# ECG-ATK-GAN: Robustness against Adversarial Attacks on ECGs using Conditional Generative Adversarial Networks

Khondker Fariha Hossain<sup>1</sup>, Sharif Amit Kamran<sup>1</sup>, Alireza Tavakkoli<sup>1</sup>, and Xingjun Ma<sup>2</sup>

Dept. of Computer Science & Engineering, University of Nevada, Reno, NV, USA School of Computer Science, Fudan University, China

Abstract. Automating arrhythmia detection from ECG requires a robust and trusted system that retains high accuracy under electrical disturbances. Many machine learning approaches have reached human-level performance in classifying arrhythmia from ECGs. However, these architectures are vulnerable to adversarial attacks, which can misclassify ECG signals by decreasing the model's accuracy. Adversarial attacks are small crafted perturbations injected in the original data which manifest the out-of-distribution shifts in signal to misclassify the correct class. Thus, security concerns arise for false hospitalization and insurance fraud abusing these perturbations. To mitigate this problem, we introduce the first novel Conditional Generative Adversarial Network (GAN), robust against adversarial attacked ECG signals and retaining high accuracy. Our architecture integrates a new class-weighted objective function for adversarial perturbation identification and new blocks for discerning and combining out-of-distribution shifts in signals in the learning process for accurately classifying various arrhythmia types. Furthermore, we benchmark our architecture on six different white and black-box attacks and compare them with other recently proposed arrhythmia classification models on two publicly available ECG arrhythmia datasets. The experiment confirms that our model is more robust against such adversarial attacks for classifying arrhythmia with high accuracy.

**Keywords:**  $ECG \cdot Adversarial Attack \cdot Generative Adversarial Network \cdot Electrocardiogram \cdot Deep Learning.$ 

# 1 Introduction

ECG is a crucial clinical measurement that encodes and identifies severe electrical disturbances like cardiac arrhythmia and myocardial infractions. Many artificial intelligence and machine learning approaches have been proposed to detect different types of ECGs accurately [19,25,11]. Recently, deep convolutional neural networks (CNNs) [30,1,20,2,32] has become the norm for achieving near-human-level performance for classifying cardiac arrhythmia and other cardiac abnormalities. Popular systems such as Medtronic LINQ II ICM [27], iRhythm Zio [31], and Apple Watch Series 4 [17] use embedded DNN models to analyze cardiac irregularities by monitoring the signals. Accurately detecting arrhythmia in real-time enables an immediate referral of the patient

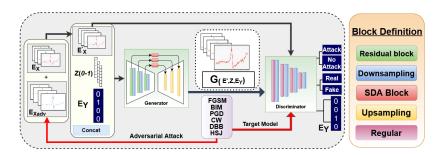


Fig. 1: Proposed ECG-ATK-GAN consisitng of a Generator and a Discriminator. The discriminator is utilized for generating the six attacked signals  $E_{x_{adv}}$ , namely FGSM, BIM, PGD, CW, HSJ, and DBB. These are then added with the non-attacked signals  $E_x$  to create the training data-set  $E_x^{'}$ . Contrarily, the Generator takes both attacked and non-attacked ECG signals,  $E_x^{'}$ , a noise vector z, and the class labels  $E_y$  as input.

to appropriate medical facilities. In addition to providing the patient with timely medical help, this will benefit the insurance companies by potentially reducing long-time consequences of delayed healthcare. Despite all these benefits, state-of-the art systems used to predict arrhythmia are vulnerable to adversarial attacks. These vulnerabilities are crucial as they can result in false hospitalization, misdiagnosis, patient data-privacy leaks, insurance fraud, and negative repercussion for healthcare companies [12,17].

Although these vulnerabilities are highly studied [9,23], a comprehensive solution is yet to be devised. Adversarial attacks misclassify ECG signals by introducing small perturbations that inject the out-of-distribution signal into the classification path. The perturbations could be introduced to the data by accessing the model parameters (Whitebox attack) or inferring the bad prediction outputs for a given set of input (Black-box attack) [7]. Current deep learning [30,1,20] and GAN-based [18,14,13,15] classifiers are not specifically designed to utilize the objective function to identify and mitigate adversarial attacked ECGs. Although recent works [17,9] illustrated the vulnerability of deep learning architectures to adversarial attacks, our work proposes a first-of-its-kind defense strategy against six different adversarial attacks for ECGs using a novel conditional generative adversarial networks. Additionally, we incorporate a class-weighted categorical loss function for identifying out-of-distribution perturbations and emphasizing the class-specific features. Both qualitative and quantitative benchmarks on two publicly available ECG datasets illustrate our proposed method's robustness.

