# CLASS-CONDITIONAL DEFENSE GAN AGAINST END-TO-END SPEECH ATTACKS

*Mohammad Esmaeilpour, Patrick Cardinal, Alessandro Lameiras Koerich*

École de Technologie Supérieure (ÉTS), Département de Génie Logiciel et des TI
1100 Notre-Dame W, Montréal, H3C 1K3, Québec, Canada

mohammad.esmaeilpour.1@ens.etsmtl.ca, {patrick.cardinal, alessandro.koerich}@etsmtl.ca

## ABSTRACT

In this paper we propose a novel defense approach against end-to-end adversarial attacks developed to fool advanced speech-to-text systems such as DeepSpeech and Lingvo. Unlike conventional defense approaches, the proposed approach does not directly employ low-level transformations such as autoencoding a given input signal aiming at removing potential adversarial perturbation. Instead of that, we find an optimal input vector for a class conditional generative adversarial network through minimizing the relative chordal distance adjustment between a given test input and the generator network. Then, we reconstruct the 1D signal from the synthesized spectrogram and the original phase information derived from the given input signal. Hence, this reconstruction does not add any extra noise to the signal and according to our experimental results, our defense-GAN considerably outperforms conventional defense algorithms both in terms of word error rate and sentence level recognition accuracy.

***Index Terms—*** Speech processing, Schur decomposition, chordal distance, adversarial subspace, adversarial defense.

## 1. INTRODUCTION

The threat of adversarial attacks has been well characterized in the domains of audio and speech recognition [1, 2]. Classifiers either trained on raw signals or their corresponding 2D representations (i.e., spectrograms) are quite vulnerable against carefully crafted adversarial examples and this poses a serious concern about safety and reliability of these models [3]. In a big picture, there are two main directions in studying adversarial attacks for speech signals: (i) generalizing strong attack algorithms developed for natural images in computer vision domain to spectrograms, taking advantage of their lower computational complexity [4, 5]; (ii) developing end-to-end attacks which require dealing directly with raw input signals [6, 7]. In this paper, we focus on the latter for defense purposes since it is closely related to a black-box attack scenario in real-life applications.

Although there are different implementations for end-to-end attacks, they unanimously use variants of the logarithmic distortion metric $l_{\text{dB}_{\vec{x}}}(\delta) = l_{\text{dB}}(\delta) - l_{\text{dB}}(\vec{x})$ [6], which measures the loudness in dB of an adversarial example $\vec{x}_{adv} = \vec{x}_{org} + \delta$ over its legitimate counterpart $\vec{x}_{org} \in \mathbb{R}^{n \times m}$, where $n$ and $m$ denote the length of the signal and the number of channels, respectively, and $\delta$ is the adversarial perturbation. Carlini and Wagner [6] have demonstrated the effectiveness of this measure as a constraint in their optimization formulation for attacking a speech-to-text model (C&W):

$$\min |\delta|_2^2 + \sum_i \vartheta_i.\mathcal{L}_i(\vec{x}_{org} + \delta_i, \pi_i) \quad \text{s.t.} \quad l_{\text{dB}_{\vec{x}}}(\delta) < \zeta \quad (1)$$

where $\pi_i$ refers to a character alignment (tokens without duplication) according to the target output phrase $\mathbf{y}_i$ in such a way that $\Pr(\pi_i|\mathbf{y}_i) = \prod_j \mathbf{y}_{\pi^j}^j$. Additionally, $\mathcal{L}_i(\cdot)$ denotes the connectionist temporal classification loss [8], and $\vartheta_i$ is a scaling factor. Finding an optimal value for $\zeta$ makes (1) brittle since it requires searching in an exponential space for a phrase $\mathbf{p}_i$, which should reduce to $\pi_i$ (after removing empty tokens). However, it has been shown that such a costly optimization formulation yields adversarial audios though sound very similar to $\vec{x}_{org}$, make the DeepSpeech system [9] generate any target phrase pre-defined by the adversary [6]. Since $\delta$ is not universal, slightly perturbing $\vec{x}_{adv}$ such as playback and recording over the air might override generating such a target phrase. In response to this issue, variants of expectation over transformation (EOT) have been developed as part of the optimization formulation inspired by [10]. Possible transformations are room impulse response, reverberation, and band-pass filters for truncating adversarial perturbation beyond human audible range [3]. However, this strong approach is more costly than (1) and it fits well for short signals with a few corresponding phrases [7]. The improved version of EOT has been recently introduced with a minor enhancement over the aforementioned distortion metric [7]:

