

Mixture GAN For Modulation Classification Resiliency Against Adversarial Attacks

Eyad Shtaiwi*, Ahmed El Ouadrhiri†, Majid Moradikia †, Salma Sultana †, Ahmed Abdelhadi†, and Zhu Han*,

*Electrical and Computer Engineering Department, University of Houston, Houston, TX, USA.

†Department of Engineering Technology, University of Houston, Houston, TX, USA.

Abstract—Automatic modulation classification (AMC) using the Deep Neural Network (DNN) approach outperforms the traditional classification techniques, even in the presence of challenging wireless channel environments. However, the adversarial attacks cause the loss of accuracy for the DNN-based AMC by injecting a well-designed perturbation to the wireless channels. In this paper, we propose a novel generative adversarial network (GAN)-based countermeasure approach to safeguard the DNN-based AMC systems against adversarial attack examples. GAN-based aims to eliminate the adversarial attack examples before feeding to the DNN-based classifier. Specifically, we have shown the resiliency of our proposed defense GAN against the Fast-Gradient Sign method (FGSM) algorithm as one of the most potent kinds of attack algorithms to craft the perturbed signals. The existing defense-GAN has been designed for image classification and does not work in our case where the above-mentioned communication system is considered. Thus, our proposed countermeasure approach deploys GANs with a mixture of generators to overcome the mode collapsing problem in a typical GAN facing radio signal classification problem. Simulation results show the effectiveness of our proposed defense GAN so that it could enhance the accuracy of the DNN-based AMC under adversarial attacks to 81%, approximately.

Index Terms—Modulation Classifications, FGSM, CNN-based classifier, GAN Countermeasure.

I. INTRODUCTION

Automatic modulation classification (AMC) is an approach that is used to automatically recognize the modulation classes of the received signals without any prior knowledge. AMC plays an important role in many military and civilian communication systems [1]–[3]. Generally, the AMC approaches fall into one of two categories [4]. The first category is based on the maximum likelihood that relies on the deployed statistical model, and their performance depends on the accuracy of the considered system model which make them susceptible to the model mismatch [5], [6]. The second category deploys the feature-based learning approach that collects the features from the received data samples and uses a classifier at the receiver to determine the modulation classes [7]. Particularly, the authors have used the convolutional neural networks (CNNs) for the classical task of AMC.

Despite their great results, deep neural networks (DNNs) have been shown to be vulnerable against adversarial attacks [8], [9]. Adversarial examples cause misclassifications of DNN-based classifier by incorporating the original input with a small and carefully designed perturbations. The authors in [9] showed that the DNN-based AMC systems are susceptible to adversarial attacks. Moreover, targeted fast gradient method

(FGM)-based adversarial examples [10] has been generated to enforce the misclassification at the receiver to the target label. However, the aforementioned adversarial examples are affected by the channel effects. To mitigate this issue, in [11], the authors have proposed some approaches to make the adversarial attack robust against realistic channel effects.

Recently, given the power of Generative Adversarial Networks (GAN) [15], the idea of defense-GAN is proposed in [16] to safeguard the DNN against adversarial attacks. The authors have shown the effectiveness of defense-GAN against both black and white-box attacks. However, the typical DNN classifier considered in [16] aimed to classify the images in the presence of an adversary. Therefore, no communication system was considered, and thus their proposed defense does not work for the radio signal classification scenario, as we intended here. Moreover, training of the GAN for the radio signals is challenging as it is easily trapped into the mode collapsing issue where the generator only concentrates to generate samples lying on a few classes instead of the whole data space [17].

To deal with this concern, motivated by the idea of Mixture Generative Adversarial Nets (MGAN) proposed in [18], we develop a novel approach where our GAN is equipped with a mixture of generators. More explicitly we deploy multiple generators, instead of using a single generator as in the typical GAN. Our simulation results validate that our proposed defense-GAN could achieve an acceptable protection, i.e., the promising accuracy of approximately 81%, against one effective kind of attacks Fast Gradient Sign Method (FGSM).

The rest of this paper is organized as follows. We describe the system in Section II. Section III details the DNN classifier. Section IV presents the FGSM-based adversarial attack. Then, we introduce the proposed MGAN in Section V. Section VI presents the simulation results. Finally, Section VII concludes the paper.