# 2 Methodology

# 2.1 Generator and Discriminator

We propose a novel GAN based on a class-conditioned generator and a robust discriminator for categorical classification of both real and adversarial attacked ECG signals as illustrated in Fig. 1. The generator concatenates both non-attacked or attacked ECG

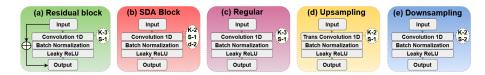


Fig. 2: Proposed (a) Residual, (b) Skip-Dilated Attention, (c) Regular, (d) Upsampling and, (e) Downsampling Blocks. Here, K=Kernel size, S=Stride, and D=Dilation rate.

signals  $E_x'$ , label  $E_y$ , and a noise vector z as input and generates  $G(E_x', Ey, z)$ . We use a Gaussian filter with  $\sigma=3$  to generate the smoothed noise vector, z. The label vector  $E_y$  in our model is utilized so that generated signal is not random. Rather it imitates class-specific ECG representing an arrhythmia. The noise vector, z ensures that the generated signal has small perturbations so that it does not fully imitate the original ECG signal and helps in overall training in extrapolation of generated signals. The generators incorporate Residual, Downsampling, Upsampling, and Skip-Dilated Attention (SDA) block as visualized in Fig. 1. The generator uses Sigmoid activation as output, so the synthesized signal is constrained within 0-1 as a continuous value.

The discriminator takes attacked/non-attacked real ECG, x and GAN synthesized ECG,  $G(E_x', E_y, z)$  signals sequentially while training. The discriminator consists of three regular blocks and three downsampling blocks (Fig. 1). The discriminator utilizes three losses: 1) Class-weighted Categorical cross-entropy for identifying adversarial attacked/non-attacked ECGs, 2) Categorical cross-entropy for normal and arrhythmia beat classification and 3) Mean-squared Error for GAN adversarial training. So we use three output activations: Sigmoid (GAN training), and two Softmax for adversarial attack and arrhythmia/normal beat classification.

- **Residual Block:** For extracting small perturbations in the attacked ECG signals, we use convolution with a small kernel, k=3 and stride, s=1 in the residual block in the Generator, illustrated in Fig. 2(a). This residual block is capable of extracting fine features that extrapolate the original signal to contain small perturbations and make it out-of-distribution. Specifically, the residual skip connection retains important signal-specific information that is added with more robust features extracted after the batch-normalization and leaky-ReLU activation.
- **Skip-Dilated Attention Block:** We use skip-dilated attention (SDA) block with kernel size, k=2, dilation rate, d=2 and stride, s=1, as illustrated in Fig. 2(b). By utilizing dilated convolution, our receptive fields become larger, covering larger areas of the attacked signals [33].
- **Regular Block:** We use the regular block for discriminators, containing convolution (k=3,s=1), batch-norm, and leaky-ReLU layers, as visualized in Fig. 2(c). Our main objective here is to encode the signals to meaningful classification outputs for two tasks, which is to 1) classify the type of arrhythmia and, 2) distinguish between non-attacked/attacked signals. Therefore, we avoid using any complex block for feature learning and extraction.
- Downsampling and Upsampling Blocks: The generator consist of both down-sampling and upsampling blocks, whereas the discriminator consist of only down-

### 4 Hossain et al.

sampling blocks to get the desired feature maps and output. The upsampling block consists of a transposed-convolution layer, batch-norm, and Leaky-ReLU activation layer successively and is given in Fig. 2(d). In contrast, The downsampling block comprises of a convolution layer, a batch-norm layer and a Leaky-ReLU activation function consecutively and is illustrated in Fig. 2(e).

