

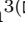





Adversarial Example Defense via Perturbation Grading Strategy

Shaowei Zhu¹, Wanli Lyu¹, Bin Li², Zhaoxia Yin³, and Bin Luo¹

¹ Anhui Provincial Key Laboratory of Multimodal Cognitive Computation,
Anhui University, Hefei, China

² Guangdong Key Laboratory of Intelligent Information Processing and Shenzhen
Key Laboratory of Media Security, Shenzhen University, Shenzhen, China
libin@szu.edu.cn

³ School of Communication and Electronic Engineering,
East China Normal University, Shanghai, China
zxyin@cee.ecnu.edu.cn

Abstract. Deep Neural Networks have been widely used in many fields. However, studies have shown that DNNs are easily attacked by adversarial examples, which have tiny perturbations and greatly mislead the correct judgment of DNNs. Furthermore, even if malicious attackers cannot obtain all the underlying model parameters, they can use adversarial examples to attack various DNN-based task systems. Researchers have proposed various defense methods to protect DNNs, such as reducing the aggressiveness of adversarial examples by preprocessing or improving the robustness of the model by adding modules. However, some defense methods are only effective for small-scale examples or small perturbations but have limited defense effects for adversarial examples with large perturbations. This paper assigns different defense strategies to adversarial perturbations of different strengths by grading the perturbations on the input examples. Experimental results show that the proposed method effectively improves defense performance. In addition, the proposed method does not modify any task model, which can be used as a preprocessing module, which significantly reduces the deployment cost in practical applications.

Keywords: Deep Neural Network · Adversarial examples · JPEG compression · Image denoising · Adversarial defense

1 Introduction

Deep neural networks (DNNs) have achieved widespread success in modern life, including image classification [14], medical image segmentation [11], and vehicle detection [27]. Research [28] has shown that attackers can add carefully crafted tiny perturbations to normal examples to mislead the model into making bad decisions. The new input generated by deliberately adding tiny perturbations to

normal examples is called adversarial examples, which can lead to misjudgment of the model and cause great harm.

However, existing studies have shown that there are also adversarial examples in real physical scenarios, so there is a significant safety problem in the practical application of DNN, such as automatic driving [22] and face recognition [7]. The safety requirements are higher in such case task scenarios.

Researchers have proposed various defense methods to reduce the impact of such adversarial perturbations, including adversarial training [8, 25], defense distillation [21], and preprocessing of input transformations [12, 17, 32]. Among them, adversarial training and defensive distillation require retraining or modification of the classifier. At the same time, input transformation-based methods focus on denoising/transforming the input before feeding it into the classifier, making it easier to deploy in practical applications. For these reasons, many input transformation-based methods have emerged in recent years. For example, ComDefend [12], Deep Image Prior (DIP) [29], and DIPDefend [3] directly reconstruct adversarial examples into normal images. Similarly, DefenseGAN [24] uses generative adversarial networks (GANs) to remove the effects of adversarial perturbations. However, these defense methods heavily depend on the dataset size and training time for training models and are computationally expensive, limiting their real-life applications. Therefore, some researchers [17, 32] turn to study how to denoise through image processing techniques.

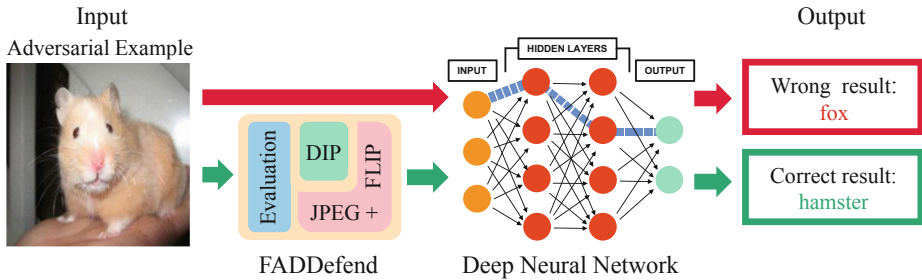


Fig. 1. An example of the proposed FADDefend. It removes the perturbations in the adversarial examples before feeding them into the classifier.

