

# MedRDF: A Robust and Retrain-Less Diagnostic Framework for Medical Pretrained Models Against Adversarial Attack

Mengting Xu, Tao Zhang, and Daoqiang Zhang

**Abstract**—Deep neural networks are discovered to be non-robust when attacked by imperceptible adversarial examples, which is dangerous for it applied into medical diagnostic system that requires high reliability. However, the defense methods that have good effect in natural images may not be suitable for medical diagnostic tasks. The pre-processing methods (e.g., random resizing, compression) may lead to the loss of the small lesions feature in the medical image. Retraining the network on the augmented data set is also not practical for medical models that have already been deployed online. Accordingly, it is necessary to design an easy-to-deploy and effective defense framework for medical diagnostic tasks. In this paper, we propose a Robust and Retrain-Less Diagnostic Framework for Medical pretrained models against adversarial attack (i.e., MedRDF). It acts on the inference time of the pertained medical model. Specifically, for each test image, MedRDF firstly creates a large number of noisy copies of it, and obtains the output labels of these copies from the pretrained medical diagnostic model. Then, based on the labels of these copies, MedRDF outputs the final robust diagnostic result by majority voting. In addition to the diagnostic result, MedRDF produces the Robust Metric (RM) as the confidence of the result. Therefore, it is convenient and reliable to utilize MedRDF to convert pre-trained non-robust diagnostic models into robust ones. The experimental results on COVID-19 and DermaMNIST datasets verify the effectiveness of our MedRDF in improving the robustness of medical diagnostic models.

**Index Terms**—Medical Image, Robust Diagnostic Framework, Adversarial Robustness, Robust Metric.

## I. INTRODUCTION

THERE are many impressive examples of deep neural networks achieving excellent performances on medical diagnostic tasks in radiology [1], dermatology [2], and ophthalmology [3], etc. However, recent studies have revealed the fact that the robustness of the state-of-the-art neural network is poor, i.e., it is easily to craft a visually imperceptible adversarial example to mislead a well-trained network with high confidence [4]–[6]. The vulnerability to adversarial examples

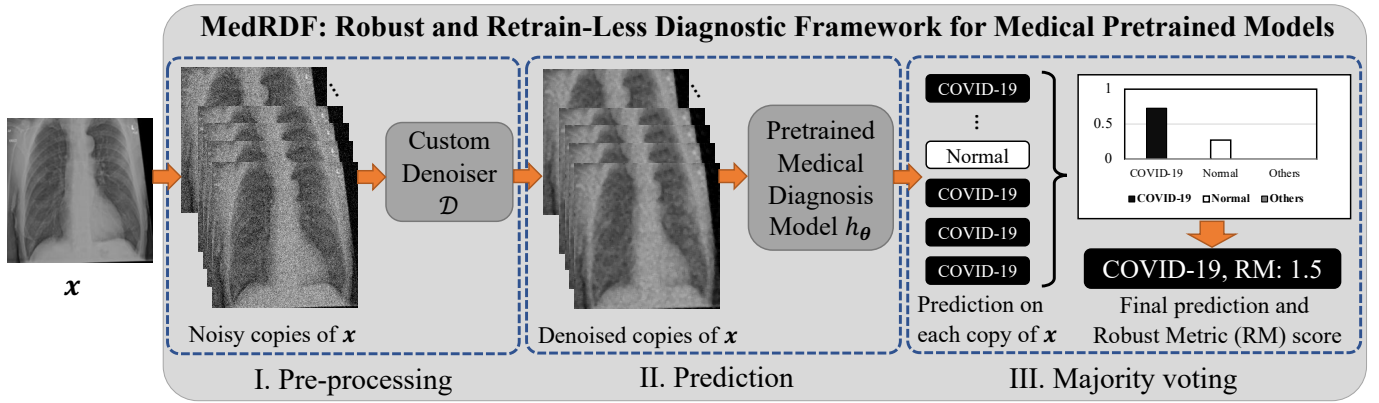
poses a huge threat to the deployment of these models to medical diagnostic tasks that require extremely high reliability [7]–[10]. For example, the misdiagnosis of COVID-19 may cause a large number of diseases to spread. Therefore, developing a robust model to defend against adversarial attacks is very crucial in medical image field.

There are many different defense strategies developed in natural image field. One of the most successful empirical defenses to date is adversarial training [5], which can be regarded as a data augmentation technique that trains neural networks on adversarial examples. However, adversarial training in medical image is problematic as it requires a large labeled training set whereas medical data sets are usually with a small amount of labeled samples. To solve this problem, Li et al. [10] propose the semi-supervised adversarial training (i.e., SSAT) which utilizes both labeled and unlabeled data to generate pseudo-labels. However, the application of SSAT is also limited, because for most medical diagnostic tasks, unlabeled data is also inaccessible, not to mention the heterogeneity between multi-site data sets acquired through different devices (i.e., data distribution difference) and the privacy of medical data. Moreover, Xue et al. [11] propose a defense mechanism which embeds an auto-encoder into the model structure and keeps high-level features invariant to general noises. However, retraining mechanism is not friendly to the medical diagnostic model that has been already deployed online. It is time-consuming and laborious to go back to the online process. Other pre-processing based-defense methods have also shown effectiveness in natural image field. For example, Xie et al. [12] use random resizing and padding (Random R-P) to pre-process the input images before feeding the images into the models. Jia et al. propose the ComDefend [13] to transform the adversarial image to its clean version by compression and reconstruction. However, these defense methods that have good effects in the field of natural images may not be suitable for medical images. For natural images, there is strong similarity and relevance between neighbor pixels in the local structure, random resizing and image compression can help reduce the redundant information of the image, while retaining the dominant information. But for medical images, medical lesions often occupy only a few pixels. Random resizing and image compression may cause the loss of lesion features, thereby affecting the classification and defense effects. To make matters worse, there is still no effective confidence indicator for doctor to evaluate the diagnostic result of the

This work was supported by the National Natural Science Foundation of China (Nos. 62136004, 61876082, 61732006), the National Key R&D Program of China (Grant Nos. 2018YFC2001600, 2018YFC2001602), and also by the CAAI-Huawei MindSpore Open Fund.

Mengting Xu, Tao Zhang, and Daoqiang Zhang are with the College of Computer Science and Technology, Nanjing University of Aeronautics and Astronautics, Nanjing 211106, China (e-mail: {xumengting, dqzhang}@nuaa.edu.cn).

Daoqiang Zhang is the corresponding author.



**Fig. 1:** The Robust and Retrain-Less Diagnostic Framework for Medical pretrained models (MedRDF). I, each test medical image  $x$  is perturbed by isotropic noises  $\eta$  to produce the noisy copies of  $x$ , then they are denoised by the pre-defined denoiser  $\mathcal{D}$ . II, the denoised copies are input to the pre-trained medical diagnostic model  $h_\theta$  to get the predictions. III, the robust diagnostic result  $g(x)$  on  $x$  and the Robust Metric (RM) of the result are obtained by the majority voting on the prediction labels of denoised ones.

model. Therefore, how to reliably improve the robustness of medical diagnostic model is still an open problem.

In this paper, we propose a novel Robust and Retrain-Less Diagnostic Framework for Medical Pretrained Models (i.e., MedRDF) to defense against adversarial attack. As shown in Fig. 1, our proposed MedRDF can easily convert the non-robust pre-trained model to robust one in inference time without retraining. Specifically, *firstly*, for each queried medical image  $x$ , MedRDF produces a large number of copies (i.e., with adding common noise and denoising) around it. *Secondly*, the denoised copies are input into the pre-trained diagnostic model to get the prediction labels. *Finally*, MedRDF outputs the robust diagnostic result of medical image  $x$  by majority voting on the prediction labels of denoised ones. What's more, MedRDF also produces the Robust Metric (RM) as the confidence of the result, which can be used to instruct the doctor to adopt the diagnostic result or re-evaluate it.

The main innovations of our MedRDF can be summarized as follows:

- A novel Robust and Retrain-Less Diagnostic Framework for medical pretrained models (i.e., MedRDF) has been proposed. The MedRDF can be applied to all medical diagnostic tasks seamlessly without retraining diagnostic models, which is very convenient for diagnostic services that are already deployed online.
- A novel Robust Metric (i.e., RM) based on MedRDF has been proposed. It can give the confidence score of the diagnostic result produced by MedRDF, so as to guide the following work of the doctor, such as adopting the result (with high RM) or re-evaluating this case (with low RM).

## II. RELATED WORK

In this section, we first briefly introduce the latest developments in deep learning in the diagnosis of coronavirus disease (COVID-19) and common pigmented skin lesions. Then, the recent adversarial attacks and defense methods on natural and medical images have been reviewed.

### A. Deep Learning for Medical Image Analysis

In the past few years, high-performance deep diagnostic classification models for disease diagnosis have emerged. Here, we are going to introduce two successful applications of deep learning models in medical image analysis.

1) *Coronavirus Disease (COVID-19)*: In recent years, the global outbreak of the Coronavirus disease (COVID-19) has caused tens of thousands of deaths and infected millions of people around the world. This undoubtedly poses a huge threat to the lives of the human beings and the national public health system. Any technical tool that can quickly screen for COVID-19 infection with high accuracy is vital to healthcare professionals. The main clinical tool currently used to diagnose COVID-19 is Reverse Transcription Polymerase Chain Reaction (RT-PCR), but it is expensive, less sensitive, and requires specialized medical personnel [14]. A clinical study of COVID-19 infected patients showed that most of these types of patients were infected by lung infections after being exposed to the virus [15]. Therefore, easy-to-use and low-cost X-ray (i.e., radiography) imaging has become an excellent alternative to COVID-19 diagnosis.

Many automatic algorithms have been proposed to diagnose COVID-19 from chest X-ray images [16]–[18]. In particular, deep learning methods have been considered the best performing methods [19], including Generative Adversarial Networks (GANs) [20], Extreme Learning Machine (ELM) [21], and Long /Short Term Memory (LSTM) [22]. Besides, Jain et al. [15] compared Inception V3 [23], Xception [24], and ResNeXt [25] models which have high performance in natural image field and examined their accuracy in diagnosis of COVID-19. Moreover, Schlemper et al. [26] proposed the Attention-Gated Sononet (AG-Sononet) model, which is carefully designed for fetal ultrasound images. It can also be used for COVID-19 disease diagnosis.