# 2.2 Objective Function and Individual Losses

To distinguish non-attacked and attacked signals with out-of-distribution perturbations and emphasize the class-specific features even under significant perturbations, we propose a class-weighted categorical cross-entropy loss. The loss function is given in Eq. 1, where m=2, for attacked/non-attacked signal and  $\kappa$  is the class weight for the ground-truth,  $E_y$  and predicted class-label,  $E_{y'}$ .

$$\mathcal{L}_{atk}(D) = -\sum_{i=0}^{m} \kappa^i E_y^i \log(E_{y'}^i) \tag{1}$$

For classification of normal and different arrhythmia signals, we use categorical cross-entropy loss. Here, k = distinct normal/arrhythmia beats, depending on the dataset.

$$\mathcal{L}_{ary}(D) = -\sum_{i=0}^{k} E_y^i \log(E_{y'}^i)$$
(2)

For ensuring that the synthesized signal contains representative features of both adversarial examples and adversarial attacks, our generator incorporates the mean-squared error (MSE) as shown in Eq. 3. This helps the generator output signals with small perturbations that guarantee the signal to misclassify. As the generator, G is class-conditioned, it takes distinct ground truth class-label  $E_y$ , along with the attacked/non-attacked ECGs,  $E_x'$  and Gaussian noise vector z as input.

$$\mathcal{L}_{mse}(G) = \frac{1}{N} \sum_{i=1}^{N} (G(E_{x}^{'}, E_{y}, z) - E_{x})^{2}$$
(3)

We use Least-squared GAN [26] for calculating the adversarial loss and training our GAN. The cost function for our adversarial loss is given in Eq. 4. The discriminator takes real ECG signal,  $E_x$  and generated ECG signal,  $G(E_x^{'}, E_y, z)$  in two iterations. The adversarial loss quadratically penalizes the error while stabilizing the min-max game between the generator and discriminator.

$$\mathcal{L}_{adv}(D) = \left[ (D(E_x', E_y) - 1)^2 \right] + \left[ (D(G(E_x', E_y, z), E_y) + 1))^2 \right]$$
(4)

By incorporating Eq. 1, 2, 3 and 4, we can formulate our final loss function as given in Eq. 5. Here,  $\lambda_{mse}$ ,  $\lambda_{atk}$ , and  $\lambda_{ary}$  denote different weights, that are multiplied with their corresponding losses. We want our generator to synthesize realistic ECGs to fool the Discriminator, while classifying the types of arrhythmia with high accuracy. So, the final goal is to maximize the adversarial loss and minimize other losses.

$$\min_{G, D_{ary}, D_{atk}} \left( \max_{D_{adv}} (\mathcal{L}_{adv}(D)) + \lambda_{mse} \left[ \mathcal{L}_{mse}(G) \right] + \lambda_{atk} \left[ \mathcal{L}_{atk}(D) \right] + \lambda_{ary} \left[ \mathcal{L}_{ary}(D) \right] \right)$$
(5)

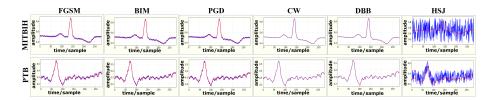


Fig. 3: The non-attacked and attacked signals (white and black-box attacks) overlapped on each other signified by Red and Blue lines.

# 2.3 Adversarial Attacks

We incorporated six established adversarial attacks (shown in Fig 3) that target our discriminator model as it is responsible for classifying different types of arrhythmia and normal beats in ECG signals. The reason for choosing these state-of-the-art attacks is to make our model more robust for intrusive perturbations in real-world applications. Four of these attacks are white-box, meaning detailed knowledge of the network architecture, the parameters, and the gradient w.r.t to the input is utilized to corrupt the data [6]. The other two are black-box attacks, meaning no knowledge of the underlying architecture or parameter is needed; instead, some output is observed for some probed inputs [6]. Moreover, the attack corrupts the data by estimating the gradient direction using the information at the decision boundary of the output [10,4]. We experimented with perturbation values,  $\epsilon$  ranging from 0.001 to 0.1, and selected the value which looked visually realistic and harder to discern. So, the visually realistic perturbations for FGSM, BIM, PGD and DBB is,  $\epsilon = 0.01$  and for CW,  $\epsilon = 0.1$ .

