

Back in Black: A Comparative Evaluation of Recent State-Of-The-Art Black-Box Attacks

Kaleel Mahmood

kaleel.mahmood@uconn.edu
Department of Computer Science and Engineering
University of Connecticut, USA

Ethan Rathbun

Department of Computer Science and Engineering
University of Connecticut, USA

Rigel Mahmood

Department of Computer Science and Engineering
University of Connecticut, USA

Marten van Dijk

CWI, The Netherlands

ABSTRACT

The field of adversarial machine learning has experienced a near exponential growth in the amount of papers being produced since 2018. This massive information output has yet to be properly processed and categorized. In this paper, we seek to help alleviate this problem by systematizing the recent advances in adversarial machine learning black-box attacks since 2019. Our survey summarizes and categorizes 20 recent black-box attacks. We also present a new analysis for understanding the attack success rate with respect to the adversarial model used in each paper. Overall, our paper surveys a wide body of literature to highlight recent attack developments and organizes them into four attack categories: score based attacks, decision based attacks, transfer attacks and non-traditional attacks. Further, we provide a new mathematical framework to show exactly how attack results can fairly be compared.

CCS CONCEPTS

• **Security and privacy**; → Computing methodologies; Machine Learning;

KEYWORDS

Adversarial machine learning; adversarial attack; black-box attack

1 INTRODUCTION

One of the first works to popularize Convolutional Neural Networks (CNN) [32] for image recognition was published in 1998. Since then, CNNs have been widely employed for tasks like image segmentation [24], object detection [49] and image classification [29]. Although CNNs are the de facto choice for machine learning tasks in the imaging domain, they have been shown to be vulnerable to adversarial examples [23]. In this paper, we discuss adversarial examples in the context of images. Specifically, an adversarial example is an input image which is visually correctly recognized by humans, but has a small noise added such that the classifier (i.e. a CNN) misclassifies the image with high confidence.

Attacks that create adversarial examples can be divided into two basic types, white-box and black-box attacks. White-box attacks require knowing the structure of the classifier as well as the associated trained model parameters [23]. In contrast to this, black-box attacks do not require directly knowing the model and trained parameters. Black-box attacks rely on alternative information like query access to the classifier [12], knowing the training dataset [46],

or transferring adversarial examples from one trained classifier to another [58].

In this paper, we survey recent advances in black-box adversarial machine learning attacks. We select this scope for two main reasons. First, we choose the black-box adversary because it represents a realistic threat model where the classifier under attack is not directly visible. It has been noted that a black-box attacker represents a more practical adversary [10] and one which corresponds to real world scenarios [46]. The second reason we focus on black-box attacks is due to the large body of recently published literature. As shown in Figure 1, many new black-box attack papers have been proposed in recent years. These attacks are not included in current surveys or systematization of knowledge papers. Hence, there is a need to categorize and survey these works, which is precisely the goal of this paper. To the best of our knowledge, the last major survey [4] on adversarial black-box attacks was done in 2020. A graphical overview of the coverage of some of the new attacks we provide (versus the old attacks previously covered) are shown in Figure 2. The complete list of important attack papers we survey are graphically shown in Figure 1 and also listed in Table 1.

While each new attack paper published contributes to the literature, they often do not compare with other state-of-art techniques, or adequately explain how they fit within the scope of the field. In this survey, we summarize 20 recent black-box attacks, categorize them into four basic groups and create a mathematical framework under which results from different papers can be compared.

1.1 Advances in Adversarial Machine Learning

In this subsection we briefly discuss the history and development of the field of adversarial machine learning. Such a perspective helps illuminate how the field went from a white-box attack like FGSM [23] in 2014 which required complete knowledge of the classifier and trained parameters, to a black-box attack in 2021 like SurFree [43] which can create an adversarial example with only query access to the classifier using 500 queries or less.

The inception point of adversarial machine learning can be traced back to several source papers. However, identifying the very first adversarial machine learning paper is a difficult task as the first paper in the field depends on how the term "adversarial machine learning" itself is defined. If one defines adversarial machine learning as exclusive to CNNs, then in [53] the vulnerability of CNNs to adversarial examples was first demonstrated in 2013. However, others [5] claim adversarial machine learning can be traced back

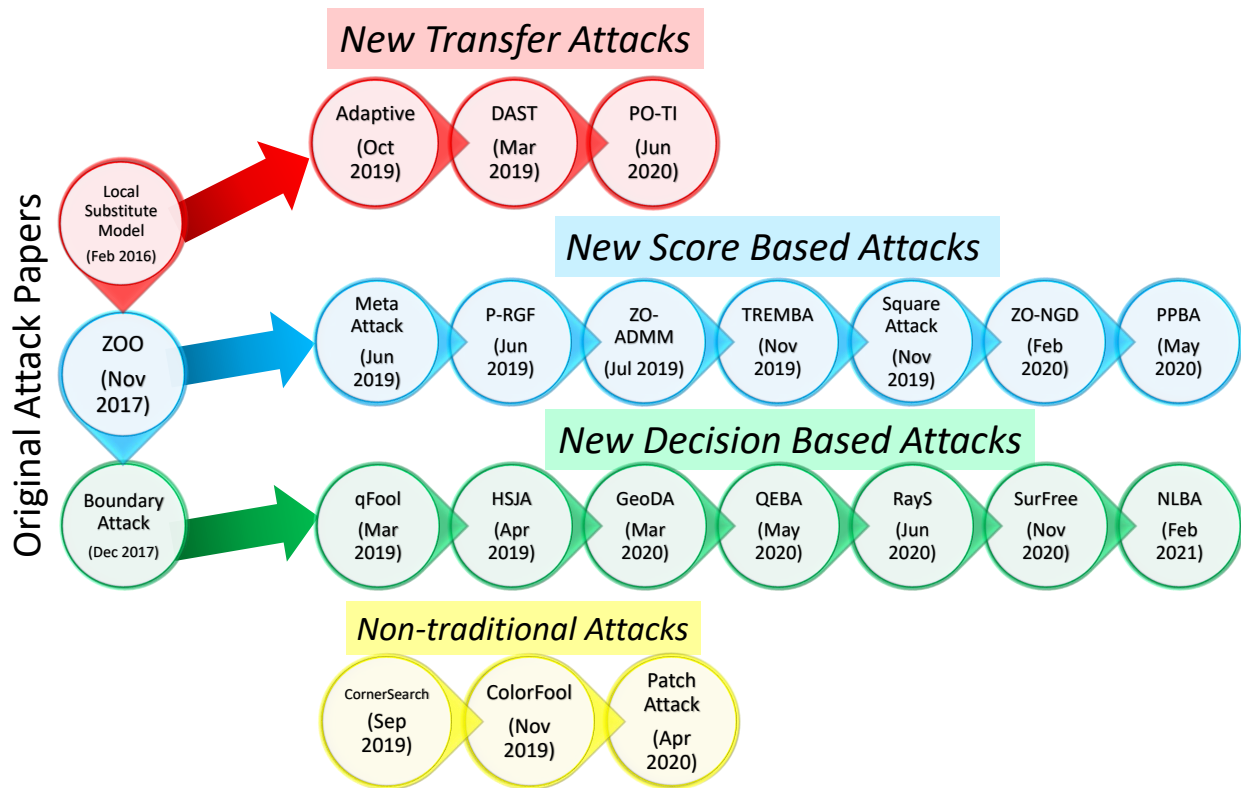


Figure 1: Timeline of recent black-box attack developments. The transfer based attacks are shown in red. The original transfer attack (Local Substitute Model) was proposed in [46]. The score based attacks are shown in blue. One of the first widely adopted score based attacks (ZOO) was proposed in [12]. The decision based attacks are shown in green. One of the first decision based attacks (Boundary Attack) was proposed in [6].

as early as 2004. In [5], the authors claim evading linear classifiers which constituted email spam detectors was one of the first examples of adversarial machine learning.

Regardless of the ambiguous starting point of adversarial examples, it remains a serious open problem which occurs across multiple machine learning domains including image recognition [23] and natural language processing [26]. Adversarial machine learning is also not just limited to neural networks. Adversarial examples have been shown to be problematic for decision trees, k-nearest neighbor classifiers and support vector machines [45].

The field of adversarial machine learning with respect to computer vision and imaging related tasks, first developed with respect to white-box adversaries. One of the first and most fundamental attacks proposed was the Fast Gradient Sign Method (FGSM) [23]. In the FGSM attack, the adversary uses the neural network model architecture F , loss function L , trained weights of the classifier w and performs a single forward and backward pass (backpropagation) on the network to obtain an adversarial example from a clean example x . Subsequent work included methods like the Projected Gradient Descent (PGD) [40] attack, which used multiple forward and backward passes to better fine tune the adversarial noise. Other attacks were developed to better determine the adversarial noise by

forming an optimization problem with respect to certain l_p norms, such as in the Carlini & Wagner [8] attack, or the Elastic Net attack [11]. Even more recent attacks [16] have focused on breaking adversarial defenses and overcoming false claims of security which are caused by a phenomena known as gradient masking [3].

All of the aforementioned attacks are considered white-box attacks. That is, the adversary requires knowledge of the network architecture F and trained weights w in order to conduct the attack. Creating a less capable adversary (i.e., one that did not know the trained model parameters) was a motivating factor in developing black-box attacks. In the next subsection, we discuss black-box attacks and the categorization system we develop in this paper.

1.2 Black-box Attack Categorization

We can divide black-box attacks according to the general adversarial model that is assumed for the attack. The four categories we use are transfer attacks, score based attacks, decision based attacks and non-traditional attacks. We next describe what defines the different categorizations and also mention the primary original attack paper in each category.

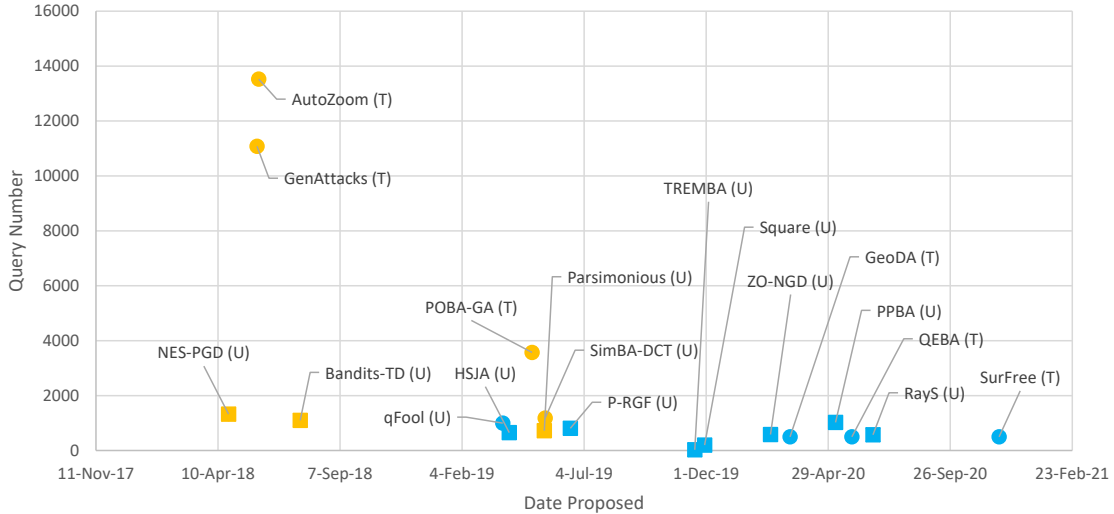


Figure 2: Graph of different black-box attacks with the respective date they were proposed (e-print made available). The query number refers to the number of queries used in the attack on an ImageNet classifier. The orange points are attacks covered in previous survey work [4]. The blue points are attacks covered in *this* work. We further denote whether the attack is targeted or untargeted by putting a U or T next to the text label in the graph. A square point represents an attack done with respect to the l_2 norm and a circular point represents attacks done with respect to the l_∞ norm.

Transfer Attacks: One of the first of black-box attacks was called the local substitute model attack [46]. In this attack, the adversary was allowed access to part of the original training data used to train the classifier, as well as query access to the classifier. The idea behind this attack was that the adversary would query the classifier to label the training data. After this was accomplished, the attacker would train their own independent classifier, which it is often referred to as the *synthetic model* [42]. Once the synthetic model was trained, the adversary could run any number of white-box attacks on the synthetic model to create adversarial examples. These examples were then submitted to the unseen classifier in the hopes the adversarial examples would *transfer* over. Here transferability is defined in the sense that adversarial examples that are misclassified by the synthetic model will also be misclassified by the unseen classifier.

Recent advances in transfer based attacks include not needing the original training data like in the DaST attack [58] and using methods that generate adversarial example with higher transferability (Adaptive [42] and PO-TI [37]).

Score Based Attacks: The zeroth order optimization based black-box attack (ZOO) [12] was one of the first accepted works to rely on a query based approach to creating adversarial examples. Unlike transfer attacks which require a synthetic model, score based attacks repeatedly query the unseen classifier to try and craft the appropriate adversarial noise. As the name implies, for score based attacks to work, they require the output from the classifier to be the score vector (either probabilities or in some cases the pre-softmax logits output).

Score based attacks represent an improvement over transfer attacks in the sense that no knowledge of the dataset is needed

since no synthetic model training is required. In very broad terms, the recent developments in score based attacks mainly focus on reducing the number of queries required to conduct the attack and/or reducing the magnitude of the noise required to generate a successful adversarial example. New score based attacks include qMeta [19], P-RGF [13], ZO-ADMM [57], TREMBA [27], Square attack [2], ZO-NGD [56] and PPBA [36].

Decision Based Attacks: We consider the type of attack that does not rely on a synthetic model and does not require the score vector output to be a decision based attack. Compared to either transfer based or score based attacks, decision based attacks represent an even more restricted adversarial model, as only the hard label output from the unseen classifier is required. The first prominent decision based attack paper was the Boundary Attack [6]. Since then, numerous decision based attacks have been proposed to improve upon the number of queries to successfully attack the unseen classifier, or reduce the noise required in the adversarial examples. The new decision attacks we cover in this paper include qFool [38], HSJA [10], GeoDA [47], QEBA [35], RayS [9], SurFree [43] and NonLinear-BA [33].

Non-traditional Attacks: The last category of attacks that we cover in this paper are called non-traditional black-box attacks. Here, we use this category to group the attacks that do not use standard black-box adversarial models. Transfer based attacks, score based attacks, and decision based attacks typically focus on designing the attack with respect to the l_2 and/or the l_∞ norm. Specifically, these attacks either directly or indirectly seek to satisfy the following condition: $\|x - x_{adv}\|_p \leq \epsilon$ where x is the original clean example, ϵ is the maximum allowed perturbation and $p = 2, \infty$.

Score based Attacks		
Attack Name	Date	Author
qMeta	6-Jun-19	Du et al. [19]
P-RGF	17-Jun-19	Cheng et al. [13]
ZO-ADMM	26-Jul-19	Zhao et al. [57]
TREMBA	17-Nov-19	Huang et al. [27]
Square	29-Nov-19	Andriushchenko et al. [2]
ZO-NGD	18-Feb-20	Zhao et al. [56]
PPBA	8-May-20	Liu et al. [36]
Decision based Attacks		
Attack Name	Date	Author
qFool	26-Mar-19	Liu et al. [38]
HSJA	3-Apr-19	Chen et al. [10]
GeoDA	13-Mar-20	Rahmati et al. [47]
QEBA	28-May-20	Li et al. [35]
RayS	23-Jun-20	Chen et al. [9]
SurFree	25-Nov-20	Maho et al. [43]
NonLinear-BA	25-Feb-21	Li et al. [33]
Transfer based Attacks		
Attack Name	Date	Author
Adaptive	3-Oct-19	Mahmood et al. [42]
DaST	28-Mar-20	Zhou et al. [58]
PO-TI	13-Jun-20	Li et al. [37]
Non-traditional Attacks		
Attack Name	Date	Author
CornerSearch	11-Sep-19	Croce et al. [15]
ColorFool	25-Nov-19	Shamsabadi et al. [52]
Patch	12-Apr-20	Yang et al. [54]

Table 1: Attacks covered in this survey, their corresponding attack categorization, publication date (when the first e-print was released) and author.

