# Universal Adversarial Defense in Remote Sensing Based on Pre-trained Denoising Diffusion Models

Weikang Yu, *Student Member, IEEE,* Yonghao Xu, *Member, IEEE,* Pedram Ghamisi, *Senior Member, IEEE*

*Abstract*—Deep neural networks (DNNs) have achieved tremendous success in many remote sensing (RS) applications, in which DNNs are vulnerable to adversarial perturbations. Unfortunately, current adversarial defense approaches in RS studies usually suffer from performance fluctuation and unnecessary re-training costs due to the need for prior knowledge of the adversarial perturbations among RS data. To circumvent these challenges, we propose a universal adversarial defense approach in RS imagery (UAD-RS) using pre-trained diffusion models to defend the common DNNs against multiple unknown adversarial attacks. Specifically, the generative diffusion models are first pre-trained on different RS datasets to learn generalized representations in various data domains. After that, a universal adversarial purification framework is developed using the forward and reverse process of the pre-trained diffusion models to purify the perturbations from adversarial samples. Furthermore, an adaptive noise level selection (ANLS) mechanism is built to capture the optimal noise level of the diffusion model that can achieve the best purification results closest to the clean samples according to their Frechet Inception Distance (FID) in deep feature space. As a result, only a single pre-trained diffusion model is needed for the universal purification of adversarial samples on each dataset, which significantly alleviates the re-training efforts and maintains high performance without prior knowledge of the adversarial perturbations. Experiments on four heterogeneous RS datasets regarding scene classification and semantic segmentation verify that UAD-RS outperforms state-of-the-art adversarial purification approaches with a universal defense against seven commonly existing adversarial perturbations. Codes and the pre-trained models are available online (https://github.com/EricYu97/UAD-RS).

*Index Terms*—Adversarial defense, adversarial purification, diffusion models, remote sensing, scene classification, semantic segmentation.

## I. INTRODUCTION

RECENT advances in artificial intelligence have significantly motivated the development of image processing techniques on remote sensing (RS) imagery [1]. In particular, deep learning algorithms have achieved promising results in various geoscience and RS applications, such as land use classification [2], change detection [3], and disaster monitoring [4]. Unfortunately, despite their tremendous successes, deep learning methods have shown vulnerability to adversarial samples [5]. By simply adding some mild perturbations to raw data, adversarial samples can be produced, which may possess imperceptible differences from the original data for

W. Yu and P. Ghamisi are with Helmholtz-Zentrum Dresden-Rossendorf, Helmholtz Institute Freiberg for Resource Technology, Machine Learning Group, 09599 Freiberg, Germany (e-mail: w.yu@hzde.de; p.ghamisi@hzdr.de).

Y. Xu is with the Institute of Advanced Research in Artificial Intelligence (IARAI), 1030 Vienna, Austria (e-mail: yonghaoxu@ieee.org).



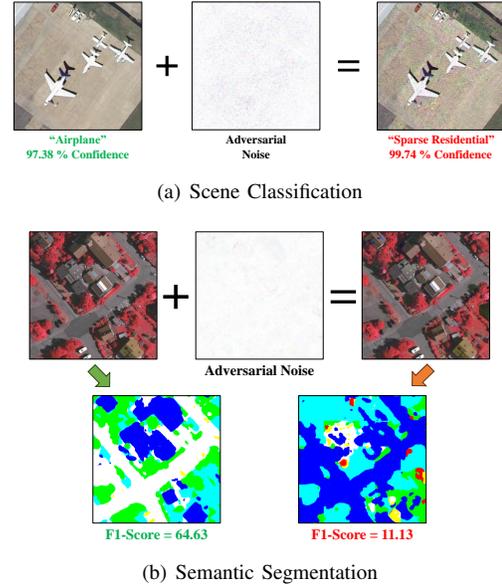(a) Scene Classification



(b) Semantic Segmentation

Fig. 1. Illustration of the adversarial attacks on scene classification and semantic segmentation of RS images. With only tiny adversarial noises added to the samples, the deep neural networks can be fooled to make completely different predictions for the images that looks nearly the same.

human observations but could fool the deep neural networks (DNNs) into making wrong predictions with high confidence, as shown in Fig. 1.

Research on adversarial samples has attracted increasing attention over the decades in computer vision society. These studies based on the game theory can be divided into two adversaries efficiently generating adversarial samples (i.e., attack) and guarding the algorithms against the attacks (i.e., defense) [6]. Debates on adversarial attack and defense have also spread to the field of geoscience and RS. In recent literature, RS researchers have also found the existence of adversarial samples which are generated based on optical data [7], LiDAR point cloud [8], or even synthetic aperture radar (SAR) data [9], [10]. Moreover, the threat of these adversarial samples in RS imagery has been emphasized that images with slight perturbations can be misclassified by the state-of-the-art DNNs into wrong categories with very high confidence [11], [12]. Since most geoscience and RS tasks are highly safety-critical, adversarial defense methods are vitally desired to improve the robustness of the deployed deep learning algorithms against adversarial samples [13]. Generally speaking, adversarial defense methods that aim to increase the model's robustness can be categorized into adversarial training, randomization,

detection, and adversarial purification. Adversarial training methods aim to improve the model's resistibility against adversarial samples by directly bringing them into the training process. For example, Goodfellow et al. [14] developed an empirical adversarial training method by generating adversarial samples with the FGSM attack and adding them back to the training set to improve the robustness of the classification model. Consequently, a min-max formulation was proposed to iteratively generate adversarial samples in the training phase, which improves the model's robustness against more powerful attack methods [15]. To reduce the efforts of labelling the additional adversarial samples in adversarial training, Xu et al. [16] developed a novel adversarial self-supervised learning framework to learn a robust pre-trained model with unlabeled adversarial samples for RS scene classification. Nevertheless, adversarial training methods failed to improve the inherent robustness of the DNNs, making the model even more fragile to newly generated adversarial samples with unseen types of adversarial attacks [17]. Even worse, generating adversarial samples in training are considered to be expensive, especially on large-scale datasets such as ImageNet [18].

Another group of adversarial defense methods incorporates randomness to defend against adversarial samples. Since adversarial perturbation can be viewed as noise, random components are introduced to improve the model's robustness in randomness methods. In particular, Cohen et al. [19] proposed a randomized smoothing technique to achieve adversarially robust classification by randomizing the input of a neural network that removes the potential adversarial perturbation. Apart from the model's input, Gaussian noise can be added to the hidden layer output or directly injected to each layer on activation and weights to introduce the randomness into the DNNs [20], [21]. Additionally, the Bayesian network is another solution for the randomness that can be leveraged into the DNNs to randomize the parameters of the model [22]. However, the performance of these methods is highly related to the specifically introduced randomness and can be significantly degraded due to the lack of both theoretical explanation and prior knowledge of the attack algorithms [23].

Instead of directly improving the robustness of the target DNNs, detection methods are designed to detect the existence of adversarial samples and remove them as disturbances from the mixed adversarial dataset. These methods are based on the assumption that adversarial samples and normal samples are of different domain distributions and, thus can be distinguished by a discriminator network. For instance, [24] proposed a soft threshold adversarial defense method incorporating the logistic regression algorithm that detects and excludes the adversarial samples to reduce the fooling rates of the attacks on RS imagery. Unfortunately, detection-based methods cannot produce correct classification results on these adversarial samples but only aims at protecting the DNNs from the mixed samples.

In contrast to the aforementioned defense methods, adversarial purification is on a different line of research that aims to purify the attacked images before sending them to target classifiers. In particular, the purification models are usually trained in an end-to-end manner independently of target DNNs without classification labels and applicable for multi kinds of adversarial attacks. As a corollary, the purification methods can combat the multi-type adversarial samples while the deployed application models remain unaffected.

With the advantage of the powerful generalization capability in specific data distribution, generative models have been the most common choice to convert adversarial samples into clean ones, which are distinguishable for the subsequent models. For example, Xue et al. [25] developed a cascaded adversarial purification method that restores the detected adversarial patches of the images by the image inpainting method under RS detection. To further improve the performance of subsequent pre-trained classifiers, Xu et al. [26] proposed a denoising network with a task-guided loss to remove the perturbations of adversarial samples for RS scene classification. Benefiting from the powerful generalization capability of generative adversarial networks (GANs), a novel perturbation-seeking GAN was developed to effectively move the adversarial samples to clean ones specified for RS scene classification [27]. Despite their success in defending against multiple kinds of adversarial attacks, current adversarial purification methods usually fall behind other defense routes [28], which is commonly due to the general shortcomings of the utilized generative models such as the mode collapse problem in GANs [29], and the low sample quality in energy-based models [30]. Therefore, it is still challenging for an adversarial purification model to remove as much of the adversarial perturbations while retaining most of the object features in the original image.

Recent studies in diffusion models have motivated the development of generative models [31]. Compared with previous image generation methods such as autoencoders and GANs, diffusion models have exhibited (1) stronger sample quality (the generated images are of higher fidelity and resemble real images more closely) and (2) mode coverage (they can generate samples that cover a broader range of possible image variations), which lead to a power-generating capability [32], [33]. In particular, diffusion models define a diffusing algorithm to first convert the data to noise by gradually adding Gaussian noise in the forward process, and then reconstructing an image by reversing the forward process with DNNs. Since the reverse process of denoising the Gaussian noise is similar to adversarial purification, diffusion models have also been introduced in adversarial defense studies. For example, Wang et al. [34] embedded purification into the diffusion-denoising process of a Denoising Diffusion Probabilistic Model (DDPM) to submerge the adversarial perturbations with gradually added Gaussian noise that is simultaneously removed following a guided denoising process. To further combat the diffusion model against strong adaptive attacks, Nie et al. [35] proposed to use the adjoint method to compute full gradients of the reverse generative process. However, the success of the existing diffusion models is highly related to the impressive generative performance of the open-source pre-trained diffusion models on several common computer vision (CV) datasets, such as ImageNet [36], CIFAR-10 [37], and CelebA-HQ [38] datasets. Apart from that, the existing diffusion model-based purification models lack the ability to automatically adapt the diffusion process for samples with different noise intensities and adversarial attack methods, which incur inevitable costs of

parameter fine-tuning for the specialists to avoid the performance fluctuating among heterogeneous adversarial samples in RS data. As a result, the implementation of these methods on RS imagery remains challenging due to the shortage of the available pre-trained model on RS datasets and the adaptive diffusion algorithms to deal with a wide variety of adversarial samples with different perturbations.

Motivated by the aforementioned challenges, this work proposes a universal adversarial defense method using pre-trained diffusion models for the purification of adversarial samples in RS imagery. In the first stage, we pre-train the diffusion models that can reconstruct clean image inputs to improve the generation capabilities on the adversarial-free image domains. After that, a universal adversarial purification framework is proposed that utilizes the forward and reverse process of the pre-trained diffusion model to transform the multiple kinds of adversarial samples into the adversarial-free domain. Consequently, an adaptive noise level selection (ANSL) algorithm is designed to automatically determine the optimal noise level of the diffusion model to adapt the purification for multi-kinds of adversarial samples, which considers the diversity of the attack perturbations and specific features of RS data, such as the heterogeneous spectrum and resolution. The ANSL can further achieve the best performance of the purification by allowing the model to remove adversarial perturbations thoroughly while preserving the distinguishable label semantics as much as possible.

Extensive experiments on four widely used optical RS datasets, including two scene classification datasets (i.e., UC-Merced [39], and Aerial Image Dataset (AID) [40]) and two semantic segmentation datasets (i.e., Vaihingen [41], and Zurich Summer dataset [42]) demonstrated the effectiveness of the proposed model for universal adversarial defense against multiple kinds of state-of-the-art adversarial attacks on the heterogeneous RS data. In addition, ablation studies are performed to explore the RS image synthesizing task, the cross-domain purification ability of the model, and the influence of model settings on defense performance.

