

Universal adversarial perturbation for remote sensing images

Qingyu Wang

*School of Computer Science and Technology
Anhui University
Heifei, China
e20301182@stu.ahu.edu.cn*

Guorui Feng

*School of Communication & Electronic Engineering
Shanghai University
Shanghai, China
grfeng@shu.edu.cn*

Zhaoxia Yin

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

Bin Luo

*School of Computer Science and Technology
Anhui University
Heifei, China
luobin@ahu.edu.cn*

Abstract—Recently, with the application of deep learning in the remote sensing image (RSI) field, the classification accuracy of the RSI has been dramatically improved compared with traditional technology. However, even the state-of-the-art object recognition convolutional neural networks are fooled by the universal adversarial perturbation (UAP). The research on UAP is mostly limited to ordinary images, and RSIs have not been studied. To explore the basic characteristics of UAPs of RSIs, this paper proposes a novel method combining an encoder-decoder network with an attention mechanism to generate the UAP of RSIs. Firstly, the former is used to generate the UAP, which can learn the distribution of perturbations better, and then the latter is used to find the sensitive regions concerned by the RSI classification model. Finally, the generated regions are used to fine-tune the perturbation making the model misclassified with fewer perturbations. The experimental results show that the UAP can make the classification model misclassify, and the attack success rate of our proposed method on the RSI data set is as high as 97.09%.

Index Terms—Remote sensing image, deep learning, universal adversarial perturbations, encoder-decoder, attention mechanism

I. INTRODUCTION

The development of surface observation instruments has promoted the rapid development of remote sensing image (RSI) technology and made the RSI scene classification widely used in coverage management, geographic spatial target detection, urban planning, and other research [1, 2].

Due to the small number, low resolution, and lack of diversity of RSIs, the classification accuracy of traditional feature extraction methods [3?] is limited. Due to the strong learning ability of deep learning, many excellent deep algorithms have been applied to remote sensing sectors such as convolutional neural networks and achieved significant results. [4] used the neural network to extract the image information, selected the consolidated full connection layer to construct the final RSI scene, and used discriminant correlation analysis for feature fusion. The problems of extracting space-spectral features and

classification model overfitting were solved by using 3D convolutional neural networks, regularization, dropout, and virtual examples [5]. [6] proposed a scene classification network architecture based on multi-objective neural evolution which can extract information of the RSI more flexible. The above classification methods based on the neural network dramatically have improved the classification accuracy and speed of the RSI classification model.

At present, many studies have shown that deep learning has serious security problems. Such as, the adversarial perturbation that could make the classification model misclassify clean examples was first discovered by [7]. [7] also proposed a method to generate adversarial perturbation by calculating the backpropagation values of gradients. [8] improved the attack success rate (ASR) by increasing the number of computing the gradient and reducing the attack step. [9] calculated the distance between clean examples and classification boundaries to improve the visual quality of adversarial examples. [10] proposed an attack method based on optimization, which considered both high ASR and minor perturbation. However, the above attack methods can only calculate a single perturbation at a time, which takes a long time and only has a high ASR on the white box model with known network parameters.

Aiming at the problems existing in the above attack methods, [11] first proposed the universal adversarial perturbation (UAP) that image-agnostic perturbations by superposing perturbation generated by [9]. UAP is generated from multiple images and is applied to more. UAP can maintain good generalization ability that can well attack black-box models with unknown parameters. [12] used spatial transformation [13] and image-to-image translation adversarial networks [14] to propose an additive and non-additive universal perturbations generation method. However, the above methods of generating UAP have the problem of poor visual quality in RSIs.

Many studies found that even RSIs that differ from ordinary images in shooting angle and resolution also face the same