2) *Pigmented Skin Lesions*: Skin cancer is one of the most commonly diagnosed cancers worldwide. According to the 2019 statistical report of the American Association, the

number of new cases and deaths of skin cancer in the United States (excluding basal cell and squamous cell skin cancer) is as high as 104,350 and 11,650, respectively [27]. Among them, melanoma accounts for the largest proportion of all lesions, and the estimated number of new cases and deaths are 92.5% and 62.1%, respectively. However, the skin cancer can be highly treated by early detection and diagnosis, thus reducing the mortality rate.

Due to the importance of early detection, many deep learning methods are used to improve the accuracy of diagnosis and expand the scale of diagnosis. For example, Li et al. [28] proposed a framework consisting of multi-scale fully-convolutional residual networks and a lesion index calculation unit (LICU) to simultaneously address lesion segmentation and lesion classification. Yan et al. [29] proposed an attention-based melanoma recognition method, which introduces an end-to-end trainable attention module regularization for melanoma recognition.

### B. Adversarial Attack

Despite the high performance of deep neural networks in medical image diagnosis, Szegedy et al. [4] first discovered that deep networks are extremely vulnerable to the adversarial examples. The so-called “adversarial example” is added carefully designed perturbation on the original example, which is invisible to the human eye, thus misleading the network output a wrong prediction with a high confidence. Even worse, due to the transferability of adversarial examples, the perturbation designed for one network can also be used to fool other networks.

In recent years, the adversarial attack methods for natural image have developed rapidly. Goodfellow et al. [30] proposed Fast Gradient Sign Method (FGSM) to generate adversarial examples. It calculates the gradient of the loss function with respect to the pixel, and modifies the pixel value of a fixed step along the direction of the gradient. Based on this work, Madry et al. [5] proposed an iterative attack method, which randomly starts a perturbation, and updates the pixel value multi-time along the direction of the gradient, which is called the Projected Gradient Descent (PGD). In addition to these gradient-based methods, Carlini and Wagner (C&W attack) [6] explored the use of maximum marginal loss and optimization method to generate adversarial examples with high fooling rate and small distortion with respect to the original image. In addition, more and more black-box attacks have been proposed. These so-called black-box attack can successfully change the model prediction without knowing the parameters and structure of the attacked model. Uesato et al. [31] proposed simultaneous perturbation stochastic approximation (SPSA) attack. It is a gradient-free query-based attack method, which minimizes the output logits of the true label and the largest logits of the rest of labels. Chen et al. [32] proposed the hard-label RayS attack, which only relies on the hard-label output of the target model and utilizes a fast check step to skip unnecessary searches. This significantly saves the number of queries needed for the hard-label attack.

Apart from the development of adversarial attack in the field of natural images, medical image domain has also paid

more and more attention to this topic. Ma et al. [33] analyzed the different behaviors of medical images and natural images when attacked by adversarial examples, and concluded that medical images are more vulnerable to adversarial attacks. Other studies [7], [8], [34] evaluated the robustness of deep diagnostic models on different tasks by adversarial attacks.

### C. Adversarial Defense

Considering the importance of network robustness, many defense methods have been proposed [35]–[37]. Among which, Adversarial Training (AT) has been demonstrated to be one of the most effective defense methods. AT can be regarded as a data augmentation technology, that trains network on adversarial examples. After that, many methods were improved based on AT and showed superior performance. TRADES [38] trades adversarial robustness off against accuracy. The objective function of TRADES is a linear combination of natural loss and regularization term. MART [39] differentiates the misclassified examples and correctly classified examples during adversarial training and adopts a regularized adversarial loss involving both adversarial and natural examples to improve the robustness of models. For medical image field, Liu et al. [40] propose the augmentation method to add adversarial synthetic nodules and adversarial attack samples to the training data to improve the generalization and the robustness of the lung nodule detection systems. However, these methods require retraining the model, which is not friendly to the medical diagnostic models that have been already deployed online.

Besides adversarial training methods, many pre-processing based-defense methods have been proposed. Xie et al. [12] use random resizing and padding (Random R-P) to pre-process the input images before feeding the images into the models to make predictions. Jia et al. propose the ComDefend [13], which consists of a compression convolutional neural network (ComCNN) and a reconstruction convolutional neural network (RecCNN) to transform the adversarial image to its clean version. However, the random resizing and compression operators may lose the lesion features of medical images.

## III. MATERIALS

In this section, we will introduce in detail the datasets and pre-trained models used in our study.

### A. Datasets

Two public datasets are used in this study, including:

- 1) *COVID-19 Radiography Database* [14]: It consists of chest X-ray images with size  $224 \times 224$  of COVID-19 positive, normal, and viral pneumonia images (i.e., 3-class diagnostic task). In the current release, there are 1200 COVID-19 positive images, 1341 normal images, and 1345 viral pneumonia images. We have split this dataset into training, validation, and test set with ratio 8 : 1 : 1.

2) *DermaMNIST* [41]: It is based on HAM10000 [42], [43], which consists of 10, 015 multi-source dermatoscopic images of common pigmented skin lesions. This dataset is labeled as 7 different categories (i.e., actinic keratoses, basal cell carcinoma, benign keratosis, dermatofibroma, melanocytic nevi, melanoma, vascular), as a 7-class classification task. We have split the images into training, validation and test set with ratio 7 : 1 : 2. The source images of  $3 \times 600 \times 450$  are resized into  $3 \times 28 \times 28$ .

### B. Pre-trained Models

In order to better explore the effect of MedRDF on different pretrained models, the base classifiers we use in experiments are natural image based ResNet-18 and ResNet-50 [44] and medical image based AG-Sononet-16 [26]. We directly train the networks on COVID-19 and DermaMNIST datasets without fine-tuning. The ResNet-18 and AG-Sononet-16 are trained for 100 epochs using stochastic gradient descent with momentum 0.9 and weight decay  $1e^{-6}$ . The initial learning rate is  $1e^{-4}$  and is decayed by 0.1 on 50 and 75 epochs. The ResNet-50 is trained for 100 epochs using stochastic gradient descent with momentum 0.9 and weight decay  $1e^{-4}$ . The initial learning rate is  $1e^{-3}$  and is decayed by 0.1 on 50 and 75 epochs. The batch size is 10.

## IV. METHODS

### A. Problem Formulations

Let  $\mathcal{X} \in \mathbb{R}^d$  denote the input space and  $\mathcal{Y} = \{1, \dots, K\}$  be a finite set consists of  $K$  possible class labels.  $D = \{(\mathbf{x}_1, y_1), \dots, (\mathbf{x}_m, y_m)\}$  is a training set with  $m$  labeled examples, where  $\mathbf{x}_i \in \mathcal{X}$  is the feature vector and  $y_i \in \mathcal{Y}$  is the label of the  $i$ -th example. Given a medical diagnostic model  $h_\theta$  with parameters  $\theta$ , it outputs the class label  $h_\theta(\mathbf{x}_i)$  for each input image  $\mathbf{x}_i \in \mathcal{X}$ :

$$h_\theta(\mathbf{x}_i) = \arg \max_{k=1, \dots, K} p_k(\mathbf{x}_i, \theta), \quad (1)$$

where  $p_k(\mathbf{x}_i, \theta)$  is the probability (softmax on logits) of  $\mathbf{x}_i$  belonging to class  $k$ . We denote  $\mathcal{A}_{h_\theta}$  as the space of adversarial examples for the pre-trained model  $h_\theta$ . Adversarial example  $\mathbf{x}' \in \mathcal{A}_{h_\theta}$  is supposed to be quasi-imperceptible to the human eye and misclassified by  $h_\theta$ , i.e.,

$$d(\mathbf{x}, \mathbf{x}') \leq \epsilon \quad \text{and} \quad h_\theta(\mathbf{x}) \neq h_\theta(\mathbf{x}'), \quad (2)$$

where  $d(\cdot)$  is the distance function,  $\epsilon$  is the maximum perturbation for adversarial attack. Here we aim to design a robust diagnostic framework  $g$  to correctly classify these adversarial examples  $\mathbf{x}' \in \mathcal{A}_{h_\theta}$  with the pre-trained model  $h_\theta$ .

### B. Framework Details

Inspired by random smoothing [45], we construct a robust and retrain-less diagnostic framework (MedRDF)  $g$  for medical pretrained model  $h_\theta$ . Fig. 1 illustrates the flowchart of proposed MedRDF  $g$ . Specifically, MedRDF returns the class label which the base classifier  $h_\theta$  is most likely to return. *Firstly*, MedRDF perturbs input  $\mathbf{x}$  by isotropic noise

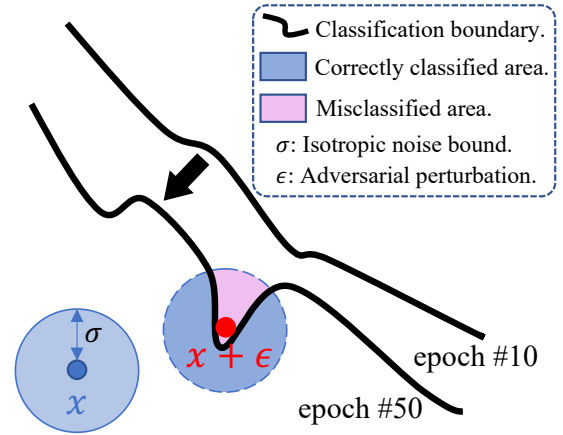


Fig. 2: The effect of isotropic noise. A significant increase in the curvature of the decision boundary during continuous training from epoch #10 to epoch #50. The hypercubes centered at  $\mathbf{x} + \epsilon$  and  $\mathbf{x}$  intersect the most with the area that examples can be correctly classified.

$\eta$  from distribution  $\mu$  to produce the noisy copies, and then denoises the noisy copies by pre-defined denoiser  $\mathcal{D}$ . *Secondly*, MedRDF inputs these copies to pre-trained medical diagnostic model  $h_\theta$  to obtain the prediction labels. *Thirdly*, the final diagnostic result is obtained by majority voting based on the labels of denoised copies. The MedRDF  $g$  is formulated as follows:

$$g(\mathbf{x}) = \arg \max_{k \in \mathcal{Y}} \mathbb{P}(h'_\theta(\mathbf{x} + \eta) = k), \quad (3)$$

$$\eta \sim \mu(0, \sigma \mathbf{I}),$$

where  $h'_\theta(\cdot)$  represents  $h_\theta(\mathcal{D}(\cdot))$ ,  $\mathcal{D}(\cdot)$  is the pre-defined denoiser. An equivalent definition is that  $g(\mathbf{x})$  returns the class  $k$  whose pre-image  $\{\mathbf{x} + \eta \in \mathbb{R}^d : h'_\theta(\mathbf{x} + \eta) = k\}$  has the largest probability measure under the distribution  $\mu(\mathbf{x}, \sigma \mathbf{I})$ . The level of noise  $\eta$  has been bounded by  $\sigma$ , where the noise level  $\sigma$  controls the tradeoff between robustness and accuracy, i.e., the robustness of the MedRDF increases with the increasing of  $\sigma$  while its standard accuracy decrease.

