e., $\mathbf{x}(i) = x_i$, and each such coordinate represents an image pixel. Now, let $C_m = \{l_1, \dots, l_m\}$ be a set of labels for m

classes, where l_i is the label for the i th class ($1 \leq i \leq m$). A deep neural network that classifies images from \mathbb{R}^n into m classes in C_m is essentially a function $\mathcal{N} : \mathbb{R}^n \rightarrow C_m$. For an input $\mathbf{x} \in \mathbb{R}^n$, $\mathcal{N}(\mathbf{x})$ is the label that the network assigns to \mathbf{x} .

Assume that input \mathbf{x} is correctly classified, and \mathbf{x}' is generated by applying subtle perturbations to \mathbf{x} . These perturbations are subtle when the distance between \mathbf{x} and \mathbf{x}' in \mathbb{R}^n is sufficiently small according to a distance metric. When this is so and $\mathcal{N}(\mathbf{x}) \neq \mathcal{N}(\mathbf{x}')$, we say that \mathbf{x}' is an *adversarial example* [72]. In other words, the network is tricked into classifying \mathbf{x}' into a different class than \mathbf{x} even though they are very similar.

In this paper, we use the L_∞ distance metric. The L_∞ distance between \mathbf{x} and \mathbf{x}' is defined as the maximum of their differences along any coordinate dimension i ($1 \leq i \leq n$):

$$\|\mathbf{x} - \mathbf{x}'\|_{L_\infty} = \max_i (|x_i - x'_i|)$$

For $d \in \mathbb{R}$, we write $\mathcal{B}(\mathbf{x}, d)$ to denote the set of images within distance d from \mathbf{x} , i.e., $\mathcal{B}(\mathbf{x}, d) = \{\mathbf{x}' \mid \|\mathbf{x} - \mathbf{x}'\|_{L_\infty} \leq d\}$, which is an n -dimensional cube. Based on the above, a deep neural network \mathcal{N} is *locally robust* for a correctly classified input \mathbf{x} with respect to distance d if it assigns the same label to all images in $\mathcal{B}(\mathbf{x}, d)$.

We mention here that numerous attacks are optimized for one distance metric (e.g., [5, 6, 13, 24, 27, 32, 40, 49, 50, 57]), just like ours. Although there exist other distance metrics, e.g., L_0 and L_2 , and extending attacks from one metric to another is possible, developing an attack that performs best in all distance metrics is not realistic. Many state-of-the-art attacks are the most effective in one metric, but their extension to other metrics performs worse than attacks specifically designed for that metric. Our paper focuses on a query-efficient L_∞ attack (as in [32, 50]), and our technique outperforms the state-of-the-art in this setting.

Query-limited blackbox threat model. We assume that attackers have no knowledge of the target network and can only query the network for its prediction scores, e.g., logits or class probabilities. Moreover, we assume that attackers have a query budget, which can be viewed as time or monetary limits in real-world settings. Thus, the blackbox attack we consider in this paper can be described as follows. Given an input \mathbf{x} , distance d , and query budget L , an attacker aims to find an adversarial example \mathbf{x}' in $\mathcal{B}(\mathbf{x}, d)$ by making at most L queries to the neural network.

3 FUZZING BINARY CLASSIFIERS

In this section, we present the technical details of how DeepSearch fuzzes (linear and non-linear) binary classifiers. We first introduce our approach for linear binary classifiers, which serves as the mathematical foundation. Then, we generalize our approach to non-linear binary classifiers through iterative linear approximations.

3.1 Linear Binary Classifiers

A *binary classifier* classifies inputs into two classes, denoted with labels $C_2 = \{l_1, l_2\}$, according to the definition below.

DEFINITION 1 (Binary Classifier). Given a classification function $f : \mathbb{R}^n \rightarrow \mathbb{R}$, a *binary classifier* $\mathcal{N}_f : \mathbb{R}^n \rightarrow C_2$ is defined as follows:

$$\mathcal{N}_f(\mathbf{x}) = \begin{cases} l_1, & \text{if } f(\mathbf{x}) > 0 \\ l_2, & \text{if } f(\mathbf{x}) < 0 \end{cases}$$

If function f is linear, then \mathcal{N}_f is a *linear* binary classifier, otherwise it is non-linear.

The set of values $\mathcal{D}_f = \{\mathbf{x} \mid f(\mathbf{x}) = 0\}$ constitute the *decision boundary* of \mathcal{N}_f , which classifies the domain \mathbb{R}^n into the two classes in \mathcal{C}_2 .

As an example, consider Fig. 1a, showing a linear binary classifier $\mathcal{N}_f : \mathbb{R}^2 \rightarrow \mathcal{C}_2$. Observe that input \mathbf{x}_0 is classified in l_1 whereas \mathbf{x}'_0 is in l_2 . Note that the decision boundary of a linear classifier $\mathcal{N}_f : \mathbb{R}^n \rightarrow \mathcal{C}_2$ is a hyperplane; it is, therefore, a straight line in Fig. 1a. Now, assume that \mathbf{x}_0 is correctly classified and that the dash-dotted square represents $\mathcal{B}(\mathbf{x}_0, d)$. Then, \mathbf{x}'_0 is adversarial because $\mathcal{N}_f(\mathbf{x}_0) \neq \mathcal{N}_f(\mathbf{x}'_0)$, which is equivalent to $f(\mathbf{x}_0)f(\mathbf{x}'_0) < 0$.

Example. We give an intuition on how DeepSearch handles linear binary classifiers using the example of Fig. 1a. Recall that \mathbf{x}_0 is a correctly classified input for which $f(\mathbf{x}_0) > 0$. To find an adversarial example, DeepSearch fuzzes \mathbf{x}_0 with the goal of generating a new input \mathbf{x}'_0 such that $f(\mathbf{x}'_0) < 0$.

Fuzzing is performed as follows. Input \mathbf{x}_0 has two coordinates x_h and x_v , for the horizontal and vertical dimensions. DeepSearch independently mutates each of these coordinates to the minimum and maximum values that are possible within $\mathcal{B}(\mathbf{x}_0, d)$, with the intention of finding the minimum value of f in $\mathcal{B}(\mathbf{x}_0, d)$. For instance, when mutating x_h , we obtain inputs $\mathbf{x}_0[l_h/x_h]$ and $\mathbf{x}_0[u_h/x_h]$ in the figure. Values l_h and u_h are, respectively, the minimum and maximum that x_h may take, and $\mathbf{x}_0[l_h/x_h]$ denotes substituting x_h with l_h (similarly for $\mathbf{x}_0[u_h/x_h]$). We then evaluate $f(\mathbf{x}_0[l_h/x_h])$ and $f(\mathbf{x}_0[u_h/x_h])$, and for x_h , we select the value (l_h or u_h) that causes function f to *decrease*. This is because, in our example, an adversarial input \mathbf{x}'_0 must make the value of f negative. Let us assume that $f(\mathbf{x}_0[u_h/x_h]) < f(\mathbf{x}_0[l_h/x_h])$; we, thus, select u_h for coordinate x_h .

DeepSearch mutates coordinate x_v in a similar way. It evaluates $f(\mathbf{x}_0[l_v/x_v])$ and $f(\mathbf{x}_0[u_v/x_v])$, and selects the value that causes f to decrease. Let us assume that $f(\mathbf{x}_0[u_v/x_v]) < f(\mathbf{x}_0[l_v/x_v])$; we, thus, select u_v for x_v .

Next, we generate input \mathbf{x}'_0 by substituting each coordinate in \mathbf{x}_0 with the boundary value that was previously selected. In other words, $\mathbf{x}'_0 = \mathbf{x}_0[u_h/x_h, u_v/x_v]$, and since $f(\mathbf{x}'_0) < 0$, DeepSearch has generated an adversarial example. Note that $f(\mathbf{x}'_0)$ is actually the minimum value of f in $\mathcal{B}(\mathbf{x}_0, d)$.

DeepSearch for linear binary classifiers. We now formalize how DeepSearch treats linear binary classifiers. Consider a linear classification function $f(\mathbf{x}) = \mathbf{w}^T \mathbf{x} + b = \sum_{i=1}^n w_i x_i + b$, where $\mathbf{w}^T = (w_1, \dots, w_n)$ and $b \in \mathbb{R}$. Note that f is monotonic with respect to all of its variables x_1, \dots, x_n . For instance, if $w_i > 0$, then f is monotonically increasing in x_i .

Recall that $\mathcal{B}(\mathbf{x}, d)$ denotes the set of inputs within distance $d \in \mathbb{R}$ of an input \mathbf{x} . $\mathcal{B}(\mathbf{x}, d)$ may be represented by an n -dimensional cube $\mathcal{I} = I_1 \times \dots \times I_n$, where $I_i = [l_i, u_i]$ is a closed interval bounded by $l_i, u_i \in \mathbb{R}$ with $l_i \leq u_i$ for $1 \leq i \leq n$. Intuitively, value l_i (resp. u_i) is the lower (resp. upper) bound on the i th dimension of \mathbf{x} . An input \mathbf{x}' is a *vertex* of \mathcal{I} if each of its coordinates $\mathbf{x}'(i)$ is an endpoint of I_i for $1 \leq i \leq n$, i.e., $\mathbf{x}'(i) = u_i$ or l_i ($1 \leq i \leq n$).

Due to the monotonicity of f , the maximum and minimum values of f on \mathcal{I} may be easily calculated by applying f to vertices of \mathcal{I} . For

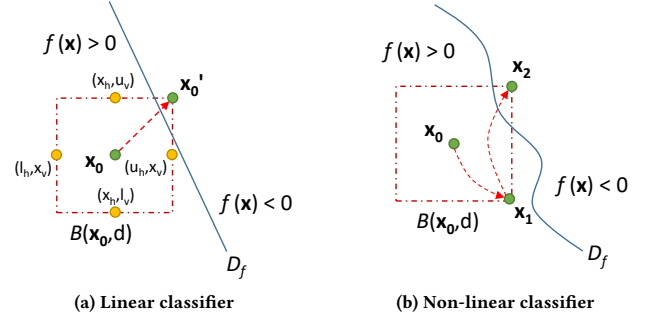


Figure 1: DeepSearch for binary classifiers.

example, consider a one-dimensional linear function $f(x) = -2x$, where $x \in [-1, 1]$, that is, -1 and 1 are the lower and upper bounds for x . Since $f(1) < f(-1)$, we get a maximum value of f at $x = -1$ and a minimum value of f at $x = 1$. N -dimensional linear functions can be treated similarly. We write $f(I)$ for the values of f on I , i.e., $f(I) = \{f(\mathbf{x}) \mid \mathbf{x} \in I\}$, and have the following theorem (whose proof can be found in [84]).