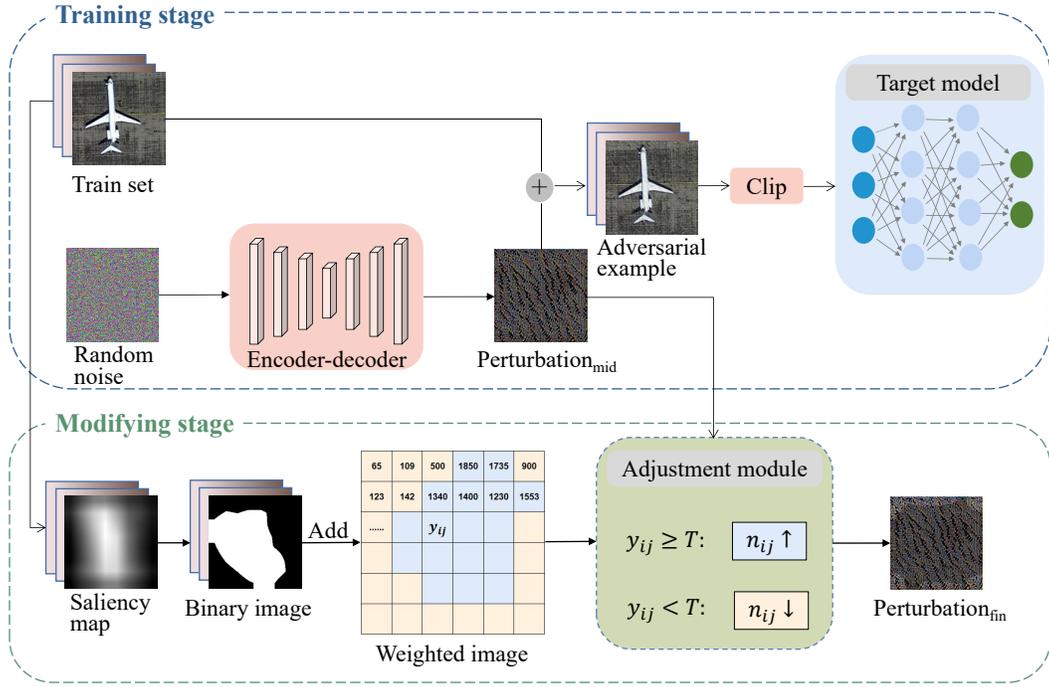


Fig. 1: Overview of generalized universal adversarial perturbation. ‘ \uparrow ’ represents that the value of n_{ij} increases and ‘ \downarrow ’ represents that the value of n_{ij} decreases.

security problems. [15] first applied the attack method on ordinary images to the RSI classification model and verified that the RSI classification models can also misclassify. [16] showed the performance of remote sensing adversarial examples under different models and data sets and found that most of the adversarial examples were wrongly divided into several specific classes, which was called attack selectivity. [17] proposed a black-box attack method using a proxy model and tested the universal of the generated adversarial examples in different models. [18] verified that synthetic aperture radar images were also affected by adversarial attacks. [19] used the method [11] to verify that synthetic aperture radar images were also affected by UAP. The above research shows that adversarial attacks can cause misclassification of the RSI classification model. However, the UAP in remote sensing images has not been studied.

To verify the effect of the UAP on the RSI classification model and whether remote sensing universal adversarial examples has the attack selectivity, a method of generating remote sensing UAP using an encoder-decoder network [20] is proposed. To better find the characteristics of the regions concerned by the RSI classification models, a saliency map [21] is used. The experimental results of our proposed method, which can achieve an ASR of 97.09%, show that the RSI is also affected by UAP.

The main contributions of our paper are:

1. An encode-decode network, which can improve the generation efficiency, is introduced to generate the universal perturbation.

2. A saliency map, which can find the sensitive region of the classification model, is used to modify the universal perturbation.
3. We verify that universal adversarial examples of RSI have the attack selectivity.

II. OUR PROPOSED METHOD

The perturbation generated needs to mislead the classification model accurately to generate the UAP with high ASR and minor perturbation. Based on this, we propose a method to attack accurately by training and modifying UAP. In the training stage, an encoder-decoder network and a classification model are used to train the perturbation. An encoder-decoder network is used to ensure that the input and output are consistent. The classification model is used to update the perturbation. To reduce the perturbation of remote sensing UAP, a saliency map is used to fine-tune the generated perturbations to improve their ASR and visual quality.