The detailed information of our MedRDF is described as follows:

1) *Isotropic Noise*  $\eta$ : Recent studies [4], [46] show that the non-robustness of deep networks against attacks is caused by the high nonlinearity of deep networks.

TABLE I: Classification accuracy (%) of the MedRDF  $g$  on COVID-19 with salt-and-pepper noise (with level as  $\sigma$ ) and median filter. The attacker is crafted by PGD with 100 steps and maximum  $L_\infty$  perturbation  $\epsilon$  on ResNet-50.

$\epsilon \backslash \sigma$	0.05	0.1	0.15	0.2	0.25	0.3
0	<b>93.4</b>	91.6	77.4	50.0	45.8	30.0
2/255	<b>93.0</b>	91.8	79.4	51.6	45.8	44.2
4/255	<b>92.4</b>	91.8	82.0	54.0	46.0	44.0
8/255	91.1	<b>91.2</b>	84.8	59.6	47.0	44.2
16/255	87.8	<b>91.4</b>	89.4	73.0	51.2	46.6

---

**Algorithm 1** MedRDF: Robust and Retrain-Less Diagnostic Framework for Medical pretrained models
 

---

**Input** base classifier  $h_\theta$ , diagnostic case  $\mathbf{x}$ , noise distribution  $\mu(0, \sigma \mathbf{I})$ , sampling numbers  $n$ , abstention threshold  $\alpha$ , denoiser operator  $\mathcal{D}$ .

Initialization array: counts[0,  $\dots$ ,  $n - 1$ ]

**for**  $i \leq n$  **do**

Sample noise  $\eta_i \sim \mu(0, \sigma \mathbf{I})$

counts[ $h_\theta(\mathcal{D}(\mathbf{x} + \eta_i))$ ] ++

**end for**

$n_A, n_B \leftarrow$  top two indices in counts

**if** Binom( $n_A, n_A + n_B, 1/2$ )  $\leq \alpha$  **Output**  $k_A$

**else return -1 (ABSTAIN)**

---

**TABLE II:** Accuracy (%) and test time (s) on each image of MedRDF on different number of copies. The base classifier is ResNet-18. The common noise is salt-and-pepper noise with  $\sigma = 0.1$ , the denoiser in median filter, and maximum  $L_\infty$  perturbation  $\epsilon = 8/255$ . The number after attack method represents the number of iteration steps. The bold number represents the result of our selection.

Datasets	$n$	Natural	I-FGSM-7	PGD-7	C&W	Time(s)
COVID-19	$1e^2$	91.2	85.4	92.2	89.2	0.1
	$1e^3$	91.2	86.2	93.0	90.0	1.1
	$1e^4$	<b>91.2</b>	<b>86.2</b>	<b>93.2</b>	<b>90.4</b>	<b>3.8</b>
	$1e^5$	91.4	86.2	93.2	90.0	87.6
DermaMNIST	$1e^2$	68.9	61.3	65.4	64.3	0.1
	$1e^3$	69.0	62.5	67.0	65.6	0.1
	$1e^4$	<b>69.0</b>	<b>63.1</b>	<b>67.3</b>	<b>66.0</b>	<b>1.1</b>
	$1e^5$	70.4	63.2	67.3	65.9	10.0

Kalimeris et al. [47] show that with the continuous training of the network, a significant increasing in the curvature of the decision boundary and loss landscape will occur, and the adversarial examples are easy to hide in these isolated regions with high curvature [48], as illustrated in Fig. 2. Based on this observation, we add the common random noise  $\eta$  bounded by  $\sigma$  to original image, which can reduce the impact of the adversarial example in isolated area on the accuracy of the model. As shown in Fig. 2, in the noise area, with  $\mathbf{x}$  and  $\mathbf{x} + \epsilon$  as the center and maximum noise  $\sigma$  as the boundary, most examples can be correctly classified. The result is also true for the adversarial example in the isolated area. Therefore, adding isotropic noise to the original image to generate noisy copies can effectively instruct the network not to be misled by adversarial examples. However, although neural networks have certain robustness to common noise, too large noise will still lead to the accuracy decrease of  $h_\theta$ , which will also affect the final prediction result of  $g$  based on  $h_\theta$ . In the following subsection, we will introduce the denoising operator to alleviate the decline of accuracy.

**2) Pre-defined Denoiser:** To alleviate the accuracy decline of the base classifier  $h_\theta$  under large isotropic noise, denoising operator has been adopted in our MedRDF. Instead of CNN-based denoiser [49], we use Gaussian Smoothing (GS) and median Filter (MF) as denoisers in our work, which have faster

inference speed and more efficient GPU memory than CNN-based denoiser.

**3) Prediction and Majority Voting:** For notational convenience, we define Equation (3) as  $P_k = \mathbb{P}(h'_\theta(\mathbf{x} + \boldsymbol{\eta}) = k)$ . Let  $\hat{k}_A = \arg \max_k P_k$ . Notice that by definition,  $g(\mathbf{x}) = \hat{k}_A$ . We draw  $n$  noise examples with Markov Chain Monte Carlo (MCMC) principle from distribution  $\mu(0, \sigma \mathbf{I})$ , and inquire  $n$  noise-corrupted copies of  $\mathbf{x}$  through the base classifier  $h'_\theta(\cdot)$ . Sample a vector of class counts  $\{n_k\}_{k \in \mathcal{Y}}$  from Multinomial( $\{P_k\}_{k \in \mathcal{Y}}, n$ ). Let  $k_A = \arg \max_k n_k$  be the class whose count is largest. Let  $n_A$  and  $n_B$  be the largest count and the second-largest count, respectively. If  $k_A$  appears much more often than any other class, then the prediction of MedRDF  $g$  returns  $k_A$ . Otherwise, it abstains from making a prediction. As Cohen et al. [45] declared, we use the hypothesis test from Hung & Fithian [50] to calibrate the abstention threshold so as to bound by  $\alpha$  the probability of returning an incorrect answer. The prediction of MedRDF  $g$  satisfies the following guarantee:

*Proposition 1:* With probability at least  $1 - \alpha$  over the randomness in the prediction of MedRDF, the probability that the prediction of MedRDF returns a class other than  $\hat{k}_A$  is at most  $\alpha$ , i.e.,

$$\mathbb{P}(g(\mathbf{x}) \neq \hat{k}_A) \leq \alpha. \quad (4)$$

We use the p-value of the two-sided hypothesis test that  $n_A$  is drawn from binomial distribution Binom( $n_A + n_B, 1/2$ ) to verify whether Equation (4) holds. If the p-value is less than  $\alpha$ , then return  $k_A$ . Else, abstain. i.e., we can adopt two-sided hypothesis test with binomial distribution (Binom) to justify the randomness in the prediction of MedRDF:

$$\text{Binom}(n_A, n_A + n_B, 1/2) \leq \alpha. \quad (5)$$

The proof is as follows:

*Proof 1:* MedRDF returns a class other than  $\hat{k}_A$  if and only if (1)  $k_A \neq \hat{k}_A$  and (2) MedRDF does not abstain. We have:

$$\begin{aligned} \mathbb{P}(g(\mathbf{x}) \neq \hat{k}_A) &= \mathbb{P}(k_A \neq \hat{k}_A, \text{MedRDF does not abstain}) \\ &= \mathbb{P}(k_A \neq \hat{k}_A) \mathbb{P}(\text{MedRDF does not abstain} | k_A \neq \hat{k}_A) \\ &\leq \mathbb{P}(\text{MedRDF does not abstain} | k_A \neq \hat{k}_A). \end{aligned} \quad (6)$$

Recall that MedRDF does not abstain if and only if the p-value of the two-sided hypothesis test that  $n_A$  is drawn from Binom( $n_A + n_B, 1/2$ ) is less than  $\alpha$ . Theorem 1 in Hung & Fithian [50] proves that the conditional probability that this event occurs given that  $k_A \neq \hat{k}_A$  is exactly  $\alpha$ . That is,

$$\mathbb{P}(\text{MedRDF does not abstain} | k_A \neq \hat{k}_A) = \alpha. \quad (7)$$

Therefore, we have:

$$\mathbb{P}(g(\mathbf{x}) \neq \hat{k}_A) \leq \alpha. \quad (8)$$

When  $\alpha$  is small, MedRDF abstains frequently but rarely returns the wrong class. When  $\alpha$  is large, MedRDF usually makes a prediction, but may often return the wrong class.  $\alpha = 0.001$  and  $n = 1e^4$  have been adopted in our framework. The complete prediction procedure of MedRDF  $g$  is described in Algorithm 1.

**TABLE III:** Accuracy (%) of different defense mechanism (rows) against white box adversarial attacks with maximum  $L_\infty$  perturbation  $\epsilon = 8/255$  (columns) on **COVID-19 and DermaMNIST datasets with ResNet-18**. The original accuracy of each defense is described in the column “Natural”. GS: gaussian smoothing, MF: median filter. The number after attack method represents the number of iteration steps.