THEOREM 1. Given a linear classification function $f(\mathbf{x}) = \mathbf{w}^T \mathbf{x} + b$, where $\mathbf{w}^T = (w_1, \dots, w_n)$ and $b \in \mathbb{R}$, an n -dimensional cube $\mathcal{I} = I_1 \times \dots \times I_n$, where $I_i = [l_i, u_i]$ for $1 \leq i \leq n$, and an input $\mathbf{x} \in \mathcal{I}$, we have:

- (1) $\min f(I) = f(\mathbf{x}')$, where $\mathbf{x}'(i) = l_i$ (resp. $\mathbf{x}'(i) = u_i$) if $f(\mathbf{x}[u_i/x_i]) > f(\mathbf{x}[l_i/x_i])$ (resp. $f(\mathbf{x}[u_i/x_i]) \leq f(\mathbf{x}[l_i/x_i])$) for $1 \leq i \leq n$
- (2) $\max f(I) = f(\mathbf{x}')$, where $\mathbf{x}'(i) = u_i$ (resp. $\mathbf{x}'(i) = l_i$) if $f(\mathbf{x}[u_i/x_i]) > f(\mathbf{x}[l_i/x_i])$ (resp. $f(\mathbf{x}[u_i/x_i]) \leq f(\mathbf{x}[l_i/x_i])$) for $1 \leq i \leq n$

According to the above theorem, we can *precisely* calculate the minimum and maximum values of f in any n -dimensional cube. In particular, assume a correctly classified input \mathbf{x} for which $f(\mathbf{x}) > 0$. For each dimension i ($1 \leq i \leq n$) of \mathbf{x} , we first construct inputs $\mathbf{x}[l_i/x_i]$ and $\mathbf{x}[u_i/x_i]$. We then compare the values of $f(\mathbf{x}[l_i/x_i])$ and $f(\mathbf{x}[u_i/x_i])$. To generate a new input \mathbf{x}' , we select the value of its i th coordinate as follows:

$$\mathbf{x}'(i) = \begin{cases} l_i, & \text{if } f(\mathbf{x}[u_i/x_i]) > f(\mathbf{x}[l_i/x_i]) \\ u_i, & \text{if } f(\mathbf{x}[u_i/x_i]) \leq f(\mathbf{x}[l_i/x_i]) \end{cases}$$

As shown here, selecting a value for $\mathbf{x}'(i)$ requires evaluating function f twice, i.e., $f(\mathbf{x}[l_i/x_i])$ and $f(\mathbf{x}[u_i/x_i])$. Therefore, for n dimensions, f must be evaluated $2n$ times. In practice however, due to the monotonicity of f , evaluating it only once per dimension is sufficient. For instance, if $f(\mathbf{x}[u_i/x_i])$ already decreases (resp. increases) the value of f in comparison to $f(\mathbf{x})$, there is no need to evaluate $f(\mathbf{x}[l_i/x_i])$. Value u_i (resp. l_i) should be selected for the i th coordinate. Hence, the minimum value of f can be computed by evaluating the function exactly n times. If, for the newly generated input \mathbf{x}' , the sign of f becomes negative, \mathbf{x}' constitutes an adversarial example.

We treat the case where $f(\mathbf{x}) < 0$ for a correctly classified input \mathbf{x} analogously. DeepSearch aims to generate a new input \mathbf{x}' such that the sign of f becomes positive. We are, therefore, selecting coordinate values that cause f to *increase*.

Algorithm 1: DeepSearch for binary classifiers.

Input: input $\mathbf{x} \in \mathbb{R}^n$, initial input $\mathbf{x}_{init} \in \mathcal{B}(\mathbf{x}, d)$,
function $f : \mathbb{R}^n \rightarrow \mathbb{R}$, distance $d \in \mathbb{R}$
Output: $\mathbf{x}' \in \mathcal{B}(\mathbf{x}, d)$

```

1 Function ApproxMax( $\mathbf{x}, f, (I_1, \dots, I_n)$ ) is
2    $\mathbf{x}' := (0, \dots, 0)$ 
3   foreach  $1 \leq i \leq n$  do
4     if  $f(\mathbf{x}[u_i/x_i]) > f(\mathbf{x}[l_i/x_i])$  then
5        $\mathbf{x}' := \mathbf{x}'[u_i/x'_i]$ 
6     else
7        $\mathbf{x}' := \mathbf{x}'[l_i/x'_i]$ 
8   return  $\mathbf{x}'$ 
9
10 Function ApproxMin( $\mathbf{x}, f, (I_1, \dots, I_n)$ ) is
11    $\mathbf{x}' := (0, \dots, 0)$ 
12   foreach  $1 \leq i \leq n$  do
13     if  $f(\mathbf{x}[u_i/x_i]) > f(\mathbf{x}[l_i/x_i])$  then
14        $\mathbf{x}' := \mathbf{x}'[l_i/x'_i]$ 
15     else
16        $\mathbf{x}' := \mathbf{x}'[u_i/x'_i]$ 
17   return  $\mathbf{x}'$ 
18
19 Function DS-Binary( $\mathbf{x}, \mathbf{x}_{init}, f, d$ ) is
20   construct intervals  $(I_1, \dots, I_n)$  such that  $\mathcal{B}(\mathbf{x}, d) = I_1 \times \dots \times I_n$ 
21   initialize  $\mathbf{x}_0 := \mathbf{x}_{init}$  and  $k := 0$ 
22   if  $f(\mathbf{x}_0) > 0$  then
23     repeat
24        $\mathbf{x}_{k+1} := \text{ApproxMin}(\mathbf{x}_k, f, (I_1, \dots, I_n))$ 
25        $k := k + 1$ 
26     until  $\mathcal{N}_f(\mathbf{x}) \neq \mathcal{N}_f(\mathbf{x}_k)$ , or  $k = \text{MaxNum}$ 
27   else
28     repeat
29        $\mathbf{x}_{k+1} := \text{ApproxMax}(\mathbf{x}_k, f, (I_1, \dots, I_n))$ 
30        $k := k + 1$ 
31     until  $\mathcal{N}_f(\mathbf{x}) \neq \mathcal{N}_f(\mathbf{x}_k)$ , or  $k = \text{MaxNum}$ 
32   return  $\mathbf{x}_k$ 

```

3.2 Non-Linear Binary Classifiers

We generalize our technique to non-linear binary classifiers. In this setting, DeepSearch iteratively *approximates* the minimum and maximum values of f in $\mathcal{B}(\mathbf{x}, d)$.

Example. As an example, consider the non-linear classification function f shown in Fig. 1b. Since f is non-linear, the decision boundary \mathcal{D}_f of the binary classifier is a curve.

Starting from correctly classified input \mathbf{x}_0 , DeepSearch treats f as linear within $\mathcal{B}(\mathbf{x}_0, d)$ and generates \mathbf{x}_1 , exactly as it would for a linear binary classifier. To explain how \mathbf{x}_1 is derived, we refer to the points in Fig. 1a. Suppose we first mutate the horizontal dimension of \mathbf{x}_0 (using the lower and upper bounds) and find that $f(l_h, x_v) > f(u_h, x_v)$. To increase the chances of crossing the decision boundary, we choose the bound for the horizontal dimension of \mathbf{x}_0 that gives us the lower value of f , i.e., we select u_h for horizontal coordinate x_h . Then, we mutate the vertical dimension of \mathbf{x}_0 and find that $f(x_h, u_v) > f(x_h, l_v)$. This means that we select l_v for vertical

coordinate x_v . Hence, we derive $\mathbf{x}_1 = (u_h, l_v)$. Observe that input \mathbf{x}_1 is not adversarial. Unlike for a linear binary classifier however, where the minimum value of f in $\mathcal{B}(\mathbf{x}_0, d)$ is precisely computed, the non-linear case is handled by iteratively approximating the minimum. In particular, after generating \mathbf{x}_1 , DeepSearch iterates starting from \mathbf{x}_1 , while again treating f as linear in $\mathcal{B}(\mathbf{x}_0, d)$. As a result, our technique generates input \mathbf{x}_2 , which is adversarial.

The reason we can treat non-linear binary classifiers as linear ones is that perturbations of pixels are only allowed in a very small n -dimensional cube, constrained by the L_∞ distance. Within such a small space, we can effectively approximate non-linear functions using iterative linear approximations.

DeepSearch for non-linear binary classifiers. Alg. 1 shows DeepSearch for binary classifiers. It uses iterative approximations to search for adversarial examples. Note that our technique is blackbox, and consequently, it cannot differentiate between linear and non-linear classifiers. Alg. 1 is, therefore, the general algorithm that DeepSearch applies to fuzz any binary classifier.

The main function in Alg. 1 is DS-Binary. Input \mathbf{x}_{init} is the input from which we start the first iteration, e.g., it corresponds to \mathbf{x}_0 in Fig. 1. Input \mathbf{x} is used to compute $\mathcal{B}(\mathbf{x}, d)$, and for now, assume that \mathbf{x} is equal to \mathbf{x}_{init} . (We will discuss why \mathbf{x} is needed in Sect. 5.) In addition to these inputs, the algorithm also takes a classification function f and the distance d .

Function DS-Binary assigns \mathbf{x}_{init} to \mathbf{x}_0 and constructs n intervals I_1, \dots, I_n to represent $\mathcal{B}(\mathbf{x}_0, d)$ (lines 20–21). Then, based on the sign of $f(\mathbf{x}_0)$, our algorithm iteratively approximates the minimum (lines 23–26) or the maximum (lines 28–31) value of f in $\mathcal{B}(\mathbf{x}_0, d)$. DS-Binary terminates when either an adversarial example is found or it has reached MaxNum iterations. To find adversarial examples in k iterations, we evaluate f at most $2n + n(k - 1)$ times.

ApproxMin and ApproxMax implement Thm. 1 to calculate the minimum and maximum values of function f in the n -dimensional cube $I_1 \times \dots \times I_n$. When f is linear, calling these functions on any input $\mathbf{x} \in I_1 \times \dots \times I_n$ does not affect the computation. In other words, the minimum and maximum values are precisely computed for any \mathbf{x} . When f is non-linear, it is still assumed to be linear within the n -dimensional cube. Given that the size of the cube is designed to be small, this assumption does not introduce too much imprecision. As a consequence of this assumption however, different inputs in the n -dimensional cube lead to computing different minimum and maximum values of f . For instance, in Fig. 1b, calling ApproxMin on \mathbf{x}_0 returns \mathbf{x}_1 , while calling it on \mathbf{x}_1 returns \mathbf{x}_2 .

4 FUZZING MULTICLASS CLASSIFIERS

In this section, we extend our technique for blackbox fuzzing of binary classifiers to multiclass classifiers.

4.1 Linear Multiclass Classifiers

A *multiclass classifier* classifies inputs in m classes according to the following definition.

DEFINITION 2 (Multiclass Classifier). For classification function $f : \mathbb{R}^n \rightarrow \mathbb{R}^m$, which returns m values each corresponding to one

class in $C_m = \{l_1, \dots, l_m\}$, a *multiclass classifier* $\mathcal{N}_f : \mathbb{R}^n \rightarrow C_m$ is defined as

$$\mathcal{N}_f(\mathbf{x}) = l_j, \quad \text{iff } j = \arg \max_i f_i(\mathbf{x}),$$

where $f_i : \mathbb{R}^n \rightarrow \mathbb{R}$ denotes the function derived by evaluating f for the i th class, i.e., $f(\mathbf{x}) = (f_1(\mathbf{x}), \dots, f_m(\mathbf{x}))^T$.

In other words, a multiclass classifier \mathcal{N}_f classifies an input \mathbf{x} in l_j if $f_j(\mathbf{x})$ evaluates to the largest value in comparison to all other functions f_i .

Function f of a multiclass classifier \mathcal{N}_f may be decomposed into multiple binary classifiers such that the original classifier can be reconstructed from the binary ones. First, to decompose a multiclass classifier into binary classifiers, for any pair of classes l_i and l_j ($1 \leq i, j \leq m$), we define a classification function $g_{ij} : \mathbb{R}^n \rightarrow \mathbb{R}$ as $g_{ij}(\mathbf{x}) = f_i(\mathbf{x}) - f_j(\mathbf{x})$. We then construct a binary classifier $\mathcal{N}_{g_{ij}} : \mathbb{R}^n \rightarrow \{l_i, l_j\}$ as follows:

$$\mathcal{N}_{g_{ij}}(\mathbf{x}) = \begin{cases} l_i, & \text{if } g_{ij}(\mathbf{x}) > 0 \\ l_j, & \text{if } g_{ij}(\mathbf{x}) < 0 \end{cases}$$

As usual, the set of values $\mathcal{D}_{g_{ij}} = \{\mathbf{x} \mid f_i(\mathbf{x}) - f_j(\mathbf{x}) = 0\}$ constitutes the pairwise decision boundary of binary classifier $\mathcal{N}_{g_{ij}}$, which classifies the domain \mathbb{R}^n into the two classes $\{l_i, l_j\}$. As an example, consider Fig. 2a depicting a multiclass classifier $\mathcal{N}_f : \mathbb{R}^2 \rightarrow C_3$, where f is linear and $C_3 = \{l_1, l_2, l_3\}$. Assume that \mathcal{N}_f correctly classifies input \mathbf{x} in l_2 . Based on the above, linear binary classifiers $\mathcal{N}_{g_{21}}$ and $\mathcal{N}_{g_{23}}$ also classify \mathbf{x} in l_2 , i.e., $g_{21}(\mathbf{x}) > 0$ and $g_{23}(\mathbf{x}) > 0$.