The main contributions of this study are fourfold:

1) We develop a universal adversarial defense framework (UAD-RS) based on the forward and reverse process of the pre-trained diffusion models for universal adversarial defense against heterogeneous attack methods and data, firstly on RS imagery.
2) The diffusion models are pre-trained on four widely used RS datasets, including each two of scene classification and semantic segmentation, which can generate high-quality RS images in different resolutions and spectra. The study is the first one that releases the pre-trained generative diffusion models as a basis implementation to motivate further studies of the diffusion models in the field of RS datasets.
3) Since only a single pre-trained diffusion model is needed for UAD-RS to purify a variety of adversarial samples generated by different attack algorithms on various victim DNNs for each dataset, the proposed method significantly alleviates the costs and efforts of retraining a new model for each attack setting targeting the same dataset in existing approaches.
4) To maximize the performance of adversarial purification, an adaptive noise level selection (ANLS) mechanism is proposed to customize the optimal inference hyper-parameters of the pre-trained diffusion model for defending the specified victim DNN against diverse RS data and adversarial attack algorithms. The ANLS explicitly considers the effects of diffusion steps and the intensity of noise added throughout the entire diffusion process to restore an adversarial-free image with the best quality.
5) Extensive experiments on four RS datasets regarding both scene classification and semantic segmentation demonstrate that the UAD-RS outperforms the state-of-the-art methods for universal adversarial defense against seven common adversarial attacks with only one pre-trained diffusion model for each dataset in RS imagery.

The article is organized as follows. Section II provides a brief review of adversarial attacks, adversarial purification, and diffusion models. Section III describes the proposed adversarial defense methodology based on pre-trained diffusion models and an adaptive noise level selection mechanism. The performance of generative diffusion models and the adversarial purification framework is experimentally evaluated in Section IV, while the experimental results are thoroughly analyzed. In Section V, some key concerns regarding this paper are discussed along with the limitations of the generative diffusion-based adversarial defense method, and several potential avenues for future research are identified. Finally, Section VI concludes this article.

## II. RELATED WORK

### A. Adversarial Attacks

*1) Fast Gradient Sign Method (FGSM):* One of the most intuitive strategies of adversarial attack is by leveraging the way they learn gradients. Based on this idea, a gradient-based attack, namely, Fast Gradient Sign Method (FGSM), was proposed to adjust the input data to maximize the objective function by adjusting the weights based on the back-propagated gradients [14]. Given an image $x$ and its true label $y$, the adversarial sample $x_{adv}$ can be generated as:

$$x_{adv} = x + \epsilon \cdot sign(\nabla_x \mathcal{L}(\theta, x, y)), \quad (1)$$

where $sign$ denotes the sign function, $\nabla_x J(\theta, x, y)$ calculates the gradients of an objective function $\mathcal{L}(\cdot, \cdot)$ with respect to each input $x$ and ground-truth $y$, and $\epsilon$ is a scalar value that restricts the norm of the perturbation. To improve the performance of the adversarial attack, Kurakin et al. [18] proposed an iterative-FGSM (IFGSM), which applies FGSM multiple times with a small step size, as follows:

$$x^{t+1} = x^t + \alpha \cdot sign(\nabla_{x_{adv}^t} \mathcal{L}(\theta, x_{adv}^t, y)), \quad (2)$$

where an adversarial sample $x_{adv}$ is iteratively calculated in $T$ steps of FGSM, and $\alpha = \frac{\epsilon}{T}$ is the step size that reduces the perturbation of each step that sums to a similar intense of attack with the FGSM.

*2) Trade-off Projected Gradient Descent (TPGD) Attack:* Since FGSM attacks aim to optimize a loss function that measures the pixel-wise difference between probability maps of adversarial and clean images, Zhang et al. [43] proposed a Trade-off Projected Gradient Descent (TPGD) Attack that can generate the adversarial perturbations according to the Kullback–Leibler (KL) divergence between probability distributions of clean and adversarial predictions as a cross-domain covariance:

$$x^{t+1} = x^t_{adv} + \nabla_{x^t_{adv}} \mathbf{KL}(\theta, (x^t_{adv}), y). \tag{3}$$

*3) Carlini and Wagner (CW) Attack:* Instead of utilizing the direct objective functions to measure the adversarial samples, Carlini and Wagner (CW) proposed to attack the DNNs by encouraging $x_{adv}$ to have a larger probability score for a wrong class than all other classes. The CW method directly optimizes the distance between the clean and adversarial samples as follows:

$$\arg\min_{x_{adv}} ||x_{adv} - x||_\infty - \mu\mathcal{L}(\theta, x_{adv}, y), \tag{4}$$

where $\mu$ is a weighting factor.

*4) Mixcut-Attack:* Since the FGSM and its variants are kinds of white-box attack that requires complete knowledge of the objective model, it is not practical for real-world scenarios in which detailed information of the deployed model is usually impossible to obtain, especially in the field of RS [17]. As a result, black-box attacks are proposed to generate adversarial samples without knowledge of the victim model [44]. One of the most advanced back-box attacks for RS imagery is the Mixup-attack, which employs a surrogate model to produce universal adversarial examples that can deceive various heterogeneous DNNs with a high success rate [45]. Given an input $x$ with groundtruth $y$ and a surrogate model $\theta_s$, the adversarial samples can be iteratively generated as follows:

$$g^{t+1} = g^t + \frac{\nabla_x\mathcal{L}(\theta_s, x^t, y)}{||\theta_s, x^t, y)||_1}, \tag{5}$$

$$x^{t+1} = clip(x^t + \alpha \cdot \frac{g^{t+1}}{||g^{t+1}||_\infty}), \tag{6}$$

$$x_{adv} = x^T, \tag{7}$$

where $g_t$ denotes the momentum term at the $t$-th iteration, and $clip(\cdot)$ clips the pixel values in the image.

### B. Adversarial Purification

The development of generative models has significantly boosted the research on adversarial purification methods. In recent literature, there are mainly two ways to purify the perturbations from an adversarial sample, namely, denoising-based methods and generative-based methods. Benefiting from the advances of denoising models in CV, denoising-based models can effectively remove the perturbations as noises by using an end-to-end convolutional neural network (CNN). Based on this assumption, Meng and Chen [46] proposed an adversarial purification framework using a reformer network that moves the adversarial samples towards the manifold of normal examples as a defense. Specifically, the reformer

employs the framework of auto-encoder that consists of two parts of, an encoder and a decoder, which extract the high-dimension features of the image input and recover a new image, respectively. With the input of an adversarial sample, the perturbations can be removed through the encoding and decoding process, while only effective features of the clean image are potentially retained. To further improve the performance of purification, a task-guided loss is introduced to align the adversarial and normal images in the perceptual domain [26].

Instead of employing an end-to-end denoising approach for purification, generative-based methods aim to convert adversarial samples into an adversarial-free domain. Empirically, GANs are introduced to fit the clean training data distribution as the target domain, where adversarial samples can be transformed, resulting in generated images that are free from the input adversarial perturbation. For instance, Samangouei et al. [47] proposed Defense-GAN, which utilizes the expressive capability of generative models to filter out perturbations for adversarial purification. It trains a generator to learn the distribution of unperturbed images and then generates an output that closely resembles a given image but lacks any adversarial alterations during the inference stage. In generative-based strategies, purification is performed indirectly, and the specific type of adversarial attack is not explicitly specified during the training process. Consequently, it has been demonstrated that generative-based methods exhibit greater robustness compared to denoising-based methods [6].

Despite the powerful generative ability of GANs, they usually suffer from difficulty in training due to common problems such as vanishing gradients, mode collapse, and failure to converge [48], [49]. To overcome these challenges, [35] developed a novel adversarial purification algorithm based on the generative diffusion model, which converts adversarial samples into noises from which adversarial-free samples are generated in a reverse process. Wang et al. [34] further improved the diffusion model by guiding the reverse process with adversarial samples for a more substantial purification effect. Compared to the generators in GANs and autoencoders, diffusion models can preserve more local details by keeping the same latent place in the transformation process [50].

### C. Diffusion Models

Diffusion models are a class of probabilistic generative models designed for unsupervised modeling, and they have demonstrated strong sample quality and diversity in image synthesis [51]. In practical scenarios, generative diffusion models typically involve two complementary processes: a diffusion process and a reverse process, which iteratively learns the distribution of the training data [52]. Specifically, the diffusion process gradually introduces noise to the input image until it becomes Gaussian noise, while the reverse process denoises the noise iteratively to reconstruct a clean image. Let $p_{data}$ represent the distribution of all input images $x$, and $p_{latent}$ denote the latent distribution. The diffusion process $q$ with $T$
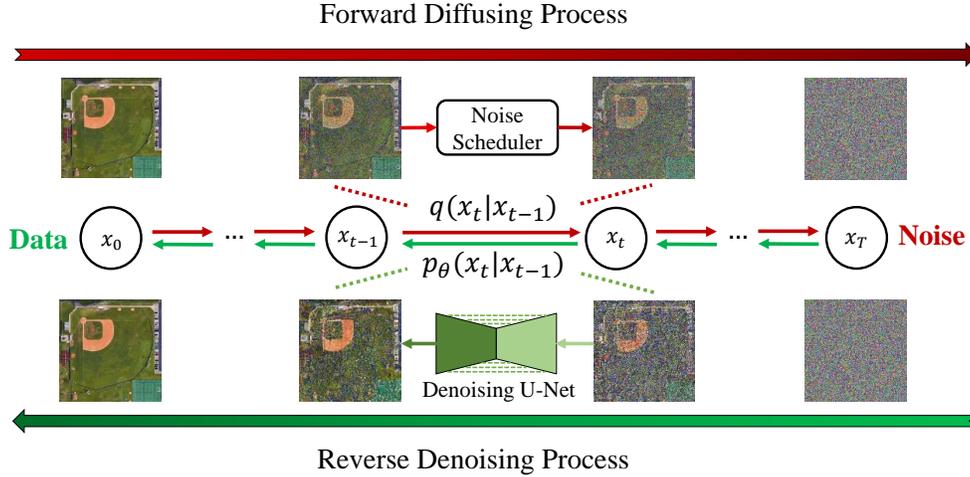
Forward Diffusing Process



Fig. 2. Illustration of the forward and reverse processes of the generative diffusion models. The forward diffusion process gradually adds Gaussian noise to the images using the noise scheduler, and finally, pure Gaussian noise is generated. After that, the reverse process progressively recovers the noise to reconstruct an image with the aid of a denoising U-Net model.

steps can be defined as follows:

$$q(x_1, ..., x_T|x_0) = \prod_{t=1}^{T} q(x_t|x_{t-1}). \quad (8)$$

In contrast, the reverse process iteratively eliminates the noise added to the diffusion process, gradually restoring a clean image. Given the latent variable $x_T$, the reverse process $p$ also consists of $T$ steps, resulting in the generation of the clean data $x_0$ as follows:

$$p_\theta(x_0, ..., x_{T-1}|x_T) = \prod_{t=1}^{T} p_\theta(x_{t-1}|x_t), \quad (9)$$

where $\theta$ parameterized the reverse diffusion process. In the literature, there are mainly two variants of diffusion models: Denoising diffusion probabilistic model (DDPM) [52] and stochastic differential equation (SDE)-based generative methods [53]. Both methods introduce Gaussian noise to the input data in the forward process but employ different denoising algorithms in the reverse process. On one hand, DDPM incorporates two Markov chains for the forward and reverse processes and aims to match the reverse transitional kernel $p_\theta (x_{t-1}|x_t)$ with the forward transitional kernel $q (x_t|x_{t-1})$ at each time step $t$ by adjusting the parameter $\theta$ in the reverse Markov chain. On the other hand, SDE-based generative methods typically utilize a noise-conditioned score network (NCSN) to estimate score functions for all noise distributions. These score functions are sequentially applied to decrease the noise levels and eventually sample a clean image. Intuitively, DDPM tackles the Gaussian Markov chain as a discrete SDE using Ancestral Sampling, while SDE-based generative models focus on solving continuous-time SDEs based on Langevin Dynamics.

### III. METHODOLOGY

#### A. Pre-training Generative Diffusion Model

Diffusion models have demonstrated an impressive ability in image synthesis within the CV community. Several pre-trained models have been made available on popular open-source datasets such as ImageNet, CelebAHQ, and Cifar-10. These pre-trained models, such as the diffusers project [54] and guided diffusion [33], have been released to facilitate further research. However, due to the heterogeneous features of RS images, these publicly available models often experience significant performance degradation when applied to RS and geoscience applications. Since there are currently few available generative diffusion models specifically pre-trained for RS datasets, research on diffusion models in the field of RS and geoscience remains limited due to the inherent domain shift between RS and CV datasets.