Dataset	Method	Denoiser	Natural	I-FGSM				PGD			C&W
				1	2	5	7	7	20	100	
COVID-19	ResNet-18	None	94.4	5.6	0.0	0.0	0.0	0.0	0.0	0.0	0.0
		GS	94.4	5.6	1.4	0.0	0.0	0.4	0.0	0.0	0.0
		MF	94.4	35.2	23.8	1.2	0.2	1.8	0.0	0.0	0.0
	MedRDF with gaussian noise	None	28.8	40.3	51.1	36.4	35.5	30.8	28.8	27.7	20.8
		GS	94.8	39.4	72.8	48.6	45.4	73.8	55.0	47.6	55.4
		MF	94.2	68.4	87.6	75.2	73.8	89.4	82.2	78.2	82.4
	MedRDF with s.p. noise	None	65.4	28.4	44.6	24.8	23.8	39.4	26.2	21.2	26.0
		GS	<b>95.6</b>	42.6	79.6	53.2	49.0	82.2	64.2	54.0	64.0
		MF	91.2	<b>84.6</b>	<b>89.8</b>	<b>86.6</b>	<b>86.2</b>	<b>93.2</b>	<b>90.2</b>	<b>89.6</b>	<b>90.4</b>
	MedRDF with poisson noise	None	32.4	28.4	29.8	28.4	28.2	29.4	28.6	28.6	28.8
		GS	94.4	45.2	75.6	54.6	50.8	78.0	62.8	55.6	62.6
		MF	92.2	71.2	86.6	77.2	75.6	89.0	81.0	79.0	81.4
DermaMNIST	ResNet-18	None	<b>74.1</b>	1.7	4.2	0.0	0.0	0.1	0.0	0.0	0.0
		GS	71.5	12.5	25.2	2.5	1.5	16.2	1.1	0.4	0.9
		MF	72.5	55.3	53.2	31.3	22.3	38.2	15.1	4.0	16.7
	MedRDF with gaussian noise	None	44.6	5.4	26.7	7.0	5.1	26.0	12.3	10.1	16.9
		GS	67.9	44.6	59.1	49.7	47.9	60.9	53.3	49.8	54.8
		MF	70.9	61.1	<b>66.8</b>	<b>63.5</b>	62.7	<b>67.4</b>	64.8	63.3	65.8
	MedRDF with s.p. noise	None	68.5	15.6	45.9	15.0	9.4	44.6	17.1	8.3	24.6
		GS	68.4	48.1	60.5	52.2	50.9	62.4	56.0	53.5	57.4
		MF	69.0	<b>62.1</b>	66.5	<b>63.5</b>	<b>63.1</b>	67.3	<b>65.1</b>	<b>64.1</b>	<b>66.0</b>
	MedRDF with poisson noise	None	58.9	13.4	38.0	16.7	11.6	37.0	18.9	11.4	26.5
		GS	68.3	47.5	59.1	51.9	51.0	60.8	55.0	52.0	57.4
		MF	69.4	60.0	65.0	62.0	61.5	66.3	63.9	62.9	65.3

### C. Robust Metric

In medical diagnostic tasks, in addition to the diagnostic results output by the model, we also hope to obtain the confidence score of the results, so as to better guide the doctor’s follow-up work, such as adopting the result or re-evaluating this case. Therefore, in this subsection, in order to provide doctors with intuitive and effective indicator, we define a Robust Metric (RM) based on MedRDF. The formulation of RM is presented as follows:

$$\text{RM} = \frac{K * (n_A - n_B)}{n}, \quad (9)$$

where  $n_A$  and  $n_B$  denote the number of classes  $k_A$  and  $k_B$  with the most and second most occurrences of  $g$ , respectively.  $K$  is number of diagnosis categories. Setting the threshold of RM, when the RM output by MedRDF is greater than the threshold, the doctor can accept this diagnostic result. Otherwise, doctor should consider re-evaluating this result. The effectiveness of RM is analyzed as follows:

From Equation (9) we can obtain:

$$(P_{k_A})_{\min} = \left(\frac{n_A}{n}\right)_{\min} = \frac{1}{K} + \frac{K-1}{K^2} \text{RM}. \quad (10)$$

Then for different classification tasks, doctors can set different thresholds to make the probability of output labels reach their expectations. Take the 3-class diagnostic task as an example, we set a threshold of RM with 1 for indicating the diagnostic result to be robust or not, that is to say,  $k_A$  should have at least

5/9 probability for robust evaluation. For 7-class diagnostic task, when setting the threshold of RM as 3, the probability of class  $k_A$  is at least 0.51.

## V. EXPERIMENTS

In this section, we first introduce the experimental settings including the common isotropic noises and adversarial attack methods we used in this study. Second, we choose the best noise boundry  $\sigma$  and the number  $n$  of the copies for our experiments by ablation study. Then, we conduct a set of experiments to evaluate the robustness of our MedRDF under different adversarial attacks. Furthermore, we confirm the necessity and effectiveness of our RM indicator and visually present the robust diagnostic results for different cases. Finally, we have shown more comparable results on our MedRDF with other augmentation strategies and other defense methods.

### A. Experimental Settings

1) *Common Isotropic Noise*: We evaluate the robustness of MedRDF under gaussian noise, salt-and-pepper (s.p.) noise and poisson noise, and utilize gaussian smoothing (GS) and median filter (MF) as denoisers in experiments.

2) *Adversarial Attack*: The adversarial examples are crafted by the most challenging “white-box” attacks (i.e., I-FGSM [51], PGD [5], and C&W [6]) and “black-box” attacks (i.e., SPSA [31] and RayS [32]). The “white-box” attacks are under maximum  $L_\infty$  perturbation  $\epsilon = 8/255$ .

**TABLE IV:** Accuracy (%) of different defense mechanism (rows) against white box adversarial attacks with maximum  $L_\infty$  perturbation  $\epsilon = 8/255$  (columns) on **COVID-19 and DermaMNIST dataset with ResNet-50**. The original accuracy of each defense is described in the column “Natural”. GS: gaussian smoothing, MF: median filter. The number after attack method represents the number of iteration steps.

Dataset	Method	Denoiser	Natural	I-FGSM				PGD			C&W
				1	2	5	7	7	20	100	
COVID-19	ResNet-50	None	92.6	62.8	7.0	0.0	0.0	0.0	0.0	0.0	0.0
		GS	92.6	62.8	7.0	0.2	0.0	0.0	0.0	0.0	0.0
		MF	92.6	70.2	17.8	1.0	0.2	0.4	0.0	0.0	0.0
	MedRDF with gaussian noise	None	82.6	68.2	77.6	73.2	72.2	78.2	76.4	73.4	73.6
		GS	91.4	74.6	84.6	73.0	67.8	87.0	80.6	71.0	70.2
		MF	<b>93.6</b>	88.2	91.8	88.8	87.2	<b>93.0</b>	91.4	89.4	89.0
	MedRDF with s.p. noise	None	89.8	78.4	86.6	82.6	81.6	87.6	85.4	82.2	83.2
		GS	90.4	83.8	88.8	85.6	84.4	88.8	88.4	87.6	86.8
		MF	91.6	<b>91.2</b>	<b>91.8</b>	<b>91.8</b>	<b>91.2</b>	91.6	<b>91.4</b>	<b>91.2</b>	<b>91.2</b>
	MedRDF with poisson noise	None	77.0	66.6	75.4	69.8	68.4	75.4	73.0	69.2	72.4
		GS	88.0	80.8	83.4	82.0	80.4	83.6	83.0	81.6	83.2
		MF	87.8	83.8	86.4	85.0	84.0	88.2	87.4	85.6	87.0
DermaMNIST	ResNet-50	None	73.0	8.9	18.8	1.1	1.0	9.3	1.0	0.9	0.1
		GS	71.5	4.9	11.9	0.8	0.4	5.8	0.4	0.4	0.1
		MF	<b>73.2</b>	35.3	40.4	11.5	4.6	26.4	4.7	1.5	4.9
	MedRDF with gaussian noise	None	69.9	13.5	34.7	6.6	4.2	32.9	7.1	5.4	6.6
		GS	67.3	29.8	52.7	27.7	25.8	54.4	31.6	30.0	32.4
		MF	72.0	54.6	73.1	55.0	53.0	75.9	60.0	57.9	62.0
	MedRDF with s.p. noise	None	71.3	12.7	32.4	6.2	4.3	31.0	8.4	6.2	8.1
		GS	65.3	30.3	53.0	30.0	28.6	56.0	35.0	33.1	35.8
		MF	68.7	<b>60.8</b>	<b>75.5</b>	<b>62.8</b>	61.3	<b>78.5</b>	67.3	66.4	69.1
	MedRDF with poisson noise	None	70.0	21.4	38.5	19.1	17.3	39.0	21.2	19.3	21.8
		GS	64.5	33.8	56.2	35.9	34.7	59.1	40.9	38.5	41.0
		MF	70.3	60.5	74.0	62.5	<b>62.7</b>	77.3	<b>68.1</b>	<b>67.2</b>	<b>70.0</b>

## B. Ablation Study

### 1) The level of common noise and adversarial perturbation:

We first explore the influence of common noise  $\sigma$  and adversarial perturbation  $\epsilon$  on original and robust accuracy. As shown in TABLE I, since we set  $\epsilon = 8/255$  for adversarial attack, we choose  $\sigma = 0.1$  for the boundry of the common noise for its high accuracy.

2) *The number of the copies:* In this part, we explore the influence of different number of copies to the final robustness of MedRDF. The defense accuracy and test time of MedRDF on different number  $n$  of copies are recorded in TABLE II. As shown in TABLE II, both on COVID-19 and DermaMNIST datasets, although the natural accuracy on  $n = 1e^5$  is higher than it on  $n = 1e^4$ , the test time on each image when  $n = 1e^5$  is much longer than it on  $n = 1e^4$ , which is not conducive to the clinical application of the MedRDF. For example, on COVID-19 dataset, the natural accuracy is 91.4% on  $n = 1e^5$ , which is little higher than 91.2% on  $n = 1e^4$ . However, its test time on each image is 87.6s, which is not easily tolerated when compared with the test time 3.8s on  $n = 1e^4$ . Besides, in terms of defense accuracy, it can be seen that  $n = 1e^4$  has a greater impact on the final defense accuracy on MedRDF, compared with that on  $n = 1e^3$  and  $n = 1e^5$ . For instance, with  $n = 1e^4$ , the accuracy on DermaMNIST attacked by C&W is 66.0%, which is higher than the accuracy (65.6% and 65.9%) of  $n = 1e^3$  and  $n = 1e^5$ , respectively. In summary, we choose the number  $n = 1e^4$  in our experiments.

## C. Robustness Evaluation and Analysis.

1) *Quantitative Results:* In this part, we present the quantitative results of base model and MedRDF under different white-box attacks and black-box attacks.