A. Training stage

Fig. 1 is flow charts for generating UAP. In the training stage, random noise $z \sim N(0, 1)$ is an input. The generator includes multi-layer convolution, pooling, and upsampling operations, ensuring better high-dimensional feature extraction. Perturbation is added to each clean example to get the adversarial example which is then clipped. Next, the target model is used to predict it. For target model $C_\theta(x)$ with parameter θ , the model can correctly identify the clean example x if $C_\theta(x) = c$, where c is the correct label of the corresponding example. When the example is added with perturbation δ , the

model will misclassify the example if $C_\theta(x + \delta) \neq c$. The UAP is to find a perturbation v that many clean examples satisfies the formula $C_\theta(x + v) \neq c$. This paper aims to find a perturbation v that misclassifies most examples. The training of universal perturbations with an encoder-decoder network can be summarized by

$$\min_{\theta_1} E(x, y) \sim D\left[\frac{\max}{\|\delta\|_\infty \leq \epsilon} L(C_\theta(x + \delta), y)\right], \quad (1)$$

where y is the example label, and $L(C_\theta(x + \delta), y)$ is the loss function, i.e., cross-entropy loss. $\max(L)$ is the optimization objective, that is, to find the perturbation bounded by an infinite norm that maximizes the loss function. The outer layer is the minimum formula for optimizing the encoder-decoder network with parameter θ_1 . When the model is fixed, we minimize the perturbation to make the model classification errors.

The loss function used to optimize the encode-decode network in this paper is

$$L = -\frac{1}{N} \sum_i \sum_{c=1}^M y_{ic} \log(P_{ic}), \quad (2)$$

where N is the total number of examples. M represents the number of categories, $y_{ic} = 1$ if the true category of example i is equal to c , otherwise 0, P_{ic} is the predicted probability that the example i belongs to category c .

B. Modifying stage

In the modifying stage, firstly, we generate the saliency map, which represents the region of concern for the model of the training set. The reason why the saliency map is used to modify the perturbation is saliency map can find the concern region of the model, which is more sensitive to the perturbation than the other region. In order to better the statistical model for each sample area of concern, the saliency

map is binarized. Next, a weighted image of the model concern region is obtained by adding all binary images. The weighted image is divided into two parts by setting a threshold of T .

An adjustment module is designed to modify the perturbation $_{mid}$ according to the weighted image, which is modified by

$$n_{ij} = \begin{cases} n_{ij} * \alpha & \text{if } y_{ij} \geq T \\ n_{ij} * \beta & \text{if } y_{ij} < T \end{cases}, \quad (3)$$

where y_{ij} represents the value of the weighted image in position (i, j) and n_{ij} represents the value of perturbation $_{mid}$ in position (i, j) . The parameter of α and β is a positive number, which is $\alpha > 1$ and $\beta < 1$. The value of the parameter is different due to the different areas of concern in different models. When the value y_{ij} greater than T , the value n_{ij} is increased; otherwise, n_{ij} is decreased. Finally, the perturbation $_{fin}$ is generated, which is used to generate the universal adversarial examples.

III. EXPERIMENTS

The experiments are carried out on the remote sensing data set PatternNet [22] with 38 classes. The data set has 38 classes and each includes 800 images with a size of 256×256 . We randomly select 50 images for each class as the training set and 9158 as the validation set. The latest classification models of VGG16 [23], VGG19 [23], ResNet34 [24] and ResNet101 [24] are trained to evaluate the attack effect of our proposed method. The classification accuracy of the four models is 94.92%, 94.18%, 99.72% and 99.75%, respectively. An infinite norm bounds the UAP in training, such as $\|\delta\|_\infty \leq \epsilon = 10$. The parameters of the learning rate and weight decay are 0.001, and the number of iterations is 50. Fig. 2 illustrates some clean and corresponding adversarial examples. For example, the label in Fig. 2 (a) changes from airplane to (Fig. 2 (e)) storage_tank when a universal perturbation is added. It can be seen that the universal adversarial perturbation is imperceptibility. The