Yin et al. [33] found that the image compression method based on image processing technology can remove the perturbation from the small-perturbed image, and combined with the mirror flip, the defense effect can be improved. However, as the adversarial perturbation increases, the defense effect becomes worse. Furthermore, DIP [29] can be obtained from a single low-dimensional robust feature extracted from the input image to reconstruct the image without additional training cost. The reconstructed images of this method can significantly reduce the aggressiveness of large-perturbed adversarial examples. Therefore, we propose a perturbation grading strategy-based defense method called FADDefend,

which classifies adversarial perturbations and performs different operations on different levels of perturbations to improve the defense performance. The proposed method's adversarial example defense process is shown in Fig. 1.

The main contributions of our paper are:

- We propose an effective adversarial defense method that takes different defensive operations on adversarial examples with different perturbations.
- The proposed method uses a perturbation grading strategy to divide adversarial examples into large and small-perturbed examples.
- Experimental results show that the proposed method not only improves the performance of the adversarial defense but also reduces the computational cost.

2 Related Works

2.1 Adversarial Attack Methods

Adversarial examples are input images to which tiny perturbations are maliciously added to fool the neural network classifier. Common adversarial attack algorithms can be divided into attack methods based on white-box settings and attack methods based on black-box settings. Under the white-box attack, the attacker can fully understand the structure and parameters of the target model, while under the black-box attack, the attacker knows nothing about the relevant information of the target model.

White-Box Attacks. Szegedy et al. generate adversarial examples using L-BFGS [28], and Goodfellow et al. proposed a method to generate adversarial examples with only a one-step attack called the Fast Gradient Sign Method (FGSM) [9]. In order to improve the attack performance of FGSM, iteration-based multi-step attack methods are proposed, such as the Basic Iterative Method (BIM) [13], Momentum-based Iterative FGSM (MIFGSM) [5], and Projected Gradient Descent (PGD) [18]. DeepFool [19] attack method proposed to generate adversarial examples by finding the minimum perturbation on the hyperplane. The C&W attack [2] is another way to find adversarial examples through optimization. JSMA [20] is a sparse attack method that only modifies a small number of pixels. Athalye et al. proposed an attack method called Backward Pass Differentiable Approximation [1] by approximating the gradient of the defense to break this stochastic defense method.

Black-Box Attacks. Black-box attacks usually generate black-box adversarial examples on surrogate models and then exploit the transferability of adversarial examples to attack the target model. Various methods are proposed to improve the adversarial transferability of black-box attacks, such as query-based attacks [7,30] and transfer-based attacks [6].

2.2 Adversarial Defense Methods

Many adversarial defense methods have been proposed, including adversarial training and input transformation methods. The first method improves the robustness of the model by adding some adversarial examples to the normal training dataset, but this method slightly reduces the accuracy of normal examples. The second method modifies the input examples by preprocessing before entering the model to eliminate the adversarial perturbations in the examples. Classic digital image processing techniques, such as color depth reduction [31], image stitching [4], and JPEG compression [17], are used to improve model accuracy. However, these methods perform unsatisfactorily in defending against large-perturbed adversarial examples.

Later, model-based image reconstruction methods are proposed, including 1) denoising adversarial examples and 2) restoring them to clean images through CNN networks. Liao et al. proposed a high-level representation-guided denoiser [15] to remove adversarial perturbation. However, it requires a large dataset and more iterations to train a denoising model that transforms adversarial examples into clean images, limiting its practical scope.

Our work focuses on improving the unsatisfactory robustness of these defenses to large-perturbed adversarial examples. Since previous defense methods have shown better defense against small-perturbed adversarial examples, this problem can be solved by transforming large-perturbed images into small-perturbed images through the DIP [29] network.

3 Proposed Method

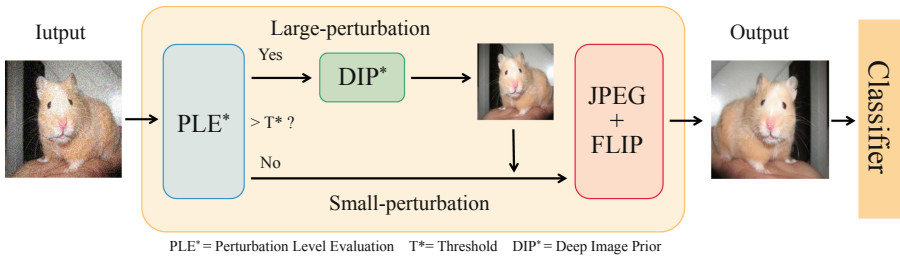


Fig. 2. FADDefend defense framework.