However, there are attacks that work outside of this traditional scheme.

CornerSearch [15] proposes a black-box attack based on finding an adversarial example with respect to the l_0 norm. Abandoning norm based constraints completely, Patch Attack [54] replaces a certain area of the image with an adversarial patch. Likewise, ColorFool [52] disregards norms and instead recolors the image to make it adversarial. While the non-traditional norm category is not strictly defined, it gives us a concise grouping that highlights the advances being made outside of the l_2 and l_∞ based black-box attacks.

1.3 Paper Organization and Major Contributions

In this paper we survey state-of-the-art black-box attacks that have recently been published. We provide three major contributions in this regard:

- (1) **In-Depth Survey:** We summarize and distill the knowledge from 20 recent significant black-box adversarial machine learning papers. For every paper, we include explanation of the mathematics necessary to conduct the attacks

and describe the corresponding adversarial model. We also provide an experimental section that brings together the results from all 20 papers, reported on three datasets (MNIST, CIFAR-10 and ImageNet).

- (2) **Attack Categorization:** We organize the attacks into four different categories based on the underlying adversarial model used in each attack. We present this organization so the reader can clearly see where advances are being made under each of the four adversarial threat models. Our breakdown concisely helps new researchers interpret the rapidly evolving field of black-box adversarial machine learning.
- (3) **Attack Analysis Framework:** We analyze how the attack success rate is computed based on different adversarial models and their corresponding constraints. Based on this analysis, we develop an intuitive way to define the threat model used to compute the attack success rate. Using this framework, it can clearly be seen when attack results reported in different papers can be compared, and when such evaluations are invalid.

The rest of our paper is organized as follows: in Section 2, we summarize score based attacks. In Section 3, we cover the papers that propose new decision based attacks. In Section 4, we discuss transfer attacks. The last type of attack, non-traditional attacks are described in Section 5. After covering all the new attacks, we turn our attention to analyzing the attack success rate in Section 6. Based on this analysis, we compile the experimental results for all the attacks in Section 7, and give the corresponding threat model developed from our new adversarial model framework. Finally, we offer concluding remarks in Section 8.

2 SCORE BASED ATTACKS

In this section we summarize recent advances in adversarial machine learning with respect to attacks that are score based or logit based. The adversarial model for these attacks allow the attacker to query the defense with input x and receive the corresponding probability outputs $p_1(x), \dots, p_k(x)$, where k is the number of classes. We also include logit based black-box attacks in this section. The logits are the pre-softmax outputs from the model, $l_1(x), \dots, l_k(x)$.

We cover 7 recently proposed score type attacks. These attacks include the square attack [2], the Zeroth-Order Natural Gradient Descent attack (ZO-NGD) [56], the Projection and Policy Driven Attack (PPBA) [36], the Zeroth-order Optimization Alternating Direction Method of Multipliers (ZO-ADMM) attack [57], the prior-guided random gradient-free (P-RGF) attack [13], the TRansferable Embedding based Black-box Attack (TREMBA) [27] and the qMeta attack [19].

2.1 Square Attack

The Square attack is a score based, black-box adversarial attack proposed in [2] that focuses primarily on being query efficient while maintaining a high attack success rate. The novelty of the attack comes in the usage of square shaped image perturbations which have a particularly strong impact on the predicted outputs of CNNs. This works in tandem with the implementation of the randomized search optimization protocol. The protocol is independent of model

gradients and greedily adds squares to the current image perturbation if they lead to an increase in the target model’s error. The attack solves the following optimization problem:

$$\min_{\hat{x} \in [0,1]^d} L(f(\hat{x}), y), \quad \text{s.t.} \quad \|\hat{x} - x\|_p \leq \epsilon \quad (1)$$

Where f is the classifier function, K is the number of classes, \hat{x} is the adversarial input, x is the clean input, y is the ground truth label, and ϵ is the maximum perturbation.

$$\begin{aligned} \text{Untargeted :} \quad & L(f(\hat{x}), y) = f_y(\hat{x}) - \max_{k \neq y} f_k(\hat{x}) \\ \text{Targeted :} \quad & L(f(\hat{x}), t) = -f_t(\hat{x}) + \log(\sum_{i=1}^K e^{f_i(\hat{x})}) \end{aligned} \quad (2)$$

The attack algorithm begins by first applying random noise to the clean image. Then an image perturbation, δ , is generated according to a perturbation generating algorithm defined by the attacker. If $L(f(\hat{x} + \delta), y) < L(f(\hat{x}), y)$ δ is applied to the current \hat{x} . This step is done iteratively until the targeted model outputs the desired label or until the max number of iterations are reached.

The distributions used for the iterative and initial image perturbations are chosen by the attacker. In [2] two different initial and iterative perturbation algorithms are proposed for the l_2 and l_∞ norm attacks.

For the l_∞ norm the perturbation is initialized by applying one pixel wide vertical stripes to the clean image. The color of each stripe is sampled uniformly from $\{-\epsilon, \epsilon\}^c$ where c is the number of color channels. The distribution used in the iterative step generates a square of a given size at a random location such that the magnitude of the perturbation in each color channel is chosen randomly from $\{-2\epsilon, 2\epsilon\}$. The resulting, clipped adversarial image will then differ from the clean image by either ϵ or $-\epsilon$ at each modified point.

The l_2 norm attack is initialized by generating a grid-like tiling of squares on the clean image. The perturbation is then rescaled to have l_2 norm ϵ and is clipped to $[0, 1]^d$. The iterative perturbation is motivated by the realization that classifiers are particularly susceptible to large, localized perturbations rather than smaller, more sparse ones. Thus the iterative attack places two squares of opposite sign either vertically or horizontally in line with each other, where each square has a large magnitude at its center that swiftly drops off but never reaches zero. After each iteration of the attack the current x_{adv} is clipped such that $\|x_{adv} - x\|_p < \epsilon$ and $x_{adv} \in [0, 1]^d$, where d is the dimensionality of the clean image.

The attack is tested on contemporary models like ResNet-50, Inception v3, and VGG-16-BN which are trained on ImageNet. It achieves a lower attack failure rate while requiring significantly less queries to complete than attacks like Bandits, Parsimonious, DFO-MCA, and SignHunter. Similarly the square attack is compared to the white box Projected Gradient Descent (PGD) attacks on the MNIST and CIFAR-10 datasets where it performs similarly to PGD in terms of attack success rate despite operating within a more difficult threat model.

2.2 Zeroth-Order Natural Gradient Descent Attack

The Zeroth-Order Natural Gradient Descent (ZO-NGD) attack is a score-based, black box attack proposed in [56] as a query efficient

attack utilizing a novel attack optimization technique. In particular the attack approximates a Fisher information matrix over the distribution of inputs and subsequent outputs of the classifier. The attack solves the following optimization problem:

$$\min_{\delta} f(x + \delta, t), \quad \|\delta\|_\infty \leq \epsilon \quad (3)$$

$$f(x + \delta, t) = \max\{\log p(t|x + \delta) - \max_{i \neq t} \{\log p(i|x + \delta)\}, -k\} \quad (4)$$

Where x is the clean image, δ is an image perturbation, ϵ is the maximum allowed image perturbation, t is the clean image’s ground truth label, $p(i|x)$ is the classifier’s predicted score for class i given input x , and f is the attack’s loss. The attack is an iterative algorithm that initializes the image perturbation, δ , as a matrix of all zeros. At each step the algorithm first approximates the gradient of the loss function, f , according to the following equation:

$$\widehat{\nabla} f(\delta) = \frac{1}{R} \sum_{j=1}^R \frac{f(\delta + \mu u_j, t) - f(\delta, t)}{\mu} u_j \quad (5)$$

Where each $u_j \sim N(0, I_d)$ is a random perturbation chosen i.i.d. from the unit sphere, μ is a smoothing parameter, and R is a hyper parameter for the number of queries used in the approximation. Next, the attack approximates the gradient of the log-likelihood function. This is necessary for calculating the Fisher information matrix and subsequently the perturbation update.

$$\widehat{\nabla} \log p(t|x + \delta) = \frac{1}{R\mu} \sum_{j=1}^R (\log p(t|x + \delta + \mu u_j) - \log p(t|x + \delta)) u_j \quad (6)$$

Here the notation is consistent with the notation seen in Equation 5. This can be calculated using the same queries that were used in Equation 5. The Fisher information matrix is approximated and δ is updated according to the following equations:

$$\widehat{F} = \widehat{\nabla} \log p(t|x + \delta) \widehat{\nabla} \log p(t|x + \delta)^T + \gamma I \quad (7)$$

$$\delta_{k+1} = \prod (\delta_k - \lambda \widehat{F}^{-1} \widehat{\nabla} f(\delta_k)) \quad (8)$$

Where γ is a constant and λ is the attack learning rate. \prod is the projection function which projects its input onto the set $S = \{\delta | (x + \delta) \in [0, 1]^d, \|\delta\|_\infty \leq \epsilon\}$. It is also worth recognizing that δ is represented as a matrix since images, like x , are also represented as matrices. This makes the addition seen in Equation 8 valid. The iterative process can be continued for a predetermined number of iterations or until the perturbation yields a satisfactory result. The Fisher information matrix is a powerful tool, however its size can prove it impractical for use on datasets with larger inputs, thus an approximation of δ_{k+1} may be necessary.

The attack is tested on the MNIST, CIFAR-10, and ImageNet datasets where it achieves a similar attack success rate to the ZOO, Bandits, and NES-PGD attacks while requiring less queries to be successful. The attack is then also shown to have an extremely high attack success rate within 1200 queries on all three aforementioned datasets.

2.3 Projection and Probability Driven Attack

The Projection and Probability-driven Black-box Attack (PPBA) proposed in [36] is a score based, black box attack that achieves high attack success rates while being query efficient. It achieves this by shrinking the solution space of possible adversarial inputs to those which contain low-frequency perturbations. This is motivated by an observation that contemporary neural networks are particularly susceptible to low frequency perturbations. The attack solves the following optimization problem:

$$\min_{\delta} L(\delta) = [f(x + \delta)_t - \max_{j \neq t} f(x + \delta)_j]^+ \quad (9)$$

Where $f(x)_j$ is the model's predicted probability that the input is of class j , x is the clean image, t is the ground truth label, δ is the adversarial perturbation, and $[\cdot]^+$ is shorthand for $\max(\cdot, 0)$. The attack utilizes a sensing matrix, A , which is composed of a Discrete Cosine Transform matrix, Ψ , and a Measurement matrix, Φ , along with the corresponding measurement vector, z . The exact design of the measurement matrix varies according to practice [1] [48]. The relationship between all these variables is as follows: $A = \Psi\Phi$, $z = A\delta$, $\delta \approx A^T z$.

One point to note is that Φ should be an orthonormal matrix which allows $\delta \approx A^T z$ to be true. Once A is calculated the attack utilizes a query efficient version of the random walk algorithm. In particular, the attack stores a Confusion matrix C_j for each dimension j of Δz , which is the change in z at each iteration. C_j can be seen below:

	$-\rho$	0	ρ
# effective steps	$e_{-\rho}$	e_0	e_{ρ}
# ineffective steps	$i_{-\rho}$	i_0	i_{ρ}

Where ρ is a predefined step size, e_v is the number of times the loss function descended when $\Delta z_j = v$, and i_v is the number of times the loss function increased or remained the same when $\Delta z_j = v$ for $v \in \{-\rho, 0, \rho\}$. The algorithm then uses C to determine its sampling probability for Δz_j as seen below:

$$P(a|\Delta z_j = v) = \frac{e_v}{e_v + i_v}, v \in \{-\rho, 0, \rho\} \quad (10)$$

$$P(\Delta z_j = v) = \frac{P(a|\Delta z_j = v)}{\sum_u P(a|\Delta z_j = u)}, u, v \in \{-\rho, 0, \rho\} \quad (11)$$

Where a is a probabilistic variable that is true when the step is determined to be effective. The attack algorithm begins by first calculating A and then initializing all values of C to be 1. The iterative part of the algorithm then begins, at each step the algorithm generates a new Δz according to the probability distribution described in Equation 11. If $L(A^T(z + \Delta z)) < L(A^T z)$ then z is updated as $z = \text{clip}(z + \Delta z)$. Here the clip function forces $x + z$ to remain within the clean image's input space, $[0, 1]^d$. If at any point the perturbation generated causes the model to output an incorrect class label the attack terminates and returns the penultimate perturbation.

PPBA is tested on the ImageNet dataset with the classifiers ResNet50, Inception v3 and VGG-16. PPBA achieves high attack success rates while maintaining a low query count. It is also tested on Google Cloud Vision API where it achieves a high attack success rate in this more realistic setting.

2.4 Alternating Direction Method of Multipliers Based Black-Box Attacks

A new black-box attack framework is proposed in [57] based on the distributed convex optimization technique, the Alternating Direction Method of Multipliers (ADMM). The advantage of using the ADMM technique is that it can be directly combined with the zeroth-order optimization attack (ZOO-ADMM) or Bayesian optimization (BO-ADMM) to create a query-efficient, gradient free black-box attack. The attack can be run with score based or decision based output from the defense.

The main concept presented in [57] is the conversion of the black-box attack optimization problem from a traditional constrained optimization problem, into an unconstrained objective function that can be iteratively solved using ADMM. The original formulation of the black-box attack optimization problem can be written as:

$$\begin{aligned} & \underset{\delta}{\text{minimize}} && f(x_0 + \delta, t) + \gamma D(\delta) \\ & \text{subject to} && (x_0 + \delta) \in [0, 1]^d, \|\delta\|_{\infty} \leq \epsilon \end{aligned} \quad (12)$$

where $f(\cdot)$ is the loss function of the classifier, δ is the perturbation added to the original input x_0 , t is the target class that the adversarial example $(x_0 + \delta)$ should be misclassified as and D is a distortion function to limit the difference between the adversarial example and x_0 . In Equation 12, γ controls the weight given to the distortion function and ϵ specifies the maximum tolerated perturbation.

Instead of directly solving Equation 12, the constraints can be moved into the objective function and an auxiliary variable z can be introduced in order to write the optimization problem in an ADMM style form:

$$\begin{aligned} & \underset{\delta, z}{\text{minimize}} && f(x_0 + \delta, t) + \gamma D(\delta) + \mathcal{I}(z) \\ & \text{subject to} && z = \delta \end{aligned} \quad (13)$$

where $\mathcal{I}(z)$ is 0 if $(x_0 + z) \in [0, 1]^d, \|z\|_{\infty} \leq \epsilon$ and ∞ otherwise. The augmented Lagrangian of Equation 13 is written as:

$$\begin{aligned} \mathcal{L}(z, \delta, u) &= \gamma D(z) + \mathcal{I}(z) + f(x_0 + \delta, t) \\ &+ \frac{\rho}{2} \|z - \delta + \frac{1}{\rho} u\| - \frac{1}{2\rho} \|u\|_2^2 \end{aligned} \quad (14)$$

where u is the Lagrangian multiplier and ρ is a penalty parameter. Equation 14 can be iteratively solved using ADMM in the k^{th} step through the following update equations:

$$z^{k+1} = \arg \min_z \mathcal{L}(z, \delta^k, u^k) \quad (15)$$

$$\delta^{k+1} = \arg \min_{\delta} \mathcal{L}(z^{k+1}, \delta, u^k) \quad (16)$$

$$u^{k+1} = u^k + \rho(z^{k+1} - \delta^{k+1}) \quad (17)$$

While Equation 15 has a closed form solution, minimizing Equation 16 requires a gradient descent technique like stochastic gradient decent, as well as access to the gradient of $f(x_0 + \delta, t)$. In the black-box setting this gradient is not available to the adversary and hence must be estimated using a special approach. If the gradient is estimated using the random gradient estimation technique, then the attack is referred to as ZOO-ADMM. Similarly, if the gradient

is estimated using bayesian optimization, the attack is denoted as BO-ADMM.