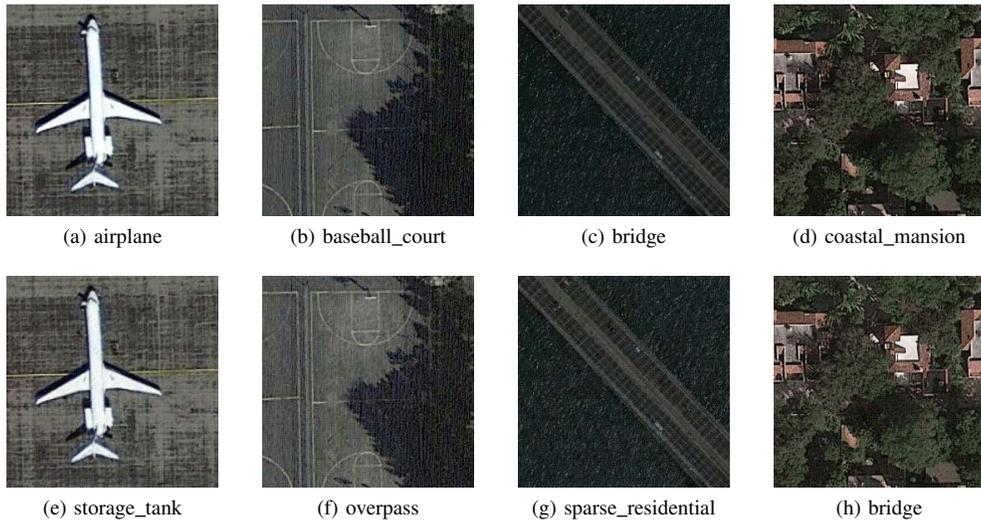


Fig. 2: Some clean and adversarial examples in VGG16. (a)~(d) are clean examples, (e)~(h) are adversarial examples.

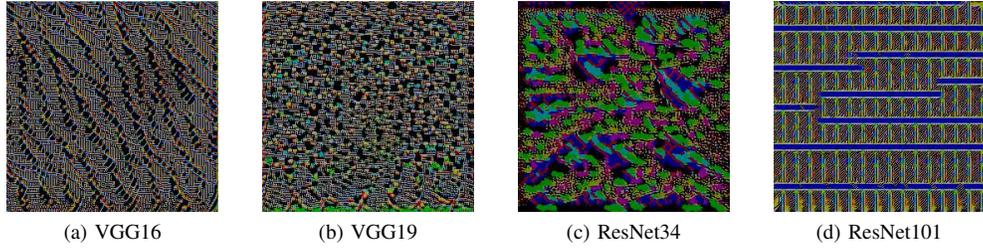


Fig. 3: Universal perturbations computed for different network architectures. The pixel values are scaled for visibility.

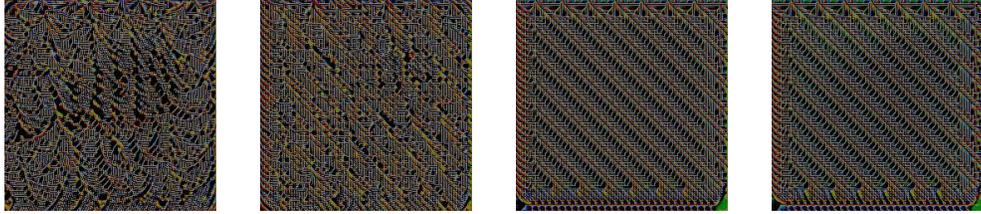


Fig. 4: Diversity of universal perturbations for the VGG16 network. The four perturbations are generated using different random shufflings of the training set.

universal perturbations of different networks are visualized in Fig. 3. It shows that the universal perturbations are not the same in different networks. In Fig. 4, four different universal perturbations obtained by using random shufflings in the training set are visualized. It can be seen that such universal perturbations are different, although they are generated by the same networks and constraints.

A. The performance of the ASR and perturbation magnitude (PM)

The evaluation criteria are the ASR, which refers to the proportion of adversarial examples that the model misclassifies, and PM, which is the difference between the adversarial and clean examples. We compare ASR and PM on the validation set with the other methods. The experiment results are shown in Table I. It can find that our proposed method improves the ASR and reduces the PM. This is because an encoder-decoder network and a saliency map are used to mislead the classification model accurately. In addition, we test the effectiveness of the saliency map used in our proposed method. The performance before and after perturbations modification is shown in Table II. It can be seen that the saliency map can improve the visual quality and ASR of the adversarial examples.

To demonstrate the effectiveness of our method, we show the ASR under different perturbations in Fig. 5. These perturbations are calculated by averaged l_2 distances in the validation set. Different perturbation norms are achieved by scaling each perturbation to have the target norm. It can be seen from Fig. 5 that the overall attack effect of our proposed method is higher compared to the other methods under the same perturbations.