Second, a multiclass classifier may be composed from multiple binary classifiers as follows. An input \mathbf{x} is classified in class l_i by multiclass classifier \mathcal{N}_f if and only if it is classified in l_i by all $m-1$ binary classifiers $\mathcal{N}_{g_{ij}}$ for $1 \leq j \leq m, i \neq j$, where $l_i \in C_m$ and $g_{ij}(\mathbf{x}) = f_i(\mathbf{x}) - f_j(\mathbf{x})$:

$$\mathcal{N}_f(\mathbf{x}) = l_i, \quad \text{iff } \forall 1 \leq j \leq m, i \neq j : \mathcal{N}_{g_{ij}}(\mathbf{x}) = l_i$$

For instance, in Fig. 2a, if both $\mathcal{N}_{g_{21}}$ and $\mathcal{N}_{g_{23}}$ classify input \mathbf{x} in class l_2 , then the multiclass classifier also classifies it in l_2 .

Based on the above, a multiclass classifier has an adversarial input if and only if this input is also adversarial for a constituent binary classifier.

COROLLARY 1. Let \mathcal{N}_f be a multiclass classifier and $\mathbf{x} \in \mathbb{R}^n$ a correctly classified input, where $\mathcal{N}_f(\mathbf{x}) = l_i$ and $l_i \in C_m$. There exists an adversarial example $\mathbf{x}' \in \mathcal{B}(\mathbf{x}, d)$ for \mathcal{N}_f , where $d \in \mathbb{R}$, if and only if \mathbf{x}' is an adversarial example for a binary classifier $\mathcal{N}_{g_{ij}}$ ($1 \leq j \leq m, i \neq j$), where $g_{ij}(\mathbf{x}) = f_i(\mathbf{x}) - f_j(\mathbf{x})$:

$$\mathcal{N}_f(\mathbf{x}') \neq l_i, \quad \text{iff } \exists 1 \leq j \leq m, i \neq j : \mathcal{N}_{g_{ij}}(\mathbf{x}') \neq l_i$$

Example. This corollary is crucial in generalizing our technique to multiclass classifiers. Assume a correctly classified input \mathbf{x} , for which $\mathcal{N}_f(\mathbf{x}) = l_i$. According to the above corollary, the robustness of \mathcal{N}_f in $\mathcal{B}(\mathbf{x}, d)$ reduces to the robustness of all $m-1$ binary classifiers $\{\mathcal{N}_{g_{ij}} \mid 1 \leq j \leq m, i \neq j\}$ in $\mathcal{B}(\mathbf{x}, d)$. We, therefore, use DeepSearch for binary classifiers to test each binary classifier in this set. If there exists an adversarial input \mathbf{x}' for one of these classifiers, i.e., $\mathcal{N}_{g_{ij}}(\mathbf{x}') \neq l_i$ for some j , then \mathbf{x}' is also an adversarial input for \mathcal{N}_f , i.e., $\mathcal{N}_f(\mathbf{x}') \neq l_i$.

Let us consider again the example of Fig. 2a. Recall that multiclass classifier \mathcal{N}_f correctly classifies input \mathbf{x} in class l_2 , and so do binary

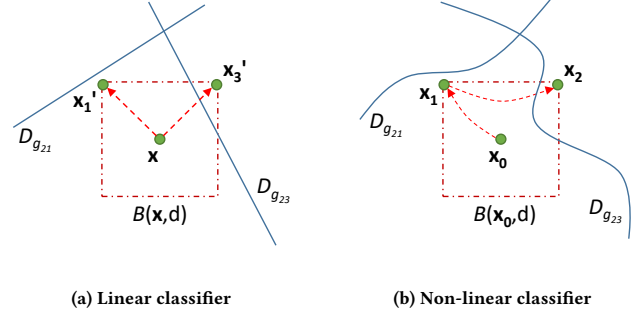


Figure 2: DeepSearch for multiclass classifiers.

classifiers $\mathcal{N}_{g_{21}}$ and $\mathcal{N}_{g_{23}}$, i.e., $g_{21}(\mathbf{x}) > 0$ and $g_{23}(\mathbf{x}) > 0$. As a result, DeepSearch tries to generate inputs that decrease the value of each of these functions in $\mathcal{B}(\mathbf{x}, d)$ in order to find adversarial examples. Function g_{21} evaluates to its minimum value in $\mathcal{B}(\mathbf{x}, d)$ for input \mathbf{x}'_1 , and function g_{23} for input \mathbf{x}'_3 . Observe that \mathbf{x}'_3 is an adversarial example for $\mathcal{N}_{g_{23}}$, and thus also for \mathcal{N}_f , whereas \mathbf{x}'_1 is not.

DeepSearch for linear multiclass classifiers. Let us assume a linear classification function $f(\mathbf{x}) = \mathbf{W}^T \mathbf{x} + \mathbf{b}$, where $\mathbf{W}^T = (\mathbf{w}_1^T, \dots, \mathbf{w}_m^T)^T$, $\mathbf{w}_i \in \mathbb{R}^n$ ($1 \leq i \leq m$), and $\mathbf{b} = (b_1, \dots, b_m)^T \in \mathbb{R}^m$. Then, f_i , which denotes the function derived by evaluating f for the i th class, is of the form $f_i(\mathbf{x}) = \mathbf{w}_i^T \mathbf{x} + b_i$ for $1 \leq i \leq m$. For any pair of class labels l_i and l_j , function g_{ij} is defined as $g_{ij}(\mathbf{x}) = f_i(\mathbf{x}) - f_j(\mathbf{x}) = (\mathbf{w}_i^T - \mathbf{w}_j^T) \mathbf{x} + (b_i - b_j)$. Hence, g_{ij} is also linear, and $\mathcal{N}_{g_{ij}}$ is a linear binary classifier.

Assume that classifier \mathcal{N}_f correctly classifies input $\mathbf{x} \in \mathbb{R}^n$ in l_i , $\mathcal{N}_f(\mathbf{x}) = l_i$ ($l_i \in C_m$). According to Cor. 1, the robustness of \mathcal{N}_f in $\mathcal{B}(\mathbf{x}, d)$ ($d \in \mathbb{R}$) reduces to the robustness of each binary classifier $\mathcal{N}_{g_{ij}}$ ($1 \leq j \leq m, i \neq j$) in $\mathcal{B}(\mathbf{x}, d)$. To find an adversarial example for a binary classifier $\mathcal{N}_{g_{ij}}$ in $\mathcal{B}(\mathbf{x}, d)$, DeepSearch must generate an input $\mathbf{x}' \in \mathcal{B}(\mathbf{x}, d)$ such that $g_{ij}(\mathbf{x}') < 0$. (Recall that by definition $g_{ij}(\mathbf{x}) > 0$.) Since all functions g_{ij} ($1 \leq j \leq m, i \neq j$) are linear, we easily find their minimum values in $\mathcal{B}(\mathbf{x}, d)$ as follows.

Let I_1, \dots, I_n be intervals such that $\mathcal{B}(\mathbf{x}, d) = I_1 \times \dots \times I_n$, where $I_k = [l_k, u_k]$ for $1 \leq k \leq n$. As in Sect. 3.1, for each dimension k , DeepSearch evaluates function f twice to compare the values of $f(\mathbf{x}[u_k/x_k])$ and $f(\mathbf{x}[l_k/x_k])$. To generate a new input \mathbf{x}'_j for which function g_{ij} evaluates to its minimum value, we select its k th coordinate as follows:

$$\mathbf{x}'_j(k) = \begin{cases} l_k, & \text{if } g_{ij}(\mathbf{x}[u_k/x_k]) > g_{ij}(\mathbf{x}[l_k/x_k]) \\ u_k, & \text{if } g_{ij}(\mathbf{x}[u_k/x_k]) \leq g_{ij}(\mathbf{x}[l_k/x_k]). \end{cases}$$

Note that, although we calculate the minimum value of $m-1$ linear functions, we still evaluate f $2n$ times. This is because a function g_{ij} is defined as $g_{ij}(\mathbf{x}) = f_i(\mathbf{x}) - f_j(\mathbf{x})$, where $f_i(\mathbf{x})$ and $f_j(\mathbf{x})$ are the values of $f(\mathbf{x})$ for the i th and j th classes, respectively. If the sign of $g_{ij}(\mathbf{x}'_j)$ becomes negative for some j , then DeepSearch has found an adversarial example for \mathcal{N}_f in $\mathcal{B}(\mathbf{x}, d)$.

4.2 Non-Linear Multiclass Classifiers

We now extend our technique to non-linear multiclass classifiers. Analogously to Sect. 3.2, DeepSearch iteratively approximates the minimum values of functions g_{ij} in $\mathcal{B}(\mathbf{x}, d)$.

Algorithm 2: DeepSearch for multiclass classifiers.

Input: input $\mathbf{x} \in \mathbb{R}^n$, initial input $\mathbf{x}_{init} \in \mathcal{B}(\mathbf{x}, d)$,
function $f : \mathbb{R}^n \rightarrow \mathbb{R}^m$, distance $d \in \mathbb{R}$
Output: $\mathbf{x}' \in \mathcal{B}(\mathbf{x}, d)$

```

1 Function DS-Multiclass( $\mathbf{x}, \mathbf{x}_{init}, f, d$ ) is
2   construct intervals  $(I_1, \dots, I_n)$  such that  $\mathcal{B}(\mathbf{x}, d) = I_1 \times \dots \times I_n$ 
3   compute  $\mathcal{N}_f(\mathbf{x})$  and assume that  $\mathcal{N}_f(\mathbf{x}) = l_i$ 
4   initialize  $\mathbf{x}_0 := \mathbf{x}_{init}$  and  $k := 0$ 
5   define  $\{g_{ij} \mid g_{ij}(\mathbf{x}) = f_i(\mathbf{x}) - f_j(\mathbf{x}), 1 \leq j \leq m, i \neq j\}$ 
6   repeat
7      $r := \arg \min_j g_{ij}(\mathbf{x}_k)$ 
8      $\mathbf{x}_{k+1} := \text{ApproxMin}(\mathbf{x}_k, g_{ir}, (I_1, \dots, I_n))$ 
9      $k := k + 1$ 
10  until  $\mathcal{N}_f(\mathbf{x}) \neq \mathcal{N}_f(\mathbf{x}_k)$ , or  $k = \text{MaxNum}$ 
11  return  $\mathbf{x}_k$ 

```

Example. As an example, consider Fig. 2b depicting a multiclass classifier $\mathcal{N}_f : \mathbb{R}^2 \rightarrow C_3$, where f is non-linear and $C_3 = \{l_1, l_2, l_3\}$. Assume that \mathcal{N}_f classifies input \mathbf{x}_0 in class l_2 , and thus, so do non-linear binary classifiers $\mathcal{N}_{g_{21}}$ and $\mathcal{N}_{g_{23}}$.