II. SYSTEM MODEL

As shown in Fig. 1, we consider a wireless communication system that consists of three nodes including transmitter (Tx), receiver (Rx), and an adversary (Ad). Next, we describe the working principle of each node as follows.

The transmitted modulated signal over the wireless channel is an L -dimensional vector denoted by $\mathbf{x} = [x[0], x[1], \dots, x[L]]^T \in \mathcal{C}^L$. The wireless channel interacts

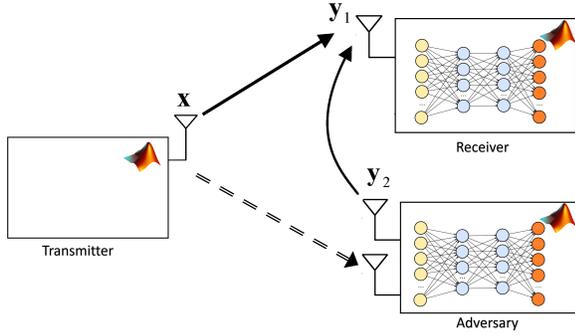


Fig. 1: System Model

with the modulated signal by introducing carrier frequency and phase shift offsets, sample rate offset (SRO), selective fading, and additive white Gaussian noise (AWGN). Let $h[k]$ be the channel's impulse response including the aforementioned radio imperfection at sample time k , and then the received signal is expressed as

$$r[k] = x[k] * h[k] + n[k], \quad (1)$$

where $n[k]$ is the complex zero mean AWGN which is modeled as $\mathcal{CN}(0, \sigma_n^2)$, where σ_n^2 is the noise power. It is worth mentioning, the received complex signal is represented using a 2-dimensional reals. Equivalently, $\mathbf{r} = [\mathbf{r}_I \ \mathbf{r}_Q] \in \mathbb{R}^{2 \times L}$, where \mathbf{r}_I and \mathbf{r}_Q represent the in-phase (I) and quadrature (Q) components of \mathbf{r} , respectively. In order to minimize the receiver sensitivity to the inter-symbol interference (ISI) and achieve a band-limited modulated signal, we low-pass filter the I/Q samples through a pulse shaping filter called a raised cosine filter. The time-domain impulse response of the shaping filter is given by

$$p[k] = \frac{\sin\left(\frac{\pi k}{\tau}\right)}{\left(\frac{\pi k}{\tau}\right)} \frac{\cos\left(\frac{\pi \rho k}{\tau}\right)}{1 - \left(\frac{2\pi \rho k}{\tau}\right)^2}, \quad (2)$$

where τ and ρ denote the pulse period and the roll-off factor which determines the modulated signal bandwidth, respectively.

III. CLASSIFIER ARCHITECTURE

Let \mathcal{X} denote the set of modulated signals where $\mathcal{X} \subset \mathbb{R}^{2 \times L}$. Each modulated signal, i.e., $\mathbf{x} \in \mathcal{X}$, belongs to of one of the modulation classes C . Therefore, the ML-based classifier which is denoted by $f(\cdot, \theta)$, where maps each frame of the modulated signals into C , where θ represents the model parameters. In other words, $f(\cdot, \theta_0) : \mathbb{R}^{L \times 2} \rightarrow \mathbb{R}^C$. After aggregation of the data in \mathcal{X} , the trained classifier assigns a label $\hat{C}(\mathbf{x}; \theta_0)$ for each input \mathbf{x} where $\hat{C}(\mathbf{x}; \theta_0)$ is given by

$$\hat{C}(\mathbf{x}; \theta_0) = \arg \max_n g_n(\mathbf{x}, \theta_0), \quad (3)$$

where $f_n(\mathbf{x}, \theta_0)$ denotes the classification probability that the signal \mathbf{x} belongs to the n^{th} class for $n = 1, 2, \dots, C$. In this paper, we consider $f(\cdot, \theta_0)$ for the CNN architecture.

The adopted CNN-classifier consists of multiple convolutional layers, where each layer is able to extract an underlying spatially correlated information from its input without need to complex a priori constraint.

