DATA-DRIVEN LIPSCHITZ CONTINUITY: A COST EFFECTIVE APPROACH TO IMPROVE ADVERSARIAL ROBUSTNESS

Anonymous authors

Paper under double-blind review

Abstract

The security and robustness of deep neural networks (DNNs) have become increasingly concerning. This paper aims to provide both a theoretical foundation and a practical solution to ensure the reliability of DNNs. We explore the concept of Lipschitz continuity to certify the robustness of DNNs against adversarial attacks, which aim to mislead the network with adding imperceptible perturbations into inputs. We propose a novel algorithm that remaps the input domain into a constrained range, reducing the Lipschitz constant and potentially enhancing robustness. Unlike existing adversarially trained models, where robustness is enhanced by introducing additional examples from other datasets or generative models, our method is almost cost-free as it can be integrated with existing models without requiring re-training. Experimental results demonstrate the generalizability of our method, as it can be combined with various models and achieve enhancements in robustness. Furthermore, our method achieves the best robust accuracy for CIFAR10, CIFAR100, and ImageNet datasets on the RobustBench leaderboard.

026 027 028

029

024

025

006

008 009 010

011

013

014

015

016

017

018

019

021

1 INTRODUCTION

Deep neural networks (DNNs) have demonstrated promising results across various tasks (Krizhevsky et al., 2012; Redmon et al., 2016), prompting concerns about AI security as these networks are increasingly deployed in our daily lives. A single erroneous prediction could lead to catastrophic consequences. For example, the Overload attack can significantly inflate the inference time of detecting objects for self-driving systems (Chen et al., 2023), while even minor typos in input prompts can cause large language models to produce unexpected responses (Zhu et al., 2023).

The focus of this paper is to design robust DNNs that can defend against adversarial attacks, which
 aim to create perturbations in inputs that are imperceptible to humans but can mislead DNNs. Pre vious studies have revealed the existence of adversarial examples in diverse domains, such as image
 pixels (Szegedy et al., 2013), audio data (Carlini & Wagner, 2018), and textual content (Li et al.,
 2018). Consequently, exploring the vulnerabilities of DNNs and developing theoretically grounded
 explainable AI is crucial for ensuring the reliability of DNN-based applications.

- 042 Adversarial training (Madry et al., 2017) has proven to be an effective strategy for enhancing the 043 robustness of DNNs. It achieves this by generating adversarial examples on the fly during the train-044 ing phase and optimizing the model's weights to minimize the losses caused by these examples. Recent studies have shown that robustness can be further improved by introducing additional examples from other datasets (Carmon et al., 2019) or using generative models (Gowal et al., 2021; 046 Wang et al., 2023) to cover low-frequency data. Despite the promising improvements in robustness, 047 training costs increase significantly due to the demand for additional data, which can be up to 20 to 048 100 times larger than the original dataset. This poses a trade-off between training cost and robust-049 ness. The concern over high computational costs becomes a significant obstacle in deploying robust 050 DNN-based applications, especially in fields like medicine, autonomous driving systems, and other 051 areas where human lives are at stake. 052
- ⁰⁵³ In this paper, we explore how robustness is certified by the theorem of Lipschitz continuity, which theoretically gauges how much outputs are amplified by the perturbations. However, we argue that



Figure 1: The empirical Lipschitz constant of specific layers that can be represented by linear systems, such as convolutional or fully connected layers, can be reduced by remapping their input domain to a constrained range.

the set of observed data is finite and cannot cover the entire real data space, leading to an overestimation of the Lipschitz constant derived from the theorem. Therefore, we propose an algorithm
that can remap the input domain into a constrained range, resulting in a Lipschitz constant for the
modified function that is less than or equal to the Lipschitz constant of the original function, thus
potentially enhancing robustness. Our key contributions are outlined as follows:

- We introduce the concept of the empirical Lipschitz constant, which can more precisely reflect the robustness of the corresponding observed data. Compared with the original definition of Lipschitz constant, the empirical value is derived from a set of observed data, thereby eliminating the influence of space that is never drawn from real data. As illustrated in Figure 1, we prove that any function that can be formulated as a linear system, when combined with our proposed function to remap the input domain of a specific layer to a constrained range, can reduce its empirical Lipschitz constant, resulting in better robustness.
- Our method is well-suited for inference-time model optimization to enhance robustness cost-effectively. The proposed function can enhance the robustness of adversarially trained models with minimal additional costs. Specifically, it introduces only one parameter, the value of which can be determined by scanning an observed data once without the need for re-training or fine-tuning.
- The experimental results suggest that our method can be combined with various existing methods and gain robustness improvements. Besides, our method achieves the best robust accuracy against adversarial examples generated by AutoAttack (Croce & Hein, 2020), a state-of-the-art ensemble attack, for CIFAR10, CIFAR100, and ImageNet datasets on the RobustBench leaderboard (Croce et al., 2020). By assessing accuracy against adaptive attacks, transfer attacks, and evaluation methods for validating obfuscated gradients (Athalye et al., 2018; Carlini et al., 2019), we believe that the proposed algorithm should not cause robustness to be overestimated.

The rest of this paper is organized as follows. Section 2 introduces the background on adversarial attacks and adversarial training. Section 3 presents the theoretical proof of how robustness is enhanced by manipulating the domain of linear functions and introduces the proposed algorithm. Section 4 shows the experimental resultsand ablation studies on various hyper-parameters, combination with different activation functions and gradient masking verification. The last section is our conclusion.

- 2 RELATED WORKS

2.1 ADVERSARIAL ATTACKS

Adversarial attacks aim to inject tiny perturbations into inputs, causing victim DNNs to output in correct predictions with high confidence (Chen & Hsieh, 2022). These attacks have been observed in numerous vision applications (Goodfellow et al., 2014; Chen et al., 2018; Wang et al., 2022; Yin

et al., 2022). Furthermore, these tiny perturbations can be embedded not only in image pixels but also in textual contexts (Kumar et al., 2023; Yao et al., 2023), audio space (Xie et al., 2021), and other fields (Ilahi et al., 2021). Some research has shown how adversarial attacks threaten real applications (Xu et al., 2020; Komkov & Petiushko, 2021; Du et al., 2022; Wei et al., 2022). Investigating the vulnerability of DNNs and theoretically avoiding adversarial examples when optimizing model weights or designing architecture is an ongoing challenge.

Adversarial attacks can be classified into two types based on the amount of information the attacker has access to: white-box attacks and black-box attacks. In the case of white-box attacks, the attacker has full access to all information about the victim model. Methods such as the PGD attack (Madry et al., 2017) and AutoAttack (Croce & Hein, 2020) generate adversarial examples by leveraging the gradient direction of the model. Although this scenario is often unrealistic in practical settings, research in this area is valuable for developing more robust models in the future.

In contrast, black-box attacks, such as the square attack (Andriushchenko et al., 2020) and ZO-NGD
 (Zhao et al., 2020), only provide the attacker with access to the model's output predictions. In these
 cases, adversarial examples can still be generated through transferability, where models with similar
 architectures are used to create adversarial examples (Wang et al., 2021; Chen et al., 2024). The
 goal of black-box attacks is to investigate the potential risks of adversarial attacks in real-world
 application scenarios, where attackers may have limited access to model internals.

126

128

127 2.2 DEFENSIVE STRATEGIES

Adversarial training is a defensive strategy that aims to find optimal weights against adversarial attacks. It achieves this by generating adversarial examples on the fly during the training phase and optimizing the model's weights to minimize the losses caused by these examples. Despite the superior robustness achieved by adversarial training, the associated training costs of adversarially trained models are generally ten times more expensive than those of models trained utilizing a standard policy. The concern over high computational costs becomes a significant obstacle in deploying DNN-based applications.

Balancing between training cost and robustness is a challenge for adversarial training. Fast adver-136 sarial training has been proposed for applications pursuing higher robustness under a limited budget 137 (Chen & Lee, 2020; Zhang et al., 2022). However, numerous adversarial examples cannot be drawn 138 from these approaches, potentially leading to catastrophic overfitting, where robust accuracy signif-139 icantly decreases without warning signs (Rice et al., 2020). On the contrary, some studies attempted 140 to refine robustness by introducing additional examples from other datasets (Carmon et al., 2019) 141 or using generative models (Gowal et al., 2021; Wang et al., 2023). Alternatively, another line of 142 research has demonstrated that the removal of partial adversarial examples does not compromise ro-143 bust accuracy, addressing the issue of unaffordable training costs (Zhang et al., 2020; Chen & Lee, 144 2024).