Let us also assume that $g_{21}(\mathbf{x}_0) < g_{23}(\mathbf{x}_0)$. Since g_{21} evaluates to a smaller value than g_{23} for input \mathbf{x}_0 , we consider it more likely to have an adversarial example. In other words, we first approximate the minimum value of g_{21} because it is closer to becoming negative for the initial input. DeepSearch treats g_{21} as linear within $\mathcal{B}(\mathbf{x}_0, d)$ and generates \mathbf{x}_1 . Observe that input \mathbf{x}_1 is not adversarial. Now, assume that $g_{21}(\mathbf{x}_1) > g_{23}(\mathbf{x}_1)$. As a result, DeepSearch tries to find the minimum of function g_{23} in $\mathcal{B}(\mathbf{x}_0, d)$, also by treating it as linear. It generates input \mathbf{x}_2 , which is an adversarial example for classifiers $\mathcal{N}_{g_{23}}$ and \mathcal{N}_f .

DeepSearch for non-linear multiclass classifiers. Alg. 2 is the general DeepSearch algorithm for multiclass classifiers. The inputs are the same as for Alg. 1. For now, assume that \mathbf{x} is equal to \mathbf{x}_{init} . Again, the algorithm executes at most MaxNum iterations, and it terminates as soon as an adversarial example is found.

Function DS-Multiclass assigns \mathbf{x}_{init} to \mathbf{x}_0 and constructs n intervals I_1, \dots, I_n to represent $\mathcal{B}(\mathbf{x}_0, d)$. It also computes the class label l_i of \mathbf{x} , and defines functions g_{ij} ($1 \leq j \leq m, i \neq j$) (lines 2–5). The rest of the algorithm uses ApproxMin from Alg. 1 to iteratively approximate the minimum of one function g_{ij} per iteration, which is selected on line 7 such that its value for input \mathbf{x}_k is smaller in comparison to all other constituent binary classification functions. Intuitively, g_{ij} corresponds to the binary classifier that is most likely to have an adversarial example near \mathbf{x}_k . This heuristic allows our algorithm to find an adversarial example faster than having to generate an input \mathbf{x}_{k+1} for all $m - 1$ functions g_{ij} per iteration.

To find an adversarial example in k iterations, we need at most $2n + n(k - 1)$ queries for the value of f .

An alternative objective function. In each iteration of Alg. 2, we construct a different objective function g_{ij} and approximate its minimum value. An alternative choice of an objective function is f_i itself. In multiclass classification, decreasing the value of f_i amounts to decreasing the score value of the i th class, which implicitly increases the score values of other classes.

We refer to the algorithm derived by substituting lines 7–8 of Alg. 2 with the following assignment as Alg. 2' :

$$\mathbf{x}_{k+1} := \text{ApproxMin}(\mathbf{x}_k, f_i, (I_1, \dots, I_n))$$

It uses ApproxMin to iteratively approximate the minimum value of f_i . We find it very effective in our experiments.

5 ITERATIVE REFINEMENT

The closer the adversarial examples are to a correctly classified input, the more subtle they are. Such adversarial examples are said to have a low distortion rate. In this section, we extend DeepSearch with an iterative-refinement approach for finding subtle adversarial examples. On a high level, given an input \mathbf{x} and a distance d , for which we have already found an adversarial example \mathbf{x}' in region $\mathcal{B}(\mathbf{x}, d)$, DeepSearch iteratively reduces distance d as long as the smaller region still contains an adversarial example. If none is found, the distance is not reduced further.

Let $\mathcal{I} = I_1 \times \dots \times I_n$ be an n -dimensional cube, where $I_i = [l_i, u_i]$ is a closed interval bounded by $l_i, u_i \in \mathbb{R}$ with $l_i \leq u_i$ for $1 \leq i \leq n$. For an input \mathbf{x} with an i th coordinate $\mathbf{x}(i) \in (-\infty, l_i] \cup [u_i, +\infty)$ for $1 \leq i \leq n$, we define a projection operator PROJ that maps \mathbf{x} to a vertex of \mathcal{I} as follows

$$\text{PROJ}(\mathcal{I}, \mathbf{x})(i) = \begin{cases} u_i, & \text{if } \mathbf{x}(i) \geq u_i \\ l_i, & \text{if } \mathbf{x}(i) \leq l_i, \end{cases}$$

where $\text{PROJ}(\mathcal{I}, \mathbf{x})(i)$ denotes the i th coordinate of $\text{PROJ}(\mathcal{I}, \mathbf{x})$. As an example, consider Fig. 3a showing a linear multiclass classifier. Input \mathbf{x}_2 is a projection of \mathbf{x}_1 .

Using this operator, the minimum and maximum values of a linear classification function f may also be projected on \mathcal{I} , and we have the following theorem (whose proof is available in [84]).

THEOREM 2. Let $f(\mathbf{x}) = \mathbf{w}^T \mathbf{x} + b$ be a linear classification function, and $\mathcal{I}_1, \mathcal{I}_2$ two n -dimensional cubes such that $\mathcal{I}_1 \subseteq \mathcal{I}_2$. Assuming that \mathbf{x} is a vertex of \mathcal{I}_2 , we have:

- (1) if $\min f(\mathcal{I}_2) = f(\mathbf{x})$, then $\min f(\mathcal{I}_1) = f(\text{PROJ}(\mathcal{I}_1, \mathbf{x}))$
- (2) if $\max f(\mathcal{I}_2) = f(\mathbf{x})$, then $\max f(\mathcal{I}_1) = f(\text{PROJ}(\mathcal{I}_1, \mathbf{x}))$

In Fig. 3a, assume that input \mathbf{x}_0 is correctly classified in class l_2 . Then, in region $\mathcal{B}(\mathbf{x}_0, d_1)$, function g_{23} obtains its minimum value for input \mathbf{x}_1 . When projecting \mathbf{x}_1 to vertex \mathbf{x}_2 of $\mathcal{B}(\mathbf{x}_0, d_2)$, notice that g_{23} evaluates to its minimum for input \mathbf{x}_2 in this smaller region.

Example. Fig. 3 shows two multiclass classifiers $\mathcal{N}_f : \mathbb{R}^2 \rightarrow C_3$, where $C_3 = \{l_1, l_2, l_3\}$. In Fig. 3a, function f is linear, whereas in Fig. 3b, it is non-linear. For correctly classified input \mathbf{x}_0 , we assume that $\mathcal{N}_f(\mathbf{x}_0) = l_2$, and thus, $\mathcal{N}_{g_{21}}(\mathbf{x}_0) = l_2$ and $\mathcal{N}_{g_{23}}(\mathbf{x}_0) = l_2$.

In both subfigures, assume that \mathbf{x}_1 is an adversarial example found by DS-Multiclass (see Alg. 2) in $\mathcal{B}(\mathbf{x}_0, d_1)$. Once such an example is found, our technique with refinement uses bisection search to find the smallest distance d' such that the projection of \mathbf{x}_1 on $\mathcal{B}(\mathbf{x}_0, d')$ is still adversarial. In Fig. 3, this distance is d_2 , and input \mathbf{x}_2 constitutes the projection of \mathbf{x}_1 on $\mathcal{B}(\mathbf{x}_0, d_2)$. So, \mathbf{x}_2 is closer to \mathbf{x}_0 , which means that it has a lower distortion rate than \mathbf{x}_1 . In fact, since we are using bisection search to determine distance d_2 , \mathbf{x}_2 is the closest adversarial input to \mathbf{x}_0 that may be generated by projecting \mathbf{x}_1 on smaller regions.

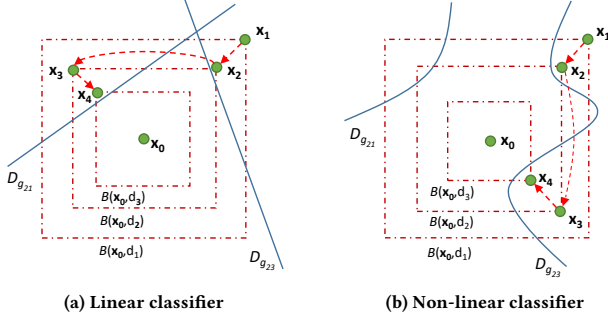


Figure 3: DeepSearch with iterative refinement.

However, in region $\mathcal{B}(\mathbf{x}_0, d_2)$, there may be other vertices that are adversarial and get us even closer to \mathbf{x}_0 with projection. To find such examples, we apply DS-Multiclass again, this time starting from input \mathbf{x}_2 and searching in region $\mathcal{B}(\mathbf{x}_0, d_2)$. As a result, we generate adversarial input \mathbf{x}_3 in the subfigures. Now, by projecting \mathbf{x}_3 to the smallest possible region around \mathbf{x}_0 , we compute \mathbf{x}_4 , which is the adversarial example with the lowest distortion rate so far.

Assume that applying DS-Multiclass for a third time, starting from \mathbf{x}_4 and searching in $\mathcal{B}(\mathbf{x}_0, d_3)$, does not generate any other adversarial examples. In this case, our technique returns \mathbf{x}_4 .

DeepSearch with iterative refinement. Alg. 3 describes our technique with iterative refinement. Each iteration consists of a refinement and a search step, which we explain next.

The refinement step (lines 3–4) first calculates the L_∞ distance d between \mathbf{x} and \mathbf{x}' . In Fig. 3, input \mathbf{x} of the algorithm is \mathbf{x}_0 , and \mathbf{x}' is \mathbf{x}_1 . So, \mathbf{x}' is an adversarial input that was generated by our technique, and consequently, a vertex of $\mathcal{B}(\mathbf{x}, d)$. On line 4, we use bisect search to find the minimum distance d' such that the input derived by projecting \mathbf{x}' on $\mathcal{B}(\mathbf{x}, d')$ is still adversarial. In Fig. 3, this is distance d_2 , and $\text{PROJ}(\mathcal{B}(\mathbf{x}, d'), \mathbf{x}')$ of the algorithm corresponds to adversarial input \mathbf{x}_2 in the figure.

Note that this refinement is possible because of Thm. 2, which guarantees that a linear function f evaluates to its minimum in $\mathcal{B}(\mathbf{x}, d')$ for the input derived with projection. When f is non-linear, it might not evaluate to its minimum for adversarial input $\text{PROJ}(\mathcal{B}(\mathbf{x}, d'), \mathbf{x}')$. However, it is still the case that this projected input is closer to \mathbf{x} , i.e., $\|\mathbf{x} - \text{PROJ}(\mathcal{B}(\mathbf{x}, d'), \mathbf{x}')\|_{L_\infty} \leq \|\mathbf{x} - \mathbf{x}'\|_{L_\infty}$, and thus, has a lower distortion rate.

After selecting an input from which to start the search (line 5), the search step (lines 6–9) calls function DS-Binary (Alg. 1) or DS-Multiclass (Alg. 2), depending on whether f is a binary classification function. The goal is to search for another adversarial example (other than $\text{PROJ}(\mathcal{B}(\mathbf{x}, d'), \mathbf{x}')$) in region $\mathcal{B}(\mathbf{x}, d')$. In Fig. 3, an adversarial input found by this step, when starting the search from \mathbf{x}_2 , is \mathbf{x}_3 , which is also a vertex of $\mathcal{B}(\mathbf{x}_0, d_2)$.

On lines 10–13, we essentially check whether the search step was successful in finding another adversarial input \mathbf{x}'' . However, DS-Binary and DS-Multiclass might not return an adversarial example. If they do, like input \mathbf{x}_3 in Fig. 3, the algorithm iterates (line 14). If not, like when starting the search from input \mathbf{x}_4 in the figure, which is a projection of \mathbf{x}_3 on $\mathcal{B}(\mathbf{x}_0, d_3)$, then we return the projected input and terminate.

Algorithm 3: DeepSearch with iterative refinement.