Fig. 2 depicts the CNN-based classifier architecture that consists of six convolutional layers, one fully connected layer, and a softmax layer. Each convolutional layer has followed by a batch normalization layer, rectified linear unit (ReLU) activation layer, and a max-pooling layer to form CNN blocks, i.e., CNN1, \dots , CNN6. In the last convolutional layer, the max-pooling layer is replaced with an average pooling layer. The output layer has softmax activation, as well. The j^{th} convolutional layer produces output, i.e., feature map, is given by

$$L_j^{(l)}(m, n) = \sigma(O_j^{(l)}(m, n)), \quad (4)$$

where $\sigma(\cdot)$ and $O_j^{(l)}(m, n)$ represent the activation function, and the weighted sum output of the previous convolutional layer which is computed by

$$O_j^{(l)}(m, n) = \sum_{i=1}^{M^{(l)}} \sum_{u,v=0}^{S-1} W_{ji}^{(l)}(u, v) L_i^{(l-1)}(m-u, n-v) + b_j^{(l)}, \quad (5)$$

where M and S represent the size of the kernel filter, while b and W are the bias and the weights parameters that are optimized using the back-propagation (BP) method. Therefore, the updated weights and bias formulas of the l^{th} convolutional layer are

$$W_{ji}^{(l)} = W_{ji}^{(l)} + \alpha \cdot \frac{\partial L}{\partial W_{ji}^{(l)}}, \quad (6)$$

$$b_j^{(l)} = b_j^{(l)} + \alpha \cdot \frac{\partial L}{\partial b_j^{(l)}}. \quad (7)$$

Consequently, the derivatives of (6) and (7) are computed using the chain rule as

$$\frac{\partial L}{\partial W_{ji}^{(l)}} = \sum_{m,n} \delta_j^{(l)}(m, n) \cdot L_j^{(L-1)}(m-u, y-v), \quad (8)$$

$$\frac{\partial L}{\partial b_j^{(l)}} = \sum_{m,n} \delta_j^{(l)}(m, n), \quad (9)$$

where $\delta_j^{(l)}$ represents the error map and is given as

$$\delta_j^{(l)} = \sum_j \sum_{u,v=0}^{S-1} \mathcal{W}_{ji}^{(l+1)}(u, v) \cdot \delta_j^{(l+1)}(m+u, n+v). \quad (10)$$

After training the model with RML2016.10a dataset [19], we tested the trained model for 4 modulation types, namely, BPSK, QPSK, 8PSK, and 16QAM. The resultant confusion matrix is shown in Fig. 3, has been obtained using 10^4 frames per modulation.

It can be seen that the accuracy approximately reached 99%, and the CNN classifier is barely confused between 8PSK and QPSK.

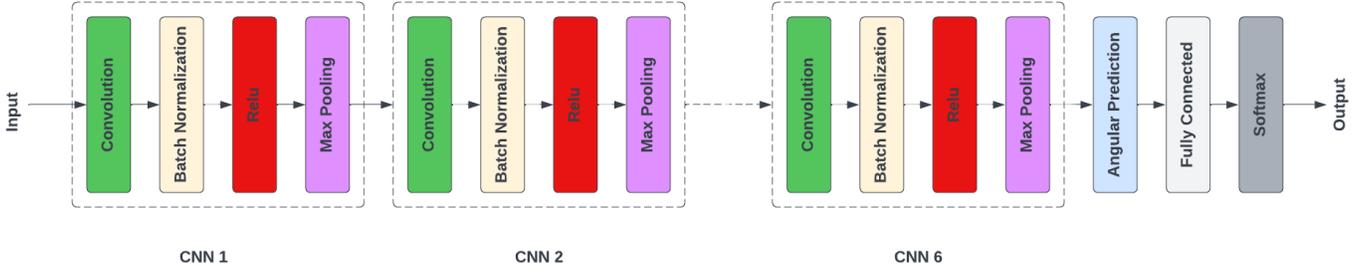


Fig. 2: The CNN architecture.

True Class \ Predicted Class	16QAM	8PSK	BPSK	QPSK
16QAM	999	1		
8PSK		979		21
BPSK			1000	
QPSK		25		975

99.9%	0.1%
97.9%	2.1%
100.0%	
97.5%	2.5%

Fig. 3: Confusion Matrix for the DNN-based AMC.