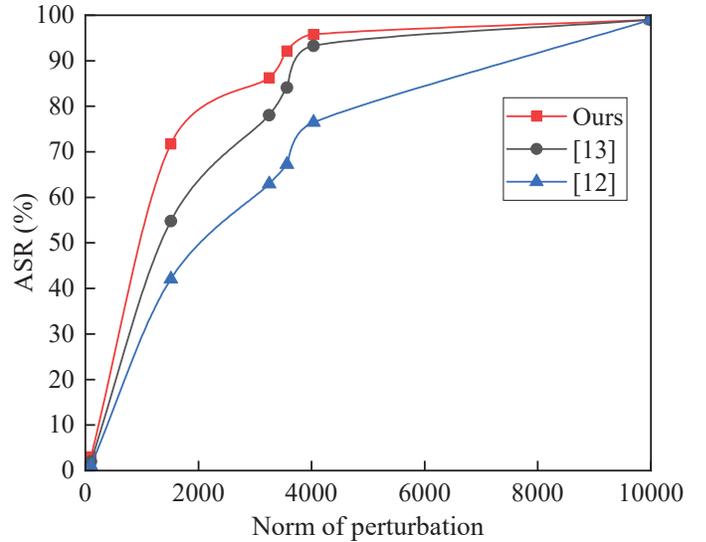


Fig. 5: Comparison of ASR (%) under different perturbations on the VGG16 network.

B. Generalization ability between models

The generalization ability of our proposed method with and without a saliency map is tested to show the attack ability of the UAP between different models, as shown in Table III. Each row represents the ASR of the perturbation generated in the target model on the other models, and each column represents a different target model. As in Table III, our proposed method preserves a good ASR between different models. Such as, the UAP generated by VGG16 can achieve an ASR of over 90% on VGG19 and over 85% on ResNet34 and ResNet101. It also can be seen that the generalization ability of the same model

TABLE I: Comparison of ASR (%) and PM (l_2 norm) with other methods on the Patternet dataset.

		VGG16	VGG19	ResNet34	ResNet101
ASR	[11]	73.27	66.35	82.91	89.29
	[12]	95.40	93.31	96.22	96.27
	Ours	95.75	94.74	96.26	97.09
PM	[11]	3348.64	3422.67	4015.34	4322.52
	[12]	4126.38	4247.44	4042.89	4209.33
	Ours	3375.28	3288.22	3845.84	3878.72

TABLE II: The ablation experiment with and without saliency map in our proposed method.

		VGG16	VGG19	ResNet34	ResNet101
ASR	Without	95.72	94.41	95.82	97.02
	With	95.75	94.74	95.96	97.09
PM	Without	4034.78	4053.13	3665.37	4159.35
	With	3375.28	3288.22	3512.71	3878.72

architecture is better, such as the perturbation generated by the VGG16 model has a higher ASR in the VGG19 model. It also can be seen that saliency maps has little effect on the generalization ability, which reduce the required PM.

C. Attack selectivity

The paper [25] proposed a soft-threshold defense against adversarial examples method based on attack selectivity. To verify that UAP generated by the proposed method has attack selectivity, we make statistics on the classification results of adversarial examples, as shown in Fig. 6. It can be seen that it will cause attack selectivity when the perturbation is added to the examples. For example, on the VGG19 model, over 88% of the validation set, which includes 38 classes, is misclassified into three classes. At the same time, under the ResNet34, VGG16, and VGG19 models, the proportion of examples in certain classes has exceeded 80%, and ResNet101 has exceeded 89%, equivalent to achieving a targeted attack.