Input: input $\mathbf{x} \in \mathbb{R}^n$, adversarial input $\mathbf{x}' \in \mathcal{B}(\mathbf{x}, d)$ ($d \in \mathbb{R}$),
function $f : \mathbb{R}^n \rightarrow \mathbb{R}^m$
Output: an adversarial input $\mathbf{x}'' \in \mathcal{B}(\mathbf{x}, d')$ ($d' \leq d$)

```

1 Function DS-Refinement( $\mathbf{x}, \mathbf{x}', f$ ) is
2   repeat
3      $d := \|\mathbf{x} - \mathbf{x}'\|_{L_\infty}$ 
4     apply bisect search to find the smallest distance  $d' \leq d$ 
       such that input  $\text{PROJ}(\mathcal{B}(\mathbf{x}, d'), \mathbf{x}')$  is an adversarial
       example.
5     choose an  $\mathbf{x}_{\text{new}} \in \mathcal{B}(\mathbf{x}, d')$ , from which to start a new
       search, e.g.  $\mathbf{x}_{\text{new}} = \text{PROJ}(\mathcal{B}(\mathbf{x}, d'), \mathbf{x}')$ .
6     if  $f$  is binary then
7        $\mathbf{x}'' := \text{DS-Binary}(\mathbf{x}, \mathbf{x}_{\text{new}}, f, d')$ 
8     else
9        $\mathbf{x}'' := \text{DS-Multiclass}(\mathbf{x}, \mathbf{x}_{\text{new}}, f, d')$ 
10    if  $\mathbf{x}''$  is an adversarial example then
11       $\mathbf{x}' := \mathbf{x}''$ 
12    else
13       $\mathbf{x}' := \text{PROJ}(\mathcal{B}(\mathbf{x}, d'), \mathbf{x}')$ 
14  until  $\mathbf{x}''$  is not an adversarial example
15  return  $\mathbf{x}'$  and  $d'$ 

```

6 HIERARCHICAL GROUPING

For an n -dimensional input \mathbf{x} , our technique makes at least n queries per iteration. For high-dimensional inputs, it could cost a significantly large number of queries to perform even one iteration of our attack. One basic strategy for query reduction is to divide pixels of an input image into different groups and mutate all pixels in a group to the same direction, e.g., all pixels in a group are moved to their upper bounds. Thus, we only need one query for all pixels in the same group. To exploit spatial regularities in images for query efficiency, we adapt hierarchical grouping [50] to our setting.

DeepSearch with hierarchical grouping. Alg. 4 summarizes our technique with hierarchical grouping, which consists of the following three main steps.

- (1) *Initial grouping* (line 8): For an n -dimensional input image \mathbf{x} , we first divide the n dimensions into $\lceil \frac{n}{k^2} \rceil$ sets $G_1, \dots, G_{\lceil \frac{n}{k^2} \rceil}$, where each set G_i ($1 \leq i \leq \lceil \frac{n}{k^2} \rceil$) contains indices corresponding to $k \times k$ neighboring pixels in \mathbf{x} . This amounts to dividing the original image \mathbf{x} into $\lceil \frac{n}{k^2} \rceil$ square blocks. The definition of Initial-Group($\{1, \dots, n\}, k$) is omitted due to space limitations.
- (2) *Fuzzing* (line 10): We extend DeepSearch to handle groups of pixels and write $\text{DeepSearch}(\mathbf{x}, \mathbf{x}', f, d, \mathcal{G})$ to mean such an extension. For each set $G_i \in \mathcal{G}$, our technique mutates all coordinates that correspond to indices in the set toward the same direction at the same time. Hence, DeepSearch only compares two values per set, namely $f[u_{i_1}/x_{i_1}, \dots, u_{i_l}/x_{i_l}]$ and $f[l_{i_1}/x_{i_1}, \dots, l_{i_l}/x_{i_l}]$, where $i_1, \dots, i_l \in G_i$ and $l = |G_i|$.
- (3) *Group splitting* (line 11–13): If the current partition of the image is still too coarse for DeepSearch to find adversarial examples, we perform DeepSearch in finer granularity. We further divide each set G_i into $m \times m$ subsets $G_{i,1}, \dots, G_{i,m \times m}$.

Algorithm 4: DeepSearch with hierarchical grouping.

Input: input $\mathbf{x} \in \mathbb{R}^n$, initial input $\mathbf{x}_{init} \in \mathcal{B}(\mathbf{x}, d)$, initial group size k , parameter m for group splitting, function $f : \mathbb{R}^n \rightarrow \mathbb{R}^m$, distance $d \in \mathbb{R}$, query budget L

Output: $\mathbf{x}' \in \mathcal{B}(\mathbf{x}, d)$

```

1 Function Divide-Group( $\mathcal{G}, m$ ) is
2   foreach  $G_i \in \mathcal{G}$  do
3     divide  $G_i$  into  $m \times m$  subset  $\{G_{i,1}, \dots, G_{i,m \times m}\}$ 
4      $\mathcal{G}' := \mathcal{G} \cup \{G_{i,1}, \dots, G_{i,m \times m}\}$ 
5   return  $\mathcal{G}'$ 
6
7 Function DS-Hierarchy( $\mathbf{x}, \mathbf{x}_{init}, f, k, m$ ) is
8    $\mathcal{G} := \text{Initial-Group}(\{1, \dots, n\}, k)$  and  $\mathbf{x}' := \mathbf{x}_{init}$ 
9   repeat
10     $\mathbf{x}' := \text{DeepSearch}(\mathbf{x}, \mathbf{x}', f, d, \mathcal{G})$ 
11    if  $1 < k/m$  then
12       $\mathcal{G} := \text{Divide-Group}(\mathcal{G}, m)$ 
13       $k := k/m$ 
14  until  $N_f(\mathbf{x}) \neq N_f(\mathbf{x}')$ , or reached query budget  $L$ 
15  return  $\mathbf{x}'$ 

```

where each set $G_{i,j}$ ($1 \leq j \leq m \times m$) contains indices corresponding to $k/m \times k/m$ neighboring pixels in \mathbf{x} . After splitting all sets, the total number of sets is multiplied by $m \times m$. This results in a more fine-grained partition of input \mathbf{x} . We then go back to step (2).

In the query-limited setting, we use single-step DeepSearch on line 10, i.e., we fix MaxNum to 1 in Alg. 2 and Alg. 2', and choose \mathbf{x}_{init} to be a vertex of $\mathcal{B}(\mathbf{x}, d)$ to avoid unnecessary queries. Hence, when there are $\lceil \frac{n}{k^2} \rceil$ sets in \mathcal{G} , the total number of queries per iteration in Alg. 4 reduces to $\lceil \frac{n}{k^2} \rceil$.

7 EXPERIMENTAL EVALUATION

We evaluate DeepSearch by using it to test the robustness of deep neural networks trained for popular datasets. We also compare its effectiveness with state-of-the-art blackbox attacks. Our experiments are designed around the following research questions:

- RQ1:** Is DeepSearch effective in finding adversarial examples?
- RQ2:** Is DeepSearch effective in finding adversarial examples with low distortion rates?
- RQ3:** Is DeepSearch a query-efficient blackbox attack?
- RQ4:** Is the hierarchical grouping of DeepSearch effective in improving query efficiency?

We make our implementation open source¹. Our experimental data, including detected adversarial examples, are also available via the provided link.

7.1 Evaluation Setup

Datasets and network models. We evaluate our approach on deep neural networks trained for three well known datasets, namely SVHN [53] (cropped digits), CIFAR-10 [38], and ImageNet [61]. For

each dataset, we randomly selected 1000 correctly classified images from the test set on which to perform blackbox attacks.

For SVHN and CIFAR-10, we attack two wide ResNet w32-10 [83] networks, where one is naturally trained (without defense) and the other is adversarially trained with a state-of-the-art defense [49]. For SVHN, the undefended network we trained has 95.96% test accuracy, and the adversarially trained network has 93.70% test accuracy. During adversarial training, we used the PGD attack [49] (that can perturb each pixel by at most 8 on the 0–255 pixel scale) to generate adversarial examples. For CIFAR-10, we trained an undefended network with 95.07% test accuracy. The defended network we attack is the pretrained network provided in Madry's challenge². For ImageNet, we attack a pretrained Inception v3 network [71], which is undefended.

Defenses for ImageNet networks are also important, and we would have attacked defended ImageNet networks if they were publicly available. In this work, we did not attack such networks (using the defense in [49]) for the following reasons. First, there are no publicly available ImageNet networks that use the defense in [49]. Second, none of the state-of-the-art attacks that we used for comparison in our paper (i.e., [27, 32, 33, 50]) have been evaluated on defended networks for this dataset. Therefore, we did not compare DeepSearch with these attacks on such networks. Third, due to the extremely high computational cost, implementing the defense in [49] for an ImageNet network is impractical.

Existing approaches. We compare DeepSearch with four state-of-the-art blackbox attacks for generating adversarial examples:

- The NES attack [32], optimized for the L_∞ distance metric, is developed for various settings, including a query-limited setting. It uses natural evolution strategies (NES) [62] for gradient estimation and performs projected gradient-descent (PGD) [49] style adversarial attacks using estimated gradients. We compare with the NES attack developed for a query-limited setting, i.e., QL-NES.
- The Bandits attack [33] extended gradient-estimation-based blackbox attacks by integrating gradient priors, e.g., time-dependent and data-dependent priors, through a bandit optimization framework. The Bandits attack can perform both L_2 and L_∞ attacks.
- The Simple Blackbox Attack [27] is optimized for the L_2 distance metric. Starting from an input image, it finds adversarial examples by repeatedly adding or subtracting a random vector sampled from a set of predefined orthogonal candidate vectors. We compare with their SimBA algorithm, which can also be easily constrained using L_∞ distance.
- The Parsimonious blackbox attack [50], optimized for the L_∞ distance metric, encodes the problem of finding adversarial perturbations as finding solutions to linear programs. For an input \mathbf{x} and distance d , it searches among the vertices of $\mathcal{B}(\mathbf{x}, d)$ and finds adversarial examples by using efficient algorithms in combinatorial optimization.

DeepSearch implementation. We briefly introduce some implementation details of DeepSearch. In our implementation, the classification function f of multiclass classifiers maps input images

¹<https://github.com/Practical-Formal-Methods/DeepSearch>

²https://github.com/MadryLab/cifar10_challenge

to the logarithm of class probabilities predicted by neural networks. In this setting, the objective function in Alg. 2 (resp. Alg. 2') corresponds to logit loss [12] (resp. cross-entropy loss [23]).

To reduce the number of queries, for input \mathbf{x} and distance d , we choose $\mathbf{x}_{init} \in \mathcal{B}(\mathbf{x}, d)$ such that it is derived from \mathbf{x} by setting the values of all its pixels to the lower bounds in $\mathcal{B}(\mathbf{x}, d)$. In the refinement step, when calculating new adversarial examples within smaller distances, we set \mathbf{x}_{new} to \mathbf{x}_{init} for convenience.

In our experiments, we used Alg. 2 to attack the undefended networks. We used Alg. 2' to attack the defended networks since they are more robust against cross-entropy model attacks. To attack the SVHN and CIFAR-10 networks, we set the initial grouping size to 4×4 . For ImageNet, we set the initial grouping size to 32×32 due to their large image size. For group splitting, we set $m = 2$ so that we always divide a group into 2×2 subgroups.

In the hierarchical-grouping setting, we mutate groups of pixels in random orders. To avoid overshooting query budgets, we mutate groups in batches, and the batch size is 64 in all our experiments.

Parameter settings. For all datasets, we set L_∞ distance $d = 8$ on the 0–255 pixel scale to perform attacks. For both SVHN and CIFAR-10 networks, we set the query budget to 20,000. For the ImageNet network, we set the query budget to 10,000 as done in related work [33, 50].

For the QL-NES attack, we set $\sigma = 0.001$, size of NES population $n = 100$, learning rate $\eta = 0.001$, and momentum $\beta = 0.9$ for SVHN and CIFAR-10. We set $\sigma = 0.01$, size of NES population $n = 100$, learning rate $\eta = 0.0001$, and momentum $\beta = 0.9$ for ImageNet.