Despite the potential of adversarial training to enhance model robustness, budgetary constraints often limit the scope of their crafting to one or two specific attack types during the training stage. This restricted approach may inadvertently render adversarially trained models susceptible to novel, unseen attacks. As an alternative, Lipschitz-based certified training offers a theoretical framework for ensuring an upper bound on prediction errors (Gowal et al., 2018; Huang et al., 2021; Müller et al., 2022). However, it is important to acknowledge that these training methods often suffer from scalability issues.

152

153 3 METHODOLOGY

155 3.1 MOTIVATION

In this paper, we approach robustness from a theoretical perspective, aiming to demonstrate that all risks posed by adversarial examples are limited while minimizing the additional costs associated with improving robustness. Our evaluation is conducted under the white-box scenario, where the target model is capable of defending against various types of known adversarial attacks, including white-box attacks (Madry et al., 2017; Chen et al., 2018; Croce & Hein, 2020), black-box attacks (Chen et al., 2017), and transfer attacks (Demontis et al., 2019; Qin et al., 2022). Additionally, we

conduct a set of experiments to verify that gradient masking (Athalye et al., 2018) does not occur in 163 our method and to ensure that robustness is not overestimated.

164 165

166

171 172

173

174 175

176

192 193 3.2 LIPSCHITZ CONTINUITY

167 To achieve our goal, we introduce a quantitative metric known as the Lipschitz constant, which gauges how much outputs are amplified by the perturbations within the input domain. The mathe-168 matical definition is as follows, a function $f: \mathbb{R}^m \to \mathbb{R}^n$ is globally Lipschitz continuous if there 169 exists an constant K > 0 such that 170

$$D_f(f(x_1), f(x_2)) \le K D_x(x_1, x_2) \quad \forall x_1, x_2 \in \mathbb{R}^m,$$

$$\tag{1}$$

where D_x is a metric on the domain of f; D_f is a metric on the range of f; and $x_1 \neq x_2$. For a DNN, it can be considered as a composite function:

$$F(x) = (f_1 \circ f_2 \circ \dots \circ f_L)(x), \tag{2}$$

where f_i is the function of *i*-th layer. If there exists a Lipschitz constant for each individual layer, we can derive an upper bound of the Lipschitz constant for the victim model as follows,

$$K_F \le \prod_{i=1}^{L} K_i,\tag{3}$$

181 where K_i is the Lipschitz constant of f_i .

182 By defining adversarial examples x^{adv} within a ϵ -ball centered at an image x as the inputs of (1), we 183 can assess the impact caused by adversarial examples. Therefore, the Lipschitz constant serves as 184 a bridge that connects the design of robust models with the measurement of risks posed by adver-185 sarial examples. A small Lipschitz constant for the victim model implies that the increase in loss is minimal, indicating a higher ability to resist adversarial attacks. Consequently, the objective of this 187 paper is to lower the upper bound of Lipschitz constant for the given models. 188

As indicated by previous studies (Yoshida & Miyato, 2017; Farnia et al., 2018), Lipschitz constant 189 of the given model defined in (3) can be minimized by reducing the output discrepancy of individual 190 linear layers. Under the L_2 norm, we have 191

$$\frac{||f(x^{\mathrm{adv}}) - f(x)||_2}{||x^{\mathrm{adv}} - x||_2} = \frac{(||Wx^{\mathrm{adv}} + b) - (Wx + b)||_2}{||\delta||_2} = \frac{||W\delta||_2}{||\delta||_2},\tag{4}$$

194 where W is the weight matrix; and δ is the distance between x^{adv} and x. Therefore, the original op-195 timization problem of minimizing Lipschitz constant is transformed into the following minimization 196 problem: 197

$$\min_{W} \max_{\delta \neq 0, \delta \in \mathbb{R}^m} \frac{||W\delta||_2}{||\delta||_2} = \min_{W} \sigma_{\max}(W), \tag{5}$$

199 where $\sigma_{\max}(W)$ represents the largest singular value of the matrix W. Notably, there is a relation to 200 eigenvalues: 201

$$\sigma_i^2(W) = \lambda_i(WW^{\dagger}) = \lambda_i(W^{\dagger}W), \tag{6}$$

202 where W^{\dagger} is the conjugate transpose of W. Each singular value of the matrix W is the square root 203 of the eigenvalue of the matrices WW^{\dagger} or $W^{\dagger}W$. In other words, minimizing $\lambda_{max}(WW^{\dagger})$, the 204 largest eigenvalue of the matrices, can achieve the same objective. 205

Rather that minimizing the objective directly, Gershgorin circle theorem provides an alternative 206 solution to estimate the robustness of the given linear system. 207

Theorem 1. (Gershgorin Circle Theorem) For an $m \times m$ matrix A with entries a_{ij} , each eigenvalue 208 of A is in at least one of the disk: 209

$$R_i = \{ z \in \mathbb{C} : |z - a_{ii}| \le \sum_{i \ne j} |a_{ij}| \} \text{ for } i = \{1, 2, \dots, m\}.$$
(7)

211 212

210

213 Theorem 1 indicates each row vector can be represented as a disk which is centered at the diagonal entry a_{ii} and whose radius is the sum of the off-diagonal entries a_{ij} . For any layer which can be 214 represented by a linear system, such as convolutional or fully connected layers, robustness can be 215 improved by shrinking the radius of the disk with the largest eigenvalue.



Figure 2: Insertion points of the forged function. In ConvNets, it is inserted into the residual blocks, while in Transformers, it is inserted into the MLP layers.

3.3 FORGED FUNCTION

232

233 234 235

237

238

239

241 242

243

251

253

267

268

We argue that the largest singular value provides a loose bound for the Lipschitz constant. To precisely reflect the robustness of the corresponding observed data, we define the empirical Lipschitz constant that eliminates the influence of space that is never drawn from real data.

Definition 1. Empirical Lipschitz constant:

$$\max_{\delta \neq 0, x \in \mathcal{S}} \frac{||Wx||_2}{||x||_2} \quad \forall x \in \mathcal{S},$$
(8)

where S is an observed dataset. As can be seen, the empirical Lipschitz constant on the finite dataset is less than or equal to its Lipschitz constant derived from the theorem.

Based on Definition 1, we can build robust models by manipulating the output ranges of individual layers, thereby restricting the input domain of the next layer. If input vectors do not align with the direction of the eigenvector with the largest eigenvalue, the empirical constant should be bounded. Therefore, we proposed a forged function defined as follows:

$$f^{\text{forge}}(x) = \begin{cases} 0 & \text{if } |x| \le c_i^{\text{th}}, \\ x & \text{otherwise,} \end{cases}$$
(9)

where c_i^{th} is a threshold for the *i*-th layer. Compared with the original functions, the range of the forged function is suppressed if its value is less than the threshold. When c_i^{th} is set to 0, the forged function degrades into the original function.

The forged function aim to reduce the empirical Lipschitz constant of the layers that can be represented as linear systems by remapping the input domain of these layers into a constrained set. Figure provides a visual representation of potential insertion points for the forged function, while maintaining the integrity of other layers. For the ResNet architecture, the forged function is placed before the convolutional layers in each residual block. Similarly, for vision transformer architectures, the structure of MLP layers is adapted to seamlessly integrate the forged function.

Here is the proof that the largest eigenvalue can be shrunk by the forged function. Let W be the weight of the target layer, which can be represented by an $m \times n$ matrix, and **t** be the input vector. Without loss of generality, we assume that $A = W^{\dagger}W$ and $f^{\text{forge}}(\mathbf{t})$ is defined as:

$$f^{\text{forge}}(t_i) = \begin{cases} 0 & i \le k \\ t_i & \text{otherwise,} \end{cases}$$
(10)

where t_i is the *i*-th element of **t** and *k* is a positive number.