**White box attack.** The accuracy of original models (i.e.,  $h_\theta = \text{ResNet-18, ResNet-50, and AG-Sononet-16}$ ) and our MedRDF (i.e.,  $g_\theta$  based on ResNet-18, ResNet-50, and AG-Sononet-16) are recorded in TABLE III, TABLE IV and TABLE V, respectively. As shown in TABLE III, the original model ResNet-18 is vulnerable to adversarial attacks both on COVID-19 and DermaMNIST diagnostic tasks (e.g., its accuracy drops to 0.0% under C&W attack on COVID-19 and DermaMNIST). The other result we can observe from TABLE III is that, MedRDF markedly improves the robustness of original model in all attack settings (e.g., when under PGD-7 attack, the accuracy of MedRDF with gaussian noise and MF denoiser is 67.4% while the original pre-trained model is 0.1% on DermaMNIST). The same result can also be found in TABLE IV on ResNet-50 and TABLE V on AG-Sononet-16. As shown in TABLE IV, MedRDF even maintains a better performance on natural accuracy (i.e., the natural accuracy on COVID-19 of MedRDF with gaussian noise and MF denoiser is 93.6%, while the original ResNet-50 is 92.6%). As shown in TABLE V, one can observe that the classification results of two datasets have been significantly improved after using MedRDF. For instance, on DermaMNIST dataset, the defense accuracy of MedRDF with gaussian noise and median filter is

**TABLE V:** Accuracy (%) of different defense mechanism (rows) against white box adversarial attacks with maximum  $L_\infty$  perturbation  $\epsilon = 8/255$  (columns) on **COVID-19 and DermaMNIST dataset with AG-Sononet-16**. The original accuracy of each defense is described in the column “Natural”. GS: gaussian smoothing, MF: median filter. The number after attack method represents the number of iteration steps.

Dataset	Method	Denoiser	Natural	I-FGSM				PGD			C&W
				1	2	5	7	7	20	100	
COVID-19	AG-Sononet-16	None	<b>93.4</b>	29.2	12.0	0.8	0.2	0.2	0.0	0.0	0.0
		GS	90.4	41.8	10.0	0.0	0.0	16.2	0.0	0.0	0.0
		MF	92.0	75.6	39.6	17.6	10.6	26.8	3.0	0.0	0.0
	MedRDF with gaussian noise	None	28.2	28.2	28.2	28.2	28.2	28.2	28.2	28.2	28.2
		GS	75.0	44.8	64.6	53.0	50.4	66.4	60.2	49.6	52.4
		MF	88.6	85.6	87.4	85.6	84.2	87.2	86.2	83.6	83.8
	MedRDF with s.p. noise	None	28.2	28.2	28.2	28.2	28.2	28.2	28.2	28.2	28.2
		GS	92.0	65.2	79.2	64.2	58.8	83.0	73.2	60.6	60.0
		MF	87.0	<b>88.0</b>	87.0	<b>86.4</b>	83.8	<b>89.2</b>	<b>88.8</b>	<b>88.6</b>	<b>88.6</b>
	MedRDF with poisson noise	None	28.2	28.2	28.2	28.2	28.2	28.2	28.2	28.2	28.2
		GS	88.8	84.2	88.0	85.6	<b>85.4</b>	88.6	87.8	86.4	87.2
		MF	90.0	86.4	<b>90.0</b>	86.0	83.4	88.7	88.4	86.8	87.3
DermaMNIST	AG-Sononet-16	None	<b>70.6</b>	38.7	31.4	11.5	5.5	17.7	2.7	0.1	3.5
		GS	68.2	42.0	38.4	24.3	18.8	31.9	19.3	6.8	19.3
		MF	71.4	55.1	49.4	35.7	30.7	36.4	24.0	9.6	25.7
	MedRDF with gaussian noise	GS	42.2	25.5	36.4	29.4	27.1	37.9	31.6	27.3	32.4
		MF	65.6	57.4	62.4	60.3	59.8	64.4	63.4	62.6	63.5
		MF	70.5	63.2	68.8	66.2	65.9	<b>70.3</b>	<b>69.5</b>	<b>68.7</b>	<b>69.5</b>
	MedRDF with s.p. noise	None	66.4	45.1	56.8	46.6	43.4	59.5	53.6	47.5	55.0
		GS	66.3	58.5	64.3	62.7	61.7	66.4	65.5	64.9	65.5
		MF	68.3	<b>66.0</b>	<b>67.4</b>	66.1	<b>66.3</b>	68.9	68.5	68.2	68.8
	MedRDF with poisson noise	None	61.9	45.4	54.5	47.6	45.1	56.6	52.1	48.3	53.4
		GS	68.5	58.9	65.7	62.9	62.8	66.7	66.3	65.7	66.2
		MF	70.3	64.3	67.2	<b>66.5</b>	<b>66.3</b>	69.2	68.8	68.0	69.0

**TABLE VI:** Accuracy (%) of original model ResNet-50 and AG-Sononet-16 on COVID-19 under different settings. The common noise  $\sigma = 0.1$ , GS: gaussian smoothing, MF: median filter.

Method	Natural	Gaussian Noise			Salt-and-Pepper Noise		
		None	GS	MF	None	GS	MF
ResNet-50	92.6	81.8	91.2	93.4	90.0	90.6	91.2
AG-Sononet-16	93.4	<b>28.2</b>	87.4	88.0	<b>28.2</b>	85.0	85.6

70.3% when attacked by PGD-7, which is much better than the accuracy of base model AG-Sononet-16 (i.e., 17.7%), and even is comparable with that without any attack (i.e., natural accuracy 70.6%). These results indicate the effectiveness of our framework to convert non-robust models to robust ones.

Moreover, TABLE VI records the accuracy of COVID-19 with original models ResNet-50 and AG-Sononet-16 under different noise settings. The result we can obtain from TABLE VI is that, since original AG-Sononet-16 model is not robust to common noise (i.e., the natural accuracy is 28.2% after adding noise without denoiser), the MedRDF without denoiser will lose its discrimination ability (i.e., the accuracy is 28.2% under all attacks with random guess in TABLE V). This result has attracted our attention that the robustness of base model  $h_\theta$  under common noise will affect the final robustness of MedRDF under adversarial attack.

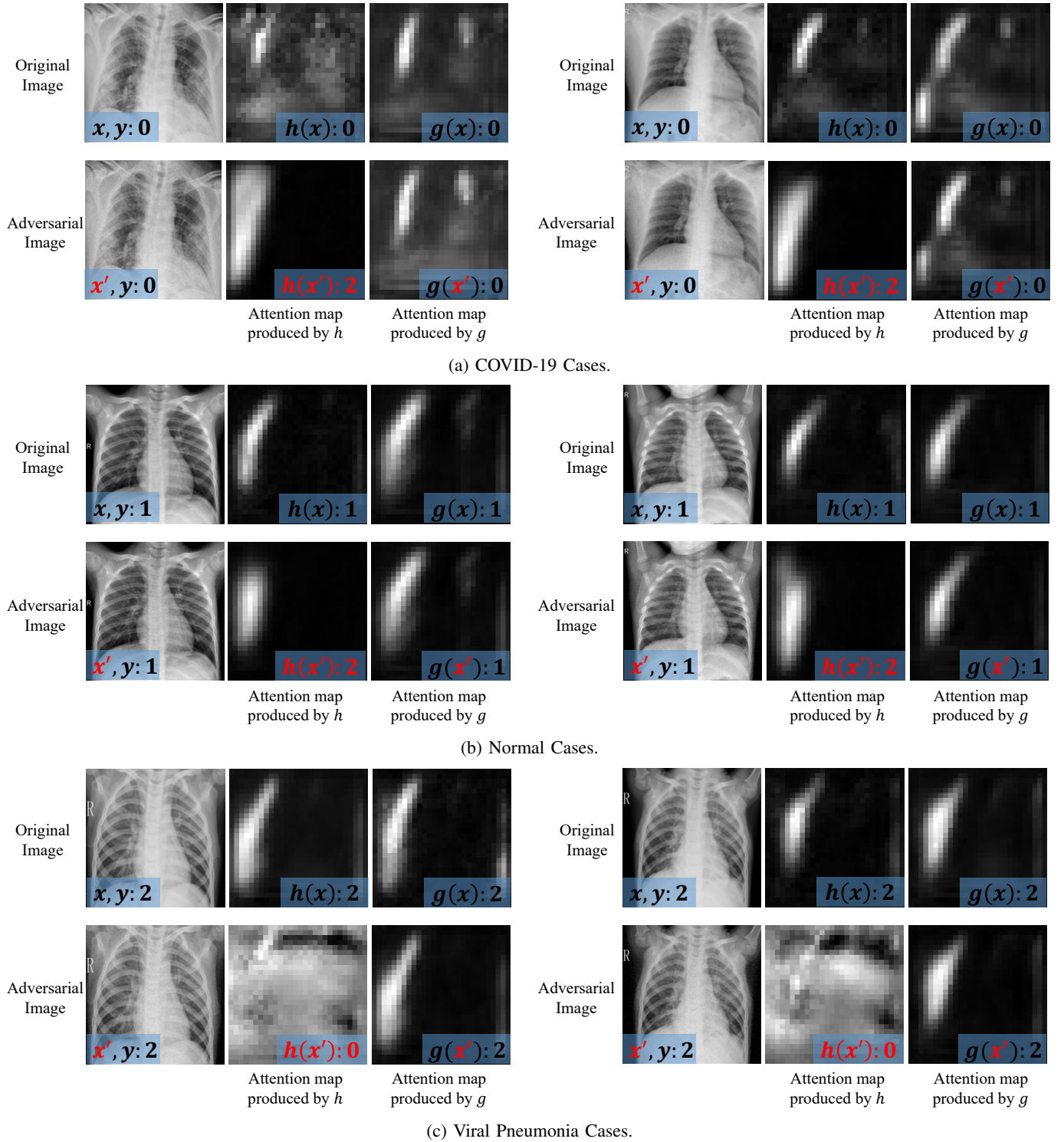
**Black box attack.** For SPSA attack, to estimate the gradi-

**TABLE VII:** Accuracy (%) of different defense mechanism (rows) against **black box adversarial attacks** on COVID-19 and DermaMNIST dataset with **ResNet-18**. GS: gaussian smoothing, MF: median filter, s.p. noise: salt-and-pepper noise.