In this paper, we propose to pre-train DDPM as a generative diffusion model on a series of RS datasets for image synthesis and adversarial defense tasks, as shown in Fig. 2. In DDPM, the transitional kernel $q(x_t|x_{t-1})$ in the forward diffusion process (Eq. (8)) is handcrafted using Gaussian perturbation. This allows for the incremental transformation of the input data distribution into a tractable Gaussian noise distribution, defined as follows:

$$q(x_t|x_{t-1}) = \mathcal{N}(x_t; \sqrt{1 - \beta_t}x_{t-1}, \beta_t I)), \quad (10)$$

where $\mathcal{N}$ is the Gaussian distribution with variance $I$, and $\beta_t$ represents the variance schedule for noise addition in the $t$-th diffusion step, which is predetermined before model training. Concerning the reverse diffusion process in Eq. (9), DDPM incorporates a learnable Markov chain parameterized jointly by a prior distribution $p(X_T) = \mathcal{N}(x_T; 0, I)$, approximately the same as $q(x_T)$, and a reverse transitional kernel $p(x_{t-1}|x_t)$ as follows:

$$p_\theta(x_{t-1}|x_t) = \mathcal{N}(x_{t-1}; \mu_\theta(x_t, t), \sum_\theta(x_t, t)), \quad (11)$$

where the mean $\mu_\theta(x_t, t)$ and variance $\sum_\theta(x_t, t)$ are usually calculated from DNNs parameterized by $\theta$. To simplify the implementation of the forward process, DDPM further pro-
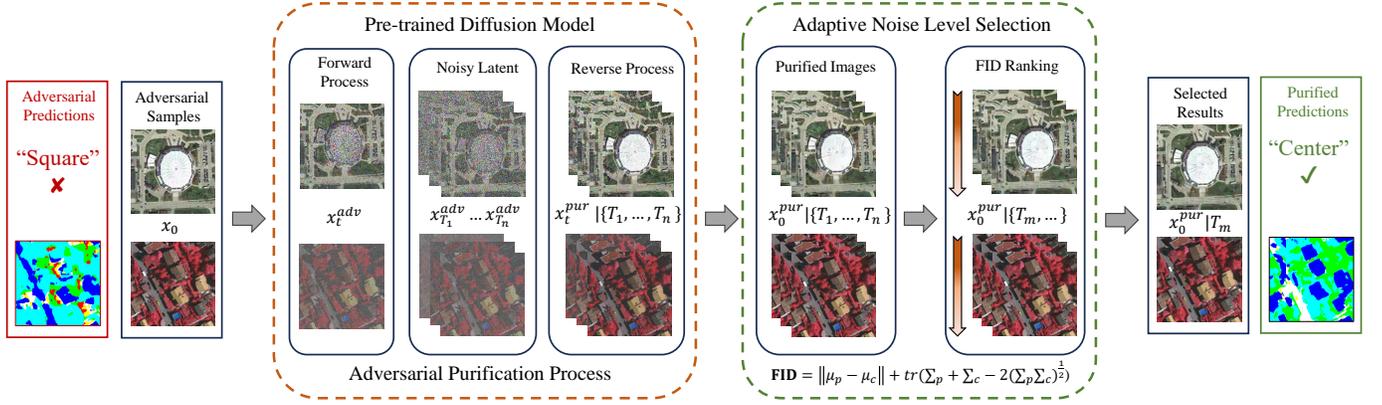
Fig. 3. Overview of the proposed UAD-RS adversarial purification framework: The adversarial samples are first diffused with Gaussian noises for $T_m$ steps into a noisy latent. The noise latent is then denoised by the reverse process, reconstructing a new image. The adversarial purification process of UAD-RS can purify the perturbations from adversarial samples by first mixing them with Gaussian noises and then denoising the mixture into clean images. After that, the adaptive noise level selection (ANLS) algorithm is utilized to find the optimal noise level $T_m$ that can generate the best-purified results compared to other noise levels $\{T_1, ..., T_n\}$, accomplished by calculating and ranking the FID scores of their results. Finally, the purified predictions can avoid the influence of the original perturbations and deliver the correct results.

posed a closed form of perturbed representations $x_t$ sampled from $x_0$ as follows:

$$q(x_t|x_0) = \mathcal{N}(x_t; \sqrt{\overline{\alpha}_t}x_0, (1 - \overline{\alpha}_t)I), \quad (12)$$

$$x_t = \sqrt{\overline{\alpha}_t}x_0 + \sqrt{1 - \overline{\alpha}_t}\epsilon, \quad (13)$$

where $\alpha_t := 1 - \beta_t$, $\overline{\alpha}_t := \prod_{s=0}^{t} \alpha_s$, and $\epsilon$ is a standard Gaussian noise. After determining the two opposite diffusion processes, DDPM utilizes a Kullback-Leibler (KL) divergence to optimize the reverse process, aiming to approximately match the actual time reversal of the forward Markov chain [55].

$$\mathbf{KL}(q(x_0, x_1, ..., x_T)|p_\theta(x_T, x_{T-1}, ..., x_0)), \quad (14)$$

$$= -\mathbb{E}_{q(x_0, x_1, ..., x_T)}[\log p_\theta(x_0, x_1, ..., x_T)] + const, \quad (15)$$

$$= -\mathbb{E}_{q(x_0, x_1, ..., x_T)}[-\log p(x_T) - \sum_{t=1}^{T} \log \frac{p_\theta(x_{t-1}|x_t)}{q(x_t|x_{t-1})}] + const, \quad (16)$$

where the first term in Eq. (16) represents the variational lower bound of the log-likelihood of the data $x_0$, which is commonly considered as a training objective for probabilistic generative models. By further incorporating a step-wise weight schedule $\lambda(t)$, the training objective function can be further derived from KL divergence as:

$$\mathbb{E}_{t \sim \mathcal{U}[\![1, T]\!], x_0 \sim q(x_0), \epsilon \sim \mathcal{N}(0, I)}[\lambda(t)\|\epsilon - \epsilon_\theta(x_t, t)\|^2], \quad (17)$$

where $x_t$ can be calculated from $x_0$, $\epsilon \sim (0, I)$ is the Gaussian vector for sampling, $\mathcal{U}[\![1, T]\!]$ is a uniform set over $\{0, 1, 2, ..., T\}$, and $\epsilon_\theta$ is a DNN that predicts the noise vector $\epsilon$ given $x_t$ and $t$.

### B. Adversarial Purification

Taking inspiration from the powerful generative ability of diffusion models within marginalized domains, we propose the use of a pre-trained DDPM model to purify adversarial samples on RS imagery. The DDPM model consists of a

forward process that injects Gaussian noise into a clean image, and a reverse process that eliminates the noise in the image. These processes can be considered as means to remove the perturbations from the adversarial samples and reconstruct images in the clean domain, respectively. Specifically, as shown in Fig. 3, the forward diffusion gradually adds Gaussian noise to the adversarial samples, progressively submerging the adversarial perturbations. Conversely, the reverse diffusion simultaneously eliminates both the Gaussian noise and the adversarial samples, resulting in predictions within an adversarial-free domain using the pre-trained models.

Given an adversarial sample $x^{adv}$ consisting of the original image $x$ and the adversarial perturbation $\delta$, the forward diffusion process can be computed according to Eq. (13) as:

$$x_{T_m}^{adv} = \sqrt{\overline{\alpha}_{T_m}}x^{adv} + \sqrt{1 - \overline{\alpha}_{T_m}}\epsilon, \quad (18)$$

$$x^{adv} = x + \delta, \quad (19)$$

where $T_m$ represents the noise level that controls the total steps of the forward diffusion. In this process, the first term of the adversarial sample becomes smaller, while the second term of Gaussian noise grows larger as $\overline{\alpha}T_m = \prod_{s=1}^{T_m} \alpha_s$ decreases with a larger $T_m$. Subsequently, the reverse diffusion process is applied to obtain the reconstructed image $x^{clean}$, as shown in Eq. (9), using a pre-trained model as:

$$p_{\theta^*}(x_{t-1}^{pur}|x_t^{pur}) = \mathcal{N}(x_{t-1}^{pur}; \mu_{\theta^*}(x_t^{pur}, t), \sum_{\theta^*}(x_t^{pur}, t)), \quad (20)$$

$$x_{T_m}^{pur} := x_{T_m}^{adv}, \quad (21)$$

where $\theta^*$ is the pre-trained parameters of the DNN, and $t = \{t \in N | t \leq T\}$. After that, the purified image $x_0^{pur}$ can be obtained from $T$ denoising steps and then forwarded to subsequent applications, such as a scene classification model or a semantic segmentation model.

## C. Adaptive Noise Level Selection

As the noise level $T_m$ increases, $\overline{\alpha}_{T_m} = \prod_{s=1}^{T_m} \alpha_s$ gradually decreases while $\sqrt{1 - \overline{\alpha}_{T_m}}$ gradually increases, which indicates a smaller proportion of adversarial sample and a larger proportion of injected Gaussian noise, denoted as first and second terms in Eq. (18), respectively. In regular RS adversarial samples, adversarial perturbation $\delta$ usually occupies only a small proportion of the total sample that is not clearly perceptible [56]. Therefore, selecting a proper value for $T_m$ is critical to purify the adversarial images with high quality. When the $T_m$ is set to a high level, object features in $x$ can also be destroyed and thus hard to be recovered in reverse diffusion. If the $T_m$ is set to a low level, the adversarial perturbation could still exist in latent noise and remain in the purified image. As a corollary, the key to improve the performance of adversarial purification is by optimizing the value of $T_m$ to balance the Gaussian noise intensity that can maximize the mixture of the adversarial perturbations while preserving most of the original features.

To achieve the trade-off performance of image reconstruction, we propose an unsupervised hybrid scoring function that can estimate an optimal noise level $T_m$ for each test dataset. Inspired by the advances in unsupervised domain adaptation and image generation studies, a Frechet Inception Distance (FID)-based scoring algorithm is proposed in this study to measure the domain gap between the purified and adversarial-free images from their feature overlap in DNNs. Given a victim classification model pre-trained on adversarial-free images $F_\varphi$, the FID score can be calculated as follows:

$$\mathbf{FID} = ||\mu_p - \mu_c|| + tr(\textstyle\sum_p + \sum_c - 2(\sum_p \sum_c)^{\frac{1}{2}}), \quad (22)$$

where $tr$ represents the trace linear algebra operation, $(\mu_p, \sum_p)$ and $(\mu_c, \sum_c)$ denote the mean and covariance of the adversarial-purified and clean features extracted by $F_\varphi$ from the purification results and adversarial-free samples, respectively. The FID score can measure the difference between two data distributions in the latent feature space and thus can distinguish the anomaly deep representations with the input of adversarial samples. Therefore, when the FID score is minimized, the optimal classification performance is considered to be achieved with the deep features extracted from the purified results that are with the closest distribution of clean samples. Considering the extensive volume of RS datasets, the scoring function is applied to a subset of $N$ adversarial examples that are randomly chosen from the full dataset. Consequently, the subset is purified with different noise level settings, and the groups of results are evaluated by the score function to deliver a **FID** score for each noise level $T_m$. Finally, the optimal noise level can be obtained by ranking the scores and can be consequently applied for purifying the full test dataset.

## IV. EXPERIMENTAL RESULTS

### A. Dataset Description

In this article, an adversarial attack benchmark, namely, Universal Adversarial Examples in RS (UAE-RS) dataset [45], is employed for the evaluation of the proposed UAD-RS model on defending scene classification and semantic segmentation models. The UAE-RS dataset is constructed on the basis of two benchmark RS image datasets for scene classification (i.e., UC-Merced (UCM) and Aerial Image Dataset (AID)) and the other two very high-resolution RS image datasets for semantic segmentation (i.e., Vaihingen and Zurich Summer).

*1) UCM:* consists of 2100 overhead scene images extracted from large images from the U.S. Geology Survey (USGS) National map. The data of UCM datasets are categorized into 21 land-use classes; each class contains 100 images measuring $256 \times 256$ pixels with a spatial resolution of 0.3m per pixel in the RGB color space.

*2) AID:* consists of 10000 aerial images within 30 scene types collected from Google Earth. The numbers of images vary a lot with different classes, from 220 to 420 samples. The AID dataset has multiple spatial resolutions, altering from 8m to 0.5m per pixel. The samples from all categories have the same image size of $600 \times 600$ pixels.

*3) Vaihingen:* is a subset of a semantic segmentation benchmark dataset in RS imagery provided by the International Society for Photogrammetry and RS (ISPRS). The Vaihingen dataset contains 33 aerial images annotated with 6 land cover classes, which are composed of three bands of near-infrared, red and green, with a spatial resolution of 0.09m. The average size of an image sample is about $2500 \times 1900$ pixels, covering approximately an area of $1.38 \text{ km}^2$.