Lemma 2. There exists a matrix A' whose largest eigenvalue, $\lambda_{max}(A')$, is less than or equal to the largest eigenvalue of A, $\lambda_{max}(A)$, if **Lemma 2.** There exists a matrix A' whose largest eigenvalue, $\lambda_{max}(A')$, is less than or equal to the largest eigenvalue of A, $\lambda_{max}(A)$, if (11)

$$Af^{forge}(t) = A't. \tag{11}$$

274 *Proof.* Since the first k entries of the vector \mathbf{t} are replaced with zeros, above condition can be 275 achieved by replacing the corresponding column vectors of the matrix A with zero vectors. There-276 fore, the entries of A' are formulated as

$$a_{ij}' = \begin{cases} 0 & j \le k \\ a_{ij} & \text{otherwise.} \end{cases}$$
(12)

The matrix A is a positive semidefinite matrix, implying that the diagonal entries are non-negative. Moreover, with the entry representation of A' in (12), we observe that modifications are only applied to the first k columns, while the rest remain unchanged. Combining the Gershgorin Circle Theorem, we know that the centers of the first k disks of the matrix A' are shifted towards zero. Additionally, the radii of all disks, the absolute values of the off-diagonal entries in A', are shrunk. Consequently, the upper bound of the largest eigenvalue of the matrix A' is tighter compared to that of the original matrix A.

Notably, the outputs of $f^{\text{forge}}(\mathbf{t})$ vary depending on the inputs, resulting in each input having its own A'. The upper bound of the largest eigenvalue of each matrix A' must be not greater than the largest eigenvalue of the matrix A. With Lamma 2, a precise upper bound of the largest eigenvalue can be obtained by feeding a set of observed images. On the contrary, there might be cases in which solving the minimization problem in (5) leads to the theoretical largest eigenvalue being minimized, but the empirical Lipschitz constant remains unchanged.

The choice of a proper c_i^{th} is a crucial factor in reducing the largest eigenvalue. In this paper, we propose obtaining the value of c_i^{th} through the following equation:

$$c_i^{\text{th}} = c^r \max(F_{1 \to i}(x)) \quad \forall x \in \mathcal{S},$$
(13)

where S can include all or a subset of images in the training set, c^r is a positive number and $F_{1\to i}(x)$ 298 represents the output of the *i*-th layer. Specifically, each layer has its own $c_i^{\rm th}$, but they share the 299 same hyper-parameter c^r . Algorithm 1 specifies the implementation details of the forged function. 300 The variable b is used to store the maximum value that appeared in S, as defined in (13), and is 301 initialized during construction. Similar to the implementation of the batchnorm layer, the behavior 302 is depended on the mode configuration. When the mode is set to tracking mode, the variable b is 303 updated accordingly, and the input is set to the output without any modification. Conversely, when 304 the mode is set to inference mode, the value of b is frozen, but the input is updated as defined by (9). 305 By default, the mode is set to inference, and the values of b and c^r are zero, respectively. As a result, 306 the set \mathcal{M} is empty, and the algorithm is degraded to the identical function.

307 It is worth emphasizing that by feeding all images in the set S once in track mode beforehand, the 308 value of $c_i^{\rm th}$ can be obtained. The elements satisfying the constraints are appropriately deactivated 309 during inference. Notably, this operation does not necessitate gradient computations and incurs 310 minimal time consumption, typically only a few minutes, even when executed on commonly used 311 GPUs. In comparison to adversarial training, this process is nearly cost-free. The overall procedure 312 shares many similarities with post-pruning techniques. Nevertheless, we posit that the proposed 313 function is very similar to the ReLU function, as it suppresses the output values within a specific range, but the defined range in the forged function is adaptive to the observed dataset. 314

315 316

317

273

278

279

287

296 297

4 EXPERIMENTS

318 4.1 SETUP

We evaluated the performance on CIFAR10, CIFAR100, and ImageNet datasets under the whitebox scenario with an L_{∞} norm. To ensure comparability of results, we assessed robustness using AutoAttack (Croce & Hein, 2020), For CIFAR10 and CIFAR100 datasets, ϵ is set to 8/255, while for the ImageNet dataset, ϵ is set to 4/255. The model weights are publicly accessible from Robust-Bench. The ablation study involves exploring the selection of the optimal c^r , the combination of

Alg	orithm 1 Forged Function
1:	require : Input x , Mode m , Hyper-parameter c^r
2:	if m is tracking mode then
3:	$b = \max(b, \mathbf{x})$
4:	else
5:	$\mathcal{M} = \{x abs(x) \le c^r b\}$
6:	for all $s \in \mathcal{M}$ do
7:	s = 0
8:	end for
9:	end if
10:	return x

Table 1: The results of top-3 competitors on Robustbench.

	(a) CIFAR10 data	set	(b) CIFAR100 dataset					
#	Method	acc _{nat}	acc _{AA}	#	Method	acc _{nat}	acc _{AA}	
* 1 2 3	Wang et al. (2023) + Ours Peng et al. (2023) Wang et al. (2023) Bai et al. (2024)	93.20 93.27 93.25 95.19	71.70 71.07 70.69 69.71	* 1 2 3	Wang et al. (2023) + Ours Wang et al. (2023) Bai et al. (2024) Cui et al. (2023)	74.97 75.22 83.08 73.85	44.00 42.67 41.80 39.18	

various models trained from different techniques, the verification of gradient masking, and assessments of certified adversarial robustness via randomized smoothing. Due to the page limit, the full experimental results of the ablation study are listed in Appendix A.

4.2 WHITE-BOX EVALUATION

4.2.1 PERFORMANCE ANALYSIS ON CIFAR10 AND CIFAR100 DATASETS

The model used in this study is based on WRN-70-16 architecture with SiLU function while generative data were involved during the training phase. The value of c_i^{th} is obtained by feeding all images from the training set without any augmentation, and c^r was set to 2^{-8} for this experiment. Tables la and 1b summarize the top-3 competitors on Robustbench for CIFAR10 and CIFAR100 datasets, respectively, where # represent the rankings, our results are marked by the asterisk (*), acc_{nat} and acc_{AA} denote the accuracy against clean data and adversairal examples generated by AutoAttack, respectively.

As can be seen, our method combined with WRN-70-16 with SiLU function gains improvement in robustness by at least 0.9% and achieves the best results on Robustbench for both datasets. Nevertheless, standard accuracy (acc_{nat}) is decreased. Many factors might affect the results. For example, the single additional hyper-parameter introduced in this study might not provide sufficient granularity to fit all layers in the target model.

365

335 336

347

348 349 350

351

352

4.2.2 PERFORMANCE ANALYSIS ON IMAGENET DATASET

368 In this experiment, we utilized the Swin (Liu et al., 2021) model architecture, a variant of transform-369 ers. However, scanning the approximately 1.2 million training images provided by the ImageNet dataset to determine the value of the hyper-parameter introduced in the forged function defined in 370 (13) might take a long time. Alternatively, we randomly selected about 5,000 images as the observed 371 images to determine the value of the hyper-parameter. Ideally, determining the optimal choice of c^r 372 requires conducting an ablation study to explore the relationship between the chosen c^r and robust 373 accuracy on a validation set. To accelerate this procedure, we first seek a value of c^r with the highest 374 standard accuracy. The candidate values are selected in a small range centered around this value. 375

Tables 2 lists the top-3 competitors on Robustbench for ImageNet dataset, including ranking, architecture, standard accuracy, and robust accuracy against AutoAttack. The experimental results demonstrate that the Swin-L model with GELU combined with our method can obtain improve-

379 380

381

382

383

384

385 386 387

388

389 390

#	Method	Architecutre	acc _{nat}	acc _{AA}
* Li	u et al. (2023) + Ours	Swin-L	78.88	60.04
1	Liu et al. (2023)	Swin-L	78.92	59.56

Table 2: The results of top-3 competitors for ImageNet dataset on Robustbench.

ConvNeXtV2-L + Swin-L

ConvNeXt-L

58.50

58.48

81.48

78.02

ments in robust accuracy and achieve the best result while standard accuracy has a tiny drop. This finding verifies that our method can be applied to both convolutional and fully connected layers.