$$10 \log_{10} |\rho_\delta|^2 - 10 \log_{10} |\rho_{\vec{x}_{org}}|^2 \quad (2)$$

where $\rho$ denotes the power spectral density (PSD) function. They have also introduced a new formulation for the loss function according to the configuration of the Lingvo speech-to-text system [11] :

$$\ell(\vec{x}_i, \delta_i, \mathbf{y}_i) = \mathbb{E}_{t \sim \tau} \left[ \ell_{net} \left( \mathbf{y_i}^o, \mathbf{y_i}^t \right) + \alpha \ell(\vec{x}_i, \delta_i) \right] \quad (3)$$

where $\alpha$ is a scalar and $\ell_{net}$ is the cross entropy loss which constrains over the normalized PSD function. Moreover, $\mathbf{y_i}^o$ and $\mathbf{y_i}^t$ denote the output and target phrases, respectively. This algorithm, which is known as robust attack, optimizes for the minimal $\delta_i$ over a set of $\tau$ transformations under varieties of room configurations. Similar minimization process has been implemented in a black-box scenario using a genetic algorithm (GA) [12]. Specifically, this GA-based attack (GAA) incorporates a momentum mutation approach as well as gradient estimation in order to obtain optimal candidate populations associated with a predefined target phrase.

While the fooling rate of the aforementioned adversarial attacks on DeepSpeech and Lingvo systems is almost 100%, there are few studies on defense approaches for speech-to-text systems. This might be due to the immaturity of the end-to-end attack algorithms since several playbacks of the crafted adversarial signal over the air might bypass the achieved perturbations [7]. Moreover, adversarial signals usually carry audible noises, even with $l_{\text{dB}_{\vec{x}}}(\delta) < 0$, which makes their detection easier [6]. However, reliable defense algorithms are still on demand against strong adversarial examples with less audible noises. Although there are some investigations for
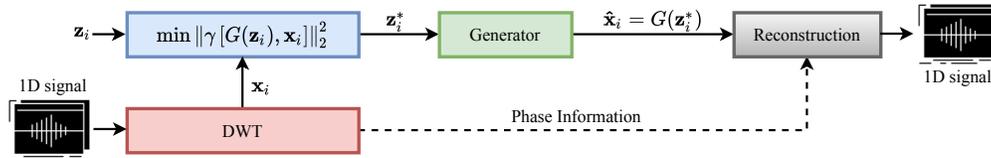
**Fig. 1**. Overview of the proposed end-to-end defense-GAN approach. The 1D signal converted to a 2D-DWT spectrogram is denoted as $\mathbf{x}_i$ and the prior $p_z$ for $\mathbf{z}_i \in \mathbb{R}^{d_z}$ is $\mathcal{N}(0, 0.4I)$. Additionally $\gamma\left[\cdot\right]$ is the chordal distance adjustment in the generalized Schur decomposition domain [2] and $\hat{\mathbf{x}}_i$ represents the synthesized spectrogram from the generator. 1D signal is reconstructed using inverse DWT.

both proactive and reactive defense approaches [13, 14], they are characterized in a small scale.

In this paper, we propose a new reactive adversarial defense using a class-conditional generative adversarial network [15]. We show that, our proposed defense scheme can be effective for large-scale systems such as DeepSpeech and Lingvo. The rest of the paper is organized as follows. In Section 2 we provide a brief introduction to GANs focused on defense strategies for speech signals. Section 3 presents our defense approach that includes three major steps for removing potential adversarial perturbations from signals. Section 4 summarizes and discusses the experiments carried out on Mozilla common voice (MCV) and LibriSpeech datasets. Conclusion and perspective of future work are presented in the last section.

## 2. GAN FOR ADVERSARIAL DEFENSE

In a typical GAN configuration organized as a two-player minimax optimization problem [16], the generator network $G(\mathbf{z}; \theta_g)$ with $\mathbf{z} \in \mathbb{R}^{d_z}$ and training parameters $\theta_g$ learns to map from the designated distribution $p_z \sim \mathcal{N}(0, I)$ to $p_g$ as:

$$\min_G \max_D \mathbb{E}_{\mathbf{x} \sim p_r(\mathbf{x})}\left[\log D(\mathbf{x})\right] + \mathbb{E}_{\mathbf{z} \sim p_z(\mathbf{z})}\left[\log\left(1 - D(G(\mathbf{z}))\right)\right]$$