This paper proposes a new defense method FADDefend. FADDefend is divided into three modules. The first module is an evaluation module used to evaluate the noise level in the example. The second module is a JPEG compression and mirror flip module used to process the small-perturbed examples after perturbation grading. The third module is the DIP reconstruction module used to process large-perturbed examples after perturbation grading. Specifically, when an example is fed into FADDefend, the evaluation module evaluates its perturbation level. If the perturbation level is less than a pre-set threshold, it is defined

as a small perturbation. It is denoised using a JPEG compression algorithm with a quality factor (QF) of 95, then mirror flipped and fed into the classifier. Otherwise, it is defined as a large perturbation and sent to the DIP image restoration module, which is sent to the JPEG compression and mirror flip module after processing. Figure 2 shows the framework of FADDefend.

The basic idea of the proposed method to resist adversarial perturbation is introduced in this section. Then the pros and cons of the JPEG compression and DIP defense methods are revealed. Finally, the proposed method combines the advantages of the two defense methods through the perturbation level evaluation module to achieve a better defense effect.

3.1 Image Perturbation Level Evaluation

Defenders assume a known perturbation level and conduct targeted defenses, but this assumption is unrealistic. When the model is deployed, it is attacked by an unknown perturbation strength, and its perturbation level cannot be judged based on the input image. Therefore, this problem can be solved by evaluating the perturbation level of the input image by blind image perturbation level evaluation.

Liu et al. proposed a blind image perturbation level assessment algorithm [16] to select low-level patches without high-frequency information from noisy images. The Principal Component Analysis technique estimates the perturbation level based on the selected patches. The eigenvalues of the image gradient covariance matrix are used as the standard to measure the texture intensity. A stable threshold is selected to distinguish the perturbation level through an iterative method.

In this section, the experimental examples are composed of 500 examples from the data set and the expected accuracy of the default is set to 50%: A choice that balances detection accuracy and defense time consumption. Therefore, the threshold is chosen by comparing the defense accuracies of adversarial examples with different perturbation levels, as shown in Fig. 3 (Threshold is 2.13). Those smaller than the threshold are defined as small-perturbed adversarial examples, and those greater than the threshold are defined as large-perturbed adversarial examples.

The perturbation strength can be divided into different intervals by selecting different thresholds, convenient for subsequent selection of different image processing modules.

3.2 Defense Methods Based on Image Processing Technology

Most existing defense methods cannot defend well against large-perturbed adversarial examples. For example, the standard JPEG compression algorithm can remove the high-frequency information in the image well to retain the low-frequency information. Just as research [15] shows that adversarial perturbations can be viewed as high-frequency information with a specific structure so that this image compression can achieve a specific denoising effect. However, as

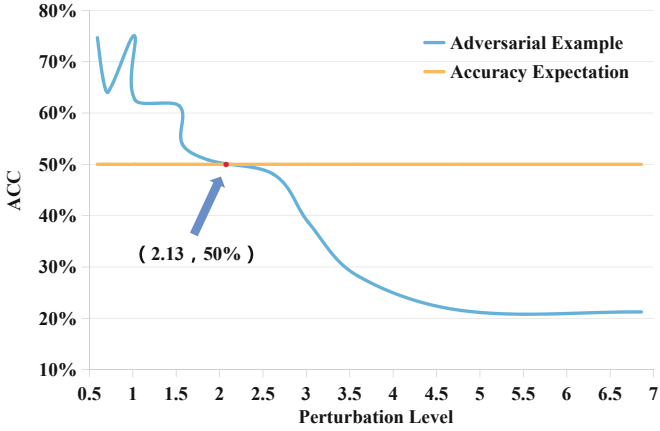


Fig. 3. Different thresholds can be chosen by the intersection of expected accuracy and adversarial example defense accuracy.