4.2.3 COMBINATION WITH VARIOUS MODELS

Bai et al. (2024)

Liu et al. (2023)

2

3

391 The experiment aims to assess the generability of the proposed function on adversarial trained mod-392 els with identical architecture but from various training strategies and to evaluate the potential cost 393 reduction of adversarial training. We integrated the proposed approach with partial models selected 394 from RobustBench, whose weights are obtained directly from the official without any modifications, and also included a model trained by TRADES (Zhang et al., 2019) as a baseline for CIFAR10 396 dataset. The selected models were trained using different techniques, such as adding perturbations 397 in internal layers, retrieving information using knowledge distillation, reducing inefficient train-398 ing data, or involving additional images from generated models or another dataset. Except for the 399 model used in RST-WAP, which is WRN-28-10, the model architecture we utilized is WRN-34-10 400 with ReLU, as it is the most popular network on the RobustBench leaderboard (Croce et al., 2020).

401 For the white-box evaluation, the value of c^r is set to 2^{-7} . Table 3a and 3b present standard and 402 robust accuracy of models integrated with our method for CIFAR10 and CIFAR100 dataset, respec-403 tively. In these tables, the column *Original* indicates the original results reported by RobustBench, 404 and the column Original+Ours demonstrates the results of the proposed method. As indicated in 405 these tables, for CIFAR10 dataset, the proposed method enhances robust accuracy by more than 2%406 for RST-AWP, DefEAT, and LTD models, while other models receive approximately 1 to 1.5% improvement in robustness. Similarly, for CIFAR100 dataset, these models meet at least a 1% increase 407 in robustness. The empirical results prove that the resilience of existing models against adversarial 408 attacks can be improved by Lemma 2. We believe that the proposed solution is general as it achieves 409 great success in models incorporating different training techniques. 410

411 Another advantage of the proposed method that we would like to highlight is that the cost of our 412 approach can almost be ignored compared to the cost of adversarial training as the cost involves only a single pass scan of a set of images to determine the hyper-parameter c^{th} . This implies that 413 these models can enhance robustness for free. Specifically, LefEAT can achieve a robust accuracy 414 of 57.30% by removing inefficient training data. By combining LefEAT model with our approach, a 415 robust accuracy of 59.55% can be achieved, which is comparable to RST-AWP (60.04%). However, 416 RST-AWP introduces more images from another dataset, resulting in a higher cost in each epoch. 417 Similarly, for the CIFAR100 dataset, DefEAT with our proposed method achieves a robust accuracy 418 of 32.11%, which is better than EffAug (31.85%), which involves more complex data augmentation 419 during the training stage. This aligns with the suggestion by DefEAT that some data can be removed 420 without hurting robustness. Holistically, we believe that our approach might provide a hint during 421 the late phase of adversarial training to drop inefficient weights, resulting in further cost savings or 422 enhanced resilience.

An interesting observation from these tables is that standard accuracy improves across all models
 for both CIFAR10 and CIFAR100 datasets. While this phenomenon is not directly explained by
 Lemma 2, we hypothesize that the output ranges of ReLU and our proposed functions are highly
 similar, allowing for the maintenance of accuracy on clean data.

428 4.2.4 GRADIENT MASKING VERIFICATION 429

Previous studies suggest that the resilience of models might be unintentionally overestimated (Athalye et al., 2018; Carlini et al., 2019). The proposed function in (9) suppresses values to zero if the condition is satisfied. One might argue that this property could unintentionally cause obfuscated

432	Tabl	e 3: Standard and robust accurac	y of mod	lels integ	rated wit	h our metl	hod.
433			D 10.1				
434		(a) CIFA	R10 datas	set			
435			Orig	ginal	Origina	l+Ours	
436		Method	acc _{nat}	acc _{AA}	acc _{nat}	acc_{AA}	
437		RST-AWP Wu et al. (2020)	88.25	60.04	89.50	62.76	
430		DefEAT Chen & Lee (2024)	86.54	57.30	87.40	59.55	
439		LTD Chen & Lee (2021)	85.21	56.94	85.98	59.25	
440		AWP Wu et al. (2020)	85.36	56.17	86.19	57.85	
442		TRADESZhang et al. (2019)	85.34	52.86	85.78	53.80	
443 444		(b) CIFA	R100 data	iset			
445		(-) -					
446		Method	Ori	ginal	Origin	al+Ours	
447			acc _{nat}	acc_{AA}	acc _{nat}	acc _{AA}	
448		EffAug Addepalli et al. (2022)	68.75	31.85	69.14	32.57	
449		DKLD Cui et al. (2023)	64.08	31.65	64.26	32.58	
450		DefEAT Chen & Lee (2024)	64.32	31.13	66.42	32.11	
451		LTD Chen & Lee (2021)	64.07	30.59	64.29	31.95	
452		AWF Wu et al. (2020)	00.38	20.00	00.05	29.12	
453							
454	gradiants resul	ting in gradient attacks being u	nable to	efficient	ly produ	ca advars	arial avar
455	Therefore to v	erify that the proposed method of	loes not	encounte	er the gra	dient mas	sking issi
456	should conduct	more experiments from the follo	wing asp	bects:	i ille gre		
457		1.	0 1				
458	1. White-	box attacks should be better than	ı black-b	ox attack	s.		
459	2. Iterativ	e attacks should have better perf	ormance	than one	-step atta	acks.	
460	3. Robus	t accuracy should gradually decre	ease to z	ero when	the radiu	1s of <i>∈</i> -bal	1 increase
462	1 The m	adified model should defense ag	inct ody	arcarial a	vomnlog	ganaratad	by the or
463	4. model	s.	anist auv		xampies	generateu	by the of
464 465	5. Certifi	ed robustness that conducted by	random s	smoothing	g (Cohen	et al., 20	19).
466	AutoAttack has	examined the first item, which i	nvolves	three whi	te-box a	ttacks and	one blac
467	attack. By com	paring robust accuracy shown in	Table 1a	a, 1b and	2, the m	odels con	ibined wi
468	proposed method	od perform better robust accurac	y than th	ne origina	al model	s. It indic	ates blac
469	attacks cannot p	produce more adversarial exampl	es.				
470	The full experi	mental results for the rest of th	e experi	ments ca	n be fou	ind in Ap	pendix B
471	results demonst	trate that the proposed algorithr	n does n	ot violate	e any of	the above	e rules a
472	certified robust	ness improves by our method ac	ross mos	st settings	s. From	the evider	ice, we b
473	that the propose	ed method does not encounter the	gradien	t masking	g probler	n among o	lifferent l
474	parameters and	various models on CIFAR10 and	I CIFAR	100 datas	ets.		
475							
476	4.3 Cost An	JALYSIS					
477	East the OIEAP	10 and CIEAD 100 data at a	a +-+-1 +		ma f.		V100 C
478	FOR the CIFAR	-10 and UFAK-100 datasets, the 20 hours In contrast other met	bods use	alning ti	me for A	wr on a	. VIUU G
479	data (Wang et a	1 2023. Peng et al 2023. Rai e	1003 use $120'$	2 larger II 24) Jeadi	no to a to	ntal traini	ng time o
480	200 hours. In c	omparison, scanning the entire (CIFAR-10	$\frac{1}{2}$ or CIFA	R -100 d	lataset. or	partial tr
481	set on ImageNe	t dataset, takes less than 5 minut	es on a V	/100 GPU	J.		1
482	Dagandin - 41	posts of hunamonometer and	og 4!	unod : (Soution -	A 1 +1	0.00 0.001
483 484	candidate parar	neters to evaluate performance.	We can	efficient	ly assess	white-bc	are only ox perform

nodels inte Table 3: Standard and robust rated with ethod f.

mples. ue, we

- e.
- riginal

ck-box ith the ck-box

3. The nd the believe hyper-

GPU is itional of over aining

three mance using partial data from the training set, which significantly reduces the computational time. In 485 practice, the total time for hyperparameter search is about 1-3 hours, depending on the model size

and the dataset used. This demonstrates that our approach incurs significantly lower computational overhead.

489

490

4.4 DISCUSSION

A common pitfall is the misconception that minimizing the magnitude of the Lipschitz constant necessarily leads to improved robustness. However, in an extreme scenario, if the weights of the linear layers are replaced with an identity matrix, the Lipschitz constant would equal 1. This change would completely alter the output distribution, resulting in zero natural accuracy and, therefore, no meaningful robustness to assess. This work demonstrates that by manipulating the input domain of a linear system, it is possible to obtain an equivalent matrix that produces the same output while having a Lipschitz constant that is equal to or less than the original.