IV. ADVERSARIAL ATTACK FOR DNN-BASED CLASSIFIER MODEL

The DNN-based AMC systems are usually vulnerable to various security attacks. For instance, different adversarial attacks cause either a misclassification or loss of the accuracy of the DNN-based modulation classifiers. Generally, the attacker in different wireless applications aims to misclassify the CNN-based classifier by transmitting a well-designed perturbation over the wireless channel [9].

The current adversarial attacks on the DNN-based AMC systems adopt the white-box attack. This type of attack setting assumes that adversarial node fully accesses to the the internal information of DNN platforms, such as trained model of $f(\cdot; \theta)$ and its parameters and architecture. However, this assumption is almost impossible in real case scenarios as the model owners do not disclose any information about their models. Hence, we develop an black-box adversarial attack against the adopted DNN-classifier.

The adversarial node affects the the classification decision process for the legitimate receiver DNN-classifier, i.e., $f(\cdot; \theta)$, by injecting adversarial interference into transmitted signals. The perturbed signals can be modeled as [11]

$$\mathbf{x}_b = \mathbf{x} + \boldsymbol{\varrho}, \quad (11)$$

where $\boldsymbol{\varrho} \in \mathbb{R}^{L \times 2}$ denotes the additive adversarial interference

which is computed as follow

$$\underset{\boldsymbol{\varrho}}{\operatorname{argmin}} \|\boldsymbol{\varrho}\|_{\infty} \quad (12)$$

$$\text{s.t.} \quad (13)$$

$$\hat{C}(\mathbf{x}; \theta_0) \neq \hat{C}(\mathbf{x} + \boldsymbol{\varrho}; \theta_0) \quad (14)$$

$$\mathbf{x}_b + \boldsymbol{\varrho} \in \mathcal{X}, \quad (15)$$

where $\|\cdot\|_{\infty}$ denotes the l_{∞} norm. It is noteworthy that solving (15) is difficult and not global optimal. Hence, the generated perturbation, i.e., $\boldsymbol{\varrho}$ is not unique and there are sub-optimal approaches to compute $\boldsymbol{\varrho}$. The most common approach among these sub-optimal methods is the FGSM [20]. Specifically, the FGSM crafted the perturbation signal as follows

$$\tilde{\mathbf{x}}_b = \mathbf{x} + \eta \cdot \operatorname{sgn}(\nabla_{\mathbf{x}} J(\mathbf{x}, y, \theta_0)), \quad (16)$$

where $\operatorname{sgn}(\cdot)$, $\nabla_{\mathbf{x}} J(\mathbf{x}, y, \theta_0)$ denote the sign operation, and the gradient of the model loss function which is function of the input sample, \mathbf{x} , the model parameters θ_0 , and the label vector $\mathbf{y} = \{0, 1\}^C$. The FSGM incorporates a small perturbed signal, η , into each feature of the input sample in the direction of the sign of the classifier's cost function, i.e., $J(\mathbf{x}, y, \theta_0)$.

It is worth mentioning that there are two variants of the FGSM attacks, namely, the targeted FGSM and untargated FGSM. In the first type, the ad node aims to generate the perturbation that causes the DNN-based classifier to have a predetermined misclassification, e.g., the DNN-based AMC classifies BPSK modulation as 8PSK modulation. While, the adversarial node does not select any specific misclassification when computing the perturbation signal. Based on the output modulation class of the CNN classifier and the attack type, the ad node apply the FGSM method to compute the needed perturbation to confuse the Rx node classifier. Since the ad node does not have full access to the DNN model, the overall accuracy of the substitute classifier is almost 92.5%.

V. GAN-BASED DEFENSE PROPOSED APPROACH

In this section, we introduce the countermeasure approach to mitigate the effect of the adversarial interference and improve the robustness of the legitimate DNN-based classifier. Specifically, we propose a GAN model as a defensive approach against the attacker. A GAN is a type of neural network framework for the generative modeling approach which uses an existing distribution of samples from a dataset to generate

new instances that follow the same distribution of the training dataset [21]. Moreover, A GAN is considered as a generative model that is trained using two neural network models: 1) The generative network model, also known as generator G , which aims to learn to generate new plausible samples similar to the training dataset. 2) The discriminative network model, also known as discriminator D , which is trained to differentiate between the generated samples by the generator and real samples of the training dataset. The generator and the discriminator are set up in a game during the training; the generator aims to fool the discriminator while the discriminator aims to distinguish between real and generated samples. The objective of the generator and the discriminator is to minimize the following min-max loss [21]

$$\min_G \max_D V(D, G) = \mathbb{E}_{d \sim p_{data}(d)} [\log D(d)] + \mathbb{E}_{h \sim p_h(h)} [\log(1 - D(G(h)))] \quad (17)$$