- Fast Gradient Sign Method (FGSM): This white-box attack creates attacked ECGs,  $E_{X_{adv}}$  by perturbing the original signal,  $E_X$ . For this, it calculates the gradients of the loss,  $\mathcal{L}_{ary}$  (Eq. 2) based on the input signal to create new adversarial signals that maximize the loss [16].
- Basic Iterative Method (BIM): This is an improved white-box attack, where the FGSM attack is iteratively updated in a smaller step size and clips the signals values of intermediate results to ensure the  $\epsilon$ -neighborhood of the original signal,  $E_X$  [22].
- **Projected Gradient Descent (PGD):** This white-box attack is considered the most decisive first-order attack. Though similar to BIM, it varies in initializing the example to a random point in the  $\epsilon$ -ball of interest (decided by the  $L_{\infty}$  norm) and does random restarts. In contrast, BIM initializes in the original point [24].
- Carlini-Wagner (CW): This is an optimization-based white-box attack [5]. It resolves the unboundedness issue by using line search to optimize the attack objective. We utilized the version with the  $L_{\infty}$  norm, i.e., for maximum perturbation applied to each point in the signal.
- Decision-based Boundary Attack (DBB): This is a decision-based black-box attack that starts from querying a large adversarial perturbation and then seeks to reduce the perturbation while staying adversarial [4]. It only requires the final class prediction of the model.

Hop Skip Jump Attack (HSJ): A powerful black-box attack that only requires
the final class prediction of the model [10]. And it is an advanced version of the
boundary attack, requiring significantly fewer model queries than Boundary Attack.

# 3 Experiments

# 3.1 Data Set Preparation

We used the PhysioNet MIT-BIH Arrhythmia dataset for our experiment [28]. We divided the dataset into four categories, N [Normal beat, Left and right bundle branch block beats, Atrial and Nodal escape beat], S [Atrial premature beat, Aberrated atrial, Supraventricular and Nodal premature beat], V [Premature ventricular contraction, Ventricular escape beat], and F [Fusion of the ventricular and regular beat]. We first find the R-peak for every signal, use a sampling rate of 280 centering on R-peak, and then normalize the amplitude between [0,1]. In the benchmarking, we combine and split the samples into 80% and 20% sets of train and test data. So we end up having train samples of N: 69958, S: 4766, V: 1965, F:617, and test samples of N: 17571, S: 1126, V: 473, and F: 157. To overcome the lack of minority class samples, we use Synthetic Minority Over-sampling Technique (SMOTE) [8] to increase the number of samples for S, V, and F to 10,000 each. We do not use SMOTE on test data. Next, we use the train and test ECG signals to create the six types of adversarial attacked ECGs (using Adversarial Robustness toolbox [29]). So we end up having same number attacked ECGs as non-attacked ones for each adversarial attacks. Next, we combine the original and adversarial ECGs to create our whole training dataset,  $E_x + E_{xadv} = E_x'$  (Fig. 1). We use 5-fold cross validation and select the model with the best validation score.

We also benchmark on PTB Diagnostic ECG Database [3], which consists of Normal and Myocardial Infraction beats. For each category, we use 10,000 samples, meaning we end up having 20,000 ECGs in total. We split them into 80% training and 20% test data. In similar manner to MITBH, we apply six adversarial attacks on these ECG signals. For training we end up having 32,000 (16,000 non-attacked and 16,0000 attacked) signals for each attack types. We use the same 5-fold cross-validation method.

# 3.2 Hyper-parameters

We chose  $\lambda_{atk}=10$  (Eq. 1),  $\lambda_{ary}=10$  (Eq. 2), and  $\lambda_{mse}=1$  (Eq. 3), to give more weight to classification losses than to adversarial loss. We give more weight to attacked signals than non-attacked ones by using  $\kappa=[1,1.05]$  (Eq. 1). We used Adam optimizer [21] with a learning rate of  $\alpha=0.0001$ ,  $\beta_1=0.5$  and  $\beta_2=0.999$ . We used Tensorflow 2.0 to train the model with batch size, b=128 for 100 epochs taking 4 hours to train on NVIDIA P100 GPU. We initialized the noise vector, z with float values between [0,1]. Code repository is provided in this link.

## 3.3 Quantitative Evaluation

We perform the quantitative evaluation by comparing our model with other state-ofthe-art architectures [30,1,20] on both attacked and non-attacked data from MITBH

Table 1: **MIT-BIH Dataset** : Comparison of architectures trained and evaluated on **non-attacked/attacked** ECGs for normal and three arrhythmia beat classification.