(4)

where $p_r$ is the real sample distribution and $D(\mathbf{x}; \theta_d)$ denotes the discriminator network with training parameters $\theta_d$. Upon carefully training $G(\mathbf{z}; \theta_g)$, it can generate seamless samples almost without recognizable perturbations compared to $\mathbf{x}_i \sim p_r$. In fact, the generator semantically learns real sample distribution and we should expect unnoticeable differences between the generated samples and random test inputs except for adversarial examples. Based on this idea, a reactive defense approach has been introduced by Samangouei et al. [17], which iteratively minimizes for $\|G(\mathbf{z}) - \mathbf{x}\|_2^2$. Since the $L_2$ distance (or any other similarity metrics such as $L_\infty$) between crafted adversarial examples and their corresponding legitimate samples is fairly small, they extended their optimization problem subject to finding the most optimal $\mathbf{z}_i$. Unfortunately, this adversarial filtration defense scheme shatters gradient information and it can be easily disrupted by running a backward pass differentiable approximation (BPDA) attack [18]. On the contrary, the generator network can be trained to minimize the similarity between adversarial and legitimate samples where the discriminator iteratively learns to span possible adversarial manifolds [19]. Training such a defense-GAN requires exploring a massive adversarial subspace since not every attack algorithm generates a universal perturbation scale [2].

Autoencoder-based GAN (A-GAN) has also been investigated for defending speech emotion recognition models using long-short term memory networks [20]. This defense-GAN configuration introduces complex architecture for transforming a feature vector into another one aiming at bypassing potential adversarial perturbation. However, with the assumption of stable training without oversmoothing, this model might not necessarily enhance adversarial robustness against translation-invariant [21] or black-box attacks. However, these attacks are robust against low-level feature reconstruction using encoder-decoder blocks. In response to this issue and to the BPDA attack, we introduce a new defense-GAN architecture in a class-conditional framework which can be effectively used to increase the robustness of large-scale speech datasets and the state-of-the-art speech-to-text systems such as DeepSpeech and Lingvo.

## 3. PROPOSED DEFENSE APPROACH: CC-DGAN

The proposed adversarial defense approach is made up of three steps, as shown in Fig. 1: (i) generating signal representation; (ii) minimizing the relative chordal distance adjustment for the given input signal relative to $G(\mathbf{z}_i)$; and (iii) signal reconstruction with the preserved phase information. We explain all these steps in detail as follows.

### 3.1. 2D Signal Representation

Due to the high dimensionality of audio and speech signals, adversarial training either on single or multi-channel waveforms is very challenging and the model often undergoes complete collapse at early iterations. Therefore, a conventional approach in speech processing is to convert a given signal into a frequency-plot representation (spectrogram). Thus, as suggested by Esmaeilpour et al. [22], we divide the input signal into smaller chunks sampled at 16 kHz using discrete wavelet transform (DWT). Additionally, we set the frame length to 50 ms and use the complex Morlet mother function. Moreover, for enhancing the quality of the resulting spectrogram ($\mathbf{x}_i$ in Fig. 1), we represent its magnitude in a logarithmic scale. It has been shown that these settings for spectrogram production outperform short-time Fourier transform both in terms of recognition accuracy and robustness against adversarial attacks [5, 23]. Since the dimensions of the generated $\mathbf{x}_i$ are not necessarily square, we bilinearly resize them to 128×128 in compliance of computing the relative chordal distance in a non-Cartesian space.

### 3.2. Chordal Distance Adjustment Minimization

The chordal distance [24] is a metric that measures subspace adjacency for two similar samples in the domain of generalized Schur decomposition [2]. This metric has been used for characterizing the existence of adversarial examples in subspaces different from legitimate and noisy samples [2]. The chordal distance between an adversarial example $G(\mathbf{z}_i)$ and $\mathbf{x}_i$ is:

$$\mathrm{chord}\left(\lambda\left[G(\mathbf{z}_i)\right], \lambda\left[\mathbf{x}_i\right]\right) = \frac{\left|\lambda\left[G(\mathbf{z}_i)\right] - \lambda\left[\mathbf{x}_i\right]\right|}{\sqrt{1 + \lambda\left[G(\mathbf{z}_i)\right]^2}\sqrt{1 + \lambda\left[\mathbf{x}_i\right]^2}}$$
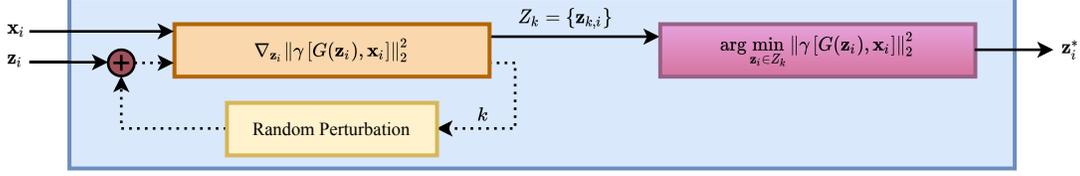
(5)