*4) Zurich Summer:* is an urban RS semantic segmentation dataset consisting of 20 satellite images annotated with 8 land-use classes. The average size of the samples is about $1000 \times 1000$ pixels with a spatial resolution of 0.62m. The images are captured from four bands of near-infrared, red, green and blue. We adopt similar settings with the UAE-RS that selects the near-infrared, red, and green channels in the experiments.

### B. Adversarial Attack Settings

In this article, a series of seven representative adversarial attack algorithms including both white-box (i.e., FGSM [14], IFGSM [18], Jitter [57], TPGD [43], and CW [58]) and black-box attack (i.e., Mixcut and Mixup Attacks [45]) from different technique routes are adopted for a comprehensive evaluation of the universal defense. In particular, the white-box attacks can generate adversarial perturbations based on the features and gradients extracted from the victim classification models, and the performance is highly related to prior knowledge of the structure of the DNNs. On the opposite, the black-box methods can directly generate adversarial samples without the information of utilized classification models by attacking the surrogate models in a standalone manner. Detailed information about these attacks has been given in Section II-A. In our experiments, the white-box attacks are performed for each pre-trained baseline classifier, while we directly incorporate the Mixcut and Mixcut attacked images from the UAE-RS dataset as the black-box attack samples.

### C. Baseline Approaches

*1) Scene Classification:* Four widely used scene classification backbones (i.e., AlexNet, ResNet-18, DenseNet-121,

TABLE I
BASELINE PERFORMANCE OF THE SCENE CLASSIFICATION MODELS.

| Datasets | DNN Classifiers | | | |
|---|---|---|---|---|
| | AlexNet | DenseNet-121 | ResNet-18 | RegNetX-400MF |
| UCM | 91.23 | 96.19 | 96.38 | 95.04 |
| AID | 91.20 | 93.76 | 96.36 | 94.32 |

TABLE II
BASELINE PERFORMANCE OF THE SEMANTIC SEGMENTATION MODELS.

| Datasets | DNN Classifiers | | | | | | | |
|---|---|---|---|---|---|---|---|---|
| | U-Net | | PSPNet | | FCN-8s | | LinkNet | |
| | OA | F1 | OA | F1 | OA | F1 | OA | F1 |
| Vaihingen | 84.19 | 68.85 | 83.59 | 67.10 | 82.61 | 64.63 | 83.19 | 67.12 |
| Zurich | 75.27 | 68.22 | 75.75 | 67.87 | 73.99 | 66.12 | 73.54 | 65.28 |

and RegNetX-400mf) are adopted as the victim DNNs to evaluate the performance of the universal adversarial defense. Concerning the complexity of the utilized RS datasets in our experiments, we chose the medium-sized variants from the backbone series, which are also mostly selected as baselines in many RS studies. Furthermore, the adopted backbones perform at different scales and, thus, can indicate a universal defense performance of the proposed adversarial purification model.

The implementation of these backbones contains a feature extractor and a classifier. The feature extractor gradually extracts the deep features, while the classifier predicts category-wise probability via softmax activation. For the training settings, we initialize the classifiers with ImageNet pre-trained weights to boost the model convergence. Afterward, each backbone model undergoes pre-training for 10 epochs, followed by testing for a baseline classification performance. The baseline performance of the DNN models for scene classification can be found in Table I.

*2) Semantic Segmentation:* Four of the most popular semantic segmentation models in RS research, namely, fully convolutional network (FCN)-8s, U-Net, pyramid scene parsing network (PSPNet), and LinkNet, are employed as baseline models in our experiments. Similar to the scene classification models, the semantic segmentation models also consist of a feature extractor and a classifier. However, the classifiers in semantic segmentation models typically utilize transposed convolutional layers to upsample the extracted latent representations into a segmentation map with the same size as the image input. With randomly initialized model parameters, we pre-train the models for 100 epochs. The baseline performance of the DNN models for semantic segmentation can be found in Table II.

*3) Adversarial Defense:* In the literature, adversarial purification with unsupervised training remains a challenging task, and we have selected three state-of-the-art studies for comparison. The first one is the Pix2Pix [59], which incorporates a GAN-based image-to-image translation model that is widely used in various image enhancement tasks and can remove the adversarial perturbations in an end-to-end manner. Apart from the Pix2Pix method, we also implement another GAN-based hybrid adversarial defense method, namely Perturbation Seeking Generative Adversarial Networks (PSGAN) [60] for comparison. The PSGAN model combines adversarial training and purification to reach a better performance. To provide

a comprehensive comparison, we have also implemented a DNN-based adversarial purification approach, namely, the task-guided denoising network (TGDN) [26], in our comparison experiments.

Unfortunately, in RS studies, most of the adversarial purification methods, like PSGAN, only focus on defending DNNs for scene classification tasks. This limitation makes it difficult for us to find more comparable methods for semantic segmentation experiments. Therefore, in our experimental settings for the semantic segmentation part, we have selected only Pix2Pix and TGDN for comparison.

*D. Evaluation Metrics*

*1) Scene Classification:* The Overall Accuracy (OA= $n_{correct}/n_{total}$) is adopted for quantitative comparison in the scene classification task, where $n_{correct}$ and $n_{total}$ represent the amount of correctly classified samples and the size of full test dataset, respectively.

*2) Semantic Segmentation:* For the quantitative evaluation of pixel-level semantic segmentation results, four commonly used metrics are employed: the OA, precision, recall, and F1-score. They are calculated according to the following equations:

$$OA = \frac{TP + TN}{TP + TN + FP + FN}, \quad (23)$$

$$Precision^{(c)} = TP^{(c)}/(TP^{(c)} + FP^{(c)}), \quad (24)$$

$$Recall^{(c)} = TP^{(c)}/(TP^{(c)} + FN^{(c)}), \quad (25)$$

$$F1^{(c)} = 2 \cdot \frac{Precision^{(c)} \cdot Recall^{(c)}}{Precision^{(c)} + Recall^{(c)}}, \quad (26)$$

where $TP^{(c)}, TN^{(c)}, FP^{(c)}, FN^{(c)}$ represents the number of pixels that are correctly classified in category $c$, correctly classified in other categories, wrongly classified in category $c$ and wrongly classified in other categories, respectively. Due to the extensive amounts of experiment settings, only the most representative evaluation metrics for semantic segmentation (i.e., OA, F1-score) are displayed. For F1-score, we calculate the average of the metrics among all the categories $C$ as the final evaluation measurements (i.e., $F1 = \sum_{c=1}^{C} F1^{(c)}/C$).

*E. Implementation Details*

For the proposed UAD-RS model, the unconditional DDPM was pre-trained under the structure of *diffusers* [54] with default training hyperparameters for 500 epochs for each dataset. The number of diffusion steps $T$ in the DDPM was set to 1000. Regarding the adversarial purification experiments, the noise level $T_m$ was set to $10, 20, 30, ..., 110, 120$ in the ANLS module, and $N = 100$ samples were tested with each $T_m$ to find the optimal one from the list. As for the adversarial attack part, the implementation of UAE-RS [45] Github repository was incorporated and utilized in our experiments. In our experiments, the parameters of the comparison methods were set following the corresponding original articles. All models were implemented using the PyTorch deep learning platform. The experiments were run on the Slurm computational system with eight NVIDIA Tesla A100 GPUs (40 GB of RAM).

TABLE III
QUANTITATIVE COMPARISONS FOR ADVERSARIAL DEFENSE OF SCENE
CLASSIFICATION ON UCM DATASET.

| Defense | Victim DNNs | Attack Algorithms | | | | | | |
|---|---|---|---|---|---|---|---|---|
| | | FGSM | IFGSM | CW | Jitter | Mixcut | Mixup | TPGD |
| None | AlexNet | 10.95 | 0.00 | 0.00 | 3.05 | 37.05 | 13.14 | 30.48 |
| | DenseNet-121 | 56.67 | 1.62 | 1.14 | 12.28 | 3.05 | 1.05 | 35.33 |
| | ResNet-18 | 33.9 | 0.00 | 0.00 | 6.28 | 8.28 | 3.62 | 35.81 |
| | RegNetX-400MF | 34.24 | 0.02 | 0.00 | 4.92 | 23.66 | 24.76 | 29.68 |
| Pix2Pix | AlexNet | 17.43 | 0.00 | 0.00 | 13.43 | 52.76 | 21.62 | 31.43 |
| | DenseNet-121 | 65.24 | 4.19 | 2.57 | 31.71 | 3.71 | 2.19 | 36.95 |
| | ResNet-18 | 42.10 | 0.00 | 0.10 | 26.00 | 11.24 | 7.52 | 35.81 |
| | RegNetX-400MF | 63.71 | 6.50 | 6.85 | 50.10 | 45.90 | 35.81 | 26.10 |
| PSGAN | AlexNet | **84.19** | 36.95 | 36.76 | **71.52** | 46.38 | 29.42 | 41.71 |
| | DenseNet-121 | **87.90** | 29.71 | 30.48 | 80.10 | 6.67 | 4.28 | 53.52 |
| | ResNet-18 | **84.10** | 38.76 | 40.19 | 80.19 | 33.14 | 25.90 | 67.03 |
| | RegNetX-400MF | **85.33** | 44.48 | 49.90 | 83.14 | 65.33 | 54.00 | 55.90 |
| TGDN | AlexNet | 46.28 | 30.48 | 31.05 | 53.14 | **52.10** | 44.28 | 48.38 |
| | DenseNet-121 | 77.81 | 64.76 | 62.57 | 79.43 | 29.62 | 30.10 | 74.48 |
| | ResNet-18 | 71.14 | 63.43 | 64.10 | 74.76 | 50.38 | 46.19 | 72.95 |
| | RegNetX-400MF | 66.19 | 60.10 | 60.19 | 66.67 | 68.10 | 61.14 | 64.95 |
| UAD-RS (Ours) | AlexNet | 65.43 | **53.33** | **54.10** | 70.10 | 51.24 | **66.28** | **64.86** |
| | DenseNet-121 | 81.71 | **68.57** | **68.19** | **82.79** | 40.48 | 39.05 | **75.33** |
| | ResNet-18 | 78.00 | **66.48** | **66.86** | 80.48 | 54.19 | 48.95 | **75.52** |
| | RegNetX-400MF | 82.48 | **75.81** | **75.81** | **84.19** | 72.57 | 62.19 | **80.28** |

TABLE IV
QUANTITATIVE COMPARISONS FOR ADVERSARIAL DEFENSE OF SCENE
CLASSIFICATION ON AID DATASET.

| Defense | Victim DNNs | Attack Algorithms | | | | | | |
|---|---|---|---|---|---|---|---|---|
| | | FGSM | IFGSM | CW | Jitter | Mixcut | Mixup | TPGD |
| None | AlexNet | 12.82 | 0.00 | 0.00 | 2.32 | 28.36 | 3.64 | 24.10 |
| | DenseNet-121 | 46.08 | 0.02 | 0.02 | 4.74 | 0.02 | 0.02 | 24.44 |
| | ResNet-18 | 15.30 | 0.00 | 0.00 | 3.68 | 2.64 | 0.14 | 34.50 |
| | RegNetX-400MF | 34.24 | 0.02 | 0.00 | 4.92 | 23.66 | 24.76 | 29.68 |
| Pix2Pix | AlexNet | 21.54 | 0.10 | 0.06 | 22.00 | 43.62 | 10.04 | 27.10 |
| | DenseNet-121 | 58.96 | 3.74 | 2.82 | 43.12 | 0.30 | 0.20 | 27.14 |
| | ResNet-18 | 31.12 | 0.28 | 0.24 | 32.62 | 3.84 | 2.70 | 34.88 |
| | RegNetX-400MF | 51.90 | 4.78 | 4.96 | 49.78 | 40.70 | 36.50 | 34.40 |
| PSGAN | AlexNet | **79.96** | 42.46 | 42.24 | 69.08 | 33.92 | 16.84 | 43.54 |
| | DenseNet-121 | **83.12** | 58.22 | 58.00 | **82.62** | 18.48 | 18.46 | **69.50** |
| | ResNet-18 | **80.52** | **61.26** | **60.86** | **81.76** | 33.16 | 29.60 | **72.60** |
| | RegNetX-400MF | **82.08** | 50.56 | 54.18 | **80.58** | 57.78 | 52.42 | 65.22 |
| TGDN | AlexNet | 26.42 | 18.24 | 18.36 | 28.78 | 22.88 | 18.28 | 27.30 |
| | DenseNet-121 | 41.96 | 34.08 | 33.86 | 41.94 | 14.46 | 14.18 | 38.74 |
| | ResNet-18 | 32.12 | 27.70 | 28.68 | 34.72 | 19.76 | 18.28 | 33.90 |
| | RegNetX-400MF | 34.44 | 30.14 | 30.66 | 35.84 | 35.24 | 32.16 | 33.78 |
| UAD-RS (Ours) | AlexNet | 68.68 | **57.12** | **57.40** | **72.56** | **47.98** | **30.12** | **59.16** |
| | DenseNet-121 | 75.02 | **59.36** | **59.08** | 75.36 | **31.54** | **32.10** | 63.04 |
| | ResNet-18 | 66.52 | 56.44 | 57.34 | 69.28 | **34.84** | **34.50** | 62.16 |
| | RegNetX-400MF | 76.88 | **67.72** | **68.20** | 78.50 | **59.42** | **53.28** | **73.00** |

### F. Experiments on UCM Dataset

*1) Quantitative Results:* The quantitative evaluation of the proposed UAD-RS method and its competitors regarding scene classification on UCM dataset is shown in Table III. UAD-RS model achieved the best adversarial defense performance with the highest OA among most of the experimental settings of classifiers and attacks. Although PSGAN achieved promising results in defending classifiers from the FGSM attacks, it suffers from a considerable performance loss in defending more complex and stronger attacks, especially the Mixcut and Mixup algorithms. The TGDN method achieves some results that beat the other two comparators, but there is still a significant performance gap compared with UAD-RS. Apart from that, the Pix2Pix methods cannot provide satisfactory defense results and can only slightly improve the performance of the classifiers from most of the attacks.