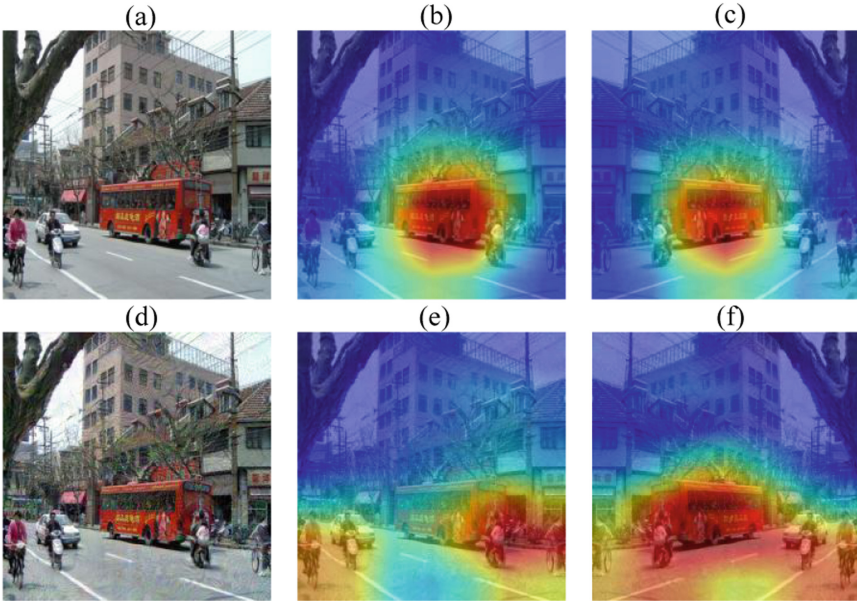


Fig. 4. (a) original example; (d) corresponding adversarial example. (b)(c)(e)(f) used class activation mapping of the images. (b) original example-bus; (c) flipped original example-bus; (e) adversarial example-bike; (f) flipped adversarial example-bus. The redder the class activation mapping of the image, the more the model pays attention to this area.

the perturbation strength increases, the adversarial perturbation may enter the low-frequency range. Even after removing some high-frequency information, the adversarial perturbation still exists.

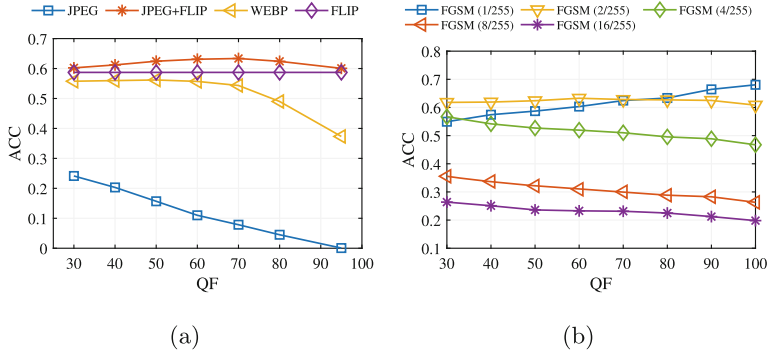


Fig. 5. Comparison of defense effects of preprocessing methods under different QFs. (a) Defense accuracy of various preprocessing methods under FGSM (2/255), (b) Defense accuracy using JPEG compression combined with mirror flip under different perturbations.

We strengthen JPEG’s defense against small perturbations by breaking the adversarial structure by mirror flipping [33]. As shown in Fig. 4, the model’s attention is observed by flipping the original and adversarial examples and using the class activation map. For the original example, the model’s attention is focused on the object’s main part before and after the flip. At the same time, for the adversarial example, the model is distracted, and the flipped adversarial example can correct the model’s attention so that it can be classified correctly.

Figure 5(a) shows the defense accuracy of different defense methods under the QFs. As the QF increases, the compression rate decreases, reducing the defense effect of JPEG compression against adversarial examples. In contrast, JPEG compression combined with mirror flip has better defense accuracy and stability. Figure 5(b) shows that the defense of JPEG compression combined with mirror flip is unsatisfactory under the large-perturbation (FGSM, $\epsilon = 8/16$).