The new attack framework is experimentally verified on the CIFAR-10 and MNIST datasets. The results of the paper [57] show ZOO-ADMM outperforms both BO-ADMM and the original boundary attack presented in [7]. This performance improvement comes in the form of smaller distortions for the l_1 , l_2 and l_∞ threat models and in terms of less queries used for the ZOO-ADMM attack.

2.5 Improving Black-box Adversarial Attacks with Transfer-based Prior

Initial adversarial machine learning black-box attacks were developed based on one of two basic principles. In query based black-box attacks [7], the gradient is directly estimated through querying. In transfer based attacks, the gradient is computed based on a trained model’s gradient that is available to the attacker [46]. In [13] they propose combining the query and transfer based attacks to create a more query efficient attack which they call the prior-guided random gradient-free method (P-RGF).

The P-RGF attack is developed around accurately and efficiently estimating the gradient of the target model f . The original random gradient-free method [44] estimates the gradient as follows:

$$\hat{g} = \frac{1}{q} \sum_{i=1}^q \frac{f(x + \sigma u_i, y) - f(x, y)}{\sigma} \cdot u_i \quad (18)$$

where q is the number of queries used in the estimate, σ is a parameter to control the sampling variance, x is the input with corresponding label y and $\{u_i\}_{i=1}^q$ are random vectors sampled from distribution \mathcal{P} . It is important to note that by selecting $\{u_i\}_{i=1}^q$ carefully (according to priors) we can create a better estimate of g . In P-RGF this choice of $\{u_i\}_{i=1}^q$ is done by biasing the sampling using a transfer gradient v . The transfer gradient v comes from a surrogate model that has been independently trained on the same data as the model whose gradient is currently being estimated. In the attack it is assumed that we have white-box access to the surrogate model such that v is known.

The overall derivation of the rest of the attack from [13] goes as follows: first we discuss the appropriate loss function $L(\cdot)$ for \hat{g} . We then discuss how to pick $\{u_i\}_{i=1}^q$ such that $L(\cdot)$ is minimized. To determine how closely \hat{g} (the estimated gradient) follows g (the true model gradient) the following loss function is used [13]:

$$\min_{b \geq 0} \mathbb{E} \|\nabla_x f(x) - b\hat{g}\|_2^2 \quad (19)$$

where b is a scaling factor included to compensate for the change in magnitude caused by \hat{g} and the expectation is taken over the randomness of the estimation algorithm. For notational convenience we write $\nabla_x f(x)$ as $\nabla f(x)$ in the remainder of this subsection. It can be proven that if x is differentiable at f then the loss function given in Equation 19 can be expressed as:

$$\lim_{\sigma \rightarrow 0} L(\hat{g}) = \|\nabla f(x)\|_2^2 - \frac{H(\mathbf{C}, x)^2}{(1 - \frac{1}{q})H(\mathbf{C}^2, x) + \frac{1}{q}H(\mathbf{C}, x)^2} \quad (20)$$

where $H(\mathbf{C}, x) = \nabla f(x)^T \mathbf{C} \nabla f(x)$ and $\mathbf{C} = \mathbb{E}[u_i u_i^T]$. Through careful choice of \mathbf{C} , $L(\hat{g})$ can be minimized to accurately estimate the

gradient, thereby making the attack query efficient. \mathbf{C} can be decomposed in terms of the transfer gradient v as:

$$\mathbf{C} = \lambda v v^T + \frac{1 - \lambda}{D - 1} (\mathbf{I} - v v^T) \quad (21)$$

where $\{\lambda_i\}_{i=1}^D$ and $\{v_i\}_{i=1}^D$ are the eigenvalues and orthonormal eigenvectors of \mathbf{C} . To exploit the gradient information of the transfer model, u_i is then randomly generated in terms of v to satisfy Equation 21:

$$u_i = \sqrt{\lambda} \cdot v + \sqrt{1 - \lambda} \cdot \overline{(\mathbf{I} - v v^T)} \xi_i \quad (22)$$

where λ controls the magnitude of the transfer gradient v and ξ_i is a random variable sampled uniformly from the unit hypersphere.

The overall P-RGF method for estimating the gradient g is as follows: First α , the cosine similarity between the transfer gradient v and the model gradient g is estimated through a specialized query based algorithm [13]. Next λ is computed as a function of α , q and the input dimension size D . Note we omitted the λ equation and explanation in our summary for brevity. After computing λ , the estimate of the gradient \hat{g} is iteratively done Q times in a two step process. In the first step of the q^{th} iteration, u_q is generated using Equation 22. In the second step \hat{g} is calculated as: $\hat{g} = \hat{g} + \frac{f(x + \sigma u_q, y) - f(x, y)}{\sigma} \cdot u_q$, where q denotes the q^{th} iteration. After Q iterations have been complete, the final gradient estimate is given as $\hat{g} \leftarrow \frac{1}{Q} \hat{g}$.

The P-RGF attack is tested on ImageNet. The surrogate model to get the transfer gradient in the attack is set as ResNet-152. Attacks are done on different ImageNet CNNs which include Inception v3, VGG-16 and ResNet50. The P-RGF attack outperforms other completing techniques in terms of having a higher attack success rate and lower number of queries for most networks.

2.6 Black-Box Adversarial Attack with Transferable Model-based Embedding

The Transferable Embedding based Black-box Attack (TREMBA) [27] is an attack that uniquely combines transfer and query based black-box attacks. In conventionally query based black-box attacks, the adversarial image is modified by iteratively fine tuning the noise that is directly added to the pixels of the original image. In TREMBA, instead of directly altering the noise, the embedding space of a pre-trained model is modified. Once the embedding space is modified, this is translated into noise for the adversarial image. The advantage of this approach is that by using the pre-trained model’s embedding as a search space, the amount of queries needed for the attack can be reduced and the attack efficiency can be increased.

The attack generates the perturbation δ for input x using a generator network \mathcal{G} . The generator network is comprised of two components, an encoder \mathcal{E} and a decoder \mathcal{D} . The encoder maps x to z , a latent space i.e., $z = \mathcal{E}(x)$. The decoder \mathcal{D} takes z as input. The outputs of the decoder \mathcal{D} is used to compute the perturbation δ which is defined as $\delta = \epsilon \tanh(\mathcal{D}(z))$. The tanh function is used to normalize the output of the decoder $\mathcal{D}(z)$ between -1 and 1 such that the final adversarial perturbation δ is bounded i.e. $\|\delta\|_\infty \leq \epsilon$.

To begin the untargeted version of the attack, the generator network \mathcal{G} is first trained. For an individual sample (x_i, y_i) , we denote the probability score associated with the correct class label

during training as:

$$P_{true}(x_i, y_i) = F_s(\epsilon \cdot \tanh(\mathcal{G}(x_i)) + x_i)_{y_i} \quad (23)$$

where ϵ is the maximum allowed perturbation, $\mathcal{G}(\cdot)$ is the output from the generator and $F_s(\cdot)_i$ is the i^{th} component of the output vector of the source model F_s . In this attack formulation the adversary is assumed to have white-box access to a pre-trained source model F_s which is different from the target model under attack. The incorrect class label with the maximum probability during training is:

$$P_{false}(x_i, y_i) = \max_{j \neq y_i} F_s(\epsilon \cdot \tanh(\mathcal{G}(x_i)) + x_i)_j \quad (24)$$

Using Equation 23 and Equation 24 the loss function for training the generator for an untargeted attack is given as:

$$\mathcal{L}_{untarget}(x_i, y_i) = \max(P_{true}(x_i, y_i) - P_{false}(x_i, y_i), -\kappa) \quad (25)$$

where (x_i, y_i) are individual training samples in the training dataset and κ is a transferability parameter (higher κ makes the adversarial examples more transferable to other models [8]).

Once \mathcal{G} is trained the perturbation δ can be calculated as a function of the embedding space z . The embedding space z is iteratively computed:

$$z_t = z_{t-1} - \frac{\eta}{b} \sum_{i=1}^b \mathcal{L}_{untarget} \nabla_{z_{t-1}} \log(\mathcal{N}(v_i | z_{t-1}, \sigma^2)) \quad (26)$$

where t is the iteration number, η is the learning rate, b is the sample size, v_i is a sample from the gaussian distribution $\mathcal{N}(z_{t-1}, \sigma^2)$ and $\nabla_{z_{t-1}}$ is the gradient of z_t estimated using the Natural Evolution Strategy (NES) [28].

Experimentally TREMBA is tested on both the MNIST and ImageNet datasets. The attack is also tested on the Google Cloud Vision API. In general, TREMBA achieves a higher attack success rate and uses less queries for MNIST and ImageNet, as compared to other attack methods. These other attack methods compared in this work include P-RGF, NES and AutoZOOM.

2.7 Query-Efficient Meta Attack

In the query-efficient meta attack [19], high query-efficiency is achieved through the use of meta-learning to observe previous attack patterns. This prior information is then leveraged to infer new attack patterns through a reduced number of queries. First, a meta attacker is trained to extract information from the gradients of various models, given specific input, with the goal being to infer the gradient of a new target model using few queries. That is, an image \mathbf{x} is input to models $\mathcal{M}_1, \dots, \mathcal{M}_n$ and a max-margin logit classification loss is used to calculate losses l_1, \dots, l_n as follows:

$$l_i(\mathbf{x}) = \max [\log[\mathcal{M}_i(\mathbf{x})]_t - \max_{j \neq t} \log[\mathcal{M}_i(\mathbf{x})]_j, 0] \quad (27)$$

where t is the true label, j is the index of other classes, $[\mathcal{M}_i(\mathbf{x})]_t$ is the probability score produced by the model \mathcal{M}_i , and $[\mathcal{M}_i(\mathbf{x})]_j$ refers to the probability scores of the subsequent classes.

After one step back-propagation is performed, n training groups for the universal meta attacker are assembled, consisting of input images $\mathbb{X} = \{\mathbf{x}\}$ and gradients $\mathbb{G}_i = \{\mathbf{g}_i\}$, $i = 1, \dots, n$ where $\mathbf{g}_i = \nabla_{\mathbf{x}} l_i(\mathbf{x})$. In each training iteration, K samples are drawn from a task $\mathcal{T}_i = (\mathbb{X}, \mathbb{G}_i)$. For meta attacker model \mathcal{A} with parameters θ ,

the updated parameters θ' are computed as: $\theta' := \theta - \alpha \nabla_{\theta} \mathcal{L}_i(\mathcal{A}_{\theta})$, where \mathcal{L}_i is the loss corresponding to task \mathcal{T}_i .

The meta attack parameters are optimized by incorporating θ'_i across all tasks $\{\mathcal{T}_i\}_{i=1, \dots, n}$ according to:

$$\theta := \theta + \epsilon \frac{1}{n} \sum_{i=1}^n (\theta'_i - \theta) \quad (28)$$

The training loss of this meta attacker \mathcal{A}_{θ} employs mean-squared error, as given below:

$$\mathcal{L}_i(\mathcal{A}_{\theta}) = \|\mathcal{A}_{\theta}(\mathbb{X}_s) - \mathbb{G}_i^s\|_2^2 \quad (29)$$

where the set $(\mathbb{X}_s, \mathbb{G}_i^s)$ refers to the K samples selected for training from $(\mathbb{X}, \mathbb{G}_i)$ for θ to θ'_i .

The high-level objective of such a meta attacker model \mathcal{A} is to produce a helpful gradient map for attacking that is adaptable to the gradient distribution of the target model. To accomplish this efficiently, a subsection q of the total p gradient map coordinates are used to fine-tune \mathcal{A} every m iterations [19], where $q \ll p$. In this manner, \mathcal{A} is trained to be able to produce the gradient distribution of various input images and learns to predict the gradient from only a few samples through this selective fine-tuning. It is of importance to note that query efficiency is further reinforced by performing the typically query-intensive zeroth-order gradient estimation only every m iterations.

Empirical results on MNIST, CIFAR-10, and tiny-ImageNet attain comparable attack success rates to other untargeted black-box attacks. However, the attack significantly outperforms prior attacks in terms of the number of queries required in the targeted setting [19].

3 DECISION BASED ATTACKS

In this section, we discuss recent developments in adversarial machine learning with respect to attacks that are decision based. The adversarial model for these attacks allows the attacker to query the defense with input x and receive the defense's final predicted output. In contrast to score based attacks, the attacker does not receive any probabilistic or logit outputs from the defense.

We cover 7 recently proposed decision based attacks. These attacks include the Geometric decision-based attack [47], Hop Skip Jump Attack [10], RayS Attack [9], Nonlinear Black-Box Attack [33], Query-Efficient Boundary-Based Black-box Attack [35], SurFree attack [43], and the qFool attack [38].

3.1 Geometric Decision-based Attacks

Geometric decision-based attacks (GeoDA) are a subset of decision based black box attacks proposed in [47] that can achieve high attack success rates while requiring a small number of queries. The attack exploits a low mean curvature in the decision boundary of most contemporary classifiers within the proximity of a data point. In particular the attack uses a hyperplane to approximate the decision boundary in the vicinity of a data point to effectively find the local normal vector of the decision boundary. The normal vector can then be used to modify the clean image in such a way that the model outputs an incorrect class label. Thus the attack solves the following optimization problem:

$$\begin{aligned} \min_{\nu} \quad & \|v\|_p \\ \text{s.t.} \quad & w^T(x+v) - w^T x_B = 0 \end{aligned} \quad (30)$$

Where w is a normal vector to the decision boundary, and x_B is point on the decision boundary and close to the clean image, x . x_B can be found by adding random noise, r , to x until the classifier’s predicted label changes, then performing a binary search in the direction of r to get x_B as close to the decision boundary as possible:

$$\begin{aligned} x_B &= x + \min_r \|r\|_2 \\ \text{s.t. } \hat{k}(x_B) &\neq \hat{k}(x) \end{aligned} \quad (31)$$

Where $\hat{k}(\cdot)$ returns the top-1 label of the target classifier. The normal vector to the decision boundary is found in the following way: N image perturbations, η_i , are randomly drawn from a multi-variate normal distribution $\eta_i \sim \mathcal{N}(0, \Sigma)$ [39]. The model is then queried on the top-1 label of each $x_B + \eta_i$ where x_b is a boundary point close to the clean image, x . Each η_i is then classified as follows:

$$S_{adv} = \{\eta_i \mid \hat{k}(x_b + \eta_i) \neq \hat{k}(x)\} \quad (32)$$

$$S_{clean} = \{\eta_i \mid \hat{k}(x_b + \eta_i) = \hat{k}(x)\} \quad (33)$$

From here the normal vector to the decision boundary can then be estimated as:

$$\hat{w}_N = \frac{\bar{\mu}_N}{\|\bar{\mu}_N\|_2} \quad (34)$$

$$\text{where } \bar{\mu}_N = \frac{1}{N} \sum_{i=1}^N \rho_i \eta_i \quad (35)$$

$$\text{and } \rho_i = \begin{cases} 1 & \eta_i \in S_{adv} \\ -1 & \eta_i \in S_{clean} \end{cases}$$

Finally the image can be modified using the following update:

$$x_{adv} = x + \hat{r} \hat{w}_N \quad (36)$$

$$\text{where } \hat{r} = \min\{r > 0 \mid \hat{k}(x + rv) \neq \hat{k}(x)\} \quad (37)$$

$$\text{and } v = \frac{1}{\|\hat{w}_N\|_a} \odot \text{sign}(\hat{w})$$

Here \odot refers to the point-wise product and $a = \frac{p}{p-1}$. This process is done iteratively, at each iteration the previous iteration’s x_{adv} is used to calculate \hat{w} which is then added to the original x to find the current iteration’s x_{adv} as seen above.

The attack is experimentally tested on the ImageNet dataset. The experiments show GeoDA outperforms the Hop Skip Jump Attack, Boundary Attack, and qFool by producing smaller image perturbations and requiring less iterations, and thus less queries, to complete.