	Model	Accuracy	N		S		V		F	
	Model		Sensitivity	Specificity	Sensitivity	Specificity	Sensitivity	Specificity	Sensitivity	Specificity
No Attack	Proposed Method	99.2	98.8	95.2	83.7	99.8	97.9	99.7	92.4	99.3
	Shaker et al. [30]	98.6	97.4	98.1	93.0	98.7	99.2	99.0	87.2	99.6
	Kachuee et al. [20]	98.1	96.8	94.5	88.7	97.6	92.5	99.6	90.4	99.3
	Acharya et al. [1]	96.4	92.8	96.2	86.2	97.0	95.9	98.8	94.2	97.1
FGSM	Proposed Method	98.7	97.9	95.0	82.6	99.2	99.2	98.7	73.2	99.8
	Shaker et al. [30]	92.6	84.7	93.3	81.8	89.9	96.5	95.6	57.9	98.7
	Kachuee et al. [20]	86.5	73.1	82.8	68.9	83.4	73.7	97.3	82.8	93.2
	Acharya et al. [1]	77.2	53.3	87.9	65.7	74.1	66.2	92.3	65.6	87.9
BIM	Proposed Method	98.1	97.1	91.2	76.1	98.6	95.0	99.4	84.1	98.9
	Shaker et al. [30]	96.2	93.1	90.1	69.1	98.2	97.6	94.9	55.4	99.8
	Kachuee et al. [20]	85.6	70.6	90.4	67.4	90.5	83.6	92.3	82.1	88.5
	Acharya et al. [1]	76.9	54.4	87.5	49.0	88.2	42.1	91.9	87.8	73.8
PGD	Proposed Method	98.4	97.0	96.1	89.0	98.0	97.5	99.3	88.5	99.6
	Shaker et al. [30]	96.5	93.4	92.8	81.3	97.9	94.5	98.2	82.8	97.4
	Kachuee et al. [20]	87.2	74.0	86.9	66.3	87.6	84.7	91.8	65.6	95.2
	Acharya et al. [1]	77.2	54.0	88.8	54.1	91.3	57.8	82.0	66.2	80.6
CW	Proposed Method	98.8	97.8	97.0	91.1	98.8	98.3	99.5	91.0	99.5
	Shaker et al. [30]	95.4	90.8	96.0	84.7	97.0	97.8	94.1	72.6	99.7
	Kachuee et al. [20]	91.9	84.5	83.3	74.4	89.3	79.7	99.3	61.1	96.3
	Acharya et al. [1]	81.2	61.8	89.1	64.6	81.5	67.7	94.1	83.4	86.8
DBB	Proposed Method	93.0	85.8	94.9	84.4	96.1	91.3	96.2	84.1	93.9
	Shaker et al. [30]	90.1	80.6	86.7	65.1	94.2	79.4	94.9	83.4	91.7
	Kachuee et al. [20]	79.7	57.9	86.4	78.0	69.3	70.9	96.6	82.1	93.6
	Acharya et al. [1]	81.9	62.6	85.8	68.1	75.6	77.0	95.8	82.8	92.7
HSJ	Proposed Method	71.9	44.8	68.2	34.6	78.6	35.7	82.6	32.4	83.8
	Shaker et al. [30]	70.5	41.2	67.4	33.0	77.1	40.4	81.1	37.5	83.4
	Kachuee et al. [20]	68.6	39.3	65.8	21.3	82.5	10.3	94.5	43.9	62.2
	Acharya et al. [1]	68.4	37.6	70.0	32.3	57.4	23.6	90.4	20.3	90.0

and PTB datasets. In the first experiment, we use either only normal or adversarial attacked test data (19,327 and 4,000 for MITBH and PTB) for benchmarking the models on normal/abnormal beat classification, which is illustrated in Table. 1 and Table. 2. We train all the models on their respective attacked and non-attacked training samples for a fair comparison. For metrics, we use Accuracy, Sensitivity, and Specificity. We can see that for 'No Attack', all models achieve comparatively good results. However, for each distinct attack, the results worsen for other models compared to ours. The architecture in [30,1,20] utilizes 1D Convolution based architecture. Out of these models, Shaker et al. [30] adopt DC-GAN, a generative network for adversarial signal generation. However, their classification architecture is trained separately, and they provide results only on real ECG signals. One reason for their model's good performance for the no-attack scenario is training with GAN-generated adversarial samples, which helps to learn out-of-distribution signals. Moreover, the two 1D CNN architectures achieve better sensitivity for minority category F for FGSM, BIM, and HSJ attacks. Similarly, our model's performance on the minority category F is best for PGD, CW, and DBB attacks and second best for FGSM and BIM. Our model performs poorly against HSJ attacks because the signals have too much high noise and no clear pattern, as illustrated in 3. Besides that, our architecture's overall performance is more robust against adversarial attacks for classifying arrhythmia and myocardial infractions, as shown in Table. 2.