The compressed image loses semantic information by using a smaller QF (the smaller the QF, the worse the image quality), making it more difficult for the model to distinguish accurately. Although such methods defend well on small-perturbed adversarial examples and low computational cost, they cannot effectively deal with adversarial examples with large perturbations. Consider using image reconstruction methods to process large perturbations in adversarial perturbations into small perturbations, followed by JPEG and flip processing, which is introduced in the next section.



Fig. 6. The above figures show the image reconstruction process of DIP and the image classification results (predicted category ID/iteration number) at different stages. The Top-5 predictions containing the correct label are marked in red (true labels: 274). PGD (8/255) is used to generate untargeted adversarial examples. (Color figure online)

3.3 Defense Methods Based on Deep Image Priors

Many model-based defense methods can achieve good defense results, such as image reconstructions [3], image denoising [15], image super-resolution [29], or GANs [24] recovering images for large-perturbed adversarial examples.

However, these methods require image reconstruction by learning a large amount of prior knowledge of external training data during the training stage. These defense methods consume a lot of computational and hardware costs. The CNN model itself has specific image prior capability. A clean image can also be reconstructed using prior internal knowledge of the image itself. Therefore, an untrained convolutional neural network is used in this module to generate denoised images through the DIP-based [29] generator network. Specifically, an adversarial example x_{adv} adds random noise z as the input example of the generator, which is solved by the following constraints by

$$\min_{\theta} \|f_{\theta}(z) - x_{adv}\|_2^2, \quad (1)$$

where f is a randomly initialized network, random but fixed noise z and adversarial examples x_{adv} as learning targets, each iteration input is a fixed noise z , and the parameter θ is updated by gradient descent. As the number of iterations increases, the output is closer to x_{adv} .

Consistent with the conclusions drawn [29], the DIP-based image reconstruction method restores the low-frequency information of the image in the early iterative stage and restores the high-frequency information (including adversarial perturbation) as the number of iterations increases. Figure 6 shows the image reconstruction process of DIP and the image classification results. In the early stage, the network recovered the image’s low-frequency information (main contours). The high-frequency information of the image will also be recovered in later iterations. It can be seen that the visual quality of the image is very high at this time. Unfortunately, adversarial perturbations are also recovered, leading to overfitting. Therefore, reconstructed images with few adversarial perturbations are obtained by stopping in an early stage. Finally, the DIP module transforms the large-perturbed example into a small-perturbed example. The small perturbation adversarial examples can increase the classification model’s defense accuracy through JPEG and flip processing.

4 Experiments

4.1 Experimental Settings

Our experiments are tested on the ImageNet dataset. 10,000 test images are selected randomly from the whole test set, and 8,850 clean example sets are screened for experiments. ResNet152 [10] is used to test robustness against various attack strengths of FGSM [9], PGD [18], BIM [13], and MIFGSM [19] ($\epsilon = 2/4/8/16$, using the L_∞ distance metric). ResNet152 [10] and VGG19 [26] are used to test cross-model defense robustness.

The DIP module adopts the U-Net [23] structure as the generative network and skip-connections to connect outputs of layers with the same spatial dimension. Specifically, LeakyReLU is used as the nonlinear activation function. Moreover, convolution operation and bilinear interpolation are used for downsampling and upsampling. An adversarial example is input as the target image. The generator network parameters are randomly initialized, random noise is used as input, and the target image is reconstructed by gradient descent training, stopping after 400 iterations.

4.2 Comparison with Image Compression Defenses

FADDefend combines image compression, mirror flip, and reconstruction modules with perturbation level evaluation modules. Since the image reconstruction module has a better processing effect on the large-perturbed adversarial examples, the image compression method is used as an auxiliary. It makes up for the disadvantage of the image compression defense method that the defense effect is unsatisfying when dealing with large-perturbed adversarial examples.