**Fig. 2**. $k$ steps minimization for the chordal distance adjustment between $G(\mathbf{z}_i)$ and $\mathbf{x}_i$. Similar to the predefined prior for $\mathbf{z}_i$, the random perturbation is also a function distributed over $\mathcal{N}(0, 0.4I)$. The inner loop is shown in dotted line.

where $\lambda[\cdot]$ denotes the vector of eigenvalues for the designated spectrograms. Achieving a valid chordal distance between two spectrograms for ensuring their subspace adjacency in the generalized Schur decomposition enquires $\|G(\mathbf{z}_i) - \mathbf{x}_i\| \simeq \xi_i$ where the threshold $\xi_i$ must be small according to the computed mean eigenvalue. For samples which lie in the same subspace, however with dissimilar spans, a minor translation is required in Eq. 5 to avoid ill-conditioned cases [24]. Specifically, for pencils $\vec{\mu}_i G(\mathbf{z}_i) - \mathbf{x}_i$ and $\vec{\mu}_i \in \text{diag}(\lambda[G(\mathbf{z}_i)]/\text{diag}(\lambda[\mathbf{x}_i])$, an adjustment $\gamma_i[\cdot] + \text{chord}(\cdot)$ is needed in (5), especially for samples with very small $L_2$ distance in Euclidean space [2].

Since the $\gamma_i[\cdot]$ adjustment is relatively large for an adversarial example $\mathbf{x}_{adv}$ [2], minimizing over $\|\gamma[G(\mathbf{z}_i)], \gamma[\mathbf{x}_{adv}]\|_2^2$ projects $\mathbf{x}_{adv}$ onto the legitimate sample subspace distribution represented by $p_g$. However, we do not filter $\mathbf{x}_{adv}$, neither by conventional encoder-decoder blocks nor by low-level transformation operations. In fact, we find an optimal $\mathbf{z}_i^* \in \mathbb{R}^{d_z}$ through an iterative approach, then pass it to the generator for crafting a spectrogram very similar to the given $\mathbf{x}_{adv}$. This approach is depicted in Fig. 2, where the number of iterations for obtaining the optimal $\mathbf{z}_i^*$ is denoted by $k$. For avoiding possible ill-conditioned pencils [24], we slightly perturb the candidate $\mathbf{z}_{k,i}$ with a random scalar and augment it with $\mathbf{z}_i$. Since $G(\mathbf{z}; \theta_g)$ is trained to support $p_g \approx p_r$, it considerably reduces the chance of generating spectrograms with adversarial perturbations. Therefore, the architectural design of both generator and discriminator has a crucial role. To this end, we propose simple yet effective class conditional architectures for reliable training.

### 3.2.1. Class-Conditional Defense GAN (CC-DGAN)

The proposed class-conditional defense GAN (CC-DGAN) is based on the vanilla GAN, where both the generator and the discriminator receive additional information on top of the noise vector $\mathbf{z}_i$ (i.e., $\mathbf{y}_i$) [15]. Unlike the baseline model (4), the CC-DGAN requires class embeddings ($c$-embeddings) mainly for the generator network: $\log(1 - D(G(\mathbf{z}|c = \mathbf{y})))$. This modification expands the learning space of the model at the risk of losing sample variety and mode collapse [25]. However, we find that $c$-embeddings provide a considerable boost in computing character probability distribution at every frame of the given signal compared to regular GANs.