3.2 Hop Skip Jump Attack

The Hop Skip Jump Attack (HSJA) is a decision based, black-box attack proposed in [10] that achieves both a high attack success rate and a low number of queries. The attack is an improvement on the previously developed Boundary Attack [6] in that it implements gradient estimation techniques at the edge of a model’s decision boundary in order to more efficiently create adversarial inputs to the classifier. Similarly to many other adversarial attacks, HSJA attempts to change the predicted class label of a given input, x , while minimizing the perturbation applied to the input. Thus the following optimization problem is proposed:

$$\min_{x'} d(x', x^*) \quad \text{s.t. } \phi_{x^*}(x') = 1 \quad (38)$$

$$\phi_{x^*}(x') = \text{sign}(S_{x^*}(x')) \quad (39)$$

$$S_{x^*}(x') = \begin{cases} \max_{c \neq c^*} F_c(x') - F_{c^*}(x') & \text{(Untargeted)} \\ F_{c^\dagger}(x') - \max_{c \neq c^\dagger} F_c(x') & \text{(Targeted)} \end{cases} \quad (40)$$

Here F_c is the predicted probability of class c , x' is the adversarial input, x^* is the clean input, and d is a distance metric. This unique optimization formulation allows HSJA to approximate the gradient of Equation 40 and thus more accurately and efficiently solve the optimization problem.

The attack algorithm starts by adding random noise, δ , to the clean image, x^* , until the model’s predicted class label changes to the desired label. Once a desired random perturbation is found the iterative process is initiated and $x^* + \delta$ is stored in x_0 which becomes an iterative parameter written as x_t for step number t . From here a binary search is performed to find the decision boundary between x^* and x_t . At the decision boundary the following operation is used to approximate the gradient of the decision boundary:

$$\widehat{\Delta S}(x_t, \delta_t) = \frac{1}{1-B} \sum_{b=1}^B (\phi_{x^*}(x_t + \delta_t u_b) - \overline{\phi_{x^*}}) u_b \quad (41)$$

$$\overline{\phi_{x^*}} = \frac{1}{B} \sum_{b=1}^B \phi_{x^*}(x_t + \delta_t u_b) \quad (42)$$

Where $\delta_t = d_t^{-1} \|x_{t-1} - x^*\|_p$ and $d_0 = \|x_0 - x^*\|$ is a small, positive parameter. Each u_b is randomly drawn i.i.d. from the uniform distribution over the d -dimensional sphere. The additional term, $\overline{\phi_{x^*}}$, is used to attempt to mitigate the bias induced into the estimation by δ . Once the gradient of the decision boundary is found an update direction is found using the following formulation:

$$v_t(x_t, \delta_t) = \begin{cases} \widehat{\Delta S}(x_t, \delta_t) / \|\widehat{\Delta S}(x_t, \delta_t)\|_2 & \text{if } p = 2 \\ \text{sign}(\widehat{\Delta S}(x_t, \delta_t)) & \text{if } p = \infty \end{cases} \quad (43)$$

Once this update direction is found a step size must be determined. The step size is initialized as $\xi_t = \|x_t - x^*\|_p / \sqrt{t}$ and is halved until $\phi_{x^*}(x_t + \xi_t v_t) \neq 0$. Then x_t is updated by $x_t = x_t + \xi_t v_t$ and d_t is updated by $d_t = \|x_t - x^*\|_p$. This process is continued for a predetermined T iterations.

In [10] HSJA is tested on the MNIST, CIFAR-10, CIFAR-100, and ImageNet datasets. HSJA outperforms the Boundary Attack and Opt Attack in terms of median perturbation magnitude and attack success rate. HSJA is also tested against multiple defenses on the MNIST dataset, where it performs better than Boundary Attack and Opt Attack when all attacks are given an equal number of queries.

3.3 RayS Attack

The RayS attack is a query efficient, decision based, black-box attack proposed in [9] as an alternative to zeroth-order gradient attacks. The attack employs an efficient search algorithm to find the nearest decision boundary that requires less queries then other contemporary decision based attacks while maintaining a high attack success rate. Specifically, the attack formulation turns the continuous problem of finding the closest decision boundary into a discrete optimization problem:

$$\min_{d \in \{-1, 1\}^n} g(d) = \arg \min_r \mathbb{1}\{f(x + \frac{rd}{\|d\|_2}) \neq y\} \quad (44)$$

Where x is the clean sample which is assumed to be a vector without loss of generality, y is the ground truth label of the clean sample,

f is the classifier’s prediction function, d is a direction vector determining the direction of the perturbation in the input space, r is a scalar projected onto d determining the magnitude of the perturbation, and n is the dimensionality of the input. This converts the continuous problem of finding the direction to the closest decision boundary into a discrete optimization problem over $d \in \{-1, 1\}^n$ which contains 2^n possible options.

The attack algorithm finds a direction, d , and a radius, r , in the input space as its final output for the attack. They can then be converted into a perturbation by projecting r onto d . The attack begins by choosing some initial direction vector, d , and setting $r = \infty$. The iterative process comes in multiple stages, s , where at each stage d is cut into 2^s equal and uniformly placed blocks. The algorithm then iterates through each of these blocks, swapping the sign of each value in the current block at a given iteration and storing the modified d into d_{temp} . If $f(x + r \cdot d_{temp}) = y$ the algorithm skips searching d_{temp} as it requires a larger perturbation than d to change the classifier’s predicted label. If $f(x + r \cdot d_{temp}) \neq y$ the algorithm performs a binary search in the direction of d_{temp} to find the smallest r such that $f(x + r \cdot d_{temp}) \neq y$ remains true. Finally d is updated to d_{temp} and r is updated to be the smallest radius found in the binary search.

The RayS attack is experimentally tested in [9] on the MNIST, CIFAR-10, and ImageNet datasets. It outperforms other black-box attacks like HSJA and SignOPT in terms of both average number of queries and attack success rate on the MNIST and CIFAR-10 datasets. On the ImageNet dataset, HSJA achieves a lower number of average queries than RayS, but attains a significantly lower attack success rate. The RayS attack is also compared to white box attacks like Projected Gradient Descent (PGD) where it outperforms the attack on the MNIST and CIFAR-10 datasets, in terms of the attack success rate.

3.4 Nonlinear Projection Based Gradient Estimation for Query Efficient Blackbox Attacks

The Nonlinear Black-box Attack (NonLinear-BA) is a query efficient, nonlinear gradient projection-based boundary blackbox attack [33]. This attack innovatively overcomes the gradient inaccessibility of blackbox attacks by utilizing vector projection for gradient estimation. AE, VAE, and GAN are used to perform efficient projection-based gradient estimation. [33] shows that NonLinear-BA can outperform the corresponding linear projections of HSJA and QEBA, as NonLinear-BA provides a higher lower bound of cosine similarity between the estimated and true gradients of the target model.

There are three components of NonLinear-BA: the first is gradient estimation at the target model’s decision boundary. While high-dimensional gradient estimation is computationally expensive, requiring numerous queries [33], projecting the gradient to lower dimensional supports greatly improves the estimation efficiency of NonLinear-BA. This desired low dimensionality is achieved through the latent space representations of generative models, e.g., AE, VAE, and GAN.

The gradient projection function \mathbf{f} is defined as $\mathbf{f} : \mathbb{R}^n \rightarrow \mathbb{R}^m$, which maps the lower-dimensional representative space \mathbb{R}^n to the

original, high-dimensional space \mathbb{R}^m , where $n \leq m$. The sample unit latent vectors v_b ’s in \mathbb{R}^n are randomly sampled to generate the perturbation vectors $u_b = \mathbf{f}(v_b) \in \mathbb{R}^m$.

Thus, the gradient estimator is as follows:

$$\widetilde{\nabla}S(x_{adv}^{(t)}) = \frac{1}{B} \sum_{b=1}^B \text{sgn}(S(x_{adv}^{(t)} + \delta \mathbf{f}(v_b))) \mathbf{f}(v_b) \quad (45)$$

where $\widetilde{\nabla}S$ is the estimated gradient, x_{adv} is the boundary image at iteration t , S is the difference function that indicates whether the image has been successfully perturbed from the original label to the malicious label, the function $\text{sgn}(S(\cdot))$ denotes the sign of this difference function, and δ is the size of the random perturbation to control the gradient estimation error.

The second component of NonLinear-BA is moving the boundary-image x_{adv} along the estimated gradient direction:

$$\hat{x}_{t+1} = x_{adv}^{(t)} + \xi_t \cdot \frac{\widetilde{\nabla}S}{\|\widetilde{\nabla}S\|_2} \quad (46)$$

where ξ_t is a step size chosen by searching with queries.

Finally, in order to enable the gradient estimation in the next iteration and move closer to the target image, the adversarial image x_{adv} is mapped back to the decision boundary through binary search. This search is aided by queries which seek to find a fitting weight α_t :

$$x_{adv}^{(t+1)} = \alpha_t \cdot x_{tgt} + (1 - \alpha_t) \cdot \hat{x}_{t+1} \quad (47)$$

where x_{tgt} is the target image, i.e., the original image whose correct label x_{adv} seeks to achieve with a crafted perturbed image.

TNonLinear-BA is evaluated on both offline model ImageNet, CelebA, CIFAR10 and MNIST datasets, as well as commercial online APIs. The nonlinear projection-based gradient estimation black-box attacks achieve better performance compared with the state-of-the-art baselines. The authors in [33] discover that when the gradient patterns are more complex, the NonLinear-BA-GAN method fails to keep reducing the MSE after a relatively small number of queries and converges to a poor local optima.

3.5 QEBA: Query-Efficient Boundary-Based Blackbox Attack

Black-box attacks can be query-free or query-based. Query-free attacks are transferability based; query access is not required, as this type of attack assumes the attacker has access to the training data such that a substitute model may be constructed. Query-based attacks can be further categorized into score-based or boundary-based attacks. In a score-based attack, the attacker can access the class probabilities of the model. In a boundary-based attack, only the final model prediction label, rather than the set of prediction confidence scores, is made accessible to the attacker. Both score-based and boundary-based attacks require a substantial number of queries.

One challenge of reducing the number of queries needed for a boundary-based attack is that it is difficult to explore the decision boundary of high-dimensional data without making many queries. The Query-Efficient Boundary-based Blackbox Attack (QEBA) seeks to reduce the queries needed by generating queries through adding perturbations to an image [35]. Thus, probing the decision boundary

is reduced to searching a smaller, representative subspace for each generated query. Three representative subspaces are studied by [35]: spatial transformed subspace, low frequency subspace, and intrinsic component subspace. The optimality analysis of gradient estimation query efficiency in these subspaces is shown in [35].

QEBA performs an iterative algorithm comprised of three steps: first, estimate the gradient at the decision boundary, which is based on the given representative subspace, second, move along the estimated gradient, and third, project to the decision boundary with the goal of moving towards the target adversarial image. These steps follow the same mathematical details as given in Equation 45 to 47 in Section 3.4. Representative subspace optimizations from spatial, frequency, and intrinsic component perspectives are then consequently explored; these subspace-based gradient estimations are shown to be optimal as compared to estimation over the original space [35].

Results for the attack are provided for models trained on ImageNet and models trained on the CelebA dataset. The results show the MSE vs the number of queries, indicating that the three proposed query efficient methods outperform HSJA significantly. The authors also show that the proposed QEBA significantly reduces the required number of queries. In addition, the attack yields high quality adversarial examples against both offline models (i.e. ImageNet) and online real-world APIs such as Face++ and Azure.

3.6 SurFree: a Fast Surrogate-free Blackbox Attack

Many black-box attacks rely on substitution, i.e., a surrogate model is used in place of the target model, the aim being that adversarial examples crafted to attack this surrogate model will effectively transfer to the target classifier. Accordingly, an accurate gradient estimate to create the substitute model requires a substantial number of queries.

By contrast, SurFree is a geometry-based black-box attack that does not query for a gradient estimate [43]. Instead, SurFree assumes that the boundary is a hyperplane and exploits subsequent geometric properties as follows. Consider the pre-trained classifier to be $f : [0, 1]^D \rightarrow \mathbb{R}^C$. A given input image \mathbf{x} produces the label $cl(\mathbf{x}) := \arg \max_k f_k(\mathbf{x})$, where $f_k(\mathbf{x})$ is the predicted probability of class k , $1 \leq k \leq C$. The goal of an untargeted attack is to find an adversarial image \mathbf{x}_a that is similar to a classified image \mathbf{x}_o such that $cl(\mathbf{x}_a) \neq cl(\mathbf{x}_o)$. Thus, an outside region is defined as $\mathcal{O} := \{\mathbf{x} \in \mathbb{R}^D : cl(\mathbf{x}) \neq cl(\mathbf{x}_o)\}$. The desired, optimal adversarial image is then:

$$\mathbf{x}_a^* = \arg \min_{\mathbf{x} \in \mathcal{O}} \|\mathbf{x} - \mathbf{x}_o\| \quad (48)$$

A key assumption of SurFree is that if a point $\mathbf{y} \in \mathcal{O}$, then there exists a point $\mathbf{x}_b \in \overline{\mathcal{O}}$ which can be found that lies on the boundary, denoted as $\partial\mathcal{O}$. Further, it is assumed that the boundary $\partial\mathcal{O}$ is an affine hyperplane that passes through $\mathbf{x}_{b,1}$ in \mathbb{R}^D with normal vector \mathbf{N} . Considering a random basis with span $(\mathbf{x}_{b,1} - \mathbf{x}_o)^\perp$ composed of $D - 1$ vectors $\{\mathbf{v}_i\}_{i=1}^{D-1}$, the inner product between \mathbf{N} and $(\mathbf{x}_{b,k} - \mathbf{x}_o) \propto \mathbf{u}_k$ can be iteratively increased by:

$$\mathbf{N}^\top \mathbf{u}_k = \prod_{i=1}^{D-k} \cos(\psi_{D-i}) \quad (49)$$

where \mathbf{u}_k is the vector that spans the plane containing \mathbf{x}_o , and $\mathbf{x}_{b,D} \in \mathcal{O}$ and $(\mathbf{x}_{b,D} - \mathbf{x}_o)$ is colinear with \mathbf{N} , which points to the projection of \mathbf{x}_o along the boundary of the hyperplane.

Additionally, restricting perturbations to a low dimensional subspace improve the estimation of the projected gradient. The low dimensional subspace is carefully chosen to incorporate meaningful, prior information about the visual content of the image. This further aids in implementing a low query budget.

It is experimentally shown that SurFree bests state-of-the-art techniques for limited query amounts (e.g., one thousand queries) while attaining competitive results in unlimited query scenarios [43]. The geometric details of approximating a hyperplane surrounding a boundary point are left to [43].

The authors present attack results using the criteria of number of queries, and the resulting distortion on the attacked image, on the MNIST and ImageNet datasets. SurFree drops significantly faster than other compared attacks (QEBA and GeoDA) to lower distortions (most notably from 1 to 750 queries).

3.7 A Geometry-Inspired Decision-Based Attack

qFool is a decision-based attack that requires few queries for both non-targeted and targeted attacks [38]. qFool relies on exploiting the locally flat decision boundary around adversarial examples. In the non-targeted attack case, the gradient direction of the decision boundary is estimated based upon the top-1 label result of each query. An adversarial example is then sought in the estimated direction from the original image. In the targeted attack case, gradient estimations are made iteratively from multiple boundary points from a starting target image. Query efficiency is further improved by seeking perturbations in low-dimensional subspace.

Prior literature [20] has shown that the decision boundary has only a small curvature near the presence of adversarial examples. This observation is thus exploited by [38] to compute an adversarial perturbation v . It conceptually follows that the direction of the smallest adversarial perturbation v for the input sample x_0 is the gradient direction of the decision boundary at x_{adv} . Due to the blackbox nature of attack, this gradient cannot be computed directly; however, from the knowledge that the boundary is relatively flat, the classifier gradient at point x_{adv} will be nearly identical to the gradient of other neighboring points along the boundary. Therefore, the direction of v can be suitably approximated by ξ , the gradient estimated at a neighbor point P . Thus, an adversarial example x_{adv} from x_0 is sought along ξ .