Table 2: **PTB Dataset**: Comparison of architectures trained and evaluated on **non-attacked/attacked** ECGs for normal and myocardial infarction beat classification.

3										
	Methods	No Attack	FGSM	BIM	PGD	CW	DBB	HSJ		
	Proposed Method	99.5	99.4	99.6	99.6	99.5	93.1	71.8		
Accuracy	Shaker et al. [30]	98.0	98.6	95.8	96.4	98.3	91.4	70.2		
Accuracy	Kachuee et al. [20]	95.2	97.1	94.4	92.2	91.3	88.6	56.5		
	Acharya et al. [1]	79.8	84.1	83.2	84.1	80.9	77.4	54.7		
	Proposed Method	99.3	99.2	99.6	99.7	99.2	92.6	79.8		
Sensitivity	Shaker et al. [30]	96.7	98.3	92.1	94.8	98.0	91.5	83.7		
Schsilivity	Kachuee et al. [20]	98.3	96.0	95.9	92.1	95.4	86.7	85.5		
	Acharya et al. [1]	82.1	93.1	93.4	90.6	90.3	88.5	90.7		
	Proposed Method	99.7	99.5	99.7	99.5	99.7	93.7	64.0		
Cnasificity	Shaker et al. [30]	99.3	98.9	99.4	98.0	98.7	91.2	56.8		
Specificity	Kachuee et al. [20]	92.1	98.2	93.0	92.2	87.3	90.3	28.1		
	Acharya et al. [1]	77.7	75.3	73.2	77.8	71.6	66.5	19.4		

Table 3: **Generator's Performance**: Similarity of adversarial and attacked / non-attacked signals.

		MITE	BIH		PTB				
	Maan Canarad Error	Structural	Cross-corelation	Normalized	Maan Sayarad Error	Structural	Cross-corelation	Normalized	
	Mean-Squared-Error	Similarity	Coefficiet	RMSE	Mean-Squared-Error	Similarity	Coefficient	RMSE	
No Attack	0.0129	99.90	99.86	3.487e-5	0.0184	99.87	99.93	8.152e-5	
FGSM	0.0117	99.81	99.89	2.890e-5	0.0001	99.84	99.89	0.02391	
BIM	0.0134	99.80	99.86	3.737e-5	0.0155	99.91	99.95	5.769e-5	
PGD	0.0122	99.84	99.89	3.115e-5	0.0179	99.88	99.92	7.722e-5	
CW	0.0065	99.95	99.97	9.038e-6	0.0188	99.87	99.92	8.498e-5	
DBB	0.0002	99.00	99.39	0.03159	0.0007	99.19	99.29	0.05532	
HSJ	0.0003	99.42	99.45	0.0393	0.0003	99.51	99.60	0.03872	

# 3.4 Qualitative Evaluation

For finding the similarity between real and synthesized attacked/non-attacked ECG signals, we benchmarked generated adversarial signals using four different metrics, i) Mean Squared Error (MSE), ii) Structural Similarity (SSIM), iii) Cross-correlation coefficient, and iv) Normalized Mean Squared Error (NRMSE). In Table. 3, We use both attacked and non-attacked signals from the test set. We score SSIM of 99.90%, 99.81% (FGSM), 99.80% (BIM), 99.84% (PGD), 99.95% (CW), 99.00% (DBB) and 99.43% (DBB) for MITBH Dataset. On the other hand we achieve SSIM of 99.87% (No Attack), 99.84% (FGSM), 99.91% (BIM), 99.84% (PGD), 99.87% (CW), 99.19% (DBB) and 99.51% (DBB) for PTB Dataset. As for cross-correlation, MSE, and NRMSE, our model generates quite realistic signals with minimal error.