The proposed generator receives $\mathbf{z}_i \in \mathbb{R}^{128} \sim \mathcal{N}(0, I)$ in the first layer followed by a linear block with dimension $50 + 128$ and shared $c - \text{embedding} = 50$ [26] including $4 \times 4 \times 16$ channels. There are two sequential residual blocks on top of the linear with $16 \rightarrow 4$ and $4 \rightarrow 1$ channels. The last hidden layer is a $128 \times 128$ non-local block with batch normalization and $\tanh$ activation function. The batch size is set 256 with orthogonal initialization [27]. Each residual block includes two linear ($128 \times 128$) and three padded convolution ($3 \times 3$ with stride 1) layers followed by upsampling, batch normalization, and ReLU activation function. In our discriminator network, the first layer requires RGB spectrogram $\mathbf{x}_i \in \mathbb{R}^{128 \times 128 \times 3}$. There is only one residual block in this network which contains two

sequential $3 \times 3$ convolution layers with concatenation, ReLU, skip-$z$ [25], and average pooling. On top of the residual block, there is a $64 \times 64$ non-local layer with 16 channels, ReLU, MaxPooling, and a linear logit layer ($\rightarrow 1$). Furthermore, we use both orthogonal regularization [28] and initialization [27] for the entire weight vectors.

### 3.3. Signal Reconstruction

This is the third step of the proposed defense approach as shown in Fig. 1. We reconstruct a given 1D signal with its own original phase information and the synthesized spectrogram $\hat{\mathbf{x}}_i$. Although synthesizing phase vectors with generative models is very challenging, there are some approaches for building them. However, they add audible hissing and whining noises to the signal. Signal reconstruction with original phase vectors often provides higher signal to noise ratio and this might help to more conveniently distinguish an adversarial example from a noisy signal [4]. The reconstruction operation only requires running an inverse DWT with basic settings such as type of mother function, sampling rate, and frame length. We use the same settings mentioned in Section 3.1 with additional quantization filter for normalizing the achieved vectors. For simplicity, we assume that signals are single-channel.

## 4. EXPERIMENTS

We have evaluated the proposed defense (CC-DGAN) against three end-to-end adversarial attacks for both Mozilla's implementation of DeepSpeech [29] and Lingvo system [11]. These speech-to-text models are trained on Mozilla common voice (MCV) [30] and LibriSpeech [31] datasets, respectively. Both these benchmarking datasets contain above 1,000 hours of voice clips with various utterances. However, as a common practice [6, 7, 12] we generate adversarial examples only for a portion of such datasets. We randomly select 11,500 and 6,000 samples from the MCV and LibriSpeech datasets for both training the CC-DGAN and running attacks, respectively. We organize these samples with their associated transcriptions into Subset-MCV and Subset-LS.

We run white-box (C&W) and black-box (GAA) adversarial attacks separately against the DeepSpeech model which uses rounds of long-short term memory blocks. We have randomly selected 1,000 samples from Subset-MCV with their original English transcriptions and we have targeted 10 different incorrect phrases (because these two attacks do not incorporate EOT) for effective attacking. Although these attacks directly optimize for achieving the minimum possible perturbation for the 1D signal, the DeepSpeech model first converts the given input into a Mel-frequency coefficient (MFC) representation. This adds more computational overhead to the attack algorithms and prohibits crafting adversarial examples for all the recordings in the dataset. The MFC layer splits the given speech signal into 50 frames per second which means the model can output up to 50 characters per second ($\mathbf{y}_i$). Therefore, this frame length is

**Table 1**. Comparison of different defense approaches against white and black-box adversarial attacks for DeepSpeech and Lingvo victim models. Better results are shown in bold face. In the robust attack, $\Delta$ is the offset scalar: $\|\delta_i\| < \zeta_i + \Delta$ [7] defined by the adversary.

| Model | Attack | Defense | Average $k$ | $\Delta$ | WER (%) | SLA (%) |
|---|---|---|---|---|---|---|
| DeepSpeech (Subset-MCV) | C&W | A-GAN | – | – | $23.98 \pm 2.14$ | $49.17 \pm 1.78$ |
| | | Compression [14] | – | – | $17.41 \pm 3.07$ | $56.96 \pm 2.38$ |
| | | Proposed CC-DGAN | 67 | – | $\mathbf{05.37 \pm 2.66}$ | $\mathbf{78.15 \pm 1.08}$ |
| | GAA | A-GAN | – | – | $18.54 \pm 5.31$ | $53.76 \pm 3.19$ |
| | | Compression [14] | – | – | $\mathbf{03.81 \pm 1.16}$ | $\mathbf{70.14 \pm 5.72}$ |
| | | Proposed CC-DGAN | 54 | – | $03.97 \pm 0.44$ | $68.35 \pm 2.51$ |
| Lingvo (Subset-LS) | Robust Attack | A-GAN | – | 300 | $21.23 \pm 4.79$ | $58.90 \pm 2.42$ |
| | | Compression [14] | – | 300 | $19.34 \pm 3.91$ | $54.88 \pm 4.52$ |
| | | Proposed CC-DGAN | 59 | 400 | $\mathbf{07.26 \pm 3.08}$ | $\mathbf{67.36 \pm 1.77}$ |