The three components of the untargeted qFool attack involve an initial point, gradient estimation, and a directional search. To begin with, the original image x_0 is perturbed by a small, random Gaussian noise to produce a starting point \mathcal{P} on the boundary:

$$\mathcal{P} := x_0 + \min_r \|r\|_2 \text{ s.t. } f_\theta(\mathcal{P}) \neq f_\theta(x_0), r \sim \mathcal{N}(0, \sigma) \quad (50)$$

Noise continues to be added ($\mathcal{P} = x_0 + r_j$) until the image is misclassified. Next, the top-1 label of the classifier is used to estimate the gradient of the boundary $\nabla f(\mathcal{P})$:

$$z_i = \begin{cases} -1 & f(\mathcal{P} + v_i) = f(x_0) \\ +1 & f(\mathcal{P} + v_i) \neq f(x_0) \end{cases}, i = 1, 2, \dots, n \quad (51)$$

where v_i are randomly generated vectors with the same norm to perturb \mathcal{P} and $f(\mathcal{P} + v_i)$ is the label produced by querying the classifier.

For the final step of qFool, the gradient direction at point x_{adv} can be approximated by the gradient direction at point \mathcal{P} , i.e., $\nabla f(\mathcal{P}) \approx \xi$. The adversarial example x_{adv} can thus be found by perturbing the decision boundary in the direction of ξ until the decision boundary is reached. Using binary search, this costs only a few queries to the classifier.

For a targeted attack, the objective becomes perturbing the input image to be classified as a particular target class, i.e., $f_\theta(x_0 + v) = t$ for a target class t . Thus, the starting point of this attack is selected to be an arbitrary image x_t that belongs to the target class t . Due to the potentially large distance between x_0 and x_t , the assumption of a flat decision boundary between the initial and targeted adversarial regions no longer holds. Instead, a linear interpolation in the direction of $(x_t - x_0)$ is utilized to find a starting point \mathcal{P}_0 :

$$\mathcal{P}_0 := \min_{\alpha} (x_0 + \alpha \cdot \frac{x_t - x_0}{\|x_t - x_0\|_2}) \text{ s.t. } f_\theta(\mathcal{P}_0) = t \quad (52)$$

The gradient direction estimation of ξ_0 at \mathcal{P}_0 follows the same method as outlined for untargeted attacks.

The qFool attack is experimentally demonstrated on the ImageNet dataset by attacking VGG-19, ResNet50 and Inception v3. The results show that qFool is able to achieve a smaller distortion in terms of MSE, as compared to the Boundary Attack when both attacks use the same number of queries. However, the overall attack success rate for qFool is not reported. The authors also test qFool on the Google Cloud Vision API.

4 TRANSFER ATTACKS

In this section, we explore recent advances in adversarial machine learning with respect to transfer attacks. The adversarial model for these attacks allows the attacker to query the target defense and/or access some of the target defense’s training dataset. The attacker then uses this information to create a synthetic model which the attacker then attacks using a white box attack. The adversarial inputs generated from the white box attack on the synthetic model are then transferred to the targeted defense.

We cover 3 recently proposed transfer attacks. These attacks include the Adaptive Black-Box Transfer attack [42], DaST attack [58] and the Transferable Targeted attack [37].

4.1 The Adaptive Black-box Attack

A new transfer based black-box attack is developed in [42] that is an extension of the original Papernot attack proposed in [46]. Under this threat model the adversary has access to the training dataset (X, Y) , and query access to the classifier under attack, C . In the original Papernot formulation of the attack, the attacker labels the training data to create a new training dataset $(X, C(X))$. The adversary is then able to train synthetic model S on $(X, C(X))$ while iteratively augmenting the dataset using a synthetic data generation technique. This results in a trained synthetic model $S(w_s)$. In the final step of the attack, a white-box attack generation method $\phi(\cdot)$ is used in conjunction with the trained synthetic model S in order to create adversarial examples X_{adv} :

$$X_{adv} = \phi(X_{clean}, S, w_s) \quad (53)$$

where X_{clean} are clean testing examples and ϕ is a white-box attack method i.e. FGSM [23].

The enhanced version of the Papernot attack is called the mixed [42] or adaptive black-box attack [41]. Where as in the original Papernot attack 0.3% of the training data is used, the adaptive version increases the strength of the adversary by using anywhere from 1% to 100% of the original training data. Beyond this, the attack generation method ϕ is varied to account for newer white-box attack generation methods that have better transferability. In general the most effective version of the attack replaces ϕ_{FGSM} with ϕ_{MIM} , the Momentum Iterative Method (MIM) [18]. The MIM attack computes an accumulated gradient [18]:

$$g_{t+1} = \mu \cdot g_t + \frac{J(x_t^{adv}, y)}{\|\nabla_x J(x_t^{adv}, y)\|_1} \quad (54)$$

where $J(\cdot)$ is the loss function, μ is the decay factor and x_t^{adv} is the adversarial sample at attack iteration t . For a L_∞ bounded attack, the adversarial example at iteration t is:

$$x_{t+1}^{adv} = x_t^{adv} + \frac{\epsilon}{T} \cdot \text{sign}(g_{t+1}) \quad (55)$$

where T represents the total number of iterations in the attack and ϵ represents the maximum allowed perturbation.

In [42], the attack is tested using the CIFAR-10 and Fashion-MNIST datasets. The adaptive black-box attack is shown to be effective against vanilla (undefended) networks, as well as a variety of adversarial machine learning defenses.

4.2 DaST: Data-free Substitute Training for Adversarial Attacks

As described in the SurFree attack in 3.6, substitute models can be difficult or unrealistic to obtain, particularly if a substantial amount of real data labeled by the target model is needed. DaST is a data-free substitute training method that utilizes generative adversarial networks (GANs) to train substitute models without the use of real data [58]. To address the potentially uneven distribution of GAN-produced samples, a multi-branch architecture and label-control loss for the GAN model is employed.

To describe the necessary context for DaST, let \mathbf{X} denote samples from the target model T , \bar{y} and y' denote the true labels and target labels of the samples \mathbf{X} , respectively, and let $T(y|\mathbf{X}, \theta)$ denote the target model parameterized by θ . Then, the objective of a targeted attack becomes:

$$\min_{\epsilon} \|\epsilon\| \text{ subject to } \underset{y_i}{\text{argmax}} T(y_i|\bar{\mathbf{X}} = \mathbf{X} + \epsilon, \theta) = y' \text{ and } \|\epsilon\| \leq r \quad (56)$$

where ϵ and r are the sample and upper bounds of the perturbation, respectively, and $\bar{\mathbf{X}} = \mathbf{X} + \epsilon$ refer to the adversarial examples that lead the target model T to misclassify a sample with a selected wrong label.

To further provide adequate context for DaST, a white-box attack under these settings would have full access to the gradient construction of the target model T and thus leverage this information to generate adversarial examples. In a black-box substitute attack under these settings, a substitute model \hat{T} would stand-in for the

target model, and the adversarial examples generated to attack \hat{T} would then be transferred to attack T . Thus, coming to the settings of a data-free black-box substitute attack, DaST utilizes a GAN to synthesize a training set for \hat{T} that is as similar as possible to the training set of the target model T .

To this end, the substitute training set crafted by the GAN aims to be evenly distributed across all categories of labels, which are produced from T . To accomplish this, for N categories, the generative network in [58] is designed to contain N upsampling deconvolutional components, which then share a post-processing convolutional network. The generative model G randomly samples a noise vector \mathbf{z} from the input space as well as the variable label n . \mathbf{z} then enters the n -th upsampling deconvolutional network and the shared convolutional network to produce the adversarial sample $\hat{\mathbf{X}} = G(\mathbf{z}, n)$. The label-control loss for G is given as:

$$\mathcal{L}_c = \text{CE}(T(G(\mathbf{z}, n)), n) \quad (57)$$

where CE is the cross-entropy.

To approximate the gradient information of T to train a label-controllable generative model, the following objective function is used:

$$\min_D d(T(\hat{\mathbf{X}}), D(\hat{\mathbf{X}})) \quad (58)$$

For the same inputs, the outputs of D will approach the outputs of T for the same inputs as training proceeds. Thus, D replaces T in Equation 57:

$$\mathcal{L}_c = \text{CE}(D(G(\mathbf{z}, n)), n) \quad (59)$$

The loss of G is then updated as:

$$\mathcal{L}_G = e^{-d(T,D)} + \alpha \mathcal{L}_c \quad (60)$$

where α is the weight of the label-control loss.

As the training stage progresses, as does the imitation quality of D , leading to a diverse set of synthetically generated samples labeled by T . These data-free substitute training-produced samples are then used to attack T .

DaST reduces the need for adversarial substitute attacks by utilizing GANs to generate synthetic samples, and thus can train substitute models without the requirement of any real data. Authors present results on using DaST to train a substitute model for adversarial attacks on the CIFAR-10 and MNIST trained models. The substitute models trained by DaST perform better than baseline models on FGSM and C&W attacks (targeted).

4.3 Towards Transferable Targeted Attack

Crafting targeted transferable examples has the dual challenges of noise curing, i.e., the decreasing gradient magnitude in iterative attacks that results in momentum accumulation, and the difficulty of moving adversarial examples toward a target class while creating distance from the true class. To this end, Li et al [37] propose a novel targeted, transferable attack that applies the Poincaré distance to combat noise curing by creating a self-adaptive gradient, and employs metric learning to improve the distance from an adversarial example's true label.

To overcome the drawback of the Poincaré distance fused logits failing to satisfy $\|l(x)\|_2 < 1$, this attack normalizes logits by the l_1 distance. To overcome the problem of potential infinite distances between a point and its target label, a constant of $\xi = 0.0001$ is

subtracted from the one-hot target label y . The Poincaré distance metric loss is given as:

$$\mathcal{L}_{P_o}(x, y) = d(u, v) = \text{arccosh}(1 + \delta(u, v)) \quad (61)$$

where d refers to the Poincaré distance, $l_k(x)$ indicates the output logits of the k -th model, $u = l_k(x) / \|l_k(x)\|_1$, $v = \max\{y - \xi, 0\}$, and $l(x)$ refer to the fused logits. By contrast, this attack formulates adversarial examples through the fusion of logits from a combination of models, as shown below:

$$l(x) = \sum_{k=1}^K w_k l_k(x) \quad (62)$$

where K is the number of ensemble models, $l_k(x)$ indicates the output logits of the k -th model, and w_k is the ensemble weight of the k -th model, with $w_k > 0$, $\sum_{k=1}^K w_k = 1$. Note that the fused logits are the average of the ensemble models.

Triplet loss is a popular targeted attack loss function that increases the distance between the adversarial example and the true label, while decreasing the distance between the adversarial example and the target label [25]. A common triplet loss function appears as:

$$\mathcal{L}_{trip}(x^a, x^p, x^n) = [D(x^a, x^p) - D(x^a, x^n) + \gamma]_+ \quad (63)$$

where x^a, x^p, x^n are the anchor, positive, and negative examples, respectively, where x^a and x^p are of the same class, while x^n is of a different class than x^a . The distance d is based on the embedding vector for the anchor, positive, and negative networks in the triplet configuration. Additionally, $\gamma \geq 0$ is a hyperparameter that regulates the margin between the distance metrics $D(x^a, x^p)$ and $D(x^a, x^n)$.

A drawback of standard triplet loss is the need to sample new data, which is often infeasible in a targeted attack. Instead, this work formulates the triplet input as the logits of clean images, $l(x_{clean})$, the one-hot target label, y_{tar} , and the true label, y_{true} :

$$\mathcal{L}_{trip}(y_{tar}, l(x_i), y_{true}) = [D(l(x_i), y_{tar}) - D(l(x_i), y_{true}) + \gamma]_+ \quad (64)$$

Due to $l(x^{adv})$ not being normalized, angular distance is used as a distance metric (note that x_i corresponds to x_{adv} below):

$$D(l(x^{adv}), y_{tar}) = 1 - \frac{|l(x^{adv}) \cdot y_{tar}|}{\|l(x^{adv})\|_2 \|y_{tar}\|_2} \quad (65)$$

However, the usage of angular loss does not account for the influence of the norm on the loss value. Thus, an additional triplet loss term appears in the final loss function, as shown below:

$$\mathcal{L}_{all} = \mathcal{L}_{P_o}(l(x), y_{tar}) + \lambda \cdot \mathcal{L}_{trip}(y_{tar}, l(x_i), y_{true}) \quad (66)$$

where x is the original, clean input image and x_i is the result of the i -th iteration of the perturbation of x , with the final result being x_{adv} .

Results on ImageNet illustrate that this attack [37] achieves improved success rates over traditional attacks for white and black-box models in the targeted setting.

5 NON-TRADITIONAL NORM ATTACKS

In this section, we discuss recent developments in black-box attacks that use non-traditional norm threat models. In all attacks covered in previous sections the focus has been on adversaries which try to create adversarial examples with respect to the l_2 or l_∞ norms. We denote the l_2 and l_∞ as "traditional" norms simply because that is what a majority of the literature (17 of the 20 attacks) thus far have focused on. While this is not a strictly technical definition, it gives us a convenient and simple way to categorize the different attacks.

We cover three non-traditional norm attacks in this section. The first attack we summarize is the sparse and imperceivable attack [15] which focuses on black-box attacks with respect to the l_0 norm. The second non-traditional norm attack we survey is Patch Attack [54]. The Patch Attack is based on completely replacing small part of the original image with an adversarial generated square (patch). The last attack we cover in this section is ColorFool [52]. This attack is based on manipulating the colors within the image as opposed to directly adding adversarial noise.

5.1 Sparse and Imperceivable Attack

The Sparse and Imperceivable attacks proposed in [15] are l_0 , black-box attacks that produce adversarial inputs while minimizing the number of perturbed pixels. The attacks come in multiple forms, but the general goal and scheme remains the same. Each attack relies on having the score based output of the network to operate, and each version of the attack attempts to solve the following optimization problem:

$$\begin{aligned} \min \gamma(x' - x) \\ \text{s.t. } \arg \max f(x') \neq \arg \max f(x) \end{aligned} \quad (67)$$

where x is the clean image, x' is an adversarial image, f is a function returning the classifier's score vector, and γ is a distance function defined as follows:

$$\gamma(x' - x) = \sum_{i=1}^d \max_j \mathbb{1}[x'_{ij} - x_{ij}] \neq 0 \quad (68)$$

Where x_{ij} refers to the i^{th} pixel of the j^{th} color channel. It is important to note that color images are typically represented as three 2-D matrices, with one matrix corresponding to each color channel. However, in the mathematical formulation for these attacks, they treat each color channel as a 1-D matrix for notational convenience.

In Equation 68, essentially γ counts the number of pixels in the adversarial image that deviate from the original image. There are three versions of the attack: l_0 , $l_0 + l_\infty$, $l_0 + \sigma$. Where $l_0 + l_\infty$ and $l_0 + \sigma$ add their own additional constraints on the optimization problem as outlined below:

attack type	Additional Constraint
$l_0 + l_\infty$	$\ x' - x\ \leq \epsilon$
$l_0 + \sigma$	Perturbations must be imperceivable

Where ϵ is the maximum allowed perturbation magnitude. Each of the attack variants follow the same scheme, with the main difference being the amount each pixel is perturbed.

The attack begins by first iterating through each pixel in the image and generating a set of pixel perturbations $\{x'_{ij}\}$ according

to the following equations:

Attack Type	Pixel Perturbation
l_0	$x'_{ij} \in \{0, 1\}$
$l_0 + l_\infty$	$x'_{ij} = x_{ij} \pm \epsilon$
$l_0 + \sigma$	$x'_{ij} = (1 \pm \kappa \sigma_{ij})x_{ij}$

Here σ is the standard deviation of the image in color channel j in proximity to pixel x_{ij} . It is calculated as follows:

$$\sigma_{ij} \sqrt{\min\{\sigma_{ij}^y, \sigma_{ij}^x\}} \quad (69)$$

Where σ_{ij}^x is the standard deviation of x_{ij} and the two pixels adjacent to it horizontally in color channel j , and σ_{ij}^y is the same but for the two pixels adjacent to x_{ij} vertically. The σ term is essential to what makes the $l_0 + \sigma$ attack imperceivable to humans. It allows the attack to avoid perturbations near edges in the image as they are more easily perceivable. It also focuses the attack on increasing the intensity of pixels rather than modifying their color. Each pixel perturbation is then clipped into the $[0, 1]$ range. After this each generated $x'_{i,j}$ is sorted in decreasing order according to the value of the following equation:

$$\pi^r(x'_{ij}) = f_r(x'_{ij}) - f_c(x'_{ij}) \quad (70)$$

Where c is the ground truth label of x , r is any class label other than c , $\pi^r(x')$ is the score of perturbation x' with respect to class r , and f_r is the model's predicted value for class r in the score vector. The perturbations $x'_{i,j}$ are then sorted by their π^r score for each r .