*2) Qualitative Results:* Fig. 4 displays six instances of the adversarial purification results of all methods on the UCM dataset. These adversarial samples are generated by strong attacks (e.g., CW, Mixup, Mixcut) to provide a more explicit comparison. It can be seen that Pix2Pix generates results that still contain some adversarial perturbations, and also suffers from color distortion over the full image. In the results of PSGAN, a serious green mask can be observed in several patches, which may be due to the guidance of the cross-
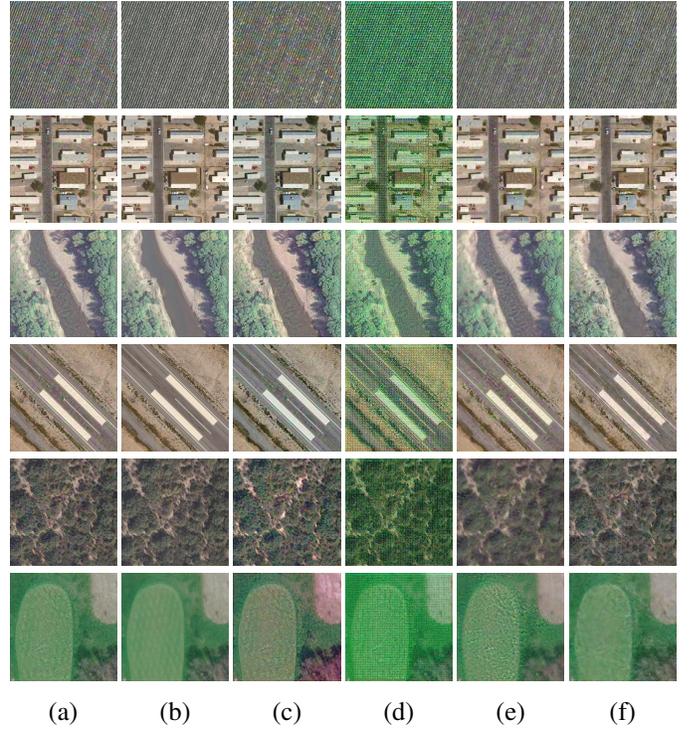
Fig. 4. Qualitative comparison for adversarial purification on the UCM dataset. (a) Adversarial samples. (b) Ground truth. (c)-(f) Purified results obtained by (c) Pix2Pix. (d) PSGAN. (e) TGDN. (f) UAD-RS.

entropy loss utilized along the auxiliary adversarial training process. Although TGDN retains most of the background features, it suffers from a loss of texture information that leads to blurs on the images, which can be especially found on the objects like rivers and plants displayed among the results. In comparison, the proposed UAD-RS achieves finer spatial information of the background and better spectral consistency with the ground truth.

### G. Experiments on AID Dataset

*1) Quantitative Results:* The proposed UAD-RS achieves the best experimental performance in most of the adversarial attacks and classifiers among all methods on the AID dataset in terms of OA values, as shown in Table IV. Although PSGAN obtained slightly better performance than the UAD-RS in a few settings, especially the ResNet-18 classifier, it remained challenging for it to yield promising results when dealing with higher-intensity attacks on the other classifiers. On the contrary, the Pix2Pix and TGDN methods can neither provide reliable results in defending an extensive RS dataset like the AID.

*2) Qualitative Results:* Qualitative results of the UAD-RS and all the comparison methods for adversarial purification on the AID dataset are shown in Fig. 5. It can be seen in the results of Pix2Pix that it failed to remove the adversarial perturbations in many patches and even fuse them into the texture features. PSGAN seems to recover the images with a color mask and cannot perceptually restore the correct images. The results of TGDN seem to have a good overall visual effect, but it generates some slight noises like black dots, especially
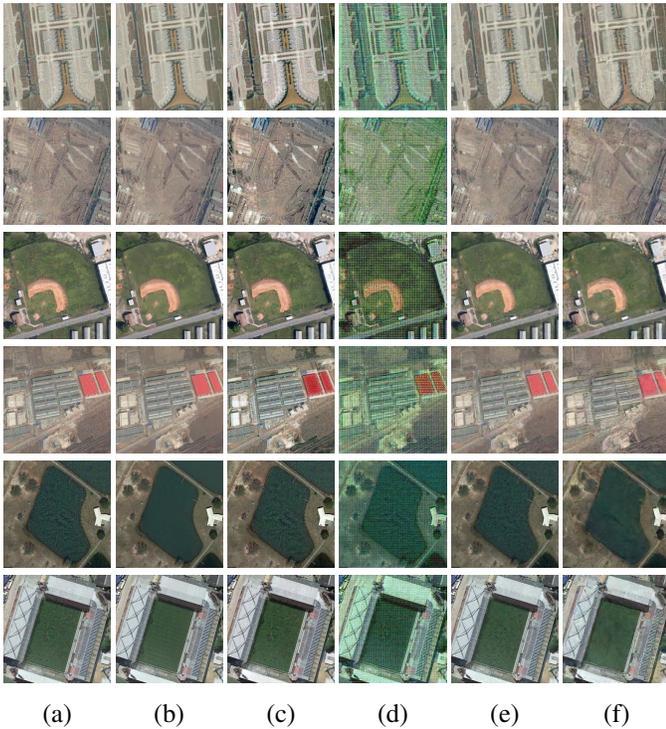
Fig. 5. Qualitative comparison for adversarial purification on the AID dataset. (a) Adversarial samples. (b) Ground truth. (c)-(f) Purified results obtained by (c) Pix2Pix. (d) PSGAN. (e) TGDN. (f) UAD-RS.

on the flattened area in the purified samples. On the contrary, UAD-RS obtained the best qualitative results that are closest to the ground truth and restored most of the crucial information about the objects.

### H. Experiments on Vaihingen Dataset

*1) Quantitative Results:* As shown in Table V, the quantitative comparison is given concerning the application of semantic segmentation on the Vaihingen dataset. Among all the attacks applied, FGSM and Jitter methods only slightly influence the performance of semantic segmentation DNNs, while the others proceed with a stronger attack that substantially degrades the segmentation. The UAD-RS model achieved to outperforms all the competitors for purifying the adversarial samples with the best quality in terms of the highest OA and F1 metrics. The Pix2Pix method successfully protects the victim DNNs from several mild attacks, but the improvement against strong attacks is very limited. The TGDN can effectively defend the victim DNNs from most of the adversarial attacks but still fails to provide satisfied purification, especially for the intensive attacks.

*2) Qualitative Results:* The qualitative comparison of adversarial purification and semantic segmentation results on Vaihingen dataset is displayed in Fig. 6. To provide a clear comparison, the samples are generated by the most powerful attack (i.e., CW) according to the quantitative results in Table V. It can be seen in the predictions of attacked adversarial samples that most of the areas are classified into wrong categories with regard to the predictions of clean samples. As for the purification results, the Pix2Pix and TGDN methods

can improve the segmentation maps for some categories, but most of the patches are still affected by the remained perturbations. On the contrary, the predictions of UAD-RS purified results successfully restore the edges of the objects and some background information, which have a closer fit to the normal segmentation maps.

### I. Experiments on Zurich Summer Dataset

*1) Quantitative Results:* The Table VI shows the quantitative comparisons of the UAD-RS and its competitors for the semantic segmentation on the Zurich Summer dataset. The adversarial attack algorithms yielded a similar success rate to the Vaihingen dataset on the Zurich Summer among four victim DNNs, where FGSM and Jitter methods generate mild perturbations and the others obtain the intensive ones. Among all the purification methods, UAD-RS achieved the best results for most of the experimental settings with the highest OA and F1 values. The Pix2Pix method could only obtain a limited enhancement for the performance of victim DNNs. The TGDN achieved marginally better than UAD-RS in a few attack settings but only gained mediocre results in most of the experiment groups.

*2) Qualitative Results:* Fig. 7 presents the qualitative comparisons of the segmentation results from all methods on the Zurich Summer dataset with the CW attacks. Among the segmentation maps obtained from different purified samples, all of them suffer from omission errors which are not well improved from the results of adversarial samples. In particular, the Pix2Pix and TGDN methods missed many pixels and classified them as the background, especially for the objects of water. In contrast, UAD-RS recovers most of the pixels that can be accurately segmented into high-precision intact targets, which is closest to the predictions of clean images.

### J. Ablation Studies

*1) RS Image Synthesis:* Diffusion models have been widely utilized to synthesize high-resolution images in CV studies and have been reported to outperform traditional GAN structures in this task [33]. Therefore, we conducted some ablation studies to explore the image-generation ability of the DDPM in the field of RS. Since we trained the DDPM in an unconditional manner, the generated images are randomly distributed within the domain of the training dataset. As shown in Fig. 8, the images were generated using a pre-trained DDPM on the AID dataset with an inference step of 1000. As seen in these results, the DDPM can yield high-quality RS images with precise textures and accurate boundaries of objects. Furthermore, the diversity of the generated images is also promising, covering a wide range of scenarios in the perspective of RS. In conclusion, the potential of RS image synthesis using diffusion models is expected, and more related studies, like conditional image generation, can be encouraged.