fairly enough for short signals with quite large transcripts. We extended these two attacks for targeting silence equivalent to generating empty tokens ($\epsilon$) for an additional 500 samples from the Subset-MCV. To this end, we updated the loss function as [6]:

$$\sum_i \max_{t \in \{\epsilon\}} \left( f(\vec{x})_t^i - \max_{\hat{t} \notin \{\epsilon\}} f(\vec{x})_{\hat{t}}^i, 0 \right), \quad f : \mathcal{X}^{50} \to [0,1]^{50 \cdot |\pi|}$$

(6)

where 50 and $\mathcal{X}$ denote the number of frames and input space, respectively. Moreover, $\hat{t}$ is the target phrase defined by the adversary in replacement of the original transcript $t$. Targeting $\epsilon$ token is easier than lexical characters and considerably reduces the computational cost. For the Lingvo victim model using the robust attack, we also randomly select 1,000 samples from Subset-LS with their associated transcripts targeting one incorrect phrase (because it incorporates EOT) with the same settings as mentioned in [12]. If the attack algorithm cannot exactly converge to a pre-defined target phrase, we replace it with another sample to keep the fooling rate at 100%.

For evaluating the proposed CC-DGAN to counteract the three adversarial attacks, we firstly train them separately on Subset-MCV and Subset-LS. In order to avoid losing sample variety and to add bias to our generative models, we exclude those nominated samples for adversarial attacks. For both generator and discriminator networks, we use Adam optimizer [32] with $\beta_1$=0, $\beta_2$=0.9, and a constant learning rate $2 \cdot 10^{-4}$. We also run an exploratory search for finding the optimal number of steps required for $G(\mathbf{z}; \theta_g)$ over $D(\mathbf{x}; \theta_d)$. We eventually opted to use two steps with decay rate 0.99 on two NVIDIA GTX-1080-Ti with $4 \times 11$GB memory in addition to a 64-bit Intel Core-i7-7700 (3.6 GHz) CPU with 64GB of RAM.

As a common issue in adversarial training, the proposed CC-DGAN configuration also undergoes collapse at about 9.3k and 6.8k iterations for Subset-MCV and Subset-LS, respectively. For improving the stability of our models, we have employed spectral normalization [33] only for $G(\mathbf{z}; \theta_g)$. However, it turns out oversmoothing the generated spectrogram. For rectifying this issue, we replaced long speech signals with shorter recordings, randomly drawn from the original datasets. The final GAN models used for further evaluations are those achieved from the checkpoints prior to potential collapse, which happens at about 10k iterations on both subsets. The $k$-step optimization algorithm for achieving $\mathbf{z}_i^*$ is depicted in Fig. 2 and finding a minimal value for it requires generalizable generative models. Regarding our experiments, for partially unstable and somewhat oversmoothed generators, $k$ never converges in less than 400 iterations.

For evaluating the performance of the proposed defense approach against the three aforementioned adversarial attacks, we use two metrics: (i) word error rate (WER), which is computed as $(I+S+D)/N \times 100$ where $I$, $S$, $D$, and $N$ are the total number of insertions, substitutions, deletions, and reference words, respec-

tively [7]; (ii) sentence level accuracy (SLA), computed as $n_c/n_{tot}$ where $n_c$ is the number of samples which could achieve the correct transcript and $n_c$ is the total number of test speech signals. Table 1 summarizes the results of our experiments, where both the SLA and WER are computed for the three defense algorithms. Specifically, these two metrics measure the performance of the defenses in producing phrases which reduce to correct transcriptions for the given adversarial signals. Note that, these two metrics while computed for the adversarial attacks, they measure fooling rates of the victim models in producing incorrect transcriptions as defined by the adversary. For consistent evaluation and in response to the raised concern of complete model vulnerability against end-to-end adversarial attacks [6], we set the SLA to 100% for all defense algorithms. Since for effective evaluations we target 10 incorrect transcriptions for every speech signal under C&W and GAA attacks, the reported results are averaged over 10 different runs. Table 1 shows that for the majority of the cases, the proposed CC-DGAN outperforms both simple compression and complex autoencoder-based GAN (A-GAN) in removing potential adversarial perturbations from speech signals and achieving lower WER and higher SLA. The only exception is for the GAA attack, which implements approximated gradient estimation, where simple compression achieves a slightly better performance than the proposed CC-DGAN. We noticed that for such attack, doubling $k$, reduces the WER in about 1.09% and increases the SLA in around 2.58% compared to $k$=54. For better investigating this issue, we attacked both victim models with the BPDA attack and measured the performance achieved by the proposed defense GAN. Our investigation on the same crafted adversarial examples uncovered the effectiveness of this attack on the CC-DGAN. More specifically, for reaching almost the same WER and SLA reported in Table 1, $k$ should be increased 2.37 and 3.12 times more for DeepSpeech and Lingvo systems, respectively.