After the sorting, an iterative process begins. At each iteration a number of maximum pixels perturbed, k , is chosen, starting at one pixel and progressing to k_{max} by the end. During each of these iterations an inner loop iterative process begins. The inner loop iterates through each of the possible class labels, r , other than c , applying k of the top N single pixel perturbations with respect to class r to the clean image. If at any point in the algorithm a perturbation leads to a changed class label, the algorithm stops and returns the penultimate perturbation. The attack is tested on the MNIST and CIFAR-10 datasets where it achieves a high attack success rate while perturbing a small amount of total pixels on average. The attack is compared to both white-box and black-box attacks where it achieves a similar attack success rate to attacks like C&W attack and Sparse Fool while having the median least pixels perturbed per attack.

5.2 Patch Attack

The Patch Attack is a black box attack proposed in [54] that utilizes textured patches and reinforcement learning to generate adversarial images. Each patch is cut out from from images in a pre-generated library of textures. Each texture is designed such that neural networks strongly associate them with a particular class label in a given dataset. The attack is perceivable to the human eye and applies a large magnitude perturbation to the clean image. However, the perturbation is localized and can be shrunk through optimization techniques. In this attack the reinforcement learning agent solves the following optimization problem:

$$\min L(y, y') = -r \cdot \ln(P) \quad (71)$$

Where L is the attack loss, y is the target model’s predicted score for the ground truth class, $y' \neq y$ is the model’s predicted score for a class other than y , r is the reinforcement learning agent’s reward, and P is the agent’s output probability of taking action a that lead to the most recent reward. $r = \ln(y')$ in a targeted attack, and $r = \ln(y' - y)$ in an untargeted attack. The reinforcement learning agent’s policy network is represented by an LSTM and a fully connected layer. An action is defined as follows:

$$a = \{u_1, v_1, i, u_2, v_2\} \quad (72)$$

Where i is a texture index in the texture library, u_1 and v_1 are corner positions used to crop the texture, and u_2 and v_2 are corner positions denoting where the cropped texture should be placed on the clean image. The attack algorithm is iterative, at each time step the environment state is determined and an agent is trained. Once trained the agent outputs a probability distribution over the possible actions. One action is sampled from the distribution and the associated patch is applied.

Before the attack the texture images must be obtained externally or generated by an algorithm. In the latter case the generation algorithm is initialized with a CNN trained on the target dataset. In our description of the attack, we denote this CNN as the texture CNN to distinguish it from the CNN being attacked. A set of data is also chosen to be used for the texture generation. The data is pre-processed using Grad-CAM [51] which masks out areas of the image that are irrelevant to the primary texture information. For each convolutional layer, j , in the i^{th} block of the texture CNN, a feature map F_i^j is generated. Furthermore, the corresponding Gram matrix, G_i^j is also calculated for each feature activation of each image. Each F_i^j and G_i^j will help encode the most important texture information in input image according to the texture CNN. The computation of the feature maps and Gram matrices are described in detail in [22].

For each image, the Gram matrices generated are then flattened and concatenated into the vector \bar{G} which encodes the texture information. From here, the \bar{G} s are organized by the original class label of the input that generated them. This is done so that the final textures can be labeled and to maximize the effectiveness of each patch. For each class label, the \bar{G} s are clustered into N clusters. In [54] N is chosen to be 30 in order to have sufficient diversity in the texture pool. In practice, the best value of N will vary based upon dataset and application. For each cluster, \bar{G}_c , is then used to generate a feature embedding of the final texture images. This is done by minimizing the following optimization problem over G_t :

$$L = \lambda(\bar{G} - G_t)^2 \quad (73)$$

Where λ is a weight constant and G_t is the feature embedding of the texture image used to generate the final texture image. We omit some details of the attack explanation for brevity, further details are given in [21].

The attack is tested on the ImageNet dataset where it achieves a high attack success rates while covering small portions of the clean image with patches. The attack also shows an ability to maintain its high attack success rate even when defense techniques are applied to the classifier.

5.3 ColorFool: Semantic Adversarial Colorization

ColorFool [52] presents a content-based black-box adversarial attack with unrestricted perturbations that selectively manipulates colors within chosen ranges to thwart classifiers, while remaining undetected by humans. ColorFool operates on the independent a and b channels of the perceptually uniform Lab color space [50]. Color modifications are implemented without changing the lightness, L , of the given image. Further, ColorFool solely selects perturbations within a defined natural-color range for particular acceptable categories [55].

ColorFool divides images into sensitive and non-sensitive regions to be considered for color modification. Sensitive regions, defined $\mathbb{S} = \{\mathbf{S}_k\}_{k=1}^S$ are separated from non-sensitive regions, defined $\bar{\mathbb{S}} = \{\bar{\mathbf{S}}_k\}_{k=1}^{\bar{S}}$, where $\mathcal{S} = \mathbb{S} \cup \bar{\mathbb{S}}$.

The color of the sensitive regions, \mathbb{S} , is modified to generate the adversarial set $\dot{\mathbb{S}}$ as follows:

$$\dot{\mathbb{S}} = \{\dot{\mathbf{S}}_k : \dot{\mathbf{S}}_k = \gamma(\mathbf{S}_k) + \alpha[0, N_k^a, N_k^b]^T\}_{k=1}^S \quad (74)$$

where color channel a ranges from green (-128) to red (+127), color channel b ranges from blue (-128) to yellow (+127), and the brightness L ranges from black (0) to white (100) within the Lab color space [50]. Further, $\gamma(\cdot)$ converts the intensities of an RGB image to the Lab colorspace, and $N_k^a \in \mathcal{N}_k^a$ and $N_k^b \in \mathcal{N}_k^b$ are randomly chosen adversarial perturbations from the set of natural color ranges \mathcal{N}_k^a and \mathcal{N}_k^b [55] within the a and b channels.

The color of the non-sensitive regions, $\bar{\mathbb{S}}$, is modified as follows to produce to the set $\dot{\bar{\mathbb{S}}}$:

$$\dot{\bar{\mathbb{S}}} = \{\dot{\bar{\mathbf{S}}}_k : \dot{\bar{\mathbf{S}}}_k = \gamma(\bar{\mathbf{S}}_k) + \alpha[0, \bar{N}_k^a, \bar{N}_k^b]^T\}_{k=1}^{\bar{S}} \quad (75)$$

where $\bar{N}_k^a, \bar{N}_k^b \in \{-127, \dots, 128\}$ are randomly chosen within the ranges of a and b , respectively. Note that the full ranges of a and b are considered, as non-sensitive regions are able to undergo greater intensity changes.

The modified sensitive and non-sensitive regions are combined to generate the adversarial image $\dot{\mathbf{X}}$, as shown below:

$$\dot{\mathbf{X}} = Q(\gamma^{-1}(\sum_{k=1}^S \dot{\mathbf{S}}_k + \sum_{k=1}^{\bar{S}} \dot{\bar{\mathbf{S}}}_k)) \quad (76)$$

where $Q(\cdot)$ is the quantization function that keeps the generated adversarial image within the dynamic range of pixel values, i.e., $\dot{\mathbf{X}} \in \mathbb{Z}^{w,h,c}$, and $\gamma^{-1}(\cdot)$ is the inverse function that converts image intensities from the Lab color space to RGB .

ColorFool provides robustness to defenses that utilize filters, adversarial training or modified training loss functions. Additionally, ColorFool is less detectable than restricted attacks, including JPEG compression. The empirical results were presented on the Private-Places365, CIFAR-10, and ImageNet datasets, indicating higher success rates in the previously mentioned categories.

6 ATTACK SUCCESS RATE ANALYSIS

In this paper we compile the experimental results from many different sources and present them together in tabular form. While it may

be tempting to directly compare attack success rates, here we give a theoretical analysis to show the fallacy of direct comparisons.

The definition of a *successful* adversarial example varies between papers based on which constraints are enforced during the execution of the attack. We can formalize this as follows: For classifier C , the associated set of clean correctly identified examples is denoted as $\mathcal{X}(C)$ such that:

$$\mathcal{X}(C) = \{(x_i, y_i) \in \mathcal{X}_t : C(x_i) = y_i\}. \quad (77)$$

where \mathcal{X}_t is the entire set of testing images. When classifier C is attacked, we can formally define the constraints on the attacker for a query-based black-box attack using the following threat vector: $W_{threat} = [w_c, w_q, w_\epsilon]$ where $w_i \in \{0, 1\}$. $w_c = 1$ corresponds to a successful attack definition where the classifier must produce the wrong class label i.e., $C(x_{adv}) \neq y_i$. Likewise, $w_q = 1$ corresponds to a successful attack where the adversarial example is generated within a fixed number of queries q_{adv} , and the number of queries are less than the query budget q : $(q - q_{adv}) \geq 0$. Lastly $w_\epsilon = 1$ ensures that the adversarial example falls within a certain $||l||_p$ distance ϵ of the original example x_i : $||x_{adv} - x_i||_p \leq \epsilon$. If any of the values in W_{threat} are 0, it simply means that the corresponding condition is not used in defining a successful adversarial example.

We denote $\phi(x_i, y_i)$ as the adversarial attack method used with respect to clean sample (x_i, y_i) returned within q_{adv} queries. Written explicitly $(x_{adv}, q_{adv}) = \phi(x_i, y_i)$ and we assume $\phi(\cdot)$ to be deterministic in nature. While this assumption may not hold true for all attacks, this simplifies the notation for the theoretical attack success rate. The attack success rate α over the clean set $\mathcal{X}(C)$ with respect to classifier C is:

$$\alpha = \frac{\left| \left\{ \begin{array}{l} (x_i, y_i) \in \mathcal{X}(C) : \\ (x_{adv}, q_{adv}) = \phi(x_i, y_i) \Rightarrow \\ (C(x_{adv}) \neq y_i \vee w_c = 0) \wedge \\ w_q(q - q_{adv}) \geq 0 \wedge w_\epsilon ||x_{adv} - x_i||_p \leq \epsilon \end{array} \right\} \right|}{|\mathcal{X}(C)|}, \quad (78)$$

From Equation 78, it can be seen that under the most restricted threat model ($W_{threat} = [1, 1, 1]$), the attack must produce an adversarial example that is misclassified i.e., $C(x_{adv}) \neq y_i$, created using limited query information i.e., $q_{adv} \leq q$ and within an acceptable $||l||_p$ norm. Such a threat model requires specification of the attack parameters q , p and ϵ . That is q the maximum allowed number of queries per sample, p the norm measurement and ϵ the maximum allowed perturbation.

Our framework for a vector defined threat model W_{threat} and corresponding attack success rate α is useful for two reasons. First, it allows us to categorize every query based black-box attack according to the three value system. Second and most importantly, this framework allows us to see where comparisons between attack success rates reported in different papers are legitimate. We illustrate this next point with an example from the literature.

Consider the Square attack [2] and the Zeroth-Order alternating direction method of multipliers attack (ZO-ADMM) [57]. The untargeted attack success rate of both attacks is reported with respect to an Inception v3 network trained on ImageNet. The Square attack reports an attack success rate α of 92.2% while the ZO-ADMM reports an attack success rate of 100%. Using ONLY these two values

Table 2: Adversarial threat models used to determine the attack success rate in each paper. $w_c = 1$ corresponds to an attack success rate where misclassification (or targeted misclassification) defines a successful adversarial attack. $w_q = 1$ corresponds to a successful adversarial attack done within a fixed query budget. $w_\epsilon = 1$ corresponds to a successful attack when the adversarial example is within a certain perturbation bound ϵ of the clean example.

Score based Attacks	
qMeta [19]	$w_c = 1, w_q = 0, w_\epsilon = 0$
P-RGF [13]	$w_c = 1, w_q = 0, w_\epsilon = 1$
ZO-ADMM [57]	$w_c = 1, w_q = 0, w_\epsilon = 0$
TREMBA [27]	$w_c = 1, w_q = 0, w_\epsilon = 1$
Square [2]	$w_c = 1, w_q = 0, w_\epsilon = 1$
ZO-NGD [57]	$w_c = 1, w_q = 0, w_\epsilon = 1$
PPBA [36]	$w_c = 1, w_q = 0, w_\epsilon = 1$
Decision based Attacks	
qFool [38]	$w_c = 0, w_q = 1, w_\epsilon = 0$
HSJA [10]	$w_c = 1, w_q = 1, w_\epsilon = 0$
GeoDA [47]	$w_c = 0, w_q = 1, w_\epsilon = 0$
QEBA [35]	$w_c = 1, w_q = 1, w_\epsilon = 1$
RayS [9]	$w_c = 1, w_q = 0, w_\epsilon = 1$
SurFree [43]	$w_c = 1, w_q = 1, w_\epsilon = 1$
NonLinear-BA [33]	$w_c = 1, w_q = 0, w_\epsilon = 0$
Transfer based Attacks	
Adaptive [42]	$w_c = 1, w_q = 0, w_\epsilon = 1$
DaST [58]	$w_c = 1, w_q = 0, w_\epsilon = 1$
PO-TI [37]	$w_c = 1, w_q = 0, w_\epsilon = 1$
Non-traditional Attacks	
CornerSearch [15]	$w_c = 1, w_q = 0, w_\epsilon = 0$
ColorFool [52]	$w_c = 1, w_q = 0, w_\epsilon = 0$
Patch [54]	$w_c = 1, w_q = 0, w_\epsilon = 0$

without our threat model framework makes it seem like the ZO-ADMM attack is much stronger than the Square attack, as it never fails. However, let us now consider the threat models. The threat model for ZO-ADMM is: $W_{threat} = [w_c = 1, w_q = 0, w_\epsilon = 0]$. The Square attack on the other hand has the following threat model: $W_{threat} = [w_c = 1, w_q = 0, w_\epsilon = 1]$. Essentially the Square attack is reporting a high attack success rate under a MORE restrictive threat model where the adversarial example must be wrongly classified *and* under a certain l_2 distance from the original clean example. The ZO-ADMM attack success rate is reported only on examples that are wrongly classified, a much weaker threat model.

7 EXPERIMENTAL RESULTS

In this section, we discuss the experimental results for all of the black-box attacks. Broadly speaking there are three common datasets that are used in measuring the attack success rate of black-box attacks.

- (1) **MNIST** - The MNIST dataset [31] consists of 60,000 training images and 10,000 test images. The dataset has 10 classes, each class is a different handwritten digit, 0-9. Each digit is a 28×28 grayscale image. In general, the maximum

allowed perturbation ϵ for MNIST is high as compared to other dataset. For example, $\epsilon = 0.2, 0.3$ as seen in table 5. This may generally be due to the fact that MNIST images can have large perturbations, while still being visually recognizable to humans.

- (2) **CIFAR-10** - The CIFAR-10 dataset [30] consists of 50,000 training images and 10,000 test images. The 10 classes in CIFAR-10 are airplane, car, bird, cat, deer, dog, frog, horse, ship and truck. Each image is $32 \times 32 \times 3$ (color images).
- (3) **ImageNet** - The ImageNet dataset [17] contains over 14 million color images that are labeled from $\approx 20,000$ categories. The images in ImageNet are color images, however the exact size of each image varies.

In the following subsection we break down the analyses according to the four different attack categories.

7.1 Score based Attack Analysis

In table 3, we show the compiled results drawn across all the papers we surveyed for the score based attacks on ImageNet classifiers. We report MNIST and CIFAR-10 results in Table 5 and Table 6. As the majority of the attacks are done with respect to the ImageNet dataset, we relegate our discussion and analysis to those results in this subsection.