# 4 Conclusions and Future Work

This paper presents ECG-ATK-GAN, a novel conditional Generative Adversarial Network for accurately predicting different types of arrhythmia from both regular and adversarially attacked ECGs. In addition, our architecture incorporates a new class-weighted categorical objective function for capturing out-of-distribution signals and robustly discerning class-specific features corrupted by adversarial perturbations. We provided an extensive benchmark on two publicly available datasets to prove the robust-

ness of our proposed architecture. One future direction is to improve our architecture by defending against other types of adversarial attacks.

**Prospect of application:** Detecting arrhythmia accurately and robustly in real-time will pave the way for better patient care and disease monitoring. In addition, insurance companies, contractors, partners, and many stakeholders will financially benefit from a trusted cardiac arrhythmia diagnostic system that is robust against adversarial attacks. This system can also help identify new attack types by distinguishing signal anomalies.

# References

- Acharya, U.R., Oh, S.L., Hagiwara, Y., Tan, J.H., Adam, M., Gertych, A., San Tan, R.: A deep convolutional neural network model to classify heartbeats. Computers in biology and medicine 89, 389–396 (2017)
- Ahmad, Z., Tabassum, A., Guan, L., Khan, N.: Ecg heart-beat classification using multi-modal image fusion. In: ICASSP 2021-2021 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP). pp. 1330–1334. IEEE (2021)
- 3. Bousseljot, R., Kreiseler, D., Schnabel, A.: Nutzung der ekg-signaldatenbank cardiodat der ptb über das internet. Biomedical Engineering / Biomedizinische Technik (1995)
- 4. Brendel, W., Rauber, J., Bethge, M.: Decision-based adversarial attacks: Reliable attacks against black-box machine learning models. arXiv preprint arXiv:1712.04248 (2017)
- 5. Carlini, N., Wagner, D.: Towards evaluating the robustness of neural networks. In: 2017 ieee symposium on security and privacy (sp). pp. 39–57. IEEE (2017)
- Chakraborty, A., Alam, M., Dey, V., Chattopadhyay, A., Mukhopadhyay, D.: Adversarial attacks and defences: A survey. arXiv preprint arXiv:1810.00069 (2018)
- Chakraborty, A., Alam, M., Dey, V., Chattopadhyay, A., Mukhopadhyay, D.: A survey on adversarial attacks and defences. CAAI Transactions on Intelligence Technology 6(1), 25– 45 (2021)
- Chawla, N.V., Bowyer, K.W., Hall, L.O., Kegelmeyer, W.P.: Smote: synthetic minority oversampling technique. Journal of artificial intelligence research 16, 321–357 (2002)
- Chen, H., Huang, C., Huang, Q., Zhang, Q., Wang, W.: Ecgadv: Generating adversarial electrocardiogram to misguide arrhythmia classification system. In: Proceedings of the AAAI Conference on Artificial Intelligence. vol. 34, pp. 3446–3453 (2020)
- Chen, J., Jordan, M.I., Wainwright, M.J.: Hopskipjumpattack: A query-efficient decisionbased attack. In: 2020 ieee symposium on security and privacy (sp). pp. 1277–1294. IEEE (2020)
- 11. Faziludeen, S., Sabiq, P.: Ecg beat classification using wavelets and svm. In: 2013 IEEE Conference on Information & Communication Technologies. pp. 815–818. IEEE (2013)
- 12. Finlayson, S.G., Bowers, J.D., Ito, J., Zittrain, J.L., Beam, A.L., Kohane, I.S.: Adversarial attacks on medical machine learning. Science **363**(6433), 1287–1289 (2019)
- Golany, T., Freedman, D., Radinsky, K.: Ecg ode-gan: Learning ordinary differential equations of ecg dynamics via generative adversarial learning. In: Proceedings of the AAAI Conference on Artificial Intelligence. vol. 35, pp. 134–141 (2021)
- Golany, T., Radinsky, K.: Pgans: Personalized generative adversarial networks for ecg synthesis to improve patient-specific deep ecg classification. In: Proceedings of the AAAI Conference on Artificial Intelligence. vol. 33, pp. 557–564 (2019)
- Golany, T., Radinsky, K., Freedman, D.: Simgans: Simulator-based generative adversarial networks for ecg synthesis to improve deep ecg classification. In: International Conference on Machine Learning. pp. 3597–3606. PMLR (2020)