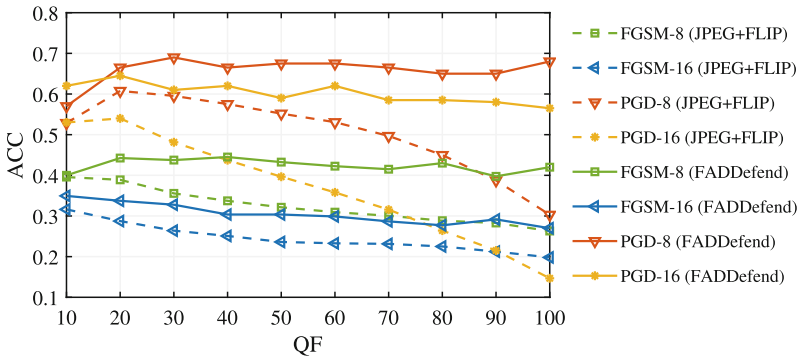


Fig. 7. The dotted part is the defense method of JPEG Compression combined with mirror flip. The solid part is the FADDefend defense method, showing the recognition accuracy curves under different QFs.

As shown in Fig. 7, compared with the defense method of JPEG image compression combined with mirror flip, FADDefend has better performance under different QFs and maintains stable defense performance with the increase of QF.

Table 1. Accuracy results (%) of defense methods on ResNet152 under attacks of different types and strengths (ImageNet, $\epsilon = 2/4/8/16$)

Method	Clean	FGSM	MIFGSM	PGD	BIM
Normal	100	0 / 0 / 0 / 0	0 / 0 / 0 / 0	0 / 0 / 0 / 0	0 / 0 / 0 / 0
JPEG [17]	86	47/45/24/13	74/64/40/14	59/64/50/37	56/68/56/63
War [32]	86	60/ 59 /36/21	79/72/53/27	65/68/59/51	61/75/63/72
ComDefend [12]	94	32/19/ 5 / 5	56/30/ 9 / 1	44/32/ 7 / 2	37/48/21/23
DIPDefend [3]	79	55/45/26/14	80/67/43/25	60/69/60/38	61/71/59/77
FADDefend	94	62/55/43/32	81/75/60/44	66/72/66/65	68/78/70/84

Table 2. Accuracy results (%) of defense methods on attacks on the VGG19, the source model is ResNet152. (ImageNet, $\epsilon = 4/8/16$)

Method	Clean	FGSM	MIFGSM	PGD	BIM
Normal	89	48/40/22	66/40/23	62/42/28	70/51/60
War [32]	74	57 /42/26	65/54/39	52/50/37	67/56/62
ComDefend [12]	86	52/38/22	70/50/31	66/47/41	75 /60/62
DIPDefend [3]	86	50/44/31	69/54/46	66/ 65 /44	73/61/72
FADDefend	86	50/ 46/31	70/57/46	66 /61/44	73/ 62/72

4.3 Results on Attacks of Different Types and Strengths

To verify the effectiveness of the proposed FADDefend method, other state-of-the-art defense methods are compared, including JPEG compression [17], War (WebP compression and resize) [32], ComDefend [12], and DIPDefend [3]. The classification accuracy on ImageNet is shown in Table 1. In contrast, the image reconstructed methods have high defense accuracies, such as DIPDefend [3] and the proposed method. Under the FGSM attack ($\epsilon = 8/16$), compared with the other four methods, the proposed method improves accuracy by at least 10%. Compared with the image compression method of ComDefend [12], the defense accuracy can even exceed 20%. The image compression methods are generally ineffective due to the large perturbation, which destroys the image structure and causes the loss of semantic information.

4.4 Results on Migration Attacks Under Different Models

Under cross-model attacks, assume that the attacker does not have access to the specific parameters of the target classifier or defense model. Attackers can only attack by exploiting the transferability of adversarial examples. The attacker acts as a surrogate model by training a model with a different structure than the target classifier. Then, adversarial examples generated by attacking the surrogate model may also lead to misclassification of the target classifier. Adversarial examples are generated under various attack algorithms ($\epsilon = 4/8/16$) using the ResNet152 model as surrogate and VGG19 as the target model. The test results are shown in Table 2. It can be seen that the recognition accuracy of the VGG19 model drops significantly on the adversarial examples generated by ResNet152. The larger the perturbation, the lower the recognition accuracy, which is the same result tested on the surrogate model and verified the attack transferability of adversarial examples. The recognition accuracy of the target model under different strengths and attack algorithms is significantly improved after adopting the proposed defense method. Compared with War [32], Comdefend [12], and DIPDefend [3] defense methods, the proposed method achieves the same or even better performance under various attacks.