For the Bandits attack, we set OCO learning rate $\eta = 0.001$, image learning rate $h = 0.0001$, bandits exploration $\delta = 0.1$, finite difference probe $\eta = 0.1$, and tile size to 16 for SVHN and CIFAR-10. We set OCO learning rate $\eta = 1$, image learning rate $h = 0.0001$, bandits exploration $\delta = 1$, finite difference probe $\eta = 0.1$, and tile size to 64 for ImageNet.

For the Parsimonious attack, we use the parameters mentioned in their paper for CIFAR-10 and ImageNet networks. For SVHN, we use the same parameters as for CIFAR-10. Moreover, their implementation offers both cross-entropy and logit loss to construct attacks. We tried both loss functions in our experiments and select the one with better performance for comparison.

7.2 Metrics

In our evaluation, we use the following metrics.

Success rate. The success rate measures the percentage of input images for which adversarial examples are found. The higher this rate, the more effective a given technique in finding adversarial examples. Assume we write $findAdv(\mathbf{x})$ to denote whether an adversarial example is found for input \mathbf{x} . If so, we have that $findAdv(\mathbf{x}) = 1$; otherwise, we have $findAdv(\mathbf{x}) = 0$. For a set of images $\mathbf{X} = \{\mathbf{x}_1, \dots, \mathbf{x}_k\}$, the success rate of a given technique is:

$$ASR(\mathbf{X}) = \frac{1}{k} \sum_{i=1}^k findAdv(\mathbf{x}_i)$$

Average distortion rate. Let sets $\mathbf{X} = \{\mathbf{x}_1, \dots, \mathbf{x}_k\}$ and $\mathbf{X}_{adv} = \{\mathbf{x}'_1, \dots, \mathbf{x}'_k\}$ be input images and adversarial examples, respectively.

Table 1: Results on SVHN networks.

Attack	Success rate	Avg. L_∞	Avg. L_2	Avg. queries	Med. queries
Undefended network					
QL-NES	62.4%	2.58%	1.80%	2157	1700
Bandits	99.2%	3.43%	2.69%	762	573
SimBA	84.7%	4.65%	3.47%	1675	1430
Parsimonious	100%	4.59%	7.63%	337	231
DeepSearch	100%	1.89%	3.17%	229	196
Defended network					
QL-NES	40.5%	4.10%	4.19%	5574	3900
Bandits	55.3%	4.38%	4.74%	2819	944
SimBA	65.9%	4.96%	3.95%	2687	2633
Parsimonious	78.9%	4.86%	8.08%	2174	423.5
DeepSearch	83.1%	3.35%	5.58%	1808	458

The average distortion rate between \mathbf{X} and \mathbf{X}_{adv} with respect to the L_∞ distance is:

$$AvgDR_{L_\infty}(\mathbf{X}, \mathbf{X}_{adv}) = \frac{1}{k} \sum_{i=1}^k \frac{\|\mathbf{x}_i - \mathbf{x}'_i\|_{L_\infty}}{\|\mathbf{x}_i\|_{L_\infty}}$$

As shown here, the lower this rate, the more subtle the adversarial examples. For approaches that achieve similar misclassification rates, we use this metric to determine which approach finds more subtle perturbations. The average distortion rate with respect to L_2 can be derived by substituting L_∞ with L_2 in the above definition. In our experimental results, we also include the average L_2 distortion for reference.

Average queries. For blackbox attacks, we use the number of queries required to find adversarial examples to measure their efficiency. An approach that requires more queries to perform attacks is more costly. For each approach, we calculate the average number of queries required to perform successful attacks. We also list their mean number of queries for reference. We point out that queries made by the refinement step of DeepSearch are not counted when calculating average queries because refinement starts after adversarial examples are already found.

7.3 Experimental Results

DeepSearch outperforms all other blackbox attacks in success rate, average queries, and average distortion rate. Experimental results are shown in Fig. 4 and Tabs. 1–5.

Results on success rate (RQ1). DeepSearch is very effective in finding adversarial examples for both undefended and defended networks. Compared to other blackbox attacks, DeepSearch has the highest attack success rate.

For undefended networks, the success rate of DeepSearch is close to 100% for all three datasets. For the SVHN defended network, DeepSearch has a success rate of 83.1%, which is 4.2% higher than that of the Parsimonious attack (the second best attack in success rate). For the CIFAR-10 defended network, DeepSearch has a success rate of 47.7%.

Results on average distortion rate (RQ2). DeepSearch has found adversarial examples with the lowest average L_∞ distortion

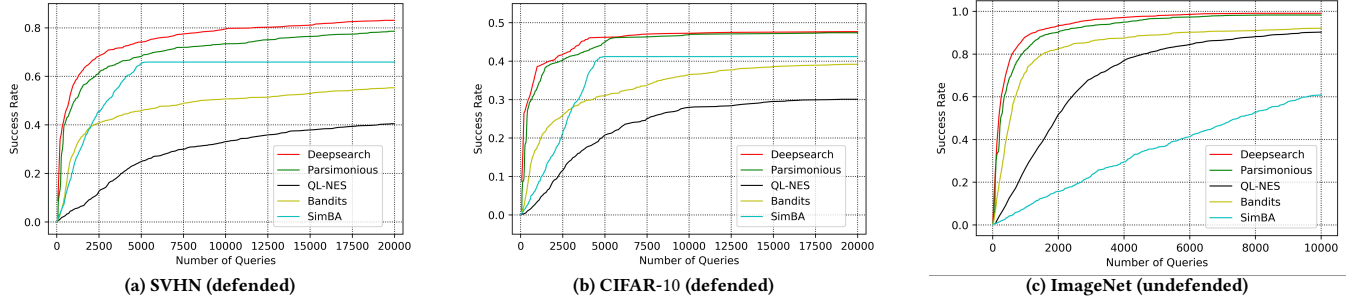


Figure 4: Results on success rate w.r.t number of queries.

Table 2: Results on CIFAR-10 networks.

Attack	Success rate	Avg. L_∞	Avg. L_2	Avg. queries	Med. queries
Undefended network					
QL-NES	52.8%	1.24%	0.99%	1360	1100
Bandits	92.6%	2.66%	2.34%	838	616
SimBA	71.6%	3.36%	2.19%	1311	1150
Parsimonious	100%	3.36%	6.36%	339	238.5
DeepSearch	100%	1.64%	3.08%	247	196
Defended network					
QL-NES	30.1%	2.71%	3.09%	4408	3200
Bandits	39.2%	2.95%	4.39%	2952	1176
SimBA	41.2%	3.46%	4.50%	2425	2424
Parsimonious	47.4%	3.45%	6.61%	1228	366
DeepSearch	47.7%	2.48%	4.70%	963	196

Table 3: Results on ImageNet undefended network.

Attack	Success rate	Avg. L_∞	Avg. L_2	Avg. queries	Med. queries
QL-NES	90.3%	1.83%	1.75%	2300	1800
Bandits	92.1%	2.15%	2.61%	930	496
SimBA	61%	3.15%	0.67%	4379	4103
Parsimonious	98.3%	3.16%	6.35%	660	241
DeepSearch	99.3%	1.50%	3.05%	561	196

rate for networks of all datasets. For the CIFAR-10 undefended network, QL-NES has a success rate of only 52.8%. To have a more fair comparison, for those 52.8% images that are successfully attacked by both DeepSearch and QL-NES, we further calculate the average distortion rate of the adversarial examples found by DeepSearch. We find it to be 1.2%, which is lower than that of QL-NES. Although DeepSearch is an L_∞ attack, adversarial examples found by DeepSearch also have low average L_2 distortion rate.

Results on query efficiency (RQ3). DeepSearch outperforms all other attacks in query efficiency. Compared to the Parsimonious attack (the second best attack in query efficiency), DeepSearch reduces the average queries by 15–32% across all datasets.

Results on query reduction (RQ4). We demonstrate the effectiveness of hierarchical grouping in DeepSearch for query reduction. We use DeepSearch to attack networks with different initial

Table 4: Query reduction for SVHN and CIFAR-10. For each dataset, success rate (resp. average queries) is shown in the first (resp. second) row.

Dataset	1	2×2	4×4	8×8	16×16
Undefended network					
SVHN	100%	100%	100%	100%	100%
	742	300	229	238	242
CIFAR-10	100%	100%	100%	100%	100%
	462	301	247	255	259
Defended network					
SVHN	81.3%	82.4%	83.1%	83.6%	83.9%
	3143	2292	1808	1565	1591
CIFAR-10	47.7%	47.4%	47.7%	47.6%	47.6%
	2292	1156	963	935	946

Table 5: Query reduction for ImageNet. Success rate (resp. average queries) is shown in the first (resp. second) row.

Dataset	8×8	16×16	32×32	64×64	128×128
ImageNet	99.1%	99.1%	99.1%	99.4%	99.4%
	765	580	533	554	580

group sizes and show their corresponding success rates and average queries in Tabs. 4 and 5.

We first notice that the initial group size can only slightly affect attack success rate. For the SVHN defended network, the success rate increases from 81.3% to 83.9% as initial group size increases from 1 to 16×16. However, the changes in success rate are negligible for all other networks.

On the other hand, we observe that hierarchical grouping improves query efficiency dramatically. We take the average queries of group size 1 and 4×4 as an example for SVHN and CIFAR-10 networks. For the SVHN undefended network, we see a 69.1% decrease of average queries from 742 to 229. For the SVHN defended network, average queries are reduced by 42.5% from 3143 to 1808. For the CIFAR-10 undefended network, the average queries are decreased by 46.5% from 462 to 247. For the CIFAR-10 defended network, average queries are decreased by 58% from 2292 to 963, and for the ImageNet network, from group size 8×8 to 32×32, we decreased the average queries by 30.3% from 765 to 533.

7.4 Threats to Validity

We have identified the following three threats to the validity of our experiments.

Datasets and network models. Our experimental results may not generalize to other datasets or network models. However, we used three of the most popular datasets for image classification, SVHN, CIFAR-10, and ImageNet. Moreover, our network models have very high test accuracy and the defense we use based on adversarial training is state of the art.

Existing approaches. The second threat is related to the choice of existing approaches with which we compare. DeepSearch uses iterative linearization of non-linear neural networks and is tailored to the L_∞ distance metric. We, thus, compare with approaches that can also perform L_∞ attacks. To our knowledge, the blackbox attacks with which we compared are all state-of-the-art L_∞ attacks.

Fairness of comparison. The selection of parameters for each approach could affect the fairness of our comparison. We tried various parameters for each attack and choose the ones yielding best performance.

8 RELATED WORK

Adversarial robustness. Szegedy et al. [72] first discovered adversarial examples in neural networks and used box-constrained L-BFGS to find them. Since then, multiple whitebox adversarial attacks have been proposed: FGSM [24], BIM [40], DeepFool [51], JSMA [57], PGD [49], and C&W [12]. Goodfellow et al. [24] first argued that the primary cause of adversarial examples is the linear nature of neural networks, and they proposed FGSM that allows fast generation of adversarial examples. BIM improved FGSM by extending it with iterative procedures. DeepFool [51] is another method that performs adversarial attacks through iterative linearization of neural networks.

Blackbox adversarial attacks are more difficult than whitebox ones, and many blackbox attacks require a large number of queries. Papernot et al. [55, 56] explored blackbox attacks based on the phenomenon of transferability [55, 72]. Chen et al. [13] and Bhagoji et al. [5] proposed blackbox attacks based on gradient estimation [41, 68]. Uesato et al. [75] used SPSA [67]. Ilyas et al. [32, 33] used NES [62] and proposed the Bandits attack. Narodytska et al. [52] performed a local-search-based attack. The boundary attack [6] only requires access to the final decision of neural networks. Guo et al. [27] further considered perturbations in low frequency space. Moon et al. [50] leveraged algorithms in combinatorial optimization.

Although research on developing adversarial attacks is moving fast, research on defending neural networks against adversarial attacks is relatively slow [1, 2, 9–12, 15, 20, 29, 46, 65]. Many defense techniques are shown to be ineffective soon after they have been developed. We refer to the work of Carlini et al. [8] for a more detailed discussion on evaluating adversarial robustness.