While it is true that the Lipschitz constant can be reduced through certified training (Mao et al., 2024) or Lipschitz-constrained methods (Zühlke & Kudenko, 2024), these approaches often introduce additional regularization terms or more complex objectives, complicating the training of robust models. In contrast, this work is specifically designed for model optimization during inference time, providing a cost-effective means to enhance robustness. Furthermore, this research offers valuable insights into identifying which weights can be eliminated with minimal impact on performance. This property could be integrated into existing training frameworks to further enhance robustness.

505 The overall procedure of the proposed algorithm shares many similarities with pruning techniques 506 (He & Xiao, 2023). However, we would like to emphasize that the proposals are distinctly different. 507 Pruning primarily aims to create a highly sparse model that accelerates inference times or reduces 508 the model size for deployment on edge devices, often without considering robustness. In particular, 509 when the pruned model exhibits extremely high sparsity without applying re-training or fine-tuning, natural accuracy drops significantly, implicitly indicating that condition (11) does not hold and that 510 Lemma 2 is not applicable in this case. On the other hand, when the pruned model has low sparsity, 511 only a small proportion of weights with values close to zero are eliminated, resulting in insignificant 512 adjustments to the center and radius of the disk. The proposed algorithm, on the other hand, can 513 identify crucial elements with minimal cost. We believe that investigating the impact of various 514 pruning techniques, such as iterative or post-training methods, on robustness, or combining these 515 techniques with our proposed approach, represents a valuable direction for future research. 516

517

5 CONCLUSION

518 519

520 In this paper, we recap how robustness is certified by the theorem of Lipschitz continuity. We 521 introduce the concept of the empirical Lipschitz constant, which minimizes the influence of the space 522 not drawn from real data, resulting in a precise estimation of the robustness of the corresponding 523 observed data. We prove that by remapping the input domain of a specific layer to a constrained range, the Lipschitz constant can be shrunk, leading to better robustness. The proposed function 524 introduces only one parameter, the value of which can be determined by scanning the training data 525 once, without re-training or fine-tuning. Compared with adversarial training, the proposed method 526 is almost cost-free. The experimental results suggest that our method can be combined with various 527 existing methods and achieve robustness improvements, and no gradient masking occurs in our 528 algorithm. Furthermore, our method can achieve the best robust accuracy for CIFAR10, CIFAR100, 529 and ImageNet datasets on the RobustBench leaderboard. 530

Numerous future directions merit exploration. Firstly, due to the property of maximization, the
 proposed function might easily be influenced by outliers. Designing a better function is an interest ing research topic. Secondly, exploring the combination with various activation functions, different
 model architectures or large-scale datasets would be beneficial. Lastly, it is worth investigating to
 understand the theoretical reasons why our proposed function improves standard accuracy.

536

537 REFERENCES

539 Sravanti Addepalli, Samyak Jain, et al. Efficient and effective augmentation strategy for adversarial training. *Advances in Neural Information Processing Systems*, 35:1488–1501, 2022.

559

585

- Maksym Andriushchenko, Francesco Croce, Nicolas Flammarion, and Matthias Hein. Square attack: a query-efficient black-box adversarial attack via random search. In *European conference on computer vision*, pp. 484–501. Springer, 2020.
- Anish Athalye, Nicholas Carlini, and David Wagner. Obfuscated gradients give a false sense of se curity: Circumventing defenses to adversarial examples. In *International conference on machine learning*, pp. 274–283. PMLR, 2018.
- Yatong Bai, Mo Zhou, Vishal M Patel, and Somayeh Sojoudi. Mixednuts: Training-free accuracy-robustness balance via nonlinearly mixed classifiers. *arXiv preprint arXiv:2402.02263*, 2024.
- Nicholas Carlini and David Wagner. Audio adversarial examples: Targeted attacks on speech-to-text. In 2018 IEEE Security and Privacy Workshops (SPW), pp. 1–7. IEEE, 2018.
- Nicholas Carlini, Anish Athalye, Nicolas Papernot, Wieland Brendel, Jonas Rauber, Dimitris
 Tsipras, Ian Goodfellow, Aleksander Madry, and Alexey Kurakin. On evaluating adversarial
 robustness. *arXiv preprint arXiv:1902.06705*, 2019.
- Yair Carmon, Aditi Raghunathan, Ludwig Schmidt, John C Duchi, and Percy S Liang. Unlabeled
 data improves adversarial robustness. *Advances in neural information processing systems*, 32, 2019.
- Erh-Chung Chen and Che-Rung Lee. Towards fast and robust adversarial training for image classification. In *Proceedings of the Asian Conference on Computer Vision*, 2020.
- Erh-Chung Chen and Che-Rung Lee. Ltd: Low temperature distillation for robust adversarial train *arXiv preprint arXiv:2111.02331*, 2021.
- Erh-Chung Chen and Che-Rung Lee. Data filtering for efficient adversarial training. *Pattern Recog- nition*, pp. 110394, 2024.
- Erh-Chung Chen, Pin-Yu Chen, I Chung, Che-rung Lee, et al. Overload: Latency attacks on object detection for edge devices. *arXiv preprint arXiv:2304.05370*, 2023.
- 570 Erh-Chung Chen, Pin-Yu Chen, I Chung, Che-Rung Lee, et al. Steal now and attack later:
 571 Evaluating robustness of object detection against black-box adversarial attacks. *arXiv preprint*572 *arXiv:2404.15881*, 2024.
- 573
 574
 575
 Pin-Yu Chen and Cho-Jui Hsieh. Adversarial robustness for machine learning. Academic Press, 2022.
- Pin-Yu Chen, Huan Zhang, Yash Sharma, Jinfeng Yi, and Cho-Jui Hsieh. 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*, pp. 15–26, 2017.
- Pin-Yu Chen, Yash Sharma, Huan Zhang, Jinfeng Yi, and Cho-Jui Hsieh. Ead: elastic-net attacks to deep neural networks via adversarial examples. In *Proceedings of the AAAI conference on artificial intelligence*, volume 32, 2018.
- Jeremy Cohen, Elan Rosenfeld, and Zico Kolter. Certified adversarial robustness via randomized
 smoothing. In *international conference on machine learning*, pp. 1310–1320. PMLR, 2019.
- Francesco Croce and Matthias Hein. Reliable evaluation of adversarial robustness with an ensemble
 of diverse parameter-free attacks. In *International conference on machine learning*, pp. 2206–2216. PMLR, 2020.
- Francesco Croce, Maksym Andriushchenko, Vikash Sehwag, Edoardo Debenedetti, Nicolas Flammarion, Mung Chiang, Prateek Mittal, and Matthias Hein. Robustbench: a standardized adversarial robustness benchmark. *arXiv preprint arXiv:2010.09670*, 2020.
- Jiequan Cui, Zhuotao Tian, Zhisheng Zhong, Xiaojuan Qi, Bei Yu, and Hanwang Zhang. Decoupled kullback-leibler divergence loss. arXiv preprint arXiv:2305.13948, 2023.