The data in Table 2 shows that the defense performance of the proposed method is almost equal to that of DIPDefend under specific attack strengths, but this does not mean that the proposed method has no advantages. At the same time, we also conduct comparative experiments on the running speed, and memory consumption with the image reconstruction method [29] and the method [3] with the defense accuracy close to the proposed method. Randomly select 1000 images from the same dataset in Table 2 to form an example dataset for experimentation. Table 3 shows that the proposed method is less than the other two methods in running time and memory consumption because the other two methods reconstruct all the examples. Due to the number of iterations required to guarantee the generation of normal examples, the DIP method takes about 900 iterations, and DIPDefend takes about 500 iterations. We allow the generated examples to contain fewer perturbations, so only 400 iterations are needed for large-perturbed examples, and image compression algorithms are used for small-perturbed examples. Therefore, compared with other methods, the proposed method has certain advantages in defense accuracy and computational cost.

Table 3. Compared with the running time and memory consumption of other existing image reconstruction methods.

Method	Run-time (s)	Memory (MB)
DIP [29]	52.896	153.9
DIPDefend [3]	31.038	146.8
FADDefend	13.616	121.1

5 Conclusion and Future Work

This paper proposes an effective adversarial defense method. Different defenses are performed according to the perturbation grading strategy by evaluating the adversarial perturbation strength. The proposed method achieves better defensive performance against the ImageNet dataset under different perturbation strengths and attack algorithms, as well as in cross-model attacks. Because this classification strategy ensures defense accuracy, it also improves the running speed and reduces the computational cost.

Since the proposed method uses mirror flip, it is well defensive against real-world adversarial examples. However, in special scenarios such as numbers, mirror flip leads to semantic information changes, which need improvement.

Acknowledgments. This research work is partly supported by National Natural Science Foundation of China No. 62172001 and No. 61860206004 Guangdong Basic and Applied Basic Research Foundation (Grant 2019B151502001) and Shenzhen R&D Program (Grant JCYJ20200109105008228).

References

1. Athalye, A., Carlini, N., Wagner, D.: Obfuscated gradients give a false sense of security: circumventing defenses to adversarial examples. In: International Conference on Machine Learning, pp. 274–283 (2018)
2. Carlini, N., Wagner, D.: Towards evaluating the robustness of neural networks. In: IEEE Symposium on Security and Privacy, pp. 39–57 (2017)
3. Dai, T., Feng, Y., Chen, B., Lu, J., Xia, S.T.: Deep image prior based defense against adversarial examples. *Pattern Recogn.* **122**, 108249 (2022)
4. Ding, L., et al.: Delving into deep image prior for adversarial defense: a novel reconstruction-based defense framework. In: Proceedings of the 29th ACM International Conference on Multimedia, pp. 4564–4572 (2021)
5. Dong, Y., et al.: Boosting adversarial attacks with momentum. In: Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition, pp. 9185–9193 (2018)
6. Dong, Y., Pang, T., Su, H., Zhu, J.: Evading defenses to transferable adversarial examples by translation-invariant attacks. In: Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition, pp. 4312–4321 (2019)
7. Dong, Y., et al.: Efficient decision-based black-box adversarial attacks on face recognition. In: Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition, pp. 7714–7722 (2019)
8. Ganin, Y.: Domain-adversarial training of neural networks. *J. Mach. Learn. Res.* **17**(1), 2030–2096 (2016)
9. Goodfellow, I.J., Shlens, J., Szegedy, C.: Explaining and harnessing adversarial examples. In: International Conference on Learning Representations (2014)
10. He, K., Zhang, X., Ren, S., Sun, J.: Identity mappings in deep residual networks. In: Leibe, B., Matas, J., Sebe, N., Welling, M. (eds.) ECCV 2016. LNCS, vol. 9908, pp. 630–645. Springer, Cham (2016). https://doi.org/10.1007/978-3-319-46493-0_38
11. Ji, W., et al.: Learning calibrated medical image segmentation via multi-rater agreement modeling. In: Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition, pp. 12341–12351 (2021)