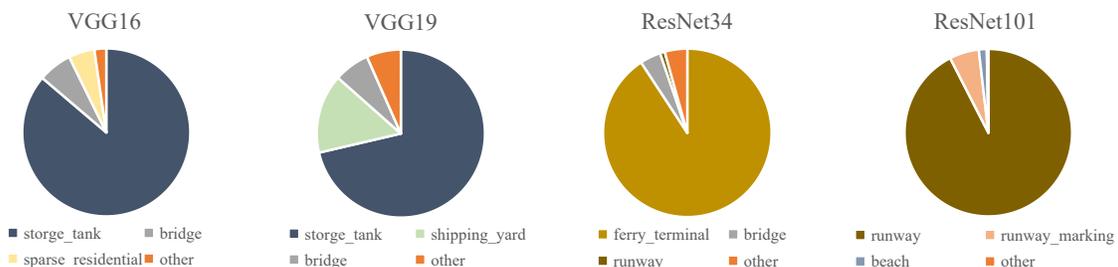


Fig. 6: A demonstration of attack selectivity on the PatternNet dataset. The pie chart represents the classification result distribution of adversarial examples.

TABLE III: Generalization ability representing the ASR (%) in different models of our proposed method with and without saliency map. The first row is attacked models, and the first column is the targeted models used to generate adversarial perturbations.

		VGG16	VGG19	ResNet34	ResNet101
VGG16	Without	95.72	92.53	86.74	86.74
	With	95.75	92.78	85.93	86.76
VGG19	Without	93.47	94.41	88.60	83.94
	With	93.46	94.74	88.03	84.23
ResNet34	Without	76.49	73.10	95.82	86.21
	With	76.10	72.30	96.26	86.31
ResNet101	Without	84.45	82.18	90.97	97.02
	With	84.22	81.67	91.13	97.09

IV. CONCLUSION

This paper proposes a UAP generation method for the first time on RSIs. The method utilizes an encoder-decoder network and saliency map to generate the universal perturbation. The former uses the training example to generate the perturbation, while the latter uses the saliency map to modify the perturbation in the sensitive area of the classification model. We apply the universal perturbation attack methods on ordinary images to RSIs and compare it with our proposed method. Experimental results show that our proposed method improves the ASR while reducing the PM. Furthermore, we also verify that it has good generalization ability and attack selectivity. Our proposed method in generating remote sensing universal perturbation has effectively achieved a high ASR and minor PM. However, the use of a saliency map to modify the perturbation is instability, so we will improve the stability of the saliency map to improve the ASR and PM in the future work.

REFERENCES

- [1] N. Kussul, M. Lavreniuk, S. Skakun, and A. Shelestov, "Deep learning classification of land cover and crop types using remote sensing data," *IEEE Geoscience and Remote Sensing Letters*, vol. 14, no. 5, pp. 778–782, 2017.