In this paper, we deploy the GAN model to generate plausible samples similar to the received frames. Then, we compare the generated frames with the perturbed received signal, i.e., x_b , to determine the true class of the modulated signal, x . Since we are dealing with multiple modulation types, we faced the mode collapse problem, where the generator fails to generate samples of all modulation types. For example, in our experiments, sometimes the generator only generates frame samples that belong to modulation types BPSK, QPSK, and 8PSK. To overcome this problem, we propose to use Mixture GAN (MGAN) [18] that allows to avoid the mode collapse problem and increases the defense GAN accuracy.

Therefore, we train a GAN model for each modulation type instead of training one GAN for multiple labels. Consequently, in our case, we will have four generative network models; a generator for each modulation type. The proposed approach using four generators does not add any complexity to the system considering the offline training. However, this approach increases the model accuracy and robustness against the ad node attacks. After obtaining the multiple generators, the actual modulation class of the received signal in the presence of adversarial attacks can be found as follows

$$h^* = \arg \min_i \arg \min_h \|G_i(h) - r\|_2^2 \quad (18)$$

Where r is the received signal after the ad node injecting the perturbation signal, and $i = 1, \dots, C$. Afterward, h^* is fed to the corresponding generator G_i to generate a signal close to the received perturbed signal. The optimization problem in (18) is a highly non-convex problem and has no optimal global solution. Therefore, we adopted the gradient descent (GD) method to find the sub-optimal solution [16]. The developed GD method is equivalent to computing a q gradients of the loss function using ζ different random initialization of h for each generator as

$$h_q^\zeta = h_{q-1}^\zeta + \eta_{q-1} \nabla_h \mathcal{L}(r, h)|_{h_q = h_{q-1}}, \forall q = 1, 2, \dots, \quad (19)$$

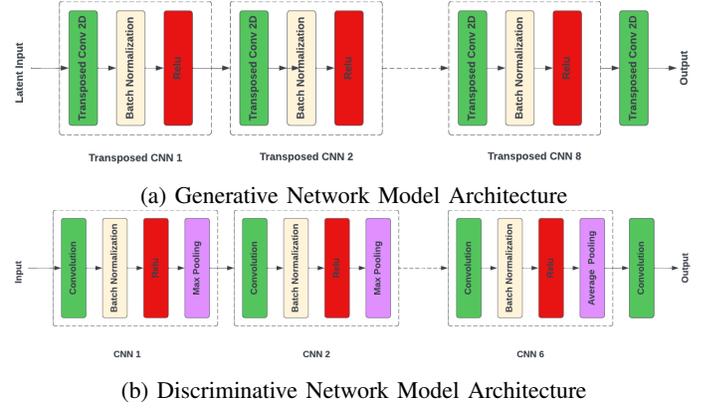


Fig. 4: GAN model architecture which consists of two neural network models: Generator and Discriminator

where $\mathcal{L}(r, h) = \|G(h) - r\|_2^2$.

Fig.4 illustrates the architecture of the GAN developed in this paper, which includes two network models:

- **Generative Network Model:** which takes as input a vector of random values, i.e., latent inputs, and generates frames similar to the training dataset. The developed generator consists of eight transposed convolution neural network layers (TCNN) as shown in Fig. 4a. The vector of latent inputs is shaped to the proper format before feeding it to the first TCNN. Each TCNN is followed by a batch normalization and a ReLU activation layer. The last layer is a TCNN that outputs the generated frame of dimension $x_{gen} \in \mathbb{R}^{L \times 2} \subset \mathcal{X}$.
- **Discriminative Network Model:** takes as inputs both the actual frames, i.e., x , and the generated frames, x_{gen} by $G_i \forall i \in 1, \dots, C$. The trained discriminators, i.e., $D_i \forall i \in 1, \dots, C$ differentiate among the real and the generated frames. Fig. 4b depicts the developed discriminator architecture which consists of six CNN layers, each CNN layer is followed by a batch normalization layer, a ReLU activation function, and a max-pooling layer, similar to the adopted CNN-classifier in Section III. The sixth CNN layer is followed by an average pooling layer instead of a max-pooling layer.