## 5. CONCLUSION

In this paper, we proposed a new defense algorithm for securing advanced DeepSpeech and Lingvo systems against three end-to-end white-box and black-box adversarial attacks. The proposed CC-DGAN uses simple architectures for both the generator and discriminator with few residual blocks and a reconstructor module. This module regenerates a test input speech with the synthesized DWT spectrogram and its original phase information for seamless reconstruction. The experimental results on subsets of MCV and LibriSpeech datasets have shown that, the proposed defense approach considerably outperforms other defense algorithms for the majority of the cases in terms of achieving lower WER and higher SLA. Since the performance of our defense approach is highly dependent on the generalizability of the CC-DGAN, we are inclined to improve its stability and increase its generalizability in our future studies.

# 6. REFERENCES

[1] L. Schönherr, K. Kohls, S. Zeiler, T. Holz, and D. Kolossa, "Adversarial attacks against automatic speech recognition systems via psychoacoustic hiding," in *Netw and Distrib Syst Secur Symposium*, 2018, pp. 1–15.

[2] M. Esmaeilpour, P. Cardinal, and A. L. Koerich, "Detection of adversarial attacks and characterization of adversarial subspace," in *IEEE Intl Conf Acoust, Speech and Signal Process (ICASSP)*, 2020, pp. 3097–3101.

[3] H. Yakura and J. Sakuma, "Robust audio adversarial example for a physical attack," in *28th Intl J Conf Artificial Intelligence*, 2018, pp. 5334–5341.

[4] K. M. Koerich, M. Esmaeilpour, S. Abdoli, A. S. Britto, and A. L. Koerich, "Cross-representation transferability of adversarial attacks: From spectrograms to audio waveforms," in *Intl Joint Conf Neural Netw*, 2020, pp. 1–7.

[5] M. Esmaeilpour, P. Cardinal, and A. L. Koerich, "A robust approach for securing audio classification against adversarial attacks," *IEEE Trans Inf Forensics Security*, vol. 15, pp. 2147–2159, 2020.

[6] N. Carlini and D. Wagner, "Audio adversarial examples: Targeted attacks on speech-to-text," in *2018 IEEE Secur Privacy Workshops*, 2018, pp. 1–7.

[7] Y. Qin, N. Carlini, G. Cottrell, I. Goodfellow, and C. Raffel, "Imperceptible, robust, and targeted adversarial examples for automatic speech recognition," in *Intl Conf Mach Learn*, 2019, pp. 5231–5240.

[8] A. Graves, S. Fernández, F. Gomez, and J. Schmidhuber, "Connectionist temporal classification: labelling unsegmented sequence data with recurrent neural networks," in *23rd Intl Conf Mach Learn*, 2006, pp. 369–376.

[9] A. Hannun, C. Case, J. Casper, B. Catanzaro, G. Diamos, E. Elsen, R. Prenger, S. Satheesh, S. Sengupta, A. Coates, et al., "Deep speech: Scaling up end-to-end speech recognition," *arXiv preprint arXiv:1412.5567*, 2014.

[10] A. Athalye, L. Engstrom, A. Ilyas, and K. Kwok, "Synthesizing robust adversarial examples," in *Intl Conf Mach Learn*, 2018, pp. 284–293.

[11] J. Shen, P. Nguyen, Y. Wu, Z. Chen, M. X. Chen, Y. Jia, A. Kannan, T. Sainath, Y. Cao, C.-C. Chiu, et al., "Lingvo: a modular and scalable framework for sequence-to-sequence modeling," *arXiv preprint arXiv:1902.08295*, 2019.