Let us consider the $l_\infty = 0.05$ norm adversary, untargeted attack with adversarial model $w_c = 1, w_q = 0, w_\epsilon = 1$. Under this adversarial threat model, three attacks have a 99% or greater attack success rate (Square, p-RGF and TREMBA). While all three attacks are done on ResNet classifiers, there is a slight difference (TREMBA and P-RGF are tested on ResNet34 and the Square attack is tested on ResNet50). Aside from this difference, if we compare results, the Square attack and TREMBA are both able to achieve a remarkable double digit query count while still maintaining a 99% or greater attack success rate. Square attack requires 73 queries on average while TREMBA requires 27.

While no attack in table 3 uses the most restrictive threat model (i.e. $w_c = 1, w_q = 1, w_\epsilon = 1$) we can see that the most common threat model is $w_c = 1, w_q = 0, w_\epsilon = 1$ making comparison between attack that use this threat model and the same classifier possible. Alternatively, one attack (ZO-ADMM) uses a highly unrestricted adversarial model $w_c = 1, w_q = 0, w_\epsilon = 0$ making it impossible to directly determine the fidelity of the ZO-ADMM attack in relation to other state-of-the-art attacks.

7.2 Decision based Attack Analysis

In Table 4 we give the results for all decision based attacks that were conducted on ImageNet classifiers. Likewise, results for decision based attacks on MNIST and CIFAR-10 can be found in Table 5 and Table 6. Due to the large number of attacks and datasets, in this subsection we specifically focus on the decision based attacks for ImageNet CNNs.

Let us first consider the l_2 norm decision based attacks that are targeted. For this setting, when looking at the most restricted threat model ($w_c = 1, w_q = 1, w_\epsilon = 1$) we can see that SurFree gives the best query efficient attack (on ResNet18) with a 90% attack success rate using only 500 queries. However, in terms of minimal distortion ($\epsilon = 0.001$), QEBA-S and NonLinear-BA can both achieve

an 80% attack success rate or higher with a query budget of 10,000. Alternatively, if we consider the l_∞ norm and an untargeted attack, it is clear the RayS attack is the best attack. RayS achieves a 98.9% attack success rate on Inception v3 using an average of 748.2 queries per sample. This same holds true for datasets like MNIST and CIFAR-10. In both cases RayS can achieve a 99% or higher attack success rate.

It is important to note that certain threat models make attack results opaque and difficult to compare. For example, the threat model ($w_c = 0, w_q = 1, w_\epsilon = 0$) is used to report attack results for GeoDA and qFool. In this case, the median distortion is considered the independent variable (i.e. the one that changes between different attacks). However, when only the median distortion is reported this does not give any information about what percent of adversarial examples are actually misclassified (which would constitute a successful attack). Reporting the median also does not give the full picture in terms of the average distortion required to create a successful adversarial example in the attack.

7.3 Transfer based Attacks and Non-traditional Attacks

In Table 7 the results for the transfer based attacks and non-traditional attacks are shown. For the transfer attacks, each attack is done under slightly different assumptions making direction comparison difficult. For example, the Adaptive attack requires all the training data to be available to the attacker, where as in DaST the attack is specifically built around not having direct access to the original training data. Overall, we can claim that the transfer based attacks, just in terms of attack success rates, are not as high as the best score based and decision based black-box attacks. For example, the Adaptive attack has a 74% attack success rate on CIFAR-10 for the l_∞ based attacker. The decision based RayS attack has a 99.8% attack success rate for CIFAR-10 (again l_∞ norm based attack).

For the non-traditional attacks, we can see several interesting trends. First, the patch attack has an extremely high attack success rate on ImageNet (greater than 99%) regardless of whether the attack is targeted or untargeted. Likewise, the l_0 based CornerSearch attack can also achieve a high untargeted attack success rate (greater than 97%) across both MNIST and CIFAR-10 datasets.

The only attack that performs relatively poorly (less than 50% attack success rate) is ColorFool. This may partially be due to the fact that the ColorFool attack can be run in both white-box and black-box form. The ColorFool black-box reported attack results are based on a transfer style attack as opposed to a query based method. As we mentioned above, the new score and decision based attacks (which use query information) that we survey have a higher attack success rates than the new transfer based attacks for the l_2 and l_∞ norms. Essentially, we conjecture there may still be room to improve the black-box ColorFool attack using a query based methodology.

8 CONCLUSION

Adversarial machine learning is advancing at a fast pace, with new attack papers being proposed every year. In light of these recent developments, we have surveyed the current state-of-the-art black-box attack and have provided three major contributions. First,

Attack Name	ASR	Avg Queries	Norm	Target	Classifier	Adv Threat Model	Source
PPBA	84.8	668	$l_2, \epsilon=5$	U	ResNet50	$w_c = 1, w_q = 0, w_e = 1$	[36]
PPBA	65.3	1051	$l_2, \epsilon=5$	U	Inception v3	$w_c = 1, w_q = 0, w_e = 1$	[36]
PPBA	96.6	481	$l_\infty, \epsilon=0.05$	U	ResNet50	$w_c = 1, w_q = 0, w_e = 1$	[36]
PPBA	67.9	1026	$l_\infty, \epsilon=0.05$	U	Inception v3	$w_c = 1, w_q = 0, w_e = 1$	[36]
ZO-NGD	97	582	$l_\infty, \epsilon=0.05$	U	Inception v3	$w_c = 1, w_q = 0, w_e = 1$	[56]
Square	99.7	197	$l_\infty, \epsilon=0.05$	U	Inception v3	$w_c = 1, w_q = 0, w_e = 1$	[2]
Square	100	73	$l_\infty, \epsilon=0.05$	U	ResNet50	$w_c = 1, w_q = 0, w_e = 1$	[2]
Square	92.2	1100	$l_2, \epsilon=5$	U	Inception v3	$w_c = 1, w_q = 0, w_e = 1$	[2]
Square	99.3	616	$l_2, \epsilon=5$	U	ResNet50	$w_c = 1, w_q = 0, w_e = 1$	[2]
TREMBA	100	27	$l_\infty, \epsilon=0.05$	U	Resnet34	$w_c = 1, w_q = 0, w_e = 1$	[27]
TREMBA	99.44	443	$l_\infty, \epsilon=0.05$	T	Resnet34	$w_c = 1, w_q = 0, w_e = 1$	[27]
ZO-ADMM	98	-	l_2	U	Inception v3	$w_c = 1, w_q = 0, w_e = 0$	[57]
ZO-ADMM	97	-	l_2	T	Inception v3	$w_c = 1, w_q = 0, w_e = 0$	[57]
P-RGF	100	328	$l_\infty, \epsilon=0.05$	U	Resnet34	$w_c = 1, w_q = 0, w_e = 1$	[27]
P-RGF	98.05	5498	$l_\infty, \epsilon=0.05$	T	Resnet34	$w_c = 1, w_q = 0, w_e = 1$	[27]
P-RGF	98.1	745	$l_2, \epsilon \approx 16.43$	U	Inception v3	$w_c = 1, w_q = 0, w_e = 1$	[13]
P-RGF	99.6	452	$l_2, \epsilon \approx 16.43$	U	ResNet50	$w_c = 1, w_q = 0, w_e = 1$	[13]
P-RGF	97.3	812	$l_\infty, \epsilon=0.05$	U	Inception v3	$w_c = 1, w_q = 0, w_e = 1$	[14]
P-RGF	99.6	388	$l_\infty, \epsilon=0.05$	U	ResNet50	$w_c = 1, w_q = 0, w_e = 1$	[14]

Table 3: Score based black-box attacks on ImageNet classifiers. The corresponding success rate (ASR) and adversarial threat model are shown for each attack along with the source paper from which the results are drawn from.

our survey covers 20 new attack papers with detailed summaries, mathematics and attack explanations. Our second contribution is a categorization of these attacks into four different types, score based, decision based, transfer based and non-traditional attacks. This organization assists new readers in comprehending the field and helps current researchers understand where each new attacks fits in the rapidly growing black-box adversarial machine learning literature. Lastly, we offer a new mathematical framework for defining the adversarial threat model. Our new framework provides a convenient and efficient way to quickly determine when attack success rates from different attacks can be compared. Without this framework, we have shown that directly comparing attack success rates from different papers with different threat models can lead to highly misleading conclusions. Overall, our work and comparative evaluations provide insight, organization and systemization to the developing field of adversarial machine learning.

REFERENCES

- [1] Valhid Abolghasemi, Saideh Ferdowsi, Bahador Makkiabadi, and Saeid Sanei. 2010. On optimization of the measurement matrix for compressive sensing. In *2010 18th European Signal Processing Conference*. 427–431.
- [2] Maksym Andriushchenko, Francesco Croce, Nicolas Flammarion, and Matthias Hein. 2020. Square attack: a query-efficient black-box adversarial attack via random search. In *European Conference on Computer Vision*. Springer, 484–501.
- [3] Anish Athalye, Nicholas Carlini, and David Wagner. 2018. Obfuscated Gradients Give a False Sense of Security: Circumventing Defenses to Adversarial Examples. In *Proceedings of the 35th International Conference on Machine Learning (Proceedings of Machine Learning Research)*, Jennifer Dy and Andreas Krause (Eds.), Vol. 80. PMLR, 274–283.
- [4] Siddhant Bhambri, Sumanyu Muku, Avinash Tulasi, and Arun Balaji Buduru. 2019. A survey of black-box adversarial attacks on computer vision models. *arXiv preprint arXiv:1912.01667* (2019).
- [5] Battista Biggio and Fabio Roli. 2018. Wild patterns: Ten years after the rise of adversarial machine learning. *Pattern Recognition* 84 (Dec 2018), 317–331. <https://doi.org/10.1016/j.patcog.2018.07.023>
- [6] Wieland Brendel, Jonas Rauber, and Matthias Bethge. 2017. Decision-based adversarial attacks: Reliable attacks against black-box machine learning models. *arXiv preprint arXiv:1712.04248* (2017).
- [7] Wieland Brendel, Jonas Rauber, and Matthias Bethge. 2018. Decision-Based Adversarial Attacks: Reliable Attacks Against Black-Box Machine Learning Models. In *International Conference on Learning Representations*.
- [8] Nicholas Carlini and David Wagner. 2017. Towards evaluating the robustness of neural networks. In *2017 IEEE Symposium on Security and Privacy (SP)*. IEEE, 39–57.
- [9] Jinghui Chen and Quanquan Gu. 2020. Rays: A ray searching method for hard-label adversarial attack. In *Proceedings of the 26th ACM SIGKDD International Conference on Knowledge Discovery & Data Mining*. 1739–1747.
- [10] Jianbo Chen, Michael I Jordan, and Martin J Wainwright. 2020. Hopskipjumpattack: A query-efficient decision-based attack. In *2020 IEEE Symposium on Security and Privacy (SP)*. IEEE, 1277–1294.
- [11] Pin-Yu Chen, Yash Sharma, Huan Zhang, Jinfeng Yi, and Cho-Jui Hsieh. 2018. Ead: elastic-net attacks to deep neural networks via adversarial examples. In *Thirty-second AAAI conference on artificial intelligence*.
- [12] 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 *Proceedings of the 10th ACM workshop on artificial intelligence and security*. 15–26.
- [13] Shuyu Cheng, Yinpeng Dong, Tianyu Pang, Hang Su, and Jun Zhu. 2019. Improving black-box adversarial attacks with a transfer-based prior. In *Proceedings of the 33rd International Conference on Neural Information Processing Systems*. 10934–10944.
- [14] Shuyu Cheng, Yinpeng Dong, Tianyu Pang, Hang Su, and Jun Zhu. 2019. Improving Black-box Adversarial Attacks with a Transfer-based Prior. (2019). [arXiv:cs.LG/1906.06919](https://arxiv.org/abs/1906.06919)
- [15] Francesco Croce and Matthias Hein. 2019. Sparse and imperceivable adversarial attacks. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*. 4724–4732.
- [16] Francesco Croce and Matthias Hein. 2020. Reliable evaluation of adversarial robustness with an ensemble of diverse parameter-free attacks. In *International conference on machine learning*. PMLR, 2206–2216.
- [17] Jia Deng, Wei Dong, Richard Socher, Li-Jia Li, Kai Li, and Li Fei-Fei. 2009. ImageNet: A large-scale hierarchical image database. In *2009 IEEE Conference on Computer Vision and Pattern Recognition*. 248–255. <https://doi.org/10.1109/CVPR.2009.5206848>
- [18] Yinpeng Dong, Fangzhou Liao, Tianyu Pang, Hang Su, Jun Zhu, Xiaolin Hu, and Jianguo Li. 2018. Boosting adversarial attacks with momentum. In *Proceedings of the IEEE conference on computer vision and pattern recognition (CVPR)*. 9185–9193.
- [19] Jiawei Du, Hu Zhang, Joey Tianyi Zhou, Yi Yang, and Jiashi Feng. 2020. Query-efficient Meta Attack to Deep Neural Networks. In *International Conference on Learning Representations*. <https://openreview.net/forum?id=Skxd6gSYDS>
- [20] Alhussein Fawzi, Seyed-Mohsen Moosavi-Dezfooli, and Pascal Frossard. 2016. Robustness of Classifiers: From Adversarial to Random Noise (NIPS’16). Curran Associates Inc., Red Hook, NY, USA, 1632–1640.