*2) Noise Level Analysis:* To visualize the influence of different noise levels selected for purification, we diffuse an adversarial sample with different noise levels $T_m$ and then denoise it, as shown in Fig. 9. When the adversarial sample is diffused with a small $T_m$, it can be seen that most of

| Impervious Surfaces | Car | Building | Tree | Low Vegetation | Background |
|---|---|---|---|---|---|



(a)                    (b)                    (c)                    (d)                    (e)                    (f)                    (g)
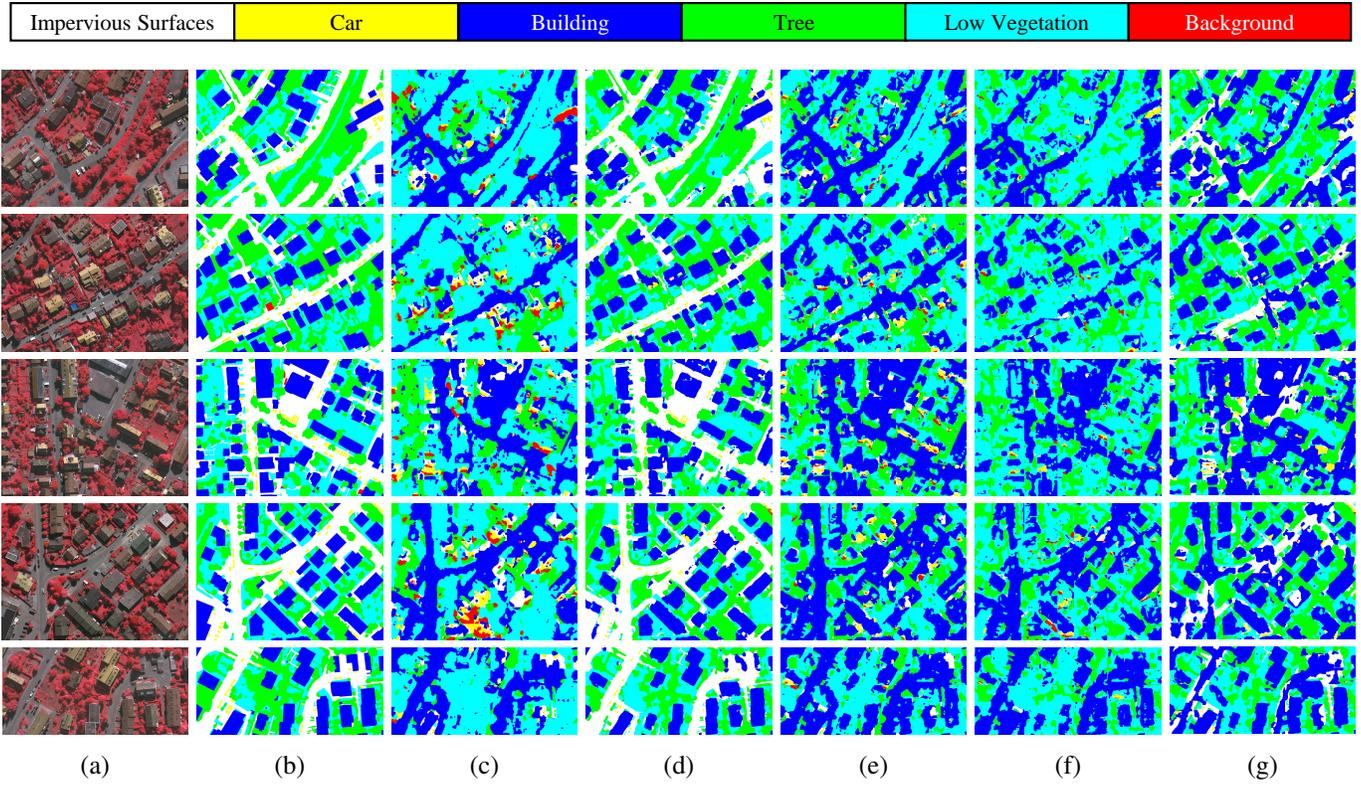
Fig. 6. Qualitative comparison for adversarial purification on Vaihingen dataset. (a) Image inputs. (b) Ground truth. (c)-(d) Segmetation maps obtained from (c) Adversarial Samples. (d) Clean images. (e)-(g) Segmentation maps of the purified results obtained by (e) Pix2Pix. (f) TGDN. (g) UAD-RS.

| Background | Building | Grass | Railway | Swimming Pool | Tree | Road | Bail Soil | Water |
|---|---|---|---|---|---|---|---|---|



(a)                    (b)                    (c)                    (d)                    (e)                    (f)                    (g)

Fig. 7. Qualitative comparison for adversarial purification on Zurich Summer dataset. (a) Image inputs. (b) Ground truth. (c)-(d) Segmetation maps obtained from (c) Adversarial samples. (d) Clean images. (e)-(g) Segmentation maps of the purified results obtained by (e) Pix2Pix. (f) TGDN. (g) UAD-RS.
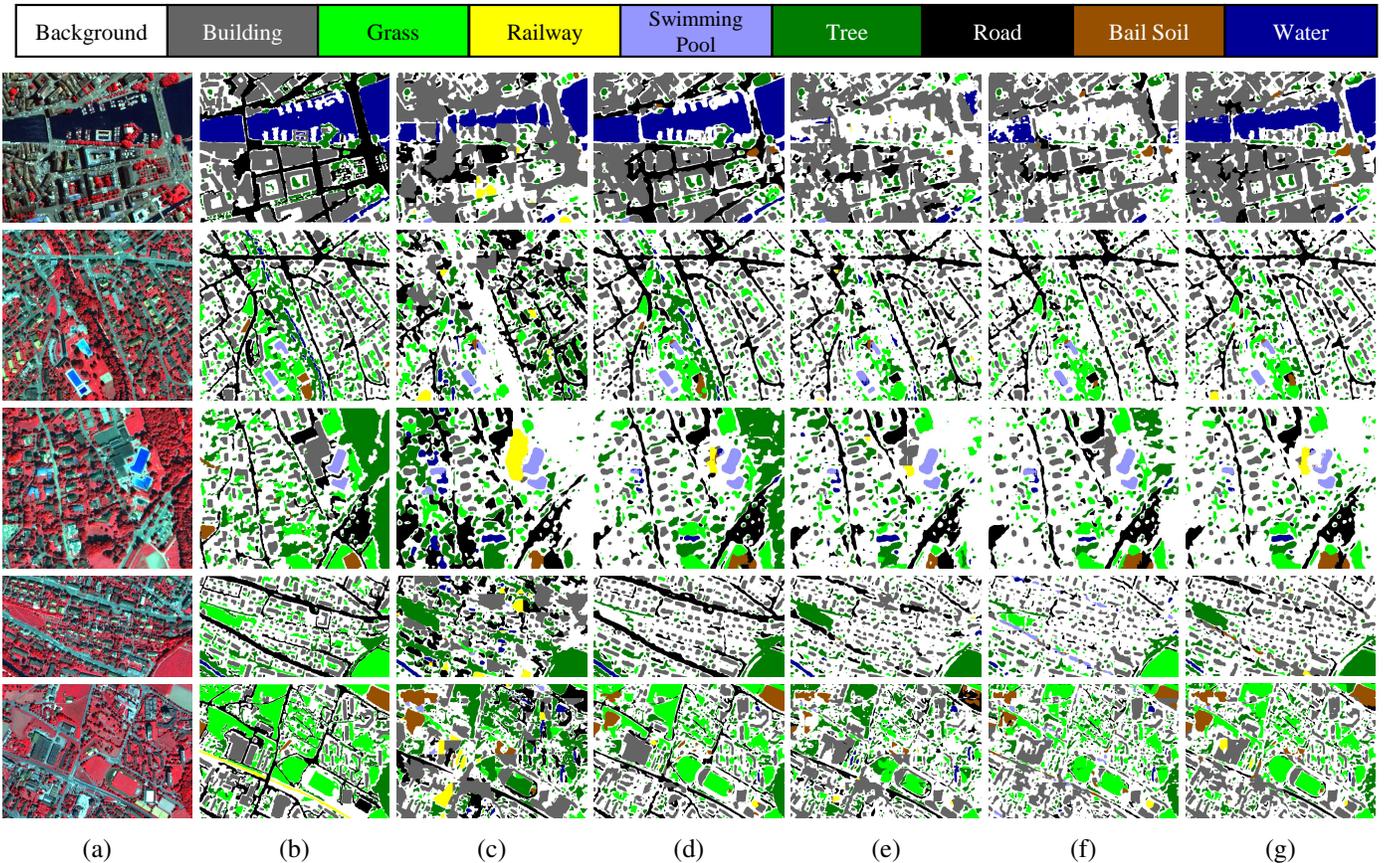
TABLE V
QUANTITATIVE COMPARISONS FOR ADVERSARIAL DEFENSE OF SEMANTIC SEGMENTATION ON VAIHINGEN DATASET.

| Datasets | Victims DNNs | Attack Algorithms | | | | | | | | | | | | |
| --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- |
| | | FGSM | | IFGSM | | CW | | Jitter | | Mixcut | | Mixup | | TPGD | |
| | | OA | F1 | OA | F1 | OA | F1 | OA | F1 | OA | F1 | OA | F1 | OA | F1 |
| No Defense | U-Net | 73.11 | 58.59 | 33.63 | 21.59 | 17.89 | 10.89 | 69.67 | 55.67 | 22.00 | 16.50 | 24.45 | 16.62 | 40.65 | 28.36 |
| | PSPNet | 78.87 | 62.83 | 64.28 | 49.36 | 59.22 | 46.01 | 77.96 | 62.05 | 47.66 | 30.53 | 48.20 | 30.88 | 70.62 | 55.13 |
| | FCN-8s | 72.45 | 54.38 | 40.20 | 27.01 | 17.93 | 11.13 | 72.31 | 54.29 | 38.12 | 25.29 | 36.66 | 24.34 | 54.06 | 38.78 |
| | LinkNet | 75.8 | 59.39 | 43.28 | 30.11 | 28.86 | 19.82 | 74.93 | 58.88 | 40.61 | 29.58 | 33.11 | 21.15 | 50.34 | 36.41 |
| Pix2Pix | U-Net | 71.43 | 55.17 | 52.30 | 37.18 | 45.43 | 30.66 | 71.20 | 54.97 | 44.79 | 30.36 | 41.92 | 27.15 | 57.85 | 42.45 |
| | PSPNet | 79.03 | 63.02 | 66.95 | 52.41 | 63.06 | 49.31 | 78.79 | 62.72 | 55.76 | 37.90 | 53.25 | 35.16 | 70.72 | 56.02 |
| | FCN-8s | 72.89 | 55.34 | 48.78 | 31.96 | 35.20 | 20.92 | 73.10 | 55.43 | 47.15 | 30.14 | 46.43 | 30.57 | 58.84 | 41.36 |
| | LinkNet | 78.56 | 62.39 | 61.27 | 47.22 | 55.49 | 42.91 | 78.35 | 62.21 | 62.32 | 47.27 | 57.00 | 41.74 | 66.42 | 51.85 |
| TGDN | U-Net | 66.67 | 52.86 | 53.00 | 39.15 | 46.54 | 32.08 | 66.29 | 52.45 | 49.21 | **40.51** | 38.28 | 24.39 | 55.61 | 41.67 |
| | PSPNet | 80.85 | 64.71 | 67.78 | 52.54 | 64.10 | 48.19 | 80.37 | 64.20 | 59.01 | 40.60 | 56.66 | 38.34 | 70.48 | 55.35 |
| | FCN-8s | 63.26 | 48.38 | 51.74 | 34.86 | 41.87 | 25.89 | 65.20 | 50.05 | 48.18 | 30.87 | 48.22 | 31.55 | 57.91 | 41.39 |
| | LinkNet | 69.85 | 54.25 | 59.65 | 44.54 | 56.46 | 42.70 | 69.69 | 54.09 | 64.42 | 50.20 | 60.53 | **46.69** | 61.50 | 46.13 |
| UAD-RS (Ours) | U-Net | **77.86** | **62.47** | **70.62** | **56.76** | **66.90** | **52.53** | **77.67** | **62.53** | 54.83 | 39.12 | **56.82** | **40.59** | **74.49** | **60.02** |
| | PSPNet | **81.80** | **65.24** | **79.54** | **63.48** | **77.81** | **61.95** | **81.75** | **65.20** | **70.60** | **51.72** | **72.15** | **53.99** | **81.68** | **65.13** |
| | FCN-8s | **75.21** | **57.33** | **68.74** | **50.77** | **58.16** | **37.03** | **74.61** | **56.88** | **61.64** | **43.57** | **62.83** | **44.76** | **75.02** | **56.67** |
| | LinkNet | **81.34** | **65.32** | **71.74** | **56.53** | **75.81** | **60.51** | **81.34** | **65.37** | **71.22** | **54.22** | 61.03 | 44.34 | **80.10** | **64.42** |

TABLE VI
QUANTITATIVE COMPARISONS FOR ADVERSARIAL DEFENSE OF SEMANTIC SEGMENTATION ON ZURICH DATASET.