[12] R. Taori, A. Kamsetty, B. Chu, and N. Vemuri, "Targeted adversarial examples for black box audio systems," in *IEEE Security and Privacy Workshops (SPW)*, 2019, pp. 15–20.

[13] J. Zhang, B. Zhang, and B. Zhang, "Defending adversarial attacks on cloud-aided automatic speech recognition systems," in *7th Intl Workshop Secur Cloud Comput*, 2019, pp. 23–31.

[14] N. Das, M. Shanbhogue, S.-T. Chen, L. Chen, M. E. Kounavis, and D. H. Chau, "Adagio: Interactive experimentation with adversarial attack and defense for audio," *arXiv preprint arXiv:1805.11852*, 2018.

[15] M. Mirza and S. Osindero, "Conditional generative adversarial nets," *arXiv preprint arXiv:1411.1784*, 2014.

[16] I. Goodfellow, J. Pouget-Abadie, M. Mirza, B. Xu, D. Warde-Farley, S. Ozair, A. Courville, and Y. Bengio, "Generative adversarial nets," in *Adv Neural Inf Process Syst*, 2014, pp. 2672–2680.

[17] P. Samangouei, M. Kabkab, and R. Chellappa, "Defense-GAN: Protecting classifiers against adversarial attacks using generative models," in *Intl Conf Learn Repres*, 2018.

[18] A. Athalye, N. Carlini, and D. Wagner, "Obfuscated gradients give a false sense of security: Circumventing defenses to adversarial examples," *arXiv preprint arXiv:1802.00420*, 2018.

[19] H. Lee, S. Han, and J. Lee, "Generative adversarial trainer: Defense to adversarial perturbations with gan," *arXiv preprint arXiv:1705.03387*, 2017.

[20] S. Latif, R. Rana, and J. Qadir, "Adversarial machine learning and speech emotion recognition: Utilizing generative adversarial networks for robustness," *arXiv preprint arXiv:1811.11402*, 2018.

[21] Y. Dong, T. Pang, H. Su, and J. Zhu, "Evading defenses to transferable adversarial examples by translation-invariant attacks," in *IEEE Conf Comput Vision Patt Recog*, 2019, pp. 4312–4321.

[22] M. Esmaeilpour, P. Cardinal, and A. L. Koerich, "Unsupervised feature learning for environmental sound classification using weighted cycle-consistent generative adversarial network," *Applied Soft Computing*, vol. 86, pp. 105912, 2020.

[23] M. Esmaeilpour, P. Cardinal, and A. L. Koerich, "From sound representation to model robustness," *CoRR*, vol. abs/2007.13703, 2020.

[24] C. F. Van Loan and G. H. Golub, *Matrix computations*, Johns Hopkins University Press, 1983.

[25] A. Brock, J. Donahue, and K. Simonyan, "Large scale GAN training for high fidelity natural image synthesis," in *Intl Conf Learn Repres*, 2019.

[26] E. Perez, F. Strub, H. de Vries, V. Dumoulin, and A. C. Courville, "Film: Visual reasoning with a general conditioning layer," in *32nd AAAI Conf Artificial Intelligence*, S. A. McIlraith and K. Q. Weinberger, Eds., 2018, pp. 3942–3951.

[27] A. M. Saxe, J. L. McClelland, and S. Ganguli, "Exact solutions to the nonlinear dynamics of learning in deep linear neural networks," in *2nd Intl Conf Learn Repres, ICLR*, Y. Bengio and Y. LeCun, Eds., 2014.

[28] A. Brock, T. Lim, J. M. Ritchie, and N. Weston, "Cneural photo editing with introspective adversarial networks," in *Intl Conf Mach Learn*, 2017.

[29] Mozilla Implementation, "Mozilla. project deepspeech," https://github.com/mozilla/DeepSpeech, 2017.

[30] Mozilla: commonvoice.mozilla.org, "Mozilla common voice dataset," https://voice.mozilla.org/en/datasets, 2019.

[31] V. Panayotov, G. Chen, D. Povey, and S. Khudanpur, "Librispeech: an asr corpus based on public domain audio books," in *IEEE Intl Conf Acoust, Speech and Signal Process (ICASSP)*, 2015, pp. 5206–5210.

[32] D. P. Kingma and J. Ba, "Adam: A method for stochastic optimization," *arXiv preprint arXiv:1412.6980*, 2014.

[33] T. Miyato, T. Kataoka, M. Koyama, and Y. Yoshida, "Spectral normalization for generative adversarial networks," in *Intl Conf Learn Repres*, 2018.