- [21] Leon Gatys, Alexander S Ecker, and Matthias Bethge. 2015. Texture synthesis using convolutional neural networks. *Advances in neural information processing systems* 28 (2015), 262–270.
- [22] Leon A. Gatys, Alexander S. Ecker, and Matthias Bethge. 2016. Image Style Transfer Using Convolutional Neural Networks. In *2016 IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*. 2414–2423. <https://doi.org/10.1109/CVPR.2016.265>
- [23] Ian J Goodfellow, Jonathon Shlens, and Christian Szegedy. 2014. Explaining and harnessing adversarial examples. *arXiv preprint arXiv:1412.6572* (2014).
- [24] Kaiming He, Georgia Gkioxari, Piotr Dollar, and Ross Girshick. 2017. Mask R-CNN. *2017 IEEE International Conference on Computer Vision (ICCV)* (Oct 2017). <https://doi.org/10.1109/iccv.2017.322>
- [25] E. Hoffer and Nir Ailon. 2015. Deep Metric Learning Using Triplet Network. In *SIMBAD*.
- [26] Yu-Lun Hsieh, Minhao Cheng, Da-Cheng Juan, Wei Wei, Wen-Lian Hsu, and Chou-Jui Hsieh. 2019. On the robustness of self-attentive models. In *Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics*. 1520–1529.
- [27] Zhichao Huang and Tong Zhang. 2019. Black-Box Adversarial Attack with Transferable Model-based Embedding. In *International Conference on Learning Representations*.
- [28] Andrew Ilyas, Logan Engstrom, Anish Athalye, and Jessy Lin. 2018. Black-box Adversarial Attacks with Limited Queries and Information. In *ICML*.
- [29] Alexander Kolesnikov, Lucas Beyer, Xiaohua Zhai, Joan Puigcerver, Jessica Yung, Sylvain Gelly, and Neil Houlsby. 2020. Big Transfer (BiT): General Visual Representation Learning. *Lecture Notes in Computer Science* (2020), 491–507. https://doi.org/10.1007/978-3-030-58558-7_29
- [30] Alex Krizhevsky, Vinod Nair, and Geoffrey Hinton. CIFAR-10 (Canadian Institute for Advanced Research). (????). <http://www.cs.toronto.edu/~kriz/cifar.html>
- [31] Y. LECUN. THE MNIST DATABASE of handwritten digits. <http://yann.lecun.com/exdb/mnist/> (????). <https://ci.nii.ac.jp/naid/10027939599/en/>
- [32] Y. LeCun, B. Boser, J. S. Denker, D. Henderson, R. E. Howard, W. Hubbard, and L. D. Jackel. 1989. Backpropagation Applied to Handwritten Zip Code Recognition. *Neural Computation* 1, 4 (1989), 541–551. <https://doi.org/10.1162/neco.1989.1.4.541>
- [33] Huichen Li, Linyi Li, Xiaojun Xu, Xiaolu Zhang, Shuang Yang, and Bo Li. 2021. Nonlinear Projection Based Gradient Estimation for Query Efficient Blackbox Attacks. In *Proceedings of The 24th International Conference on Artificial Intelligence and Statistics (Proceedings of Machine Learning Research)*, Arindam Banerjee and Kenji Fukumizu (Eds.), Vol. 130. PMLR, 3142–3150. <http://proceedings.mlr.press/v130/li21f.html>
- [34] Huichen Li, Linyi Li, Xiaojun Xu, Xiaolu Zhang, Shuang Yang, and Bo Li. 2021. Nonlinear Projection Based Gradient Estimation for Query Efficient Blackbox Attacks. (2021). [arXiv:cs.LG/2102.13184](https://arxiv.org/abs/2102.13184)
- [35] Huichen Li, Xiaojun Xu, Xiaolu Zhang, S. Yang, and B. Li. 2020. QEBA: Query-Efficient Boundary-Based Blackbox Attack. *2020 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)* (2020), 1218–1227.
- [36] Jie Li, Rongrong Ji, Hong Liu, Jianzhuang Liu, Bineng Zhong, Cheng Deng, and Qi Tian. 2020. Projection & probability-driven black-box attack. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*. 362–371.
- [37] Maosen Li, Cheng Deng, Tengjiao Li, Junchi Yan, Xinbo Gao, and Heng Huang. 2020. Towards Transferable Targeted Attack. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*.
- [38] Yujia Liu, Seyed-Mohsen Moosavi-Dezfooli, and Pascal Frossard. 2019. A Geometry-Inspired Decision-Based Attack. In *2019 IEEE/CVF International Conference on Computer Vision (ICCV)*. 4889–4897. <https://doi.org/10.1109/ICCV.2019.00499>
- [39] Yujia Liu, Seyed-Mohsen Moosavi-Dezfooli, and Pascal Frossard. 2019. A geometry-inspired decision-based attack. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*. 4890–4898.
- [40] Aleksander Madry, Aleksandar Makelov, Ludwig Schmidt, Dimitris Tsipras, and Adrian Vladu. 2018. Towards Deep Learning Models Resistant to Adversarial Attacks. In *International Conference on Learning Representations*.
- [41] Kaleel Mahmood, Deniz Gurevin, Marten van Dijk, and Phuong Ha Nguyen. 2020. Beware the Black-Box: on the Robustness of Recent Defenses to Adversarial Examples. *arXiv preprint arXiv:2006.10876* (2020).
- [42] Kaleel Mahmood, Phuong Ha Nguyen, Lam M Nguyen, Thanh Nguyen, and Marten van Dijk. 2019. BUZZ: Buffer Zones for defending adversarial examples in image classification. *arXiv preprint arXiv:1910.02785* (2019).
- [43] Thibault Maho, Teddy Furon, and Erwan Le Merrer. 2021. SurFree: A Fast Surrogate-Free Black-Box Attack. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*. 10430–10439.
- [44] Yurii Nesterov and Vladimir Spokoiny. 2017. Random gradient-free minimization of convex functions. *Foundations of Computational Mathematics* 17, 2 (2017), 527–566.
- [45] Nicolas Papernot, Patrick McDaniel, and Ian Goodfellow. 2016. Transferability in machine learning: from phenomena to black-box attacks using adversarial samples. *arXiv preprint arXiv:1605.07277* (2016).
- [46] Nicolas Papernot, Patrick McDaniel, Ian Goodfellow, Somesh Jha, Z Berkay Celik, and Ananthram Swami. 2017. Practical black-box attacks against machine learning. In *Proceedings of the 2017 ACM on Asia conference on computer and communications security*. 506–519.
- [47] Ali Rahmati, Seyed-Mohsen Moosavi-Dezfooli, Pascal Frossard, and Huaiyu Dai. 2020. GeoDA: a geometric framework for black-box adversarial attacks. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*. 8446–8455.
- [48] Andrianiaina Ravelomanantsoa, Hassan Rabah, and Amar Rouane. 2015. Compressed Sensing: A Simple Deterministic Measurement Matrix and a Fast Recovery Algorithm. *IEEE Transactions on Instrumentation and Measurement* 64, 12 (2015), 3405–3413. <https://doi.org/10.1109/TIM.2015.2459471>
- [49] Joseph Redmon, Santosh Divvala, Ross Girshick, and Ali Farhadi. 2016. You Only Look Once: Unified, Real-Time Object Detection. *2016 IEEE Conference on Computer Vision and Pattern Recognition (CVPR)* (Jun 2016). <https://doi.org/10.1109/cvpr.2016.91>
- [50] Daniel L. Ruderman, Thomas W. Cronin, and Chuan-Chin Chiao. 1998. Statistics of cone responses to natural images: implications for visual coding. *J. Opt. Soc. Am. A* 15, 8 (Aug 1998), 2036–2045. <https://doi.org/10.1364/JOSAA.15.002036>
- [51] Ramprasaath R Selvaraju, Michael Cogswell, Abhishek Das, Ramakrishna Vedantam, Devi Parikh, and Dhruv Batra. 2017. Grad-cam: Visual explanations from deep networks via gradient-based localization. In *Proceedings of the IEEE international conference on computer vision*. 618–626.
- [52] A. Shahin Shamsabadi, R. Sanchez-Matilla, and A. Cavallaro. 2020. ColorFool: Semantic Adversarial Colorization. In *2020 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*. IEEE Computer Society, Los Alamitos, CA, USA, 1148–1157. <https://doi.org/10.1109/CVPR42600.2020.00123>
- [53] Christian Szegedy, Wojciech Zaremba, Ilya Sutskever, Joan Bruna, Dumitru Erhan, Ian Goodfellow, and Rob Fergus. 2013. Intriguing properties of neural networks. (2013). [arXiv:cs.CV/1312.6199](https://arxiv.org/abs/1312.6199)
- [54] Chenglin Yang, Adam Kortylewski, Cihang Xie, Yinzhi Cao, and Alan Yuille. 2020. Patchattack: A black-box texture-based attack with reinforcement learning. In *European Conference on Computer Vision*. Springer, 681–698.
- [55] Richard Zhang, Phillip Isola, and Alexei A. Efros. 2016. Colorful Image Colorization. In *ECCV*.
- [56] Pu Zhao, Pin-Yu Chen, Siyue Wang, and Xue Lin. 2020. Towards query-efficient black-box adversary with zeroth-order natural gradient descent. In *Proceedings of the AAAI Conference on Artificial Intelligence*, Vol. 34. 6909–6916.
- [57] Pu Zhao, Sijia Liu, Pin-Yu Chen, Nghia Hoang, Kaidi Xu, Bhavya Kaikhura, and Xue Lin. 2019. On the design of black-box adversarial examples by leveraging gradient-free optimization and operator splitting method. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*. 121–130.
- [58] Mingyi Zhou, J. Wu, Y. Liu, Shuaicheng Liu, and Ce Zhu. 2020. DaST: Data-Free Substitute Training for Adversarial Attacks. *2020 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)* (2020), 231–240.

Attack Name	ASR	Avg Queries	Norm	Target	Classifier	Adv Threat Model	Source
NonLinear-BA	80	10000	$l_2, \epsilon=0.001$	T	ResNet18	$w_c = 1, w_q = 1, w_\epsilon = 1$	[34]
SurFree	90	500	$l_2, \epsilon=30$	T	ResNet18	$w_c = 1, w_q = 1, w_\epsilon = 1$	[43]
RayS	99.8	574	$l_\infty, \epsilon=0.05$	U	ResNet50	$w_c = 1, w_q = 0, w_\epsilon = 1$	[9]
RayS	98.9	748.2	$l_\infty, \epsilon=0.05$	U	Inception v3	$w_c = 1, w_q = 0, w_\epsilon = 1$	[9]
QEBA-S	82	10000	$l_2, \epsilon=0.001$	T	ResNet18	$w_c = 1, w_q = 1, w_\epsilon = 1$	[34]
QEBA-S	74	10000	$l_2, \epsilon=0.001$	T	ResNet18	$w_c = 1, w_q = 1, w_\epsilon = 1$	[35]
QEBA	71	500	$l_2, \epsilon=30$	T	ResNet18	$w_c = 1, w_q = 1, w_\epsilon = 1$	[43]
GeoDA	-	1000	$l_2=8.16$ (median)	U	ResNet50	$w_c = 0, w_q = 1, w_\epsilon = 0$	[47]
GeoDA	79	500	$l_2, \epsilon=30$	T	ResNet18	$w_c = 1, w_q = 1, w_\epsilon = 1$	[43]
HSJA	19.9	749.6	$l_\infty, \epsilon=0.05$	U	ResNet50	$w_c = 1, w_q = 0, w_\epsilon = 1$	[9]
HSJA	23.7	652.3	$l_\infty, \epsilon=0.05$	U	Inception v3	$w_c = 1, w_q = 0, w_\epsilon = 1$	[9]
HSJA	80	17000	$l_2, \epsilon=0.001$	T	ResNet18	$w_c = 1, w_q = 1, w_\epsilon = 1$	[34]
HSJA	84	20000	$l_2, \epsilon=0.001$	T	ResNet18	$w_c = 1, w_q = 1, w_\epsilon = 1$	[35]
HSJA	56	500	$l_2, \epsilon=30$	T	ResNet18	$w_c = 1, w_q = 1, w_\epsilon = 1$	[43]
qFool	-	1000	$l_2=16.05$ (median)	U	ResNet50	$w_c = 0, w_q = 1, w_\epsilon = 0$	[47]

Table 4: Decision based black-box attacks on ImageNet CNNs. The corresponding success rate (ASR) and adversarial threat model are shown for each attack along with the original source paper.

Score Based Attacks						
Attack Name	ASR	Avg Queries	Adv Threat Model	Norm	Target	Source
ZO-NGD	98.7	523	$w_c = 1, w_q = 0, w_\epsilon = 1$	$l_\infty, \epsilon=0.2$	U	[56]
TREMBBA	98	1064	$w_c = 1, w_q = 0, w_\epsilon = 1$	$l_\infty, \epsilon=0.2$	U	[27]
ZO-ADMM	98.3	-	$w_c = 1, w_q = 0, w_\epsilon = 0$	$l_2=1.975$ (avg)	T	[57]
P-RGF	68.53	16135	$w_c = 1, w_q = 0, w_\epsilon = 1$	$l_\infty, \epsilon=0.2$	U	[27]
qMeta	100	749	$w_c = 1, w_q = 1, w_\epsilon = 0$	l_2	U	[19]
qMeta	100	1299	$w_c = 1, w_q = 1, w_\epsilon = 0$	l_2	T	[19]

Decision Based Attacks						
Attack Name	ASR	Avg Queries	Adv Threat Model	Norm	Target	Source
NonLinear-BA	90	5000	$w_c = 1, w_q = 1, w_\epsilon = 1$	$l_2, \epsilon=0.005$	T	[34]
RayS	100	107	$w_c = 1, w_q = 0, w_\epsilon = 1$	$l_\infty, \epsilon=0.3$	U	[9]
QEBA-S	87	5000	$w_c = 1, w_q = 1, w_\epsilon = 1$	$l_2, \epsilon=0.005$	T	[34]
HSJA	91.2	161.6	$w_c = 1, w_q = 0, w_\epsilon = 1$	$l_\infty, \epsilon=0.3$	U	[9]

Table 5: Decision and score based black-box attacks on MNIST CNNs. The corresponding success rate (ASR) and adversarial threat model are shown for each attack along with the original source paper.

Score Based Attacks						
Attack Name	ASR	Avg Queries	Adv Threat Model	Norm	Target	Source
ZO-NGD	99.2	131	$w_c = 1, w_q = 0, w_\epsilon = 1$	$l_\infty, \epsilon=0.1$	U	[56]
ZO-ADMM	98.7	-	$w_c = 1, w_q = 0, w_\epsilon = 0$	$l_2=0.417$ (avg)	T	[57]
qMeta	92	1765	$w_c = 1, w_q = 1, w_\epsilon = 0$	l_2	U	[19]
qMeta	93	3667	$w_c = 1, w_q = 1, w_\epsilon = 0$	l_2	T	[19]

Decision Based Attacks						
Attack Name	ASR	Avg Queries	Adv Threat Model	Norm	Target	Source
NonLinear-BA	95.00	5000.00	$w_c = 1, w_q = 1, w_\epsilon = 1$	$l_2, \epsilon=0.0001$	T	[34]
RayS	99.8	792.8	$w_c = 1, w_q = 0, w_\epsilon = 1$	$l_\infty, \epsilon=0.031$	U	[9]
QEBA-S	95	5000.00	$w_c = 1, w_q = 1, w_\epsilon = 1$	$l_2, \epsilon=0.0001$	T	[34]
HSJA	99.7	1021.6	$w_c = 1, w_q = 0, w_\epsilon = 1$	$l_\infty, \epsilon=0.031$	U	[9]

Table 6: Decision and score based black-box attacks on CIFAR-10 CNNs. The corresponding success rate (ASR) and adversarial threat model are shown for each attack along with the original source paper.

Transfer Based Attacks

Attack Name	Dataset	ASR	Avg Queries	Adv Threat Model	Target	Vanilla Model	Source
Adaptive	CIFAR-10	74.1	-	$w_c = 1, w_q = 0, w_\epsilon = 1$	U	ResNet56	[41]
Adaptive	CIFAR-10	22.3	-	$w_c = 1, w_q = 0, w_\epsilon = 1$	T	ResNet56	[41]
DaST	MNIST	29.18	-	$w_c = 1, w_q = 0, w_\epsilon = 1$	U	CNN	[58]
DaST	MNIST	57.22	-	$w_c = 1, w_q = 0, w_\epsilon = 1$	U	CNN	[58]
DaST	MNIST	64.61	-	$w_c = 1, w_q = 0, w_\epsilon = 1$	U	CNN	[58]
DaST	MNIST	96.36	-	$w_c = 1, w_q = 0, w_\epsilon = 1$	U	CNN	[58]
DaST	CIFAR-10	19.78	-	$w_c = 1, w_q = 0, w_\epsilon = 1$	T	VGG-16	[58]
DaST	CIFAR-10	20.22	-	$w_c = 1, w_q = 0, w_\epsilon = 1$	T	VGG-16	[58]
DaST	CIFAR-10	28.42	-	$w_c = 1, w_q = 0, w_\epsilon = 1$	U	VGG-16	[58]
DaST	CIFAR-10	59.71	-	$w_c = 1, w_q = 0, w_\epsilon = 1$	U	VGG-16	[58]
Po+TI (Trip)	ImageNet	39.5	-	$w_c = 1, w_q = 0, w_\epsilon = 1$	T	Inception v3	[37]
Po+TI (Trip)	ImageNet	39.3	-	$w_c = 1, w_q = 0, w_\epsilon = 1$	T	ResNet50	[37]

Non-traditional Attacks

Attack Name	Dataset	ASR	Avg Queries	Adv Threat Model	Target	Vanilla Model	Source
CornerSearch	MNIST	97.38	-	$w_c = 1, w_q = 0, w_\epsilon = 0$	U	NiN	[15]
CornerSearch	CIFAR-10	99.56	-	$w_c = 1, w_q = 0, w_\epsilon = 0$	U	NiN	[15]
ColorFool	CIFAR-10	41.5	-	$w_c = 1, w_q = 0, w_\epsilon = 0$	U	ResNet50	[52]
ColorFool	ImageNet	22.3	-	$w_c = 1, w_q = 0, w_\epsilon = 0$	U	ResNet50	[52]
Patch Attack N4 4%	ImageNet	99.7	1137	$w_c = 1, w_q = 0, w_\epsilon = 0$	U	ResNet50	[54]
Patch Attack N8 2%	ImageNet	99.7	983	$w_c = 1, w_q = 0, w_\epsilon = 0$	U	ResNet50	[54]
Patch Attack N10 4%	ImageNet	99.7	8643	$w_c = 1, w_q = 0, w_\epsilon = 0$	T	ResNet50	[54]
Patch Attack N10 10%	ImageNet	100	3747	$w_c = 1, w_q = 0, w_\epsilon = 0$	T	ResNet50	[54]

Table 7: Transfer based and non-traditional attacks on various datasets (MNIST, CIFAR-10 and ImageNet). NiN stands for Network-in-Network.