VI. SIMULATION RESULTS

In this section, we investigate the effectiveness and efficiency of the proposed approach against the adversarial attacks. Therefore, we study the impact of the presence of the adversarial node on the accuracy of the developed CNN-based classifier in Section III. Subsequently, we compare the accuracy of the classifier after incorporating the multiple generators with the CNN-classifier. Specifically, we provide the confusion matrices for the classifier for the two cases.

In our simulations, we used Matlab platform for training and validating each of the classifier, the attacker, and the mixture GAN on a workstation with Nvidia Titan RTX, 2 CPUs Intel

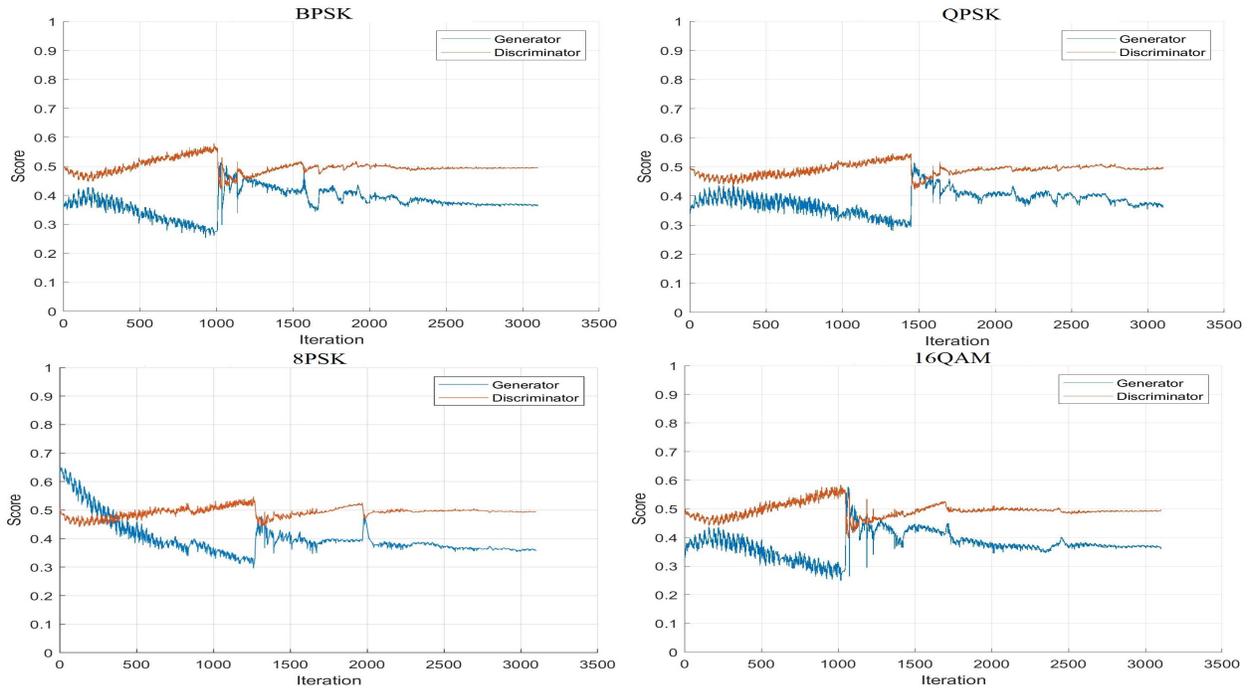


Fig. 5: The mixture GAN training for the four modulation types, i.e., BPSK, QPSK, 8PSK, 16QAM.

TABLE I: Simulation Parameters

Parameter	Value
Carrier Frequency	915 MHz
Roll-off factor	0.7
Number of frames	10^4
L	1024
C	4
Mini-batch size	12
Number of Epochs	100
Optimizer	ADAM
Learning Rate	0.02
gradient decay factor	0.5
SNR	30dB

Xeon Gold 6148, and 768GB RAM. The simulation setup parameters are summarized in Table I. For dataset generation, we consider a Rician multi-path fading channel and introduced channel impairments, such as, carrier frequency and phase offsets.