| Datasets | Victims DNNs | Attack Algorithms | | | | | | | | | | | | |
| --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- |
| | | FGSM | | IFGSM | | CW | | Jitter | | Mixcut | | Mixup | | TPGD | |
| | | OA | F1 | OA | F1 | OA | F1 | OA | F1 | OA | F1 | OA | F1 | OA | F1 |
| No Defense | U-Net | 65.63 | 59.90 | 43.26 | 40.44 | 33.53 | 31.07 | 64.95 | 59.35 | 41.69 | 26.41 | 42.63 | 34.39 | 53.76 | 50.56 |
| | PSPNet | 71.72 | 64.53 | 61.40 | 55.76 | 56.55 | 51.51 | 71.32 | 64.24 | 58.17 | 49.72 | 54.88 | 47.93 | 67.88 | 61.77 |
| | FCN-8s | 66.56 | 58.76 | 51.21 | 45.29 | 36.45 | 29.16 | 65.09 | 57.89 | 64.02 | 55.14 | 60.29 | 47.76 | 58.46 | 51.70 |
| | LinkNet | 67.11 | 59.44 | 50.42 | 44.33 | 44.94 | 38.45 | 66.93 | 59.37 | 55.62 | 40.60 | 50.99 | 35.45 | 59.39 | 53.12 |
| Pix2Pix | U-Net | 70.06 | 61.53 | 61.55 | 52.52 | 58.74 | 48.83 | 70.24 | 62.20 | 54.80 | 38.19 | 51.83 | 39.61 | 66.39 | 57.44 |
| | PSPNet | 70.35 | 61.35 | 64.79 | 55.64 | 61.74 | 51.99 | 70.51 | 61.59 | 62.36 | 48.59 | 51.71 | 41.04 | 68.56 | 58.83 |
| | FCN-8s | 69.47 | 60.23 | 58.91 | 50.08 | 49.55 | 37.20 | 69.09 | 60.77 | 57.95 | 47.13 | 57.55 | 46.20 | 65.67 | 56.21 |
| | LinkNet | 67.98 | 57.58 | 60.70 | 49.9 | 58.01 | 45.57 | 68.36 | 58.02 | 58.48 | 41.90 | 56.42 | 39.33 | 65.67 | 54.04 |
| TGDN | U-Net | 66.67 | 52.86 | 53.00 | 39.15 | 46.54 | 32.08 | 66.29 | 52.45 | 49.21 | 40.51 | 38.28 | 24.39 | 55.61 | 41.67 |
| | PSPNet | **80.85** | 64.71 | 67.78 | 52.54 | 64.10 | 48.19 | **80.37** | 64.20 | 59.01 | 40.60 | **56.66** | 38.34 | 70.48 | 55.35 |
| | FCN-8s | 63.26 | 48.38 | 51.74 | 34.86 | 41.87 | 25.89 | 65.20 | 50.05 | 64.18 | 30.87 | 48.22 | 31.55 | 57.91 | 41.39 |
| | LinkNet | 69.85 | 54.25 | 59.65 | 44.54 | 56.46 | 42.70 | 69.69 | 54.09 | **64.42** | **50.2** | 60.53 | 46.69 | 61.50 | 46.13 |
| UAD-RS (Ours) | U-Net | **71.82** | **63.80** | **66.24** | **58.98** | **64.35** | **56.26** | **71.93** | **64.21** | 59.00 | 48.13 | **59.19** | **51.04** | **68.98** | **61.84** |
| | PSPNet | 73.33 | **65.54** | **69.21** | **61.84** | **66.97** | **60.04** | 73.39 | **65.65** | 61.51 | **51.12** | 53.12 | **42.72** | **72.44** | **64.42** |
| | FCN-8s | **71.78** | **63.13** | **65.22** | **56.67** | **62.59** | **53.51** | **71.56** | **63.32** | 62.08 | **50.54** | **62.06** | **50.26** | **71.12** | **62.06** |
| | LinkNet | **71.09** | **63.01** | **65.06** | **57.73** | **63.78** | **56.65** | **71.90** | **63.84** | 59.70 | 48.77 | 56.54 | 44.94 | **69.78** | **62.02** |



Fig. 8. RS images synthesized by the pre-trained DDPM with random Guassian noises.

the texture and background of the image is retained, but there are still some adversarial perturbations remaining in the denoised result. On the contrary, when a larger noise level is selected, the diffused image is perceptually like a complete noise, and the contents seem to be destroyed in the denoised images though there is neither perturbation present. These two selections will lead to the unsatisfied performance of the classifier, and only if $T_m$ is correctly selected can the denoised image be best distinguished by the victim DNNs.

*3) Cross-Domain Purification:* Regarding the heterogeneity of the RS datasets, DNNs are considered to suffer a performance loss when tested on data from different domains. Therefore, as part of an ablation study, we also evaluate the cross-domain adversarial purification ability of the UAD-RS model. In particular, we first pre-trained the diffusion model on the UCM dataset and then utilized it to purify the adversarial samples generated from the AID dataset, using different classifiers and attacking algorithms. The adaptive noise level selection mechanism is also applied in this ablation study. Table VII reports the classification results of the cross-domain purified samples and the difference from the intra-domain performance. It can be observed that although the UAD-RS model is trained on the UCM dataset, it still exhibits some capability to purify adversarial samples from the AID dataset, particularly for some cross-domain results upon FGSM and Jitter attacks, which closely resemble the original intra-domain

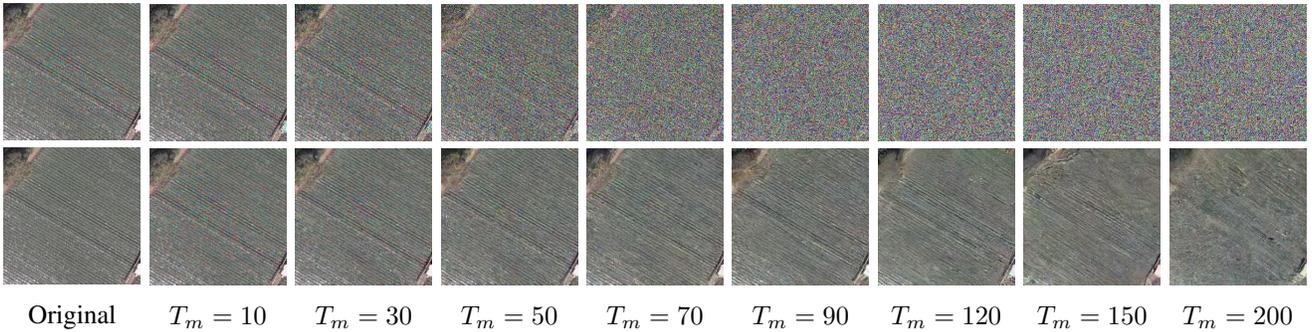| Original | $T_m = 10$ | $T_m = 30$ | $T_m = 50$ | $T_m = 70$ | $T_m = 90$ | $T_m = 120$ | $T_m = 150$ | $T_m = 200$ |

Fig. 9. Purification results with different noise levels. The first column displays the adversarial sample (top row) and the original image (bottom row). The subsequent columns exhibit the diffused images with noise level $T_m$ (top row) and the corresponding purified results (bottom row).

TABLE VII
ABLATION STUDIES FOR CROSS-DOMAIN ADVERSARIAL PURIFICATION USING UAD-RS. THE MODEL WAS FIRST PRE-TRAINED ON THE UCM DATASET AND THEN USED TO PURIFY THE ADVERSARIAL SAMPLES GENERATED FROM THE AID DATASET.

| Victim DNNs | Attack Methods | | | | | | |
| | FGSM | IFGSM | CW | Jitter | Mixcut | Mixup | TPGD |
|---|---|---|---|---|---|---|---|
| AlexNet | 53.72 (-14.96) | 35.76 (-21.36) | 36.12 (-21.28) | 60.26 (-12.30) | 45.12 (-2.86) | 22.52 (-7.60) | 41.16 (-18.00) |
| DenseNet-121 | 68.10 (-6.92) | 43.42 (-15.94) | 44.00 (-15.08) | 69.48 (-5.88) | 16.00 (-15.54) | 15.50 (-16.60) | 52.74 (-10.30) |
| ResNet-18 | 54.92 (-11.60) | 37.40 (-19.04) | 38.56 (-18.78) | 61.10 (-8.18) | 23.84 (-11.00) | 23.24 (-11.26) | 50.40 (-11.76) |
| RegNetX-400MF | 64.40 (12.48) | 50.68 (-17.04) | 50.02 (-18.18) | 69.66 (-8.84) | 51.40 (-8.02) | 41.22 (-12.06) | 60.52 (-12.48) |

results. However, a significant performance loss can also be seen compared with the model pre-trained on the AID dataset, as indicated by the values in the round brackets. In conclusion, the effectiveness of the UAD-RS is restricted to a certain extent in cross-domain adversarial purification.

## V. DISCUSSION

### A. Diffusion and GAN Models

Diffusion and GAN models are both generative models that can predict an image with a noise input, while they are frequently compared in different CV tasks. Both models can generalize the representations of the domain of the training dataset, and thus, they can be used to purify adversarial samples by transferring them to the clean domain. In particular, the defense-GAN and the diffusion methods both utilize a noise latent to reconstruct an image in the clean domain. The difference is that the defense-GAN chooses to optimize a series of random noises and find the best one that can be generated to a similar image to the exact adversarial sample, while the diffusion-based method can proactively generate a latent noise by progressively diffusing the adversarial sample. As a result, the performance of the defense-GAN is usually restricted due to the randomness of the noise optimization and the difficulty of training a GAN model, and it is neither cost-effective due to the extensive inference for hundreds of steps of noise optimization for each sample. In conclusion, the diffusion models outperforms GAN models owing to the stable conditional noise latent generation.

### B. Universal Adversarial Defense

The universal adversarial defense algorithm proposed in this paper aims to purify the adversarial examples generated by different attack algorithms on various classifiers from a single dataset using only one pre-trained diffusion model.

In the literature, most of the adversarial defense methods are performed heterogeneously based on different settings, meaning that the results can only be obtained based on a newly trained model and thus, the cost of defense is incalculable. For example, when dealing with multiple classifiers and attacking methods, these approaches would have to train a new defense model for each classifier specified with each attack. As a result, the universal adversarial defense is more cost-efficient and can perform in different settings without the efforts of unnecessary training.

### C. Unknown Adversarial Threat

The unknown adversarial threat is widely considered a challenge in adversarial defense studies, which means that the prior information of the adversarial perturbation is unknown, including the attack algorithms and intensity. Unfortunately, the lack of this prior knowledge has led to a considerable performance loss for many defense approaches, especially adversarial training methods. This is because these schemes can only achieve promising results under the assumption that the attack settings of the training and testing sessions are the same, which is not feasible in real-world applications. To overcome this difficulty, the proposed UAD-RS provides a novel universal adversarial defense paradigm that can purify the samples solely based on the pre-trained diffusion model without any information from the attacks.

### D. Limitations of the UAD-RS Model

The main limitations of the UAD-RS model are twofold. First, compared with the traditional DNN-based purification models, the UAD-RS model may take a longer time to purify the adversarial examples due to the multi-diffuse and denoising steps. Furthermore, the great amount of trainable parameters in the diffusion models induces a higher computational cost

in the pre-training process. Second, the diffusing process of the UAD-RS may be affected in some high-contrast RS data containing exotic pixels with extreme brightness, which often appears in the Zurich dataset. In these cases, the Gaussian noise added in the diffusing process may be more prominent in white pixels and less visible in dark ones, leading to an ambiguous and uneven denoising results in the consequent steps.

## VI. CONCLUSION

In this article, we have proposed a UAD-RS model for universal adversarial defense against multiple adversarial attacks on commonly used DNNs regarding scene classification and semantic segmentation in RS imagery. In the first stage, the generative diffusion models are pre-trained on some widely used scene classification and semantic segmentation datasets in the RS imagery to obtain the generalized representation ability among the data domain. Consequently, an adversarial purification framework is developed based on the pre-trained models that utilize the forward and reverse processes of them to purify the perturbation from the adversarial samples. After that, an ANLS algorithm is developed to find the optimal noise level setting of the diffusion model to have the best purification results that are closest to the clean samples according to their FID distance captured from victim DNNs. As a result, the UAD-RS model can automatically purify multiple kinds of adversarial samples using only one pre-trained model for each dataset by adapting the hyper-parameters to achieve optimal performance.

Experiments on multiple tasks, including scene classification and semantic segmentation in complex scenes, show that the UAD-RS exhibits obvious advantages in purifying the adversarial samples from a diversity of adversarial attacks targeting various DNNs. Compared with the state-of-the-art adversarial purification methods, UAD-RS is shown to achieve better performance with accurate purified results that can be consequently classified into promising predictions like scene categories and segmentation maps. To encourage further research in similar topics in the field of RS, the pre-trained diffusion models will be open-sourced, which can alleviate future researchers from the efforts and difficulties of retraining huge diffusion models on these datasets.

In the future, we will focus on applying the purification model to defend the vulnerable DNNs in more RS applications which are also potentially threatened by adversarial attacks. Moreover, the UAD-RS model can be potentially applied to other RS tasks, such as domain adaptation and image enhancement, which also deserves further studies.

## VII. ACKNOWLEDGMENT

## REFERENCES

[1] 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.

[2] F. P. Luus, B. P. Salmon, F. Van den Bergh, and B. T. J. Maharaj, "Multiview deep learning for land-use classification," *IEEE Geoscience and Remote Sensing Letters*, vol. 12, no. 12, pp. 2448–2452, 2015.