Method	Denoiser	COVID-19		DermaMNIST	
		SPSA	RayS	SPSA	RayS
ResNet-18	None	0.7	0.0	0.0	0.0
	GS	2.7	1.0	0.9	3.7
	MF	36.0	1.4	15.8	6.9
MedRDF with gaussian noise	None	52.7	64.6	63.2	70.7
	GS	86.7	85.6	68.8	69.9
	MF	86.7	86.2	<b>71.4</b>	<b>72.7</b>
MedRDF with s.p. noise	None	66.0	66.0	56.5	71.5
	GS	86.5	86.4	68.7	70.2
	MF	86.7	86.6	71.3	72.3
MedRDF with poisson noise	None	60.5	60.6	60.7	70.0
	GS	86.6	86.6	68.8	69.5
	MF	<b>86.7</b>	<b>86.6</b>	70.8	71.7

ents, we set the batch size as 128, the perturbation as  $8/255$ , and the learning rate as 0.01. We run SPSA attack for 100 iterations, and early-stop when we cause misclassification. For RayS attack, we set the  $L_\infty$  perturbation  $\epsilon$  as  $8/255$ . The accuracy of original models (i.e.,  $h_\theta = \text{ResNet-18}$  and AG-Sononet-16) and our MedRDF (i.e.,  $g_\theta$  based on ResNet-18 and AG-Sononet-16) attacked by SPSA and RayS are recorded





**Fig. 3:** Attention maps of base model  $h$  and MedRDF  $g$  on original images and adversarial images. The first row of each subfigure contains original image and corresponding attentions maps of base model  $h$  and MedRDF  $g$ , respectively. The second row of each subfigure contains adversarial image generated by PGD with 100 steps and corresponding attentions maps of base model  $h$  and MedRDF  $g$ , respectively. Red denotes the adversarial images and wrong labels. Base model  $h$  is AG-Sononet-16, MedRDF  $g$  is based on AG-Sononet-16 architecture with salt-and-pepper noise and median filter.

in TABLE VII and TABLE VIII, respectively. As shown in TABLE VII, the base model ResNet-18 with our MedRDF obtain better robustness in each attack. For instance, on COVID-

19 dataset, MedRDF with poisson noise and median filter achieves 86.7% accuracy when attacked by SPSA, which is much higher than that (0.7%) on base ResNet-18. Meanwhile,

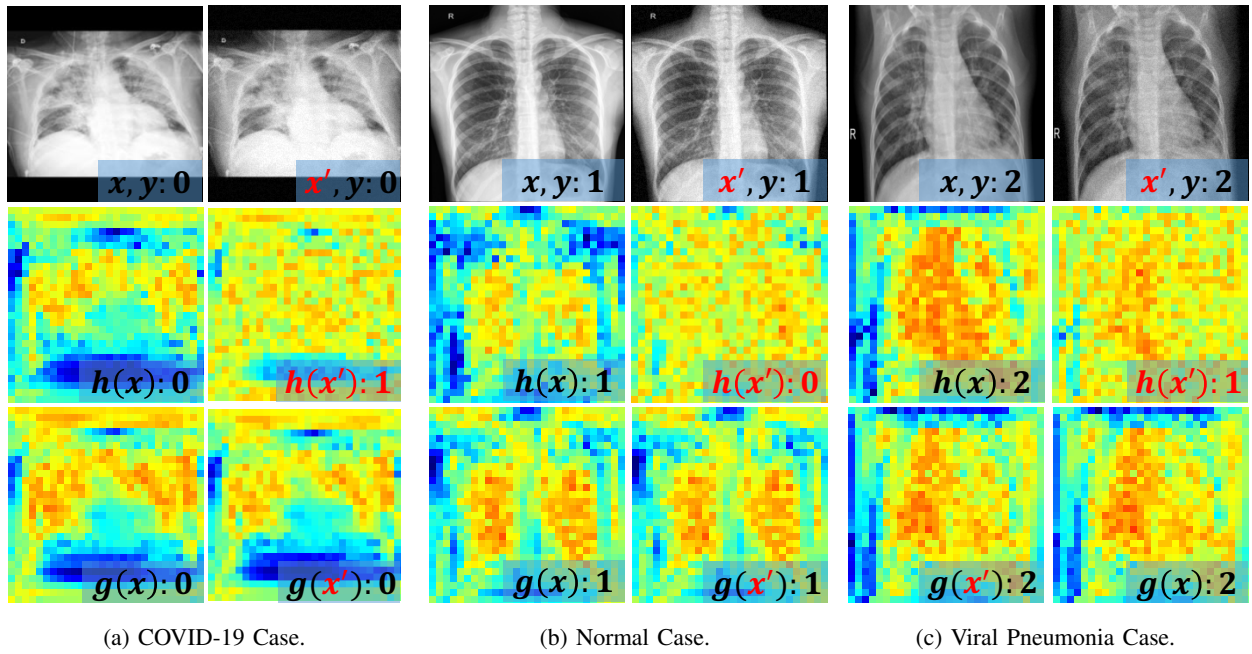


Fig. 4: Feature maps of base model  $h$  and MedRDF  $g$  on original images and adversarial images. The first row of each subfigure contains original image and its corresponding adversarial image. The second row of each subfigure shows the feature maps of original image and adversarial image on base model  $h$ , respectively. The third row of each subfigure shows the feature maps of original image and adversarial image on MedRDF  $g$ , respectively. The adversarial image is attacked by C&W attack. The feature is at the second “BasicBlock” layer of ResNet-18. MedRDF  $g$  is based on ResNet-18 with salt-and-pepper noise and median filter.

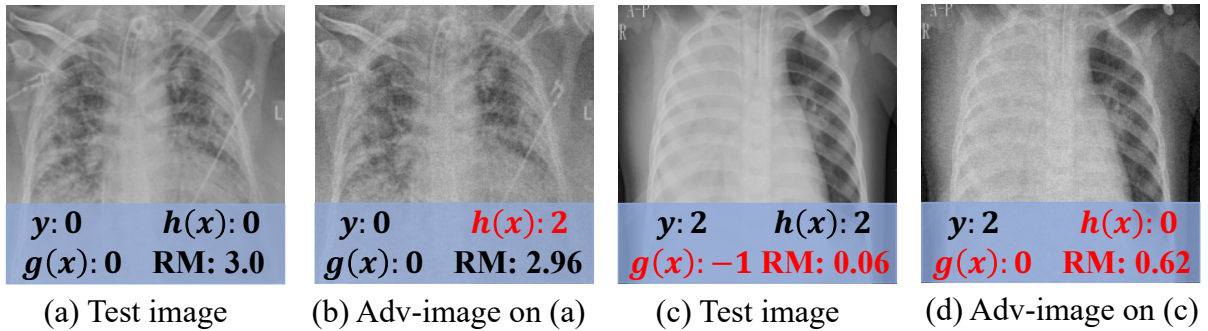


Fig. 5: Robust Metric (RM) of several cases. (a) and (c) are selected from COVID-19 test set. (b) and (d) are generated by PGD with 100 steps on (a) and (c), respectively. The base model  $h$  is ResNet-50, the MedRDF  $g$  is based on ResNet-50 with gaussian noise and median filter. Red represents the wrong label. Low RM indicates that (c) and (d) should be re-evaluated by professional doctor.

the same improvement can also be found on AG-Sononet-16 in TABLE VIII. For example, on DermaMNIST, MedRDF with poisson noise and median filter achieves 72.2% accuracy when attacked by RayS, while the accuracy of AG-Sononet-16 is 0.0%.

2) *Visualization Results*: In this part, we illustrate the superiority of our proposed MedRDF by visualizing the changes in the internal feature of each model.

**The change of attention maps.** As shown in Fig. 3, in each subfigure, the first column shows the original image and its corresponding adversarial image attacked by PGD-100 with their labels, respectively. The second column shows the

attention maps and output labels of base model AG-Sononet-16 on original image and adversarial image, respectively. And the third column denotes the attention maps and output labels of MedRDF  $g$  on original image and adversarial image, respectively. From Fig. 3 we can observe, the attention maps of base model on original image (i.e.,  $h(x)$ ) and adversarial image (i.e.,  $h(x')$ ) are extremely different. From this we can infer that, due to the changes of feature the base model focuses on, the base model can be easily fooled by adversarial example. On the contrary, we notice that MedRDF do not significantly change the attention map of the original image and the adversarial example, which shows that MedRDF are

**TABLE VIII:** Accuracy (%) of different defense mechanism (rows) against **black box adversarial attacks** on COVID-19 and DermaMNIST dataset with **AG-Sononet-16**. GS: gaussian smoothing, MF: median filter, s.p. noise: salt-and-pepper noise.

Method	Denoiser	COVID-19		DermaMNIST	
		SPSA	RayS	SPSA	RayS
AG-Sononet-16	None	9.3	0.0	1.8	0.0
	GS	1.3	0.1	13.9	3.9
	MF	31.3	13.5	20.1	5.0
MedRDF with gaussian noise	None	32.5	32.2	64.2	69.4
	GS	80.6	80.6	66.6	66.7
	MF	82.2	82.0	70.8	71.0
MedRDF with s.p. noise	None	68.6	68.7	62.1	70.8
	GS	84.0	84.0	66.7	66.8
	MF	81.9	82.0	71.4	71.9
MedRDF with poisson noise	None	55.3	57.3	63.9	70.5
	GS	<b>84.6</b>	<b>84.0</b>	67.2	67.9
	MF	82.6	82.0	<b>71.5</b>	<b>72.2</b>

**TABLE IX:** The accuracy (%) produced by base classifiers (ResNet-50 and AG-Sononet-16) and MedRDF on original COVID-19 test set and adversarial test sets. Row “Original” describes the original test set and rows “PGD-100” and “C&W” represent the test sets attacked by PGD and C&W. The abbreviation “C.” and “R.” represent Correctly classification and Robust evaluation ( $RM \geq 1$ ), respectively.

Network	Dataset	Natural Acc $h_{\theta}$	MedRDF $g$			
			C.&R.	C.&Not-R.	Not-C.&Not-R.	Not-C.&R.
ResNet-50	Original	92.6	90.8	0.6	7.2	1.4
	PGD-100	0.0	90.2	1.2	7.8	0.8
	C&W	0.0	90.4	1.0	7.8	0.8
AG-Sononet-16	Original	93.4	84.6	2.4	11.0	2.0
	PGD-100	0.0	87.4	1.8	9.2	1.6
	C&W	0.0	87.0	2.2	9.2	1.6

more robust than the base model. It can effectively improve the robustness of the original model.

