

SirenAttack: Generating Adversarial Audio for End-to-End Acoustic Systems

Tianyu Du, Shouling Ji, Jinfeng Li, Qinchun Gu, Ting Wang and Raheem Beyah

Abstract—Despite their immense popularity, deep learning-based acoustic systems are inherently vulnerable to adversarial attacks, wherein maliciously crafted audios trigger target systems to misbehave. In this paper, we present SIRENATTACK, a new class of attacks to generate adversarial audios. Compared with existing attacks, SIRENATTACK highlights with a set of significant features: (i) versatile – it is able to deceive a range of end-to-end acoustic systems under both white-box and black-box settings; (ii) effective – it is able to generate adversarial audios that can be recognized as specific phrases by target acoustic systems; and (iii) stealthy – it is able to generate adversarial audios indistinguishable from their benign counterparts to human perception. We empirically evaluate SIRENATTACK on a set of state-of-the-art deep learning-based acoustic systems (including speech command recognition, speaker recognition and sound event classification), with results showing the versatility, effectiveness, and stealthiness of SIRENATTACK. For instance, it achieves 99.45% attack success rate on the IEMOCAP dataset against the ResNet18 model, while the generated adversarial audios are also misinterpreted by multiple popular ASR platforms, including Google Cloud Speech, Microsoft Bing Voice, and IBM Speech-to-Text. We further evaluate three potential defense methods to mitigate such attacks, including adversarial training, audio downsampling, and moving average filtering, which leads to promising directions for further research.

I. INTRODUCTION

Nowadays machine learning-powered acoustic systems are ubiquitous in our everyday lives, ranging from smart locks on mobiles to speech assistants on smart home devices and to machine translation services on clouds. In general, acoustic systems can be categorized into two types according to application scenarios: classification-oriented systems and recognition-oriented systems. A classification-oriented acoustic system typically first transforms the audios from time domain to frequency domain and then performs classification on the corresponding spectrograms. As an example, a sound event classification system, which is often integrated into acoustic surveillance systems [1], [2], recognizes physical events such as glass breaking and gunshot. Compared with classification-oriented acoustic systems, a recognition-oriented acoustic system is often more complicated since it needs to first segment audios into frames, perform prediction on

each frame, and then derive the recognition results based on Connectionist-Temporal-Classification (CTC) loss [3] or attention [4]. The most typical example is the Automatic Speech Recognition (ASR) system, which is widely integrated into various popular speech assistants (e.g., Siri, Google Now, and Cortana).

Despite their immense popularity, acoustic systems based on classical models (e.g., Gaussian Mixture Model-Hidden Markov Model (GMM-HMM)) are shown vulnerable to various types of attacks: hidden voice command attack [5], in which the generated sounds are meaningless to humans while being interpreted as malicious commands (e.g., opening and unlocking doors, making unauthorized purchases, controlling sensitive home appliances) by speech recognition systems and DolphinAttack [6], in which the generated commands are inaudible to humans while audible to speech assistants. Such attacks can often be mitigated by mechanisms that differentiate the source (i.e., live speakers or synthesized replay) and nature (legitimate or malicious) of the received signal [7], [8].

Due to their superior performance, most of today's acoustic systems are built upon deep neural network models. However, such models are inherently vulnerable to adversarial inputs, which are maliciously crafted samples (typically by adding human-imperceptible noise to legitimate samples) to trigger target models to misbehave [9], [10]. Despite the plethora of work on the image domain (e.g., [11]), the research of adversarial attacks on the audio domain is still limited, due to a number of non-trivial challenges. First, the acoustic systems need to deal with information changes in the time dimension, which is more complex than image classification systems. Second, the audio sampling rate is usually very high (e.g., 16kHz, which means sampling 16,000 point per second), but images only have hundreds/thousands of pixels in total (e.g., the size of the images in the most popular datasets, i.e., MNIST and CIFAR-10, is 28×28 and 32×32 respectively). Therefore, it is harder to craft adversarial audios than images since adding slight noise to audios are less likely to impact the local features.

Recently, several mechanisms were proposed to generate adversarial audios [12], [13], [14]. They are all based on the gradient information, thereby having slight difference from each other. Even though these works against acoustic systems are seminal, they are limited in practice due to at least one of the following reasons: (i) they are designed only for a particular acoustic model under the white-box setting; (ii) they can only conduct untargeted attacks, with the goal of simply making the target systems misbehave; and (iii) they can only generate adversarial audios targeting phonetically similar

T. Du, S. Ji, J. Li are with the College of Computer Science and Technology, Zhejiang University, Hangzhou, Zhejiang 310027, China. S. Ji. is also with Alibaba-Zhejiang University Joint Research Institute of Frontier Technologies, Hangzhou, Zhejiang, China. E-mail: {zjradty, sj, lijinfeng0713}@zju.edu.cn

Q. Gu and R. Beyah are with the Department of Electrical and Computer Engineering, Georgia Institute of Technology, Atlanta, GA 30332, USA. Email: guqinchen@gatech.edu, rbeyah@ece.gatech.edu

T. Wang is with the Department of Computer Science, Lehigh University, Bethlehem, PA 18015, USA. Email: inbox.ting@gmail.com

phrases. In addition, none of them was comprehensively evaluated in end-to-end settings (more detailed analysis in Section II).