[3] X. Zhang, W. Yu, and M.-O. Pun, "Multilevel deformable attention-aggregated networks for change detection in bitemporal remote sensing imagery," *IEEE Transactions on Geoscience and Remote Sensing*, vol. 60, pp. 1–18, 2022.

[4] X. Zhang, W. Yu, M.-O. Pun, and W. Shi, "Cross-domain landslide mapping from large-scale remote sensing images using prototype-guided domain-aware progressive representation learning," *ISPRS Journal of Photogrammetry and Remote Sensing*, vol. 197, pp. 1–17, 2023.

[5] C. Szegedy, W. Zaremba, I. Sutskever, J. Bruna, D. Erhan, I. Goodfellow, and R. Fergus, "Intriguing properties of neural networks," *arXiv preprint arXiv:1312.6199*, 2013.

[6] Y. Li, M. Cheng, C.-J. Hsieh, and T. C. Lee, "A review of adversarial attack and defense for classification methods," *The American Statistician*, vol. 76, no. 4, pp. 329–345, 2022.

[7] 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.

[8] Y. Cao, C. Xiao, B. Cyr, Y. Zhou, W. Park, S. Rampazzi, Q. A. Chen, K. Fu, and Z. M. Mao, "Adversarial sensor attack on lidar-based perception in autonomous driving," in *Proceedings of the 2019 ACM SIGSAC conference on computer and communications security*, 2019, pp. 2267–2281.

[9] H. Li, H. Huang, L. Chen, J. Peng, H. Huang, Z. Cui, X. Mei, and G. Wu, "Adversarial examples for cnn-based sar image classification: An experience study," *IEEE Journal of Selected Topics in Applied Earth Observations and Remote Sensing*, vol. 14, pp. 1333–1347, 2020.

[10] B. Peng, B. Peng, J. Zhou, J. Xie, and L. Liu, "Scattering model guided adversarial examples for sar target recognition: Attack and defense," *IEEE Transactions on Geoscience and Remote Sensing*, vol. 60, pp. 1–17, 2022.

[11] 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.

[12] S. Mei, J. Lian, X. Wang, Y. Su, M. Ma, and L.-P. Chau, "A comprehensive study on the robustness of image classification and object detection in remote sensing: Surveying and benchmarking," *arXiv preprint arXiv:2306.12111*, 2023.

[13] Y. Xu, T. Bai, W. Yu, S. Chang, P. M. Atkinson, and P. Ghamisi, "AI security for geoscience and remote sensing: Challenges and future trends," *IEEE Geoscience and Remote Sensing Magazine*, vol. 11, no. 2, pp. 60–85, 2023.

[14] I. J. Goodfellow, J. Shlens, and C. Szegedy, "Explaining and harnessing adversarial examples," *arXiv preprint arXiv:1412.6572*, 2014.

[15] A. Madry, A. Makelov, L. Schmidt, D. Tsipras, and A. Vladu, "Towards deep learning models resistant to adversarial attacks," *arXiv preprint arXiv:1706.06083*, 2017.

[16] Y. Xu, H. Sun, J. Chen, L. Lei, G. Kuang, and K. Ji, "Robust remote sensing scene classification by adversarial self-supervised learning," in *2021 IEEE International Geoscience and Remote Sensing Symposium IGARSS*. IEEE, 2021, pp. 4936–4939.

[17] N. Akhtar, A. Mian, N. Kardan, and M. Shah, "Advances in adversarial attacks and defenses in computer vision: A survey," *IEEE Access*, vol. 9, pp. 155 161–155 196, 2021.

[18] A. Kurakin, I. Goodfellow, and S. Bengio, "Adversarial machine learning at scale," *arXiv preprint arXiv:1611.01236*, 2016.

[19] J. Cohen, E. Rosenfeld, and Z. Kolter, "Certified adversarial robustness via randomized smoothing," in *international conference on machine learning*. PMLR, 2019, pp. 1310–1320.

[20] X. Liu, M. Cheng, H. Zhang, and C.-J. Hsieh, "Towards robust neural networks via random self-ensemble," in *Proceedings of the European Conference on Computer Vision (ECCV)*, 2018, pp. 369–385.

[21] Z. He, A. S. Rakin, and D. Fan, "Parametric noise injection: Trainable randomness to improve deep neural network robustness against adversarial attack," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2019, pp. 588–597.

[22] X. Liu, Y. Li, C. Wu, and C.-J. Hsieh, "Adv-bnn: Improved adversarial defense through robust bayesian neural network," *arXiv preprint arXiv:1810.01279*, 2018.

[23] K. Lucas, M. Jagielski, F. Tramèr, L. Bauer, and N. Carlini, "Randomness in ml defenses helps persistent attackers and hinders evaluators," *arXiv preprint arXiv:2302.13464*, 2023.

[24] 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.

[25] W. Xue, Z. Chen, W. Tian, Y. Wu, and B. Hua, "A cascade defense method for multidomain adversarial attacks under remote sensing detection," *Remote Sensing*, vol. 14, no. 15, p. 3559, 2022.

[26] Y. Xu, W. Yu, and P. Ghamisi, "Task-guided denoising network for adversarial defense of remote sensing scene classification," in *Proc. Int. Joint Conf. Artif. Intell. Workshop*, 2022.

[27] G. Cheng, X. Sun, K. Li, L. Guo, and J. Han, "Perturbation-seeking generative adversarial networks: A defense framework for remote sensing image scene classification," *IEEE Transactions on Geoscience and Remote Sensing*, vol. 60, pp. 1–11, 2021.

[28] F. Croce and M. Hein, "Reliable evaluation of adversarial robustness with an ensemble of diverse parameter-free attacks," in *Proceedings of the 37th International Conference on Machine Learning*, ser. Proceedings of Machine Learning Research, H. D. III and A. Singh, Eds., vol. 119.   PMLR, 13–18 Jul 2020, pp. 2206–2216. [Online]. Available: https://proceedings.mlr.press/v119/croce20b.html

[29] T. Salimans, I. Goodfellow, W. Zaremba, V. Cheung, A. Radford, and X. Chen, "Improved techniques for training gans," *Advances in neural information processing systems*, vol. 29, 2016.

[30] Y. LeCun, S. Chopra, R. Hadsell, M. Ranzato, and F. Huang, "A tutorial on energy-based learning," *Predicting structured data*, vol. 1, no. 0, 2006.

[31] H. Cao, C. Tan, Z. Gao, G. Chen, P.-A. Heng, and S. Z. Li, "A survey on generative diffusion model," *arXiv preprint arXiv:2209.02646*, 2022.

[32] A. Vahdat, K. Kreis, and J. Kautz, "Score-based generative modeling in latent space," *Advances in Neural Information Processing Systems*, vol. 34, pp. 11 287–11 302, 2021.

[33] P. Dhariwal and A. Nichol, "Diffusion models beat gans on image synthesis," *Advances in neural information processing systems*, vol. 34, pp. 8780–8794, 2021.

[34] J. Wang, Z. Lyu, D. Lin, B. Dai, and H. Fu, "Guided diffusion model for adversarial purification," *arXiv preprint arXiv:2205.14969*, 2022.

[35] W. Nie, B. Guo, Y. Huang, C. Xiao, A. Vahdat, and A. Anandkumar, "Diffusion models for adversarial purification," *arXiv preprint arXiv:2205.07460*, 2022.

[36] J. Deng, W. Dong, R. Socher, L.-J. Li, K. Li, and L. Fei-Fei, "Imagenet: A large-scale hierarchical image database," in *2009 IEEE conference on computer vision and pattern recognition*.   Ieee, 2009, pp. 248–255.

[37] A. Krizhevsky, G. Hinton *et al.*, "Learning multiple layers of features from tiny images," 2009.

[38] T. Karras, T. Aila, S. Laine, and J. Lehtinen, "Progressive growing of gans for improved quality, stability, and variation," *arXiv preprint arXiv:1710.10196*, 2017.

[39] Y. Yang and S. Newsam, "Bag-of-visual-words and spatial extensions for land-use classification," in *Proceedings of the 18th SIGSPATIAL international conference on advances in geographic information systems*, 2010, pp. 270–279.

[40] G.-S. Xia, J. Hu, F. Hu, B. Shi, X. Bai, Y. Zhong, L. Zhang, and X. Lu, "Aid: A benchmark data set for performance evaluation of aerial scene classification," *IEEE Transactions on Geoscience and Remote Sensing*, vol. 55, no. 7, pp. 3965–3981, 2017.

[41] M. Cramer, "The dgpf-test on digital airborne camera evaluation overview and test design," *Photogrammetrie-Fernerkundung-Geoinformation*, pp. 73–82, 2010.

[42] M. Volpi and V. Ferrari, "Semantic segmentation of urban scenes by learning local class interactions," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition Workshops*, 2015, pp. 1–9.

[43] H. Zhang, Y. Yu, J. Jiao, E. Xing, L. El Ghaoui, and M. Jordan, "Theoretically principled trade-off between robustness and accuracy," in *International conference on machine learning*.   PMLR, 2019, pp. 7472–7482.

[44] N. Narodytska and S. P. Kasiviswanathan, "Simple black-box adversarial attacks on deep neural networks." in *CVPR Workshops*, vol. 2, 2017, p. 2.

[45] 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.

[46] D. Meng and H. Chen, "Magnet: a two-pronged defense against adversarial examples," in *Proceedings of the 2017 ACM SIGSAC conference on computer and communications security*, 2017, pp. 135–147.

[47] P. Samangouei, M. Kabkab, and R. Chellappa, "Defense-gan: Protecting classifiers against adversarial attacks using generative models," *arXiv preprint arXiv:1805.06605*, 2018.

[48] L. Weng, "From gan to wgan," *arXiv preprint arXiv:1904.08994*, 2019.

[49] Z. Zhang, M. Li, and J. Yu, "On the convergence and mode collapse of gan," in *SIGGRAPH Asia 2018 Technical Briefs*, 2018, pp. 1–4.

[50] M. Lee and D. Kim, "Robust evaluation of diffusion-based adversarial purification," *arXiv preprint arXiv:2303.09051*, 2023.

[51] Y. Song, C. Durkan, I. Murray, and S. Ermon, "Maximum likelihood training of score-based diffusion models," *Advances in Neural Information Processing Systems*, vol. 34, pp. 1415–1428, 2021.

[52] J. Ho, A. Jain, and P. Abbeel, "Denoising diffusion probabilistic models," *Advances in Neural Information Processing Systems*, vol. 33, pp. 6840–6851, 2020.

[53] Y. Song, J. Sohl-Dickstein, D. P. Kingma, A. Kumar, S. Ermon, and B. Poole, "Score-based generative modeling through stochastic differential equations," *arXiv preprint arXiv:2011.13456*, 2020.

[54] P. von Platen, S. Patil, A. Lozhkov, P. Cuenca, N. Lambert, K. Rasul, M. Davaadorj, and T. Wolf, "Diffusers: State-of-the-art diffusion models," https://github.com/huggingface/diffusers, 2022.

[55] L. Yang, Z. Zhang, Y. Song, S. Hong, R. Xu, Y. Zhao, Y. Shao, W. Zhang, B. Cui, and M.-H. Yang, "Diffusion models: A comprehensive survey of methods and applications," *arXiv preprint arXiv:2209.00796*, 2022.

[56] L. Chen, H. Li, G. Zhu, Q. Li, J. Zhu, H. Huang, J. Peng, and L. Zhao, "Attack selectivity of adversarial examples in remote sensing image scene classification," *IEEE Access*, vol. 8, pp. 137 477–137 489, 2020.

[57] L. Schwinn, R. Raab, A. Nguyen, D. Zanca, and B. Eskofier, "Exploring misclassifications of robust neural networks to enhance adversarial attacks," *Applied Intelligence*, pp. 1–17, 2023.

[58] 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.

[59] P. Isola, J.-Y. Zhu, T. Zhou, and A. A. Efros, "Image-to-image translation with conditional adversarial networks," in *Proceedings of the IEEE conference on computer vision and pattern recognition*, 2017, pp. 1125–1134.

[60] G. Cheng, X. Sun, K. Li, L. Guo, and J. Han, "Perturbation-seeking generative adversarial networks: A defense framework for remote sensing image scene classification," *IEEE Transactions on Geoscience and Remote Sensing*, vol. 60, pp. 1–11, 2021.