**The change of feature maps.** We have also shown the feature maps produced by base model ResNet-18 and MedRDF based on ResNet-18 in Fig. 4. In each subfigure, the first row shows the original image and adversarial image attacked by C&W. The second row shows the feature maps on original image and adversarial image at the second “BasicBlock” layer of the base model ResNet-18, respectively. And the feature maps at the third row are produced by MedRDF. From the second rows at Fig. 4 (a)-(c), one can observe that the learned features of base model  $h$  for the clean image focus on semantically informative regions (represented in red), while the features of the adversarial images are activated globally (without any specific focus). However, this problem can be effectively solved by MedRDF. From the third row of each subfigure, we can see that the feature map of the adversarial image generated by MedRDF is consistent with the clean image. These visualization results indicate that our MedRDF is not susceptible to adversarial perturbations, thus improving the robustness effectively.

#### D. Performance of Robust Metric.

1) *Quantitative Results:* We report the accuracy produced by base models and MedRDF on original test set and adversarial test set in TABLE IX, where the abbreviation “C.” and “R.” represent Correctly classification and Robust evaluation ( $RM \geq 1$ ), respectively. From TABLE IX one can observe that, MedRDF can obtain robust and reliable accuracy both on original test dataset and adversarial test dataset (e.g., 90.8% C.&R. accuracy on original dataset and 90.4% C.&R. accuracy on adversarial dataset attacked by C&W based on ResNet-50), even if the accuracy of original model on adversarial dataset drops to 0.0%. In addition, we can find that for examples that were misclassified by MedRDF (i.e., Not-C in TABLE IX), most examples’ RM are below the threshold (i.e., Not-C & Not-R in TABLE IX), which can effectively instruct the doctor to re-diagnose this case. These results confirm the necessity and effectiveness of the RM indicator for medical diagnostic tasks.

2) *Visualization Results:* In order to illustrate the effectiveness of RM more intuitively, several cases are presented in Fig. 5, where the last two cases should be re-evaluated with doctor due to the low RM of the result.

#### E. Comparison with Other Methods

1) *Comparison with Other Augmentation Strategies:* For each test image, MedRDF first creates a large number of noisy copies. To illustrate the effectiveness of this operator, we compare our operator of creating noisy copies with other augmentation strategies. Specifically, we use random resizing and random rotating to replace the noise in this experiment. The resizing range is [200, 224], the rotating angle is [10, 100]. The experimental results can be found in TABLE X. From TABLE X we can obtain, compared with random rotating and resizing, our proposed MedRDF with noisy copies achieves best accuracy under all attacks.

2) *Comparison with Other Defense Methods.:* To further verify the superior performance of our method, we compare MedRDF with other defense mechanisms in this section, including pre-processing based-defenses (i.e., Random R-P [12], ComDefend [13]), and retraining-based defenses (i.e., adversarial training (AT) [5], TRADES [38], and MART [39]). The accuracy and training time of each method can be found in TABLE XI. From TABLE XI we can obtain, when the dataset is COVID-19, whether the base model is ResNet-18 or ResNet-50, MedRDF not only has the highest defense accuracy (e.g. the accuracy of MedRDF based ResNet-50 attacked by C&W is 91.2% while the Random R-P is 51.0%), but its training time is much shorter than other retraining defense methods (e.g., the training time of MedRDF based ResNet-18 is 0.51 hrs while the TRADES is 1.97 hrs). For DermaMNIST, MedRDF still maintains the best defense accuracy under many attacks (e.g., PGD-20, PGD-100, and C&W on ResNet-50). Compared with Random R-P, the pre-processing defense method, MedRDF has better defense accuracy on medical images. Besides, compared with retraining methods, MedRDF which is employed in the inference phase can greatly reduce the training time and training burden. The above

**TABLE X:** Accuracy (%) of different augmentation strategies (rows) against white box adversarial attacks with maximum  $L_\infty$  perturbation  $\epsilon = 8/255$  (columns) on COVID-19 with ResNet-18. The original accuracy of each defense is described in the column “Natural”. The denoiser is MF after adding s.p. noise. The number after attack method represents the number of iteration steps.

Method	Natural	I-FGSM-1	I-FGSM-7	PGD-20	PGD-100	C&W
MedRDF with random rotating	49.6	28.6	29.2	28.8	28.7	29.0
MedRDF with random resizing	90.8	51.4	57.0	62.2	55.6	60.0
MedRDF with s.p. noise (ours)	<b>91.2</b>	<b>84.6</b>	<b>86.2</b>	<b>90.2</b>	<b>89.6</b>	<b>90.4</b>

**TABLE XI:** Accuracy (%) and training time (hrs) of MedRDF compared with other defense methods under different adversarial attacks. The maximum  $L_\infty$  adversarial perturbation is  $\epsilon = 8/255$ , the numbers after the attack methods represent the number of iterative steps. MedRDF is based on salt-and-pepper noise and median filter denoiser.

Dataset	Network	Method	Natural	I-FGSM-7	PGD-20	PGD-100	C&W	Training time (hrs)
COVID-19	ResNet-50	Random R-P [12]	53.0	49.8	51.2	50.6	51.0	0.66
		ComDefend [13]	82.3	58.7	55.6	52.1	51.9	2.53
		AT [5]	90.8	65.6	65.0	61.4	62.0	4.17
		TRADES [38]	87.8	72.8	72.4	71.6	71.4	5.57
		MART [39]	89.2	75.0	74.8	74.0	74.4	4.57
		<b>MedRDF</b>	<b>91.6</b>	<b>91.2</b>	<b>91.4</b>	<b>91.2</b>	<b>91.2</b>	<b>0.66</b>
	ResNet-18	Random R-P [12]	66.2	56.2	56.2	54.8	56.4	0.51
		ComDefend [13]	88.3	64.3	64.8	63.2	63.1	1.30
		AT [5]	95.2	67.6	67.8	65.0	67.0	1.50
		TRADES [38]	87.6	72.0	72.0	71.4	71.6	1.97
		MART [39]	94.0	75.4	75.2	74.6	75.2	1.58
		<b>MedRDF</b>	<b>95.6</b>	<b>86.2</b>	<b>90.2</b>	<b>89.6</b>	<b>90.4</b>	<b>0.51</b>
DermaMNIST	ResNet-50	Random R-P [12]	63.0	54.1	56.7	55.4	55.8	0.44
		ComDefend [13]	70.3	65.5	64.3	62.4	61.5	5.01
		AT [5]	<b>72.5</b>	57.9	58.0	56.9	57.4	6.03
		TRADES [38]	66.9	<b>66.8</b>	66.7	66.8	66.7	7.47
		MART [39]	70.6	64.0	64.1	63.5	62.9	5.14
		<b>MedRDF</b>	72.0	62.7	<b>68.1</b>	<b>67.2</b>	<b>70.0</b>	<b>0.44</b>
	ResNet-18	Random R-P [12]	50.3	40.4	41.4	40.1	41.9	0.44
		ComDefend [13]	68.3	62.4	61.8	55.2	55.1	1.50
		AT [5]	<b>71.7</b>	56.5	56.7	55.1	55.6	1.54
		TRADES [38]	68.7	62.8	64.5	<b>64.4</b>	63.9	2.03
		MART [39]	70.4	59.8	59.8	59.3	58.3	1.55
		<b>MedRDF</b>	69.0	<b>63.1</b>	<b>65.1</b>	64.1	<b>66.0</b>	<b>0.44</b>

analyses confirm that our MedRDF is effective and suitable for defending against adversarial attack on medical diagnostic tasks.

## VI. CONCLUSION

We propose a Robust and Retrain-Less Diagnostic Framework for Medical pre-trained models against adversarial attack (i.e., MedRDF). MedRDF allows users to seamlessly convert the pre-trained non-robust medical diagnostic model into robust one in inference phase, which is very convenient for diagnostic services that are already deployed online. Moreover, we also propose an effective Robustness Metric (RM) based on MedRDF, which gives the confidence score of the diagnostic result. Experimental results demonstrate a superior performance of MedRDF on COVID-19 and dermaMNIST datasets in both white-box and black-box adversarial settings. In the future, we plan to study the robustness of base medical models to common noise which plays an important role in our robust framework, as well as the trade-off between the natural accuracy and defense accuracy. In addition, we will extend our research to the field of medical image segmentation.

## REFERENCES

- [1] P. Rajpurkar, J. Irvin, K. Zhu, B. Yang, H. Mehta, T. Duan, D. Ding, A. Bagul, C. Langlotz, K. Shpanskaya *et al.*, “Chexnet: Radiologist-level pneumonia detection on chest x-rays with deep learning,” *arXiv preprint arXiv:1711.05225*, 2017.
- [2] A. Esteva, B. Kuprel, R. A. Novoa, J. Ko, S. M. Swetter, H. M. Blau, and S. Thrun, “Dermatologist-level classification of skin cancer with deep neural networks,” *Nature*, vol. 542, no. 7639, p. 115, 2017.
- [3] V. Gulshan, L. Peng, M. Coram, M. C. Stumpe, D. Wu, A. Narayanaswamy, S. Venugopalan, K. Widner, T. Madams, J. Cuadros *et al.*, “Development and validation of a deep learning algorithm for detection of diabetic retinopathy in retinal fundus photographs,” *JAMA*, vol. 316, no. 22, pp. 2402–2410, 2016.
- [4] C. Szegedy, W. Zaremba, I. Sutskever, J. Bruna, D. Erhan, I. Goodfellow, and R. Fergus, “Intriguing properties of neural networks,” in *International Conference on Learning Representations*, 2014.
- [5] A. Madry, A. Makelov, L. Schmidt, D. Tsipras, and A. Vladu, “Towards deep learning models resistant to adversarial attacks,” in *International Conference on Learning Representations*, 2018.
- [6] N. Carlini and D. Wagner, “Towards evaluating the robustness of neural networks,” in *2017 IEEE Symposium on Security and Privacy (sp)*, 2017, pp. 39–57.
- [7] M. Paschali, S. Conjeti, F. Navarro, and N. Navab, “Generalizability vs. robustness: investigating medical imaging networks using adversarial examples,” in *International Conference on Medical Image Computing and Computer-Assisted Intervention*. Springer, 2018, pp. 493–501.
- [8] M. Xu, T. Zhang, Z. Li, M. Liu, and D. Zhang, “Towards evaluating the