A. Training of MGAN

In this part, we present the training phase for the Mixture GAN with 4 generators. In Fig. 5, we study the convergence performance for both networks, i.e., G's and D's, of each class. It can be observed that the training process succeeded for each modulation class. Specifically, the discriminators reach the optimal score 0.5 after 1,000 iterations. On the other hands, the generators converge to 0.4 score. Which means that although the discriminators provide the optimal performance, the generators are able to produce modulated signals similar to the actual dataset. Fig. 5 depicts the benefits of deploying

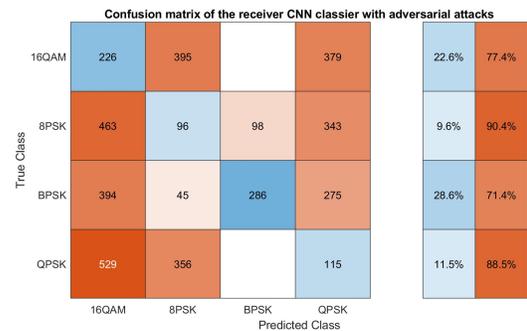


Fig. 6: Accuracy of the receiver CNN classifier in the presence of adversarial attacks.

multiple generators to overcome the collapsing mode issue of the traditional GAN.

B. Adversarial attack

Fig. 6 shows the confusion matrix for the DNN-classifier in the presence of an adversarial attack. In this paper, we only considered the untargeted FGSM variant. Compared with Fig. 3, it can be noticed that the designed perturbation signals using the FGSM approach reduce the accuracy of the adopted CNN-classifier. It can be shown that the attacker achieves its best performance when the modulated signals are 8 PSK. Generally, the accuracy performance of the classifier is significantly dropped which proves the effectiveness of the designed attack model.

True Class	Predicted Class				Accuracy	
	16QAM	8PSK	BPSK	QPSK		
16QAM	869	58		73	86.9%	13.1%
8PSK	24	694	188	94	69.4%	30.6%
BPSK	96		812	92	81.2%	18.8%
QPSK	134			866	86.6%	13.4%

Fig. 7: Accuracy of the GAN-based countermeasure.

C. GAN-based countermeasure

Fig. 7 shows the confusion matrix of the classifier in the presence of attack operations after incorporating of MGAN. The average accuracy for the classifier with the proposed defense approach reaches 81%. It can be shown that mixture GAN approach improves the DNN-based classifier accuracy. Specifically, the accuracy for 8 PSK case has been improved from 9% to approximately 70%. It is worth mentioning, that the error of this confusion matrix includes the error of the GAN model and the CNN-classifier.

VII. CONCLUSION

In this paper, we presented a countermeasure approach against adversarial attacks in AMC systems using mixture GAN. The proposed approach reduces the effect of the adversarial perturbed signals before feeding to the DNN-classifier. To prove the effectiveness of the proposed approach, we alleviated the FGSM approach to craft the perturbation. We utilized the black-box model for the attack node. The adversarial attacks reduce the DNN-classifier accuracy by injecting pre-designed perturbation signals. In addition, we addressed the collapsing mode issue of the traditional GAN by developing multiple generators. The idea of deploying multiple generators helps to capture different modes that exist in the dataset distribution. Finally, through the simulations, we demonstrated the enhancement of accuracy of the CNN-based classifier after incorporating the MGAN method.

REFERENCES

- [1] Popoola, J. & Olst, R. Automatic classification of combined analog and digital modulation schemes using feedforward neural network. *IEEE Africon'11*, Victoria Falls, Zambia, Sep. 2011.
- [2] Zhang, G. Neural networks for classification: a survey. *IEEE Transactions On Systems, Man, And Cybernetics, Part C (Applications And Reviews)*, vol.30, no. 4, 451-462, Nov. 2000.
- [3] Wu, H., Saquib, M. & Yun, Z. Novel automatic modulation classification using cumulant features for communications via multipath channels. *IEEE Transactions On Wireless Communications*, vol.7, no. 48, 3098-3105, Aug. 2008.
- [4] Zhu, Z. & Nandi, A. Automatic modulation classification: principles, algorithms and applications, *John Wiley & Sons*, 2015.