594 Ambra Demontis, Marco Melis, Maura Pintor, Matthew Jagielski, Battista Biggio, Alina Oprea, 595 Cristina Nita-Rotaru, and Fabio Roli. Why do adversarial attacks transfer? explaining transfer-596 ability of evasion and poisoning attacks. In 28th USENIX security symposium (USENIX security 597 19), pp. 321–338, 2019. 598 Andrew Du, Bo Chen, Tat-Jun Chin, Yee Wei Law, Michele Sasdelli, Ramesh Rajasegaran, and Dillon Campbell. Physical adversarial attacks on an aerial imagery object detector. In Proceedings 600 of the IEEE/CVF Winter Conference on Applications of Computer Vision, pp. 1796–1806, 2022. 601 602 Farzan Farnia, Jesse M Zhang, and David Tse. Generalizable adversarial training via spectral nor-603 malization. arXiv preprint arXiv:1811.07457, 2018. 604 Ian J Goodfellow, Jonathon Shlens, and Christian Szegedy. Explaining and harnessing adversarial 605 examples. arXiv preprint arXiv:1412.6572, 2014. 606 607 Sven Gowal, Krishnamurthy Dvijotham, Robert Stanforth, Rudy Bunel, Chongli Oin, Jonathan Ue-608 sato, Relja Arandjelovic, Timothy Mann, and Pushmeet Kohli. On the effectiveness of interval 609 bound propagation for training verifiably robust models. arXiv preprint arXiv:1810.12715, 2018. 610 611 Sven Gowal, Sylvestre-Alvise Rebuffi, Olivia Wiles, Florian Stimberg, Dan Andrei Calian, and Timothy A Mann. Improving robustness using generated data. Advances in Neural Information 612 Processing Systems, 34:4218–4233, 2021. 613 614 Yang He and Lingao Xiao. Structured pruning for deep convolutional neural networks: A survey. 615 IEEE transactions on pattern analysis and machine intelligence, 2023. 616 617 Yujia Huang, Huan Zhang, Yuanyuan Shi, J Zico Kolter, and Anima Anandkumar. Training certifiably robust neural networks with efficient local lipschitz bounds. Advances in Neural Information 618 Processing Systems, 34:22745–22757, 2021. 619 620 Inaam Ilahi, Muhammad Usama, Junaid Qadir, Muhammad Umar Janjua, Ala Al-Fuqaha, Dinh Thai 621 Hoang, and Dusit Niyato. Challenges and countermeasures for adversarial attacks on deep rein-622 forcement learning. IEEE Transactions on Artificial Intelligence, 3(2):90–109, 2021. 623 624 Stepan Komkov and Aleksandr Petiushko. Advhat: Real-world adversarial attack on arcface face 625 id system. In 2020 25th international conference on pattern recognition (ICPR), pp. 819-826. IEEE, 2021. 626 627 Alex Krizhevsky, Ilya Sutskever, and Geoffrey E Hinton. Imagenet classification with deep convo-628 lutional neural networks. Advances in neural information processing systems, 25, 2012. 629 630 Aounon Kumar, Chirag Agarwal, Suraj Srinivas, Soheil Feizi, and Hima Lakkaraju. Certifying llm 631 safety against adversarial prompting. arXiv preprint arXiv:2309.02705, 2023. 632 Jinfeng Li, Shouling Ji, Tianyu Du, Bo Li, and Ting Wang. Textbugger: Generating adversarial text 633 against real-world applications. arXiv preprint arXiv:1812.05271, 2018. 634 635 Chang Liu, Yinpeng Dong, Wenzhao Xiang, Xiao Yang, Hang Su, Jun Zhu, Yuefeng Chen, Yuan 636 He, Hui Xue, and Shibao Zheng. A comprehensive study on robustness of image classification 637 models: Benchmarking and rethinking. arXiv preprint arXiv:2302.14301, 2023. 638 Ze Liu, Yutong Lin, Yue Cao, Han Hu, Yixuan Wei, Zheng Zhang, Stephen Lin, and Baining Guo. 639 Swin transformer: Hierarchical vision transformer using shifted windows. In Proceedings of the 640 IEEE/CVF international conference on computer vision, pp. 10012–10022, 2021. 641 642 Aleksander Madry, Aleksandar Makelov, Ludwig Schmidt, Dimitris Tsipras, and Adrian Vladu. 643 Towards deep learning models resistant to adversarial attacks. arXiv preprint arXiv:1706.06083, 644 2017. 645 Yuhao Mao, Mark Niklas Mueller, Marc Fischer, and Martin Vechev. Understanding certified train-646 ing with interval bound propagation. In The Twelfth International Conference on Learning Rep-647 resentations, 2024. URL https://openreview.net/forum?id=h05eQniJsQ.

679

682

683

684

685

- Mark Niklas Müller, Franziska Eckert, Marc Fischer, and Martin Vechev. Certified training: Small boxes are all you need. *arXiv preprint arXiv:2210.04871*, 2022.
- ShengYun Peng, Weilin Xu, Cory Cornelius, Matthew Hull, Kevin Li, Rahul Duggal, Mansi Phute,
 Jason Martin, and Duen Horng Chau. Robust principles: Architectural design principles for
 adversarially robust cnns. *arXiv preprint arXiv:2308.16258*, 2023.
- Zeyu Qin, Yanbo Fan, Yi Liu, Li Shen, Yong Zhang, Jue Wang, and Baoyuan Wu. Boosting the transferability of adversarial attacks with reverse adversarial perturbation. *Advances in neural information processing systems*, 35:29845–29858, 2022.
- Joseph Redmon, Santosh Divvala, Ross Girshick, and Ali Farhadi. You only look once: Unified,
 real-time object detection. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pp. 779–788, 2016.
- Leslie Rice, Eric Wong, and Zico Kolter. Overfitting in adversarially robust deep learning. In International conference on machine learning, pp. 8093–8104. PMLR, 2020.
- Christian Szegedy, Wojciech Zaremba, Ilya Sutskever, Joan Bruna, Dumitru Erhan, Ian Goodfellow,
 and Rob Fergus. Intriguing properties of neural networks. *arXiv preprint arXiv:1312.6199*, 2013.
- Kiaosen Wang, Xuanran He, Jingdong Wang, and Kun He. Admix: Enhancing the transferability of adversarial attacks. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pp. 16158–16167, 2021.
- Yixiang Wang, Jiqiang Liu, Xiaolin Chang, Ricardo J Rodríguez, and Jianhua Wang. Di-aa: An
 interpretable white-box attack for fooling deep neural networks. *Information Sciences*, 610:14–32, 2022.
- Zekai Wang, Tianyu Pang, Chao Du, Min Lin, Weiwei Liu, and Shuicheng Yan. Better diffusion models further improve adversarial training. In *International Conference on Machine Learning*, pp. 36246–36263. PMLR, 2023.
- Kingxing Wei, Ying Guo, and Jie Yu. Adversarial sticker: A stealthy attack method in the physical
 world. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 45(3):2711–2725, 2022.
- Dongxian Wu, Shu-Tao Xia, and Yisen Wang. Adversarial weight perturbation helps robust generalization. *Advances in neural information processing systems*, 33:2958–2969, 2020.
 - Yi Xie, Zhuohang Li, Cong Shi, Jian Liu, Yingying Chen, and Bo Yuan. Enabling fast and universal audio adversarial attack using generative model. In *Proceedings of the AAAI conference on artificial intelligence*, volume 35, pp. 14129–14137, 2021.
- Kaidi Xu, Gaoyuan Zhang, Sijia Liu, Quanfu Fan, Mengshu Sun, Hongge Chen, Pin-Yu Chen, Yanzhi Wang, and Xue Lin. Adversarial t-shirt! evading person detectors in a physical world. In *Computer Vision–ECCV 2020: 16th European Conference, Glasgow, UK, August 23–28, 2020, Proceedings, Part V 16*, pp. 665–681. Springer, 2020.
- Jia-Yu Yao, Kun-Peng Ning, Zhen-Hui Liu, Mu-Nan Ning, and Li Yuan. Llm lies: Hallucinations are not bugs, but features as adversarial examples. *arXiv preprint arXiv:2310.01469*, 2023.
- Mingjun Yin, Shasha Li, Chengyu Song, M Salman Asif, Amit K Roy-Chowdhury, and Srikanth V Krishnamurthy. Adc: Adversarial attacks against object detection that evade context consistency checks. In *Proceedings of the IEEE/CVF Winter Conference on Applications of Computer Vision*, pp. 3278–3287, 2022.
- ⁶⁹⁷ Yuichi Yoshida and Takeru Miyato. Spectral norm regularization for improving the generalizability
 ⁶⁹⁸ of deep learning. *arXiv preprint arXiv:1705.10941*, 2017.
- Hongyang Zhang, Yaodong Yu, Jiantao Jiao, Eric Xing, Laurent El Ghaoui, and Michael Jordan.
 Theoretically principled trade-off between robustness and accuracy. In *International conference* on machine learning, pp. 7472–7482. PMLR, 2019.