- Goodfellow, I.J., Shlens, J., Szegedy, C.: Explaining and harnessing adversarial examples. arXiv preprint arXiv:1412.6572 (2014)
- Han, X., Hu, Y., Foschini, L., Chinitz, L., Jankelson, L., Ranganath, R.: Deep learning models for electrocardiograms are susceptible to adversarial attack. Nature medicine 26(3), 360–363 (2020)
- Hossain, K.F., Kamran, S.A., Tavakkoli, A., Pan, L., Ma, X., Rajasegarar, S., Karmaker, C.: Ecg-adv-gan: Detecting ecg adversarial examples with conditional generative adversarial networks. In: 2021 20th IEEE International Conference on Machine Learning and Applications (ICMLA). pp. 50–56. IEEE (2021)
- Jambukia, S.H., Dabhi, V.K., Prajapati, H.B.: Classification of ecg signals using machine learning techniques: A survey. In: 2015 International Conference on Advances in Computer Engineering and Applications. pp. 714–721. IEEE (2015)
- Kachuee, M., Fazeli, S., Sarrafzadeh, M.: Ecg heartbeat classification: A deep transferable representation. In: 2018 IEEE International Conference on Healthcare Informatics (ICHI). pp. 443–444. IEEE (2018)
- Kingma, D.P., Ba, J.: Adam: A method for stochastic optimization. arXiv preprint arXiv:1412.6980 (2014)
- 22. Kurakin, A., Goodfellow, I., Bengio, S., et al.: Adversarial examples in the physical world (2016)
- 23. Lam, J., Quan, P., Xu, J., Jeyakumar, J.V., Srivastava, M.: Hard-label black-box adversarial attack on deep electrocardiogram classifier. In: Proceedings of the 1st ACM International Workshop on Security and Safety for Intelligent Cyber-Physical Systems. pp. 6–12 (2020)
- 24. Madry, A., Makelov, A., Schmidt, L., Tsipras, D., Vladu, A.: Towards deep learning models resistant to adversarial attacks. arXiv preprint arXiv:1706.06083 (2017)
- 25. Mahajan, R., Kamaleswaran, R., Howe, J.A., Akbilgic, O.: Cardiac rhythm classification from a short single lead ecg recording via random forest. In: 2017 Computing in Cardiology (CinC). pp. 1–4. IEEE (2017)
- Mao, X., Li, Q., Xie, H., Lau, R.Y., Wang, Z., Paul Smolley, S.: Least squares generative adversarial networks. In: Proceedings of the IEEE international conference on computer vision. pp. 2794–2802 (2017)
- 27. Medtronic: Linq ii, cardiac monitors, https://www.medtronic.com/us-en/healthcare-professionals/products/cardiac-rhythm/cardiac-monitors/linq-ii.html
- 28. Moody, G.B., Mark, R.G.: The impact of the mit-bih arrhythmia database. IEEE Engineering in Medicine and Biology Magazine **20**(3), 45–50 (2001)
- 29. Nicolae, M.I., Sinn, M., Tran, M.N., Buesser, B., Rawat, A., Wistuba, M., Zantedeschi, V., Baracaldo, N., Chen, B., Ludwig, H., Molloy, I., Edwards, B.: Adversarial robustness toolbox v1.2.0. CoRR 1807.01069 (2018), https://arxiv.org/pdf/1807.01069
- 30. Shaker, A.M., Tantawi, M., Shedeed, H.A., Tolba, M.F.: Generalization of convolutional neural networks for ecg classification using generative adversarial networks. IEEE Access **8**, 35592–35605 (2020)
- 31. iRhythm Technologies: How could using ai impact cardiology?, https://www.irhythmtech.com/providers/evidence/ai
- 32. Wang, B., Liu, C., Hu, C., Liu, X., Cao, J.: Arrhythmia classification with heartbeat-aware transformer. In: ICASSP 2021 2021 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP). pp. 1025–1029 (2021)
- 33. Yu, F., Koltun, V.: Multi-scale context aggregation by dilated convolutions. arXiv preprint arXiv:1511.07122 (2015)