- [5] Sills, J. Maximum-likelihood modulation classification for PSKQAM. *IEEE Military Communications. Conference Proceedings*, vol. 1, pp. 217-220, Atlantic City, NJ, USA, Aug. 1999.
- [6] Panagiotou, P., Anastasopoulos, A. & Polydoros, A. Likelihood ratio tests for modulation classification. *MILCOM 2000 Proceedings. 21st Century Military Communications. Architectures And Technologies For Information Superiority*, vol. 2, Los Angeles, CA, USA, Aug. 2000.
- [7] O'Shea, T., Roy, T. & Clancy, T. Over-the-air deep learning based radio signal classification. *IEEE Journal Of Selected Topics In Signal Processing*, vol. 12, 168-179, Jan. 2018.
- [8] Szegedy, C., Zaremba, W., Sutskever, I., Bruna, J., Erhan, D., Goodfellow, I. & Fergus, R. Intriguing properties of neural networks. *ArXiv Preprint arXiv:1312.6199*, 2013.
- [9] Sadeghi, M. & Larsson, E. Adversarial attacks on deep-learning based radio signal classification. *IEEE Wireless Communications Letters*, vol. 8, 213-216, Aug. 2018.
- [10] Kurakin, A., Goodfellow, I., Bengio, S. & Others Adversarial examples in the physical world, 2016.
- [11] Kim, B., Sagduyu, Y., Davaslioglu, K., Erpek, T. & Ulukus, S. Over-the-air adversarial attacks on deep learning based modulation classifier over wireless channels. *2020 54th Annual Conference On Information Sciences And Systems (CISIS)*, Princeton, NJ, USA, May 2020.
- [12] Papernot, N., McDaniel, P., Wu, X., Jha, S. & Swami, A. Distillation as a defense to adversarial perturbations against deep neural networks. *IEEE Symposium On Security And Privacy (SP)*, pp. 582-597, 2016.
- [13] Tramèr, F., Kurakin, A., Papernot, N., Goodfellow, I., Boneh, D. & McDaniel, P. Ensemble adversarial training: Attacks and defenses. *International Conference on Learning Representations*. Vancouver, BC, Canada, Feb. 2018.
- [14] Hendrycks, D. & Gimpel, K. Early methods for detecting adversarial images. *International Conference on Learning Representations*, Toulon, France, Feb. 2017.
- [15] Goodfellow, I., Pouget-Abadie, J., Mirza, M., Xu, B., Warde-Farley, D., Ozair, S., Courville, A. & Bengio, Y. Generative adversarial nets. *Advances In Neural Information Processing Systems*, Vol. 27. 2014.
- [16] Samangouei, P., Kabkab, M. & Chellappa, R. Defense-gan: Protecting classifiers against adversarial attacks using generative models. *International Conference on Learning Representations*, Vancouver, BC, Canada, May, 2018.
- [17] Goodfellow, I., Pouget-Abadie, J., Mirza, M., Xu, B., Warde-Farley, D., Ozair, S., Courville, A. & Bengio, Y. Generative adversarial nets. *Advances In Neural Information Processing Systems*. vol. 27. Montreal, Canada, Dec. 2014.
- [18] Hoang, Q., Nguyen, T., Le, T. & Phung, D. MGAN: Training generative adversarial nets with multiple generators. *International Conference On Learning Representations*. San Diego, CA, USA, Feb. 2018.
- [19] O'Shea, T. & West, N. Radio machine learning dataset generation with gnu radio. *Proceedings Of The GNU Radio Conference*, vol. 1, University of Colorado, Boulder, Sep. 2016.
- [20] Goodfellow, I., Shlens, J. & Szegedy, C. Explaining and harnessing adversarial examples. *ArXiv Preprint arXiv:1412.6572*. 2014.
- [21] Goodfellow, I., Pouget-Abadie, J., Mirza, M., Xu, B., Warde-Farley, D., Ozair, S., Courville, A. & Bengio, Y. Generative adversarial nets. *Advances In Neural Information Processing Systems*. Curran Associates, Inc. vol. 27, 2014.