Testing deep neural networks. Recently, significant progress has been made on testing neural networks. Several useful test coverage criteria have been proposed to guide test case generation: DeepXplore [58] proposed neuron coverage and the first whitebox testing framework for neural networks; DeepGauge [47] proposed

a set of finer-grained test coverage criteria; DeepCT [48] further proposed combinatorial test coverage for neural networks; Sun et al. [69] proposed coverage criteria inspired by MC/DC; Kim et al. [36] proposed surprise adequacy for deep learning systems. Sekhon et al. [63] and Li et al. [43] pointed out the limitation of existing structural coverage criteria for neural networks. Li et al. [63] also discussed improvements for better coverage criteria.

Moreover, Sun et al. [70] proposed the first concolic testing [22, 64] approach for neural networks. DeepCheck [26] tests neural networks based on symbolic execution [14, 37]. TensorFuzz [54] proposed the first framework of coverage-guided fuzzing for neural networks. DeepHunter [82] considered various mutation strategies for their fuzzing framework. Wicker et al. [78] extracted features from images and computed adversarial examples using a two-player turn-based stochastic game. DLFuzz [28] proposed the first differential fuzzing framework for deep learning systems. DeepTest [73] and DeepRoad [86] proposed testing tools for autonomous driving systems based on deep neural networks. For more on testing neural networks, we refer to the work of Zhang et al. [85] that surveys testing of machine-learning systems.

Formal verification of deep neural networks. Verification of neural networks is more challenging than testing. Early work [59] used abstract interpretation [16] to verify small-sized neural networks. Recent work [25, 35, 66] used SMT [3] techniques and considered new abstract domains.

Liu et al. [44] classified recent work in the area into five categories: Reachability-analysis based approaches include MaxSens [81], ExactReach [80], and AI² [21]; NSVerify [45], MIPVerify [74] and ILP [4] are based on primal optimization; Duality [18], ConvDual [79] and Certify [60] use dual optimization; Fast-Lin and Fast-Lip [77], ReluVal [76] and DLV [31] combine reachability with search; Sherlock [17], Reluplex [34], Planet [19] and BaB [7] combine search with optimization. Lie et al. [44] provide a more detailed comparison and discussion of the above mentioned work.

9 CONCLUSION AND FUTURE WORK

We proposed and implemented DeepSearch, a novel blackbox-fuzzing technique for attacking deep neural networks. DeepSearch is simple and effective in finding adversarial examples with low distortion, and it outperforms state-of-the-art blackbox attacks in a query-limited setting. In our future work, we will continue improving the effectiveness of DeepSearch for an even more query-efficient L_∞ attack. We are also interested in extending DeepSearch to construct query-efficient L_2 attacks.

Designing effective defenses to secure deep neural networks against adversarial attacks is non-trivial. In this paper, we did not focus on proposing defenses against blackbox attacks. Instead, we attacked neural networks with adversarial training-based defenses [49]. Another interesting direction for future work is to develop defense techniques that specifically target blackbox attacks, for instance by identifying patterns in their sequences of queries.

ACKNOWLEDGMENTS

We are grateful to the reviewers for their valuable feedback.

This work was supported by DFG grant 389792660 as part of TRR 248 (see <https://perspicuous-computing.science>).

REFERENCES

- [1] Anish Athalye and Nicholas Carlini. 2018. On the Robustness of the CVPR 2018 White-Box Adversarial Example Defenses. *CoRR* abs/1804.03286 (2018).
- [2] Anish Athalye, Nicholas Carlini, and David A. Wagner. 2018. Obfuscated Gradients Give a False Sense of Security: Circumventing Defenses to Adversarial Examples. In *ICML (PMLR, Vol. 80)*. PMLR, 274–283.
- [3] Clark W. Barrett and Cesare Tinelli. 2018. Satisfiability Modulo Theories. In *Handbook of Model Checking*. Springer, 305–343.
- [4] Osbert Bastani, Yani Ioannou, Leonidas Lampropoulos, Dimitrios Vytiniotis, Aditya V. Nori, and Antonio Criminisi. 2016. Measuring Neural Net Robustness with Constraints. In *NIPS*. 2613–2621.
- [5] Arjun Nitin Bhagoji, Warren He, Bo Li, and Dawn Song. 2018. Practical Black-Box Attacks on Deep Neural Networks Using Efficient Query Mechanisms. In *ECCV (LNCS, Vol. 11216)*. Springer, 158–174.
- [6] Wieland Brendel, Jonas Rauber, and Matthias Bethge. 2018. Decision-Based Adversarial Attacks: Reliable Attacks Against Black-Box Machine Learning Models. In *ICLR*. OpenReview.net.
- [7] Rudy Bunel, Ilker Turkaslan, Philip H. S. Torr, Pushmeet Kohli, and Pawan Kumar Mudigonda. 2018. A Unified View of Piecewise Linear Neural Network Verification. In *NeurIPS*. 4795–4804.
- [8] Nicholas Carlini, Anish Athalye, Nicolas Papernot, Wieland Brendel, Jonas Rauber, Dimitris Tsipras, Ian J. Goodfellow, Aleksander Madry, and Alexey Kurakin. 2019. On Evaluating Adversarial Robustness. *CoRR* abs/1902.06705 (2019).
- [9] Nicholas Carlini and David A. Wagner. 2016. Defensive Distillation is Not Robust to Adversarial Examples. *CoRR* abs/1607.04311 (2016).
- [10] Nicholas Carlini and David A. Wagner. 2017. Adversarial Examples Are Not Easily Detected: Bypassing Ten Detection Methods. In *AISeC@CCS*. ACM, 3–14.
- [11] Nicholas Carlini and David A. Wagner. 2017. MagNet and “Efficient Defenses Against Adversarial Attacks” Are Not Robust to Adversarial Examples. *CoRR* abs/1711.08478 (2017).
- [12] Nicholas Carlini and David A. Wagner. 2017. Towards Evaluating the Robustness of Neural Networks. In *S&P*. IEEE Computer Society, 39–57.
- [13] Pin-Yu Chen, Huan Zhang, Yash Sharma, Jinfeng Yi, and Cho-Jui Hsieh. 2017. ZOO: Zeroth Order Optimization Based Black-box Attacks to Deep Neural Networks Without Training Substitute Models. In *AISeC@CCS*. ACM, 15–26.
- [14] Lori A. Clarke. 1976. A System to Generate Test Data and Symbolically Execute Programs. *TSE* 2 (1976), 215–222. Issue 3.
- [15] Cory Cornelius. 2019. The Efficacy of SHIELD Under Different Threat Models. *CoRR* abs/1902.00541 (2019).
- [16] Patrick Cousot and Radhia Cousot. 1977. Abstract Interpretation: A Unified Lattice Model for Static Analysis of Programs by Construction or Approximation of Fixpoints. In *POPL*. ACM, 238–252.
- [17] Souradeep Dutta, Susmit Jha, Sriram Sankaranarayanan, and Ashish Tiwari. 2018. Output Range Analysis for Deep Feedforward Neural Networks. In *NFM (LNCS, Vol. 10811)*. Springer, 121–138.
- [18] Krishnamurthy Dvijotham, Robert Stanforth, Sven Gowal, Timothy A. Mann, and Pushmeet Kohli. 2018. A Dual Approach to Scalable Verification of Deep Networks. In *UAI*. AUAI Press, 550–559.
- [19] Rüdiger Ehlers. 2017. Formal Verification of Piece-Wise Linear Feed-Forward Neural Networks. In *ATVA (LNCS, Vol. 10482)*. Springer, 269–286.
- [20] Logan Engstrom, Andrew Ilyas, and Anish Athalye. 2018. Evaluating and Understanding the Robustness of Adversarial Logit Pairing. *CoRR* abs/1807.10272 (2018).
- [21] Timon Gehr, Matthew Mirman, Dana Drachler-Cohen, Petar Tsankov, Swarat Chaudhuri, and Martin Vechev. 2018. AI2: Safety and Robustness Certification of Neural Networks with Abstract Interpretation. In *S&P*. IEEE Computer Society, 3–18.
- [22] Patrice Godefroid, Nils Klarlund, and Koushik Sen. 2005. DART: Directed Automated Random Testing. In *PLDI*. ACM, 213–223.
- [23] Ian Goodfellow, Yoshua Bengio, and Aaron Courville. 2016. *Deep Learning*. MIT Press.
- [24] Ian J. Goodfellow, Jonathon Shlens, and Christian Szegedy. 2015. Explaining and Harnessing Adversarial Examples. In *ICLR*.
- [25] Divya Gopinath, Guy Katz, Corina S. Pasareanu, and Clark W. Barrett. 2018. DeepSafe: A Data-Driven Approach for Assessing Robustness of Neural Networks. In *ATVA (LNCS, Vol. 11138)*. Springer, 3–19.
- [26] Divya Gopinath, Kaiyuan Wang, Mengshi Zhang, Corina S. Pasareanu, and Saffraz Khurshid. 2018. Symbolic Execution for Deep Neural Networks. *CoRR* abs/1807.10439 (2018).
- [27] Chuan Guo, Jacob R. Gardner, Yurong You, Andrew Gordon Wilson, and Kilian Q. Weinberger. 2019. Simple Black-box Adversarial Attacks. In *ICML*. PMLR.
- [28] Jianmin Guo, Yu Jiang, Yue Zhao, Quan Chen, and Jianguang Sun. 2018. DLFuzz: Differential Fuzzing Testing of Deep Learning Systems. In *ESEC/FSE*. ACM, 739–743.
- [29] Warren He, James Wei, Xinyun Chen, Nicholas Carlini, and Dawn Song. 2017. Adversarial Example Defense: Ensembles of Weak Defenses Are Not Strong. In *WOOT*. USENIX.
- [30] Geoffrey Hinton, Li Deng, Dong Yu, George Dahl, Abdel-rahman Mohamed, Navdeep Jaitly, Andrew Senior, Vincent Vanhoucke, Patrick Nguyen, Tara Sainath, and Brian Kingsbury. 2012. Deep Neural Networks for Acoustic Modeling in Speech Recognition. *Signal Process. Mag.* 29 (2012), 82–97. Issue 6.
- [31] Xiaowei Huang, Marta Kwiatkowska, Sen Wang, and Min Wu. 2017. Safety Verification of Deep Neural Networks. In *CAV (LNCS, Vol. 10426)*. Springer, 3–29.
- [32] Andrew Ilyas, Logan Engstrom, Anish Athalye, and Jessy Lin. 2018. Black-Box Adversarial Attacks with Limited Queries and Information. In *ICML (PMLR, Vol. 80)*. PMLR, 2142–2151.
- [33] Andrew Ilyas, Logan Engstrom, and Aleksander Madry. 2019. Prior Convictions: Black-box Adversarial Attacks with Bandits and Priors. In *ICLR*.
- [34] Guy Katz, Clark W. Barrett, David L. Dill, Kyle Julian, and Mykel J. Kochenderfer. 2017. Reluplex: An Efficient SMT Solver for Verifying Deep Neural Networks. In *CAV (LNCS, Vol. 10426)*. Springer, 97–117.
- [35] Guy Katz, Derek A. Huang, Duligur Ibeling, Kyle Julian, Christopher Lazarus, Rachel Lim, Parth Shah, Shantanu Thakoor, Haoze Wu, Aleksandar Zeljic, David L. Dill, Mykel J. Kochenderfer, and Clark W. Barrett. 2019. The Marabou Framework for Verification and Analysis of Deep Neural Networks. In *CAV (LNCS, Vol. 11561)*. Springer, 443–452.
- [36] Jinhan Kim, Robert Feldt, and Shin Yoo. 2019. Guiding Deep Learning System Testing Using Surprise Adequacy. In *ICSE*. IEEE Computer Society/ACM, 1039–1049.
- [37] James C. King. 1976. Symbolic Execution and Program Testing. *CACM* 19 (1976), 385–394. Issue 7.
- [38] Alex Krizhevsky. 2009. *Learning Multiple Layers of Features from Tiny Images*. Technical Report. University of Toronto.
- [39] Alex Krizhevsky, Ilya Sutskever, and Geoffrey E. Hinton. 2017. ImageNet Classification with Deep Convolutional Neural Networks. *CACM* 60 (2017), 84–90. Issue 6.
- [40] Alexey Kurakin, Ian J. Goodfellow, and Samy Bengio. 2017. Adversarial Examples in the Physical World. In *ICLR*. OpenReview.net.
- [41] Peter D. Lax and Maria Shea Terrell. 2014. *Calculus with Applications*. Springer.
- [42] Yann LeCun, Léon Bottou, Yoshua Bengio, and Patrick Haffner. 1998. Gradient-Based Learning Applied to Document Recognition. In *Proc. IEEE*. IEEE Computer Society, 2278–2324.
- [43] Zenan Li, Xiaoxing Ma, Chang Xu, and Chun Cao. 2019. Structural Coverage Criteria for Neural Networks Could Be Misleading. In *ICSE (NIER)*. IEEE Computer Society/ACM, 89–92.
- [44] Changliu Liu, Tomer Arnon, Christopher Lazarus, Clark W. Barrett, and Mykel J. Kochenderfer. 2019. Algorithms for Verifying Deep Neural Networks. *CoRR* abs/1903.06758 (2019).
- [45] Alessio Lomuscio and Lalit Maganti. 2017. An Approach to Reachability Analysis for Feed-Forward ReLU Neural Networks. *CoRR* abs/1706.07351 (2017).
- [46] Pei-Hsuan Lu, Pin-Yu Chen, Kang-Cheng Chen, and Chia-Mu Yu. 2018. On the Limitation of MagNet Defense Against L1-Based Adversarial Examples. In *DSN Workshops*. IEEE Computer Society, 200–214.
- [47] Lei Ma, Felix Juefei-Xu, Fuyuan Zhang, Jiyuan Sun, Minhui Xue, Bo Li, Chunyang Chen, Ting Su, Li Li, Yang Liu, Jianjun Zhao, and Yadong Wang. 2018. DeepGauge: Multi-Granularity Testing Criteria for Deep Learning Systems. In *ASE*. ACM, 120–131.
- [48] Lei Ma, Fuyuan Zhang, Minhui Xue, Bo Li, Yang Liu, Jianjun Zhao, and Yadong Wang. 2018. Combinatorial Testing for Deep Learning Systems. *CoRR* abs/1806.07723 (2018).
- [49] Aleksander Madry, Aleksandar Makelov, Ludwig Schmidt, Dimitris Tsipras, and Adrian Vladu. 2018. Towards Deep Learning Models Resistant to Adversarial Attacks. In *ICLR*. OpenReview.net.
- [50] Seunghyung Moon, Gaon An, and Hyun Oh Song. 2019. Parsimonious Black-Box Adversarial Attacks via Efficient Combinatorial Optimization. In *ICML*. PMLR.
- [51] Seyed-Mohsen Moosavi-Dezfooli, Alhussein Fawzi, and Pascal Frossard. 2016. DeepFool: A Simple and Accurate Method to Fool Deep Neural Networks. In *CVPR*. IEEE Computer Society, 2574–2582.
- [52] Nina Narodytska and Shiva Prasad Kasiviswanathan. 2017. Simple Black-Box Adversarial Attacks on Deep Neural Networks. In *CVPR Workshops*. IEEE Computer Society, 1310–1318.
- [53] Yuval Netzer, Tao Wang, Adam Coates, Alessandro Bissacco, Bo Wu, and Andrew Y. Ng. 2011. Reading Digits in Natural Images with Unsupervised Feature Learning. In *NIPS Workshop on Deep Learning and Unsupervised Feature Learning*.
- [54] Augustus Odena, Catherine Olsson, David Andersen, and Ian J. Goodfellow. 2019. TensorFuzz: Debugging Neural Networks with Coverage-Guided Fuzzing. In *ICML (PMLR, Vol. 97)*. PMLR, 4901–4911.
- [55] Nicolas Papernot, Patrick D. McDaniel, Ian J. Goodfellow, and Ian J. Goodfellow. 2016. Transferability in Machine Learning: From Phenomena to Black-Box Attacks Using Adversarial Samples. *CoRR* abs/1605.07277 (2016).
- [56] Nicolas Papernot, Patrick D. McDaniel, Ian J. Goodfellow, Somesh Jha, Z. Berkay Celik, and Ananthram Swami. 2017. Practical Black-Box Attacks Against Machine Learning. In *AsiaCCS*. ACM, 506–519.