- robustness of deep diagnostic models by adversarial attack,” *Medical Image Analysis*, p. 101977, 2021.
- [9] U. Ozbulak, A. Van Messe, and W. De Neve, “Impact of adversarial examples on deep learning models for biomedical image segmentation,” in *International Conference on Medical Image Computing and Computer-Assisted Intervention*. Springer, 2019, pp. 300–308.
  - [10] X. Li, D. Pan, and D. Zhu, “Defending against adversarial attacks on medical imaging ai system, classification or detection?” *arXiv preprint arXiv:2006.13555*, 2020.
  - [11] F.-F. Xue, J. Peng, R. Wang, Q. Zhang, and W.-S. Zheng, “Improving robustness of medical image diagnosis with denoising convolutional neural networks,” in *International Conference on Medical Image Computing and Computer-Assisted Intervention*. Springer, 2019, pp. 846–854.
  - [12] C. Xie, J. Wang, Z. Zhang, Z. Ren, and A. Yuille, “Mitigating adversarial effects through randomization,” in *International Conference on Learning Representations*, 2018.
  - [13] X. Jia, X. Wei, X. Cao, and H. Foroosh, “Comdefend: An efficient image compression model to defend adversarial examples,” in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2019, pp. 6084–6092.
  - [14] M. E. Chowdhury, T. Rahman, A. Khandakar, R. Mazhar, M. A. Kadir, Z. B. Mahbub, K. R. Islam, M. S. Khan, A. Iqbal, N. Al Emadi et al., “Can ai help in screening viral and covid-19 pneumonia?” *IEEE Access*, vol. 8, pp. 132665–132676, 2020.
  - [15] R. Jain, M. Gupta, S. Taneja, and D. J. Hemanth, “Deep learning based detection and analysis of covid-19 on chest x-ray images,” *Applied Intelligence*, vol. 51, no. 3, pp. 1690–1700, 2021.
  - [16] I. D. Apostolopoulos and T. A. Mpesiana, “Covid-19: automatic detection from x-ray images utilizing transfer learning with convolutional neural networks,” *Physical and Engineering Sciences in Medicine*, vol. 43, no. 2, pp. 635–640, 2020.
  - [17] T. Ozturk, M. Talo, E. A. Yildirim, U. B. Baloglu, O. Yildirim, and U. R. Acharya, “Automated detection of covid-19 cases using deep neural networks with x-ray images,” *Computers in biology and medicine*, vol. 121, p. 103792, 2020.
  - [18] F. Ucar and D. Korkmaz, “Covidagnosis-net: Deep bayes-squeezeenet based diagnosis of the coronavirus disease 2019 (covid-19) from x-ray images,” *Medical hypotheses*, vol. 140, p. 109761, 2020.
  - [19] M. Jamshidi, A. Lalbakhsh, J. Talla, Z. Peroutka, F. Hadjilooei, P. Lalbakhsh, M. Jamshidi, L. La Spada, M. Mirmozafari, M. Dehghani et al., “Artificial intelligence and covid-19: deep learning approaches for diagnosis and treatment,” *Ieee Access*, vol. 8, pp. 109581–109595, 2020.
  - [20] A. Creswell, T. White, V. Dumoulin, K. Arulkumaran, B. Sengupta, and A. A. Bharath, “Generative adversarial networks: An overview,” *IEEE Signal Processing Magazine*, vol. 35, no. 1, pp. 53–65, 2018.
  - [21] G.-B. Huang, Q.-Y. Zhu, and C.-K. Siew, “Extreme learning machine: theory and applications,” *Neurocomputing*, vol. 70, no. 1-3, pp. 489–501, 2006.
  - [22] S. Hochreiter and J. Schmidhuber, “Long short-term memory,” *Neural computation*, vol. 9, no. 8, pp. 1735–1780, 1997.
  - [23] C. Szegedy, V. Vanhoucke, S. Ioffe, J. Shlens, and Z. Wojna, “Rethinking the inception architecture for computer vision,” in *Proceedings of the IEEE conference on computer vision and pattern recognition*, 2016, pp. 2818–2826.
  - [24] F. Chollet, “Xception: Deep learning with depthwise separable convolutions,” in *Proceedings of the IEEE conference on computer vision and pattern recognition*, 2017, pp. 1251–1258.
  - [25] S. Xie, R. Girshick, P. Dollár, Z. Tu, and K. He, “Aggregated residual transformations for deep neural networks,” in *Proceedings of the IEEE conference on computer vision and pattern recognition*, 2017, pp. 1492–1500.
  - [26] J. Schlemper, O. Oktay, L. Chen, J. Matthew, C. Knight, B. Kainz, B. Glocker, and D. Rueckert, “Attention-gated networks for improving ultrasound scan plane detection,” in *Medical Imaging with Deep Learning*, 2018.
  - [27] R. L. Siegel, K. D. Miller, and A. Jemal, “Cancer statistics, 2019,” *CA: a cancer journal for clinicians*, vol. 69, no. 1, pp. 7–34, 2019.
  - [28] Y. Li and L. Shen, “Skin lesion analysis towards melanoma detection using deep learning network,” *Sensors*, vol. 18, no. 2, p. 556, 2018.
  - [29] Y. Yan, J. Kawahara, and G. Hamarneh, “Melanoma recognition via visual attention,” in *International Conference on Information Processing in Medical Imaging*. Springer, 2019, pp. 793–804.
  - [30] I. J. Goodfellow, J. Shlens, and C. Szegedy, “Explaining and harnessing adversarial examples,” *arXiv preprint arXiv:1412.6572*, 2014.
  - [31] J. Uesato, B. O’donoghue, P. Kohli, and A. Oord, “Adversarial risk and the dangers of evaluating against weak attacks,” in *International Conference on Machine Learning*. PMLR, 2018, pp. 5025–5034.
  - [32] J. Chen and Q. Gu, “Rays: A ray searching method for hard-label adversarial attack,” in *Proceedings of the 26th ACM SIGKDD International Conference on Knowledge Discovery & Data Mining*, 2020, pp. 1739–1747.
  - [33] X. Ma, Y. Niu, L. Gu, Y. Wang, Y. Zhao, J. Bailey, and F. Lu, “Understanding adversarial attacks on deep learning based medical image analysis systems,” *Pattern Recognition*, vol. 110, p. 107332, 2021.
  - [34] S. A. Taghanaki, A. Das, and G. Hamarneh, “Vulnerability analysis of chest x-ray image classification against adversarial attacks,” in *Understanding and interpreting machine learning in medical image computing applications*. Springer, 2018, pp. 87–94.
  - [35] N. Papernot, P. McDaniel, X. Wu, S. Jha, and A. Swami, “Distillation as a defense to adversarial perturbations against deep neural networks,” in *2016 IEEE symposium on security and privacy (SP)*. IEEE, 2016, pp. 582–597.
  - [36] J. Buckman, A. Roy, C. Raffel, and I. Goodfellow, “Thermometer encoding: One hot way to resist adversarial examples,” in *International Conference on Learning Representations*, 2018.
  - [37] A. Shafahi, M. Najibi, A. Ghiasi, Z. Xu, J. Dickerson, C. Studer, L. S. Davis, G. Taylor, and T. Goldstein, “Adversarial training for free!” in *Proceedings of the 33rd International Conference on Neural Information Processing Systems*, 2019, pp. 3358–3369.
  - [38] 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*, 2019, pp. 7472–7482.
  - [39] Y. Wang, D. Zou, J. Yi, J. Bailey, X. Ma, and Q. Gu, “Improving adversarial robustness requires revisiting misclassified examples,” in *International Conference on Learning Representations*, 2019.
  - [40] S. Liu, A. A. A. Setio, F. C. Ghesu, E. Gibson, S. Grbic, B. Georgescu, and D. Comaniciu, “No surprises: Training robust lung nodule detection for low-dose ct scans by augmenting with adversarial attacks,” *IEEE Transactions on Medical Imaging*, vol. 40, no. 1, pp. 335–345, 2020.
  - [41] J. Yang, R. Shi, and B. Ni, “Medmnist classification decathlon: A lightweight automl benchmark for medical image analysis,” in *2021 IEEE 18th International Symposium on Biomedical Imaging (ISBI)*. IEEE, 2021, pp. 191–195.
  - [42] P. Tschandl, C. Rosendahl, and H. Kittler, “The ham10000 dataset, a large collection of multi-source dermatoscopic images of common pigmented skin lesions,” *Scientific data*, vol. 5, no. 1, pp. 1–9, 2018.
  - [43] N. Codella, V. Rotemberg, P. Tschandl, M. E. Celebi, S. Dusza, D. Gutman, B. Helba, A. Kalloo, K. Liopyris, M. Marchetti et al., “Skin lesion analysis toward melanoma detection 2018: A challenge hosted by the international skin imaging collaboration (isic),” *arXiv preprint arXiv:1902.03368*, 2019.
  - [44] K. He, X. Zhang, S. Ren, and J. Sun, “Deep residual learning for image recognition,” in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2016, pp. 770–778.
  - [45] J. Cohen, E. Rosenfeld, and Z. Kolter, “Certified adversarial robustness via randomized smoothing,” in *International Conference on Machine Learning*, 2019, pp. 1310–1320.
  - [46] X. Cao and N. Z. Gong, “Mitigating evasion attacks to deep neural networks via region-based classification,” in *Proceedings of the 33rd Annual Computer Security Applications Conference*, 2017, pp. 278–287.
  - [47] D. Kalimeris, G. Kaplun, P. Nakkiran, B. L. Edelman, T. Yang, B. Barak, and H. Zhang, “{SGD} on neural networks learns functions of increasing complexity,” in *Advances in Neural Information Processing Systems 32: Annual Conference on Neural Information Processing Systems 2019*, 2019.
  - [48] S.-M. Moosavi-Dezfooli, A. Fawzi, J. Uesato, and P. Frossard, “Robustness via curvature regularization, and vice versa,” in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2019, pp. 9078–9086.
  - [49] H. Salman, M. Sun, G. Yang, A. Kapoor, and J. Z. Kolter, “Denoised smoothing: A provable defense for pretrained classifiers,” *arXiv preprint arXiv:2003.01908*, 2020.
  - [50] K. Hung and W. Fithian, “Rank verification for exponential families,” *The Annals of Statistics*, vol. 47, no. 2, pp. 758–782, 2019.
  - [51] A. Kurakin, I. Goodfellow, and S. Bengio, “Adversarial machine learning at scale,” in *International Conference on Learning Representations*, 2017.