- Jingfeng Zhang, Xilie Xu, Bo Han, Gang Niu, Lizhen Cui, Masashi Sugiyama, and Mohan Kankan halli. Attacks which do not kill training make adversarial learning stronger. In *International conference on machine learning*, pp. 11278–11287. PMLR, 2020.
- Yihua Zhang, Guanhua Zhang, Prashant Khanduri, Mingyi Hong, Shiyu Chang, and Sijia Liu. Revisiting and advancing fast adversarial training through the lens of bi-level optimization. In *International Conference on Machine Learning*, pp. 26693–26712. PMLR, 2022.
- Pu Zhao, Pin-Yu Chen, Siyue Wang, and Xue Lin. Towards query-efficient black-box adversary with zeroth-order natural gradient descent. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 34, pp. 6909–6916, 2020.
- Kaijie Zhu, Jindong Wang, Jiaheng Zhou, Zichen Wang, Hao Chen, Yidong Wang, Linyi Yang, Wei
 Ye, Neil Zhenqiang Gong, Yue Zhang, et al. Promptbench: Towards evaluating the robustness of
 large language models on adversarial prompts. *arXiv preprint arXiv:2306.04528*, 2023.
- Monty-Maximilian Zühlke and Daniel Kudenko. Adversarial robustness of neural networks from the perspective of lipschitz calculus: A survey. *ACM Computing Surveys*, 2024.
- 717 718 719

722

- A ABLATION STUDY
- 721 A.1 HYPER-PARAMETER SELECTION

This experiment investigates how the choice of hyper-parameter c^r influences standard accuracy and robust accuracy. Since most models are represented in 16 bit format, and the widths of fraction bit for FP16 format defined by IEEE-754 standard and BFloat are 10 and 7 bits, respectively, truncated errors might easily occur when performing addition on two numbers with a magnitude difference of 2^8 or higher. On the other hand, when c^r is set to 2^{-5} , all models experience a significant drop in standard accuracy, and there is meaningless in evaluating robustness at this configuration. We suggest that the candidates of c^r are 2^{-8} , 2^{-7} and 2^{-6} .

The results on CIFAR10 and CIFAR100 are presented in Table 4 and Table 5, respectively. Moreover, the results of accuracy against CW attack on L_{∞} norm for CIFAR10 and CIFAR100 datasets are presented in Tables 6a and 6b, respectively. As can be seen, when c^r is set to 2^{-8} , all models achieve better standard accuracy and robust accuracy. Additionally, the results for all models with $c^r = 2^{-7}$ are surpassed by those when c^r is set to 2^{-8} . Robust accuracy can be further enhanced by setting 2^{-6} , while standard accuracy might drop compared to the original. The results suggest that $c^r = 2^{-7}$ is a solution that balances standard accuracy and robustness. Nevertheless, when robustness is a major concern, $c^r = 2^{-6}$ is a better choice.

⁷³⁷Intuitively, we expect that standard accuracy gradually decreases when the value of c^r increases. The phenomenon can be observed when c^r is 2^{-6} or higher but two counterexamples are reported in the ablation study when setting c^r to 2^{-7} and 2^{-8} . A possible explanation is that the optimizer becomes stuck in a saddle area, as ReLU is non-differentiable at the zero point. This might cause the gradient direction to become stuck in an oscillation when values are close to zero. By shifting those values to zero, antagonistic effects among different feature maps, filters, or channels are accidentally mitigated. However, further investigation and evidence are needed to support this conjecture.

We argue that any function that satisfies the conditions defined in (11) can shrink the largest eigenvalue. There might be another function that can perform better than the proposed one. Besides, the hyper-parameter is determined by choosing the maximum value appearing in the dataset.

747 748 749

B FULL EXPERIMENTAL RESULTS OF GRADIENT MASKING VERIFICATION

Table 7a and 7b present the robust accuracy against adversarial examples generated by the original models on CIFAR10 and CIFAR100 datasets, respectively. As observed, none of the models showed lower robust accuracy than the original model. It indicates that adversarial examples can be efficiently crafted by utilizing the gradients from the victim models.

Table 8 and 9 presents the robust accuracy against FGSM and PGD attacks among different radii of the ϵ -ball on the CIFAR10 and CIFAR100 datasets, respectively. As observed, the robust accuracy

7	5	8
7	5	9
7	6	0
7	6	1
7	6	2
7	6	3
7	6	4
7	6	5
7	6	6
7	6	7
7	6	8
7	6	9
7	7	0
7	7	1
7	7	2
7	7	3
7	7	4
7	7	5
7	7	6
7	7	7
7	7	8
7	7	9
7	8	0
7	8	1
7	8	2
7	8	3
7	8	4
7	8	5
7	8	6
7	8	7
7	8	8
7	8	9
7	9	0
7	9	1
-	0	0

Table 4: Ablation study of selecting optimal c^r for CIFAR10 dataset.

Mathod	Robus	tBench	$c^r =$	$=2^{-8}$	$c^r =$	$= 2^{-7}$	$c^r =$	$= 2^{-6}$	
Methou	acc _{nat}	acc_{AA}	acc _{nat}	acc_{AA}	acc _{nat}	acc_{AA}	acc _{nat}	acc_{AA}	
RST-AWP	88.25	60.04	88.82	60.96	89.50	62.76	87.88	61.96	
DefEAT	86.54	57.30	86.88	57.81	87.40	59.55	84.59	61.08	
LTD	85.21	56.94	85.28	57.28	85.98	59.25	85.59	60.63	
AWP	85.36	56.17	85.80	56.53	86.19	57.85	84.55	59.21	
TRADES	85.34	52.86	85.57	52.97	85.78	53.80	85.49	55.37	
Table	5: Ablati	on study	of selec	ting optin	mal c^r fo	r CIFAR	100 datas	set.	
Method	Robust	Bench	$c^r =$	2^{-8}	$c^r =$	2^{-7}	$c^r =$	2^{-6}	
Witchiou	acc _{nat}	acc _{AA}	acc _{nat}	acc _{AA}	acc _{nat}	acc _{AA}	acc _{nat}	acc_{AA}	
EffAug	68.75	31.85	68.81	32.00	69.14	32.57	68.44	33.64	
DKLD	64.08	31.65	64.10	31.77	64.26	32.58	63.50	33.87	
DefEAT	65.89	30.57	66.12	31.11	66.42	32.46	65.06	34.07	
LTD	64.07	30.59	64.29	31.13	64.29	31.95	64.18	34.04	
AWP	60.38	28.86	60.18	29.10	60.63	29.72	60.71	30.82	

Table 6: The robust accuracy against CW attack on L_{∞} norm.

(a) CIFAR10 dataset (b) CIFAR100 dataset c^r c^{r} Method Origin Method Origin 2^{-8} 2^{-7} 2^{-6} 2^{-8} 2^{-7} 2^{-6} **RST-AWP** 58.98 61.84 68.24 80.92 EffAug 37.40 37.70 38.70 43.00 DefEAT 56.92 58.02 61.06 65.56 DKLD 37.50 38.06 39.38 44.20 LTD 58.12 58.56 60.50 64.86 DefEAT 36.90 37.56 39.82 44.30 AWP 57.34 60.58 66.50 LTD 36.66 37.32 38.86 43.44 56.84 AWP 34.56 35.20 35.94 40.40 TRADES 56.10 63.62 56.52 58.18

against FGSM, a one-step attack, is always higher than the robust accuracy against PGD, an iterative attack. This implies that the gradient is reliable, allowing the PGD attack to adjust the gradient direction multiple times to find adversarial examples. Additionally, we observe that the robust accuracy against PGD attacks for all models gradually decreases to zero as the radius of the ϵ -ball increases. This indicates that the quality of gradients is preserved, enabling PGD attacks to move the gradient 793 toward examples not in the observed distribution. 794

Figure 3 illustrates the certified robustness achieved by random smoothing for various models on the CIFAR10 dataset, where *Original* refers to the certified robustness of the original model, while 796 Ours denotes the robustness of the model combined with the proposed method. As can be seen, our method brings slight improvements in robustness, except for the AWP model. These results 798 demonstrate that our algorithm does not suffer from the gradient masking issue. However, the em-799 pirical Lipschitz constant is derived from the observed data. As the input distribution drawn from 800 random smoothing and the observed data might have discrepancies, this could result in fluctuations 801 in robustness.

- 802 803
- 804 805

QUANTITATIVE ANALYSIS OF EMPIRICAL LIPSCHITZ CONSTANT С

806 As mentioned in Definition 1, the Lipschitz constant in this paper is estimated based on the observed 807 data. The sparsity of the forged vectors is a crucial factor influencing the magnitude of the Lipschitz constant for the corresponding layers, although it is not the only factor. Figure 4 illustrates the 808 average proportion of pruned activations, which varies depending on the location of each linear layer. FC1 and FC2 refer to the first and second fully connected layers in the MLP blocks, respectively.



Figure 3: Certified robustness that conducted by random smoothing.