- [57] Nicolas Papernot, Patrick D. McDaniel, Somesh Jha, Matt Fredrikson, Z. Berkay Celik, and Ananthram Swami. 2016. The Limitations of Deep Learning in Adversarial Settings. In *EuroS&P*. IEEE Computer Society, 372–387.
- [58] Kexin Pei, Yinzhao Cao, Junfeng Yang, and Suman Jana. 2017. DeepXplore: Automated Whitebox Testing of Deep Learning Systems. In *SOSP*. ACM, 1–18.
- [59] Luca Pulina and Armando Tacchella. 2010. An Abstraction-Refinement Approach to Verification of Artificial Neural Networks. In *CAV (LNCS, Vol. 6174)*. Springer, 243–257.
- [60] Aditi Raghunathan, Jacob Steinhardt, and Percy Liang. 2018. Certified Defenses Against Adversarial Examples. In *ICLR*. OpenReview.net.
- [61] Olga Russakovsky, Jia Deng, Hao Su, Jonathan Krause, Sanjeev Satheesh, Sean Ma, Zhiheng Huang, Andrej Karpathy, Aditya Khosla, Michael Bernstein, Alexander C. Berg, and Fei-Fei Li. 2015. ImageNet Large Scale Visual Recognition Challenge. *IJCV* (2015), 211–252.
- [62] Tim Salimans, Jonathan Ho, Xi Chen, and Ilya Sutskever. 2017. Evolution Strategies as a Scalable Alternative to Reinforcement Learning. *CoRR* abs/1703.03864 (2017).
- [63] Jasmine Sekhon and Cody Fleming. 2019. Towards Improved Testing for Deep Learning. In *ICSE (NIER)*. IEEE Computer Society/ACM, 85–88.
- [64] Koushik Sen, Darko Marinov, and Gul Agha. 2005. CUTE: A Concolic Unit Testing Engine for C. In *ESEC/FSE*. ACM, 263–272.
- [65] Yash Sharma and Pin-Yu Chen. 2018. Bypassing Feature Squeezing by Increasing Adversary Strength. *CoRR* abs/1803.09868 (2018).
- [66] Gagandeep Singh, Timon Gehr, Markus Püschel, and Martin T. Vechev. 2019. An Abstract Domain for Certifying Neural Networks. *PACMPL* 3 (2019), 41:1–41:30. Issue POPL.
- [67] James C. Spall. 1992. Multivariate Stochastic Approximation Using a Simultaneous Perturbation Gradient Approximation. *TAC* 37 (1992), 332–341. Issue 3.
- [68] James C. Spall. 2003. *Introduction to Stochastic Search and Optimization*. John Wiley and Sons.
- [69] Yousheng Sun, Xiaowei Huang, and Daniel Kroening. 2018. Testing Deep Neural Networks. *CoRR* abs/1803.04792 (2018).
- [70] Yousheng Sun, Min Wu, Wenjie Ruan, Xiaowei Huang, Marta Kwiatkowska, and Daniel Kroening. 2018. Concolic Testing for Deep Neural Networks. In *ASE*. ACM, 109–119.
- [71] Christian Szegedy, Vincent Vanhoucke, Sergey Ioffe, Jonathon Shlens, and Zbigniew Wojna. 2016. Rethinking the Inception Architecture for Computer Vision. In *CVPR*. IEEE Computer Society, 2818–2826.
- [72] Christian Szegedy, Wojciech Zaremba, Ilya Sutskever, Joan Bruna, Dumitru Erhan, Ian J. Goodfellow, and Rob Fergus. 2014. Intriguing Properties of Neural Networks. In *ICLR*.
- [73] Yuchi Tian, Kexin Pei, Suman Jana, and Baishakhi Ray. 2018. DeepTest: Automated Testing of Deep-Neural-Network-Driven Autonomous Cars. In *ICSE*. ACM, 303–314.
- [74] Vincent Tjeng, Kai Y. Xiao, and Russ Tedrake. 2019. Evaluating Robustness of Neural Networks with Mixed Integer Programming. In *ICLR*. OpenReview.net.
- [75] Jonathan Uesato, Brendan O'Donoghue, Pushmeet Kohli, and Aaron van den Oord. 2018. Adversarial Risk and the Dangers of Evaluating Against Weak Attacks. In *ICML (PMLR, Vol. 80)*. PMLR, 5032–5041.
- [76] Shiqi Wang, Kexin Pei, Justin Whitehouse, Junfeng Yang, and Suman Jana. 2018. Formal Security Analysis of Neural Networks Using Symbolic Intervals. In *Security*. USENIX, 1599–1614.
- [77] Tsui-Wei Weng, Huan Zhang, Hongge Chen, Zhao Song, Cho-Jui Hsieh, Luca Daniel, Duane S. Boning, and Inderjit S. Dhillon. 2018. Towards Fast Computation of Certified Robustness for ReLU Networks. In *ICML (PMLR, Vol. 80)*. PMLR, 5273–5282.
- [78] Matthew Wicker, Xiaowei Huang, and Marta Kwiatkowska. 2018. Feature-Guided Black-Box Safety Testing of Deep Neural Networks. In *TACAS (LNCS, Vol. 10805)*. Springer, 408–426.
- [79] Eric Wong and J. Zico Kolter. 2018. Provable Defenses Against Adversarial Examples via the Convex Outer Adversarial Polytope. In *ICML (PMLR, Vol. 80)*. PMLR, 5283–5292.
- [80] Weiming Xiang, Hoang-Dung Tran, and Taylor T. Johnson. 2017. Reachable Set Computation and Safety Verification for Neural Networks with ReLU Activations. *CoRR* abs/1712.08163 (2017).
- [81] Weiming Xiang, Hoang-Dung Tran, and Taylor T. Johnson. 2018. Output Reachable Set Estimation and Verification for Multilayer Neural Networks. *TNNLS* 29 (2018), 5777–5783. Issue 11.
- [82] Xiaofei Xie, Lei Ma, Felix Juefei-Xu, Minhui Xue, Hongxu Chen, Yang Liu, Jianjun Zhao, Bo Li, Jianxiong Yin, and Simon See. 2019. DeepHunter: A Coverage-Guided Fuzz Testing Framework for Deep Neural Networks. In *ISSTA*. ACM, 146–157.
- [83] Sergey Zagoruyko and Nikos Komodakis. 2016. Wide Residual Networks. In *BMVC*.
- [84] Fuyuan Zhang, Sankalan Pal Chowdhury, and Maria Christakis. 2020. DeepSearch: A Simple and Effective Blackbox Attack for Deep Neural Networks. *CoRR* abs/1910.06296 (2020).
- [85] Jie M. Zhang, Mark Harman, Lei Ma, and Yang Liu. 2019. Machine Learning Testing: Survey, Landscapes and Horizons. *CoRR* abs/1906.10742 (2019).
- [86] Mengshi Zhang, Yuqun Zhang, Lingming Zhang, Cong Liu, and Sarfraz Khurshid. 2018. DeepRoad: GAN-Based Metamorphic Testing and Input Validation Framework for Autonomous Driving Systems. In *ASE*. ACM, 132–142.