- [2] A. Manno-Kovacs and T. Sziranyi, "Orientation-selective building detection in aerial images," *ISPRS Journal of Photogrammetry and Remote Sensing*, vol. 108, pp. 94–112, 2015.
- [3] A. Oliva and A. Torralba, "Modeling the shape of the scene: A holistic representation of the spatial envelope," *International journal of computer vision*, vol. 42, no. 3, pp. 145–175, 2001.
- [4] S. Chaib, H. Liu, Y. Gu, and H. Yao, "Deep feature fusion for vhr remote sensing scene classification," *IEEE Transactions on Geoscience and Remote Sensing*, vol. 55, no. 8, pp. 4775–4784, 2017.
- [5] Y. Chen, H. Jiang, C. Li, X. Jia, and P. Ghamisi, "Deep feature extraction and classification of hyperspectral images based on convolutional neural networks," *IEEE Transactions on Geoscience and Remote Sensing*, vol. 54, no. 10, pp. 6232–6251, 2016.
- [6] A. Ma, Y. Wan, Y. Zhong, J. Wang, and L. Zhang, "Scenenet: Remote sensing scene classification deep learning network using multi-objective neural evolution architecture search," *ISPRS Journal of Photogrammetry and Remote Sensing*, vol. 172, pp. 171–188, 2021.
- [7] I. J. Goodfellow, J. Shlens, and C. Szegedy, "Explaining and harnessing adversarial examples," *arXiv preprint arXiv:1412.6572*, 2014.
- [8] A. Kurakin, I. J. Goodfellow, and S. Bengio, "Adversarial examples in the physical world," in *Artificial intelligence safety and security*. Chapman and Hall/CRC, 2018, pp. 99–112.
- [9] S.-M. Moosavi-Dezfooli, A. Fawzi, and P. Frossard, "Deepfool: a simple and accurate method to fool deep neural networks," in *Proceedings of the IEEE conference on computer vision and pattern recognition*, 2016, pp. 2574–2582.
- [10] N. Carlini and D. Wagner, "Towards evaluating the robustness of neural networks," in *2017 IEEE Symposium on Security and Privacy (SP)*. IEEE, 2017, pp. 39–57.
- [11] S.-M. Moosavi-Dezfooli, A. Fawzi, O. Fawzi, and P. Frossard, "Universal adversarial perturbations," in *Proceedings of the IEEE conference on computer vision and pattern recognition*, 2017, pp. 1765–1773.
- [12] Y. Zhang, W. Ruan, F. Wang, and X. Huang, "Generalizing universal adversarial attacks beyond additive perturbations," in *2020 IEEE International Conference on Data Mining (ICDM)*. IEEE, 2020, pp. 1412–1417.
- [13] C. Xiao, J. Y. Zhu, B. Li, W. He, M. Liu, and D. Song, "Spatially transformed adversarial examples," in *6th International Conference on Learning Representations, ICLR 2018*, 2018.
- [14] C. Xiao, B. Li, J. Y. Zhu, W. He, M. Liu, and D. Song, "Generating adversarial examples with adversarial networks," in *27th International Joint Conference on Artificial Intelligence, IJCAI 2018*. International Joint Conferences on Artificial Intelligence, 2018, pp. 3905–3911.
- [15] W. Czaja, N. Fendley, M. Pekala, C. Ratto, and I.-J. Wang, "Adversarial examples in remote sensing," in *Proceedings of the 26th ACM SIGSPATIAL International Conference on Advances in Geographic Information Systems*, 2018, pp. 408–411.
- [16] Y. Xu, B. Du, and L. Zhang, "Assessing the threat of adversarial examples on deep neural networks for remote sensing scene classification: Attacks and defenses," *IEEE Transactions on Geoscience and Remote Sensing*, vol. 59, no. 2, pp. 1604–1617, 2020.
- [17] Y. Xu and P. Ghamisi, "Universal adversarial examples in remote sensing: Methodology and benchmark," *IEEE Transactions on Geoscience and Remote Sensing*, vol. 60, pp. 1–15, 2022.
- [18] L. Chen, Z. Xu, Q. Li, J. Peng, S. Wang, and H. Li, "An empirical study of adversarial examples on remote sensing image scene classification," *IEEE Transactions on Geoscience and Remote Sensing*, vol. 59, no. 9, pp. 7419–7433, 2021.
- [19] L. Wang, X. Wang, S. Ma, and Y. Zhang, "Universal adversarial perturbation of sar images for deep learning based target classification," in *2021 IEEE 4th International Conference on Electronics Technology (ICET)*. IEEE, 2021, pp. 1272–1276.
- [20] V. Badrinarayanan, A. Kendall, and R. Cipolla, "Segnet: A deep convolutional encoder-decoder architecture for image segmentation," *IEEE transactions on pattern analysis and machine intelligence*, vol. 39, no. 12, pp. 2481–2495, 2017.
- [21] R. R. Selvaraju, M. Cogswell, A. Das, R. Vedantam, D. Parikh, and D. Batra, "Grad-cam: Visual explanations from deep networks via gradient-based localization," in *Proceedings of the IEEE international conference on computer vision*, 2017, pp. 618–626.
- [22] W. Zhou, S. Newsam, C. Li, and Z. Shao, "Patternnet: A benchmark dataset for performance evaluation of remote sensing image retrieval," *ISPRS journal of photogrammetry and remote sensing*, vol. 145, pp. 197–209, 2018.
- [23] K. Simonyan and A. Zisserman, "Very deep convolutional networks for large-scale image recognition," *arXiv preprint arXiv:1409.1556*, 2014.
- [24] K. He, X. Zhang, S. Ren, and J. Sun, "Deep residual learning for image recognition," in *Proceedings of the IEEE conference on computer vision and pattern recognition*, 2016, pp. 770–778.
- [25] L. Chen, J. Xiao, P. Zou, and H. Li, "Lie to me: A soft threshold defense method for adversarial examples of remote sensing images," *IEEE Geoscience and Remote Sensing Letters*, vol. 19, pp. 1–5, 2021.