Table 7: The robust accuracy against adversarial examples generated by the original models.

(a) CIFAR10 dataset						(b) CIFAR100 dataset					
Method	Origin	$ 2^{-8}$	2^{-7}	2^{-6}		Method	Origin	2^{-8}	2^{-7}	2^{-6}	
RST-AWP	60.04	62.10	65.10	70.53		EffAug	31.85	32.87	35.08	40.04	
DefEAT	57.30	58.37	60.39	66.10		DKLD	31.65	32.91	35.04	40.58	
LTD	56.94	58.71	61.63	66.47		DefEAT	30.57	31.82	33.94	40.67	
AWP	56.17	57.49	59.74	65.58		LTD	30.59	32.05	34.07	39.11	
TRADES	52.86	55.55	55.09	58.68		AWP	28.86	29.88	32.18	36.67	

852

853

837 838

839

As shown, the proportion of pruned activations is approximately 20% for the FC1 layers, except for the last layer. This suggests that the forged function alone cannot significantly reduce the magnitude of the Lipschitz constant.

On the other hand, the pruned rates for the FC2 layer range from 30% to 95%, depending on the layer's location. However, we emphasize that a higher pruning rate does not necessarily lead to a substantial reduction in the empirical Lipschitz constant. This is because, although the Lipschitz constant estimated from most of the data may decline, the empirical Lipschitz constant is based on the worst-case scenario across the entire observed dataset. Experimental results show that, after applying our method, the eigenvalue of the worst-case scenario is approximately 95% of the original eigenvalue.

861

862

Table 8: The robust accuracy against FGSM and PGD attacks among different radii of ϵ -ball on CIFAR10 dataset.

Method	c^r	Attack	1	0	А	R	E 16	30	64	
			$\frac{1}{255}$	$\frac{2}{255}$	$\frac{4}{255}$	$\frac{3}{255}$	$\frac{10}{255}$	$\frac{32}{255}$	$\frac{04}{255}$	
	0^{-8}	FGSM	88.28	86.94	83.80	75.12	57.23	34.04	18.80	
	2	PGD	86.78	84.47	79.03	66.03	34.02	2.01	0.01	
RST-AWP	2^{-7}	FGSM	89.46	88.28	85.64	77.62	60.60	35.91	19.39	
		PGD	88.03	85.88	80.89	69.24	38.27	3.08	0.1	
	2^{-6}	FGSM	87.70	86.63	84.38	77.91	60.93	33.48	16.12	
Method RST-AWP DefEAT LTD AWP TRADES		PGD	86.38	84.69	81.19	73.72	52.24	11.89	0.19	
	2^{-8}	FGSM	86.38	85.40	81.98	72.73	53.28	30.34	18.13	
		PGD	84.52	82.07	76.51	63.71	33.87	1.76	0.0	
DefEAT	2^{-7}	FGSM	86.69	85.70	83.05	74.35	56.36	32.04	18.88	
		PGD	85.11	82.87	78.00	66.54	38.70	3.12	0.0	
	2^{-6}	FGSM	84.14	83.57	81.03	74.63	57.67	28.11	13.17	
		PGD	82.96	81.35	77.37	69.52	50.70	10.05	0.2	
LTD	2^{-8}	FGSM	84.94	83.87	81.15	72.80	55.45	33.06	18.44	
		PGD	83.13	80.68	75.53	63.52	34.81	2.64	0.0	
	2^{-7}	FGSM	85.48	84.67	82.24	74.10	57.41	35.53	17.57	
		PGD	83.88	81.88	77.00	65.38	28.57	3.95	0.0	
	2^{-6}	FGSM	85.06	84.28	82.21	75.77	60.79	34.39	14.34	
	-	PGD	83.91	82.00	78.33	69.82	49.95	11.63	0.21	
	2^{-8}	FGSM	85.11	83.90	80.68	71.28	53.78	33.21	20.86	
		PGD	83.34	80.34	75.08	61.53	30.50	1.89	0.03	
AWP	2^{-7}	FGSM	85.50	84.68	81.75	73.56	57.19	35.50	20.63	
		PGD	83.94	81.62	76.34	65.57	34.16	2.89	0.02	
	2^{-6}	FGSM	83.87	83.19	81.08	74.97	61.49	35.13	16.12	
		PGD	83.00	81.42	78.13	71.50	54.95	14.78	0.41	
	2^{-8}	FGSM	84.74	83.51	70.58	70.50	54.23	36.53	23.97	
		PGD	82.62	79.72	72.81	57.30	24.21	1.21	0.01	
TRADES	2^{-7}	FGSM	85.00	84.02	80.57	71.61	55.94	25.67	22.31	
		PGD	82.96	80.31	73.75	58.91	26.22	1.31	0.02	
	2^{-6}	FGSM	85.05	84.19	81.40	75.07	60.24	36.99	21.81	
	4	PGD	83 23	81 33	75 81	64 56	36 76	3 37	0.02	

Table 9: The robust accuracy against FGSM and PGD attacks among different radii of ϵ -ball on CIFAR100 dataset.

Method	c^r	Attack	1	2	4	ϵ_{8}	16	32	64	
			255	$\overline{255}$	$\overline{255}$	$\overline{255}$	$\overline{255}$	$\overline{255}$	$\overline{255}$	
	2^{-8}	FGSM	68.02	65.96	60.36	49.65	33.90	17.39	7.11	
	-	PGD	64.92	61.04	52.82	39.37	17.53	1.81	0.0	
EffAug	2^{-7}	FGSM	68.41	66.56	61.84	51.83	36.59	18.34	7.02	
-	2	PGD	65.66	62.02	54.58	41.55	19.70	2.26	0.0	
	2^{-6}	FGSM	67.84	67.25	64.65	57.97	44.23	22.20	8.26	
	2	PGD	66.38	64.20	$\begin{array}{c ccccccccccccccccccccccccccccccccccc$					
	2^{-8}	FGSM	63.47	61.94	58.15	48.86	34.39	17.73	6.26	
	2	PGD	60.71	57.14	50.44	38.14	17.38	1.97	0.0	
Method EffAug DKLD DefEAT LTD AWP	0^{-7}	FGSM	63.55	62.31	58.65	50.37	36.49	18.69	6.36	
	Ζ.	PGD	61.06	57.84	51.51	39.99	19.71	2.37	0.0	
	n^{-6}	FGSM	63.26	62.77	60.41	55.18	43.86	21.17	5.50	
	Ζ ~	PGD	61.67	59.62	55.83	$\begin{array}{c ccccccccccccccccccccccccccccccccccc$				
	2^{-8}	FGSM	65.57	64.36	59.88	49.38	32.48	15.38	5.50	
	2	PGD	62.39	58.96	51.85	38.59	17.18	1.45	0.0	
DefEAT	2^{-7}	FGSM	65.97	64.96	60.67	51.44	34.84	16.36	5.33	
	2	PGD	62.94	59.83	52.98	41.05	20.07	2.02	0.0	
	2^{-6}	FGSM	64.69	63.58	61.31	54.47	40.06	16.91	4.64	
	2	PGD	62.85	60.50	56.29	47.54	30.76	5.84	0.09	
	2^{-8}	FGSM	63.59	62.35	58.10	48.86	33.11	16.78	5.92	
	2	PGD	61.06	57.65	50.70	38.21	18.21	1.98	0.0	
LTD	2^{-7}	FGSM	64.05	62.87	59.02	50.16	34.90	17.27	5.54	
	2	PGD	61.51	58.32	51.84	39.89	20.32	2.26	0.0	
	2^{-6}	FGSM	63.62	62.98	60.96	54.96	41.51	19.78	4.69	
	2	PGD	61.90	59.68	55.23	46.49	29.13	5.74	0.05	
	n^{-8}	FGSM	59.77	58.06	54.21	45.54	30.98	16.69	6.49	
	2	PGD	56.72	52.92	46.53	34.90	16.03	2.11	0.0	
AWP	-7	FGSM	60.00	58.52	54.93	46.65	32.66	17.42	6.04	
	ے 	PGD	57.19	53.75	47.46	36.22	17.48	2.58	0.0	
	2^{-6}	FGSM	60.20	59.71	57.47	52.05	39.78	21.10	5.70	
	4	PGD	58.19	55.66	50.91	42.37	25.53	5.68	0.09	