12. Jia, X., Wei, X., Cao, X., Foroosh, H.: Comdefend: an efficient image compression model to defend adversarial examples. In: Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition, pp. 6084–6092 (2019)
13. Kurakin, A., Goodfellow, I.J., Bengio, S.: Adversarial examples in the physical world. In: Artificial Intelligence Safety and Security, pp. 99–112. Chapman and Hall/CRC (2018)
14. Lanchantin, J., Wang, T., Ordonez, V., Qi, Y.: General multi-label image classification with transformers. In: Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition, pp. 16478–16488 (2021)
15. Liao, F., Liang, M., Dong, Y., Pang, T., Hu, X., Zhu, J.: Defense against adversarial attacks using high-level representation guided denoiser. In: Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition, pp. 1778–1787 (2018)
16. Liu, X., Tanaka, M., Okutomi, M.: Single-image noise level estimation for blind denoising. *IEEE Trans. Image Process.* **22**(12), 5226–5237 (2013)
17. Liu, Z., et al.: Feature distillation: DNN-oriented jpeg compression against adversarial examples. In: Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition, pp. 860–868 (2019)
18. Madry, A., Makelov, A., Schmidt, L., Tsipras, D., Vladu, A.: Towards deep learning models resistant to adversarial attacks. In: International Conference on Learning Representations (2018)
19. Moosavi-Dezfooli, S.M., Fawzi, A., Frossard, P.: Deepfool: a simple and accurate method to fool deep neural networks. In: Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition, pp. 2574–2582 (2016)
20. Papernot, N., McDaniel, P., Jha, S., Fredrikson, M., Celik, Z.B., Swami, A.: The limitations of deep learning in adversarial settings. In: IEEE European Symposium on Security and Privacy, pp. 372–387 (2016)
21. Papernot, N., McDaniel, P., Wu, X., Jha, S., Swami, A.: Distillation as a defense to adversarial perturbations against deep neural networks. In: IEEE Symposium on Security and Privacy, pp. 582–597 (2016)
22. Quinonez, R., Safaoui, S., Summers, T., Thuraisingham, B., Cardenas, A.A.: Shared reality: detecting stealthy attacks against autonomous vehicles. In: Proceedings of the 2th Workshop on CPS&IoT Security and Privacy, pp. 15–26 (2021)
23. Ronneberger, O., Fischer, P., Brox, T.: U-net: convolutional networks for biomedical image segmentation. In: International Conference on Medical Image Computing and Computer-Assisted Intervention, pp. 234–241 (2015)
24. Samangouei, P., Kabkab, M., Chellappa, R.: Defense-GAN: protecting classifiers against adversarial attacks using generative models. In: International Conference on Learning Representations (2018)
25. Shafahi, A., et al.: Adversarial training for free! In: Advances in Neural Information Processing Systems, vol. 32 (2019)
26. Simonyan, K., Zisserman, A.: Very deep convolutional networks for large-scale image recognition. In: International Conference on Learning Representations (2015)
27. Srivastava, S., Narayan, S., Mittal, S.: A survey of deep learning techniques for vehicle detection from UAV images. *J. Syst. Architect.* **117**, 102–152 (2021)
28. Szegedy, C., et al.: Intriguing properties of neural networks. In: International Conference on Learning Representations (2014)
29. Ulyanov, D., Vedaldi, A., Lempitsky, V.: Deep image prior. In: Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition, pp. 9446–9454 (2018)

30. Wang, J., Yin, Z., Jiang, J., Du, Y.: Attention-guided black-box adversarial attacks with large-scale multiobjective evolutionary optimization. *Int. J. Intell. Syst.* **37**(10), 7526–7547 (2022)
31. Wang, L.Y.: Adversarial perturbation suppression using adaptive gaussian smoothing and color reduction. In: *International Symposium on Multimedia*, pp. 158–165 (2021)
32. Yin, Z., Wang, H., Wang, J.: War: An efficient pre-processing method for defending adversarial attacks. In: *International Conference on Machine Learning for Cyber Security*, pp. 514–524 (2020)
33. Yin, Z., Wang, H., Wang, J., Tang, J., Wang, W.: Defense against adversarial attacks by low-level image transformations. *Int. J. Intell. Syst.* **35**(10), 1453–1466 (2020)