In this paper, we present SIRENATTACK, a new class of adversarial attacks against deep neural network-based acoustic systems. Compared with prior work, SIRENATTACK departs in significant ways: *versatile* – SIRENATTACK is applicable to a range of end-to-end acoustic systems under both white-box and black-box settings; *targeted* – SIRENATTACK generates adversarial audio that trigger target systems to misbehave in a highly predictable manner (e.g., misclassifying the adversarial audio into a specific class); and *evasive* – SIRENATTACK generates adversarial audios by injecting a small amount of noise to legitimate audios while having negligible impact to human perception.

Our Contribution. To the best of our knowledge, this work represents the first systematical study on generating adversarial audios for various end-to-end acoustic systems. Our main contributions can be summarized as follows.

- We present SIRENATTACK, a new class of adversarial attacks against deep neural network-based acoustic systems under both white-box and black-box settings. For the white-box scenario, we combine a heuristic algorithm with a gradient-based method to conduct targeted/untargeted adversarial attacks, which is more effective and efficient than previous work [12] as demonstrated by experimental results. For the black-box scenario, we propose a new approach to conduct targeted/untargeted adversarial attacks by making use of a strong, iterative, and gradient-free algorithm.
- We evaluate SIRENATTACK on a range of state-of-the-art deep neural network models used in popular acoustic systems, including speech command recognition, speaker recognition and sound event classification systems. Experimental results show that SIRENATTACK is highly effective. For instance, it achieves 99.45% success rate on the IEMOCAP dataset against the ResNet18 model. Further, the generated adversarial audios can also be misinterpreted by multiple popular ASR platforms, including Google Cloud Speech Recognition, Microsoft Bing Voice Recognition, and IBM Speech-to-Text.
- We propose three potential defense strategies to mitigate the attacks of SIRENATTACK and conduct preliminary evaluation. Our results shed light on building more robust deep neural network-based acoustic systems, and lead to promising directions for further research.

II. RELATED WORK

A. Traditional Attacks on Acoustic Systems

In [5], Carlini *et al.* generated sounds that are unintelligible to humans while can be interpreted as commands by machine learning models. This attack targets at GMM-HMM systems rather than the advanced end-to-end neural networks used in most modern speech recognition systems as we focus on in this paper. In [6], Zhang *et al.* proposed DolphinAttack, which exploits the non-linearity of the microphones to create commands inaudible to humans while audible to speech assistants.

From the defence perspective, such attack can be eliminated by an enhanced microphone that can suppress acoustic signals on the ultrasound carrier. In [15], Yuan *et al.* embedded voice commands into songs, which can be recognized by ASR systems over the air while being imperceptible to a human listener. However, this kind of attacks can be defended by audio turbulence and audio squeezing in practice.

B. Adversarial Attacks on Acoustic Systems

Inspired by adversarial attacks on images, adversarial audios have also drawn researchers' attention. In [12], Carlini *et al.* proposed a method that can produce an adversarial audio that could be transcribed as the desired text by DeepSpeech [16] under white-box settings. Nevertheless, their method would take more than one hour to generate an adversarial audio, and thus is very inefficient. In [13], Cisse *et al.* proposed the Houdini attack that is transferable to different unknown ASR models. However, it can only construct adversarial audios targeting phonetically similar phrases. In [14], Iter *et al.* generated adversarial audios by adding perturbations to the Mel-Frequency Cepstral Coefficients (MFCC) features and then rebuilt the speech from the perturbed MFCC features. Nevertheless, the noise introduced by the inverse-MFCC process makes their adversarial audios sound strange to human. In [17], Gong *et al.* demonstrated that a 2% distortion of speech can make a Deep Neural Networks (DNNs) based model fail to recognize the identity of the speaker. However, it is an untargeted attack that is difficult to pose a real threat.

C. Defense for Acoustic Systems

As traditional attacks on acoustic systems have been extensively studied, there are many defense methods to eliminate the effects of them. These defense methods are based on the similar ideas, i.e., determining whether the received signal is from a live speaker. In [7], the authors proposed a virtual button that leverages Wi-Fi to detect human motions, and voice commands are only accepted when human motion is detected. In [8], the authors proposed VAuth, which collects the body-surface vibration of the user through a wearable device and verifies that the voice command is from the user. However, these methods are limited since voice commands are not necessarily accompanied with detectable motion, and the need for wearable devices (e.g., eyeglasses) may be inconvenient. Other defence schemes [18], [5], [19] mention the possibility of using Speaker Verification (SV) systems for defense. Nevertheless, this is not very useful since the SV system itself is vulnerable to previously recorded user speech [5]. As for adversarial attacks on acoustic systems, there are few defense schemes in published literature. Therefore, in this paper, we propose three potential defenses against such attacks. More in-depth dedicated defense research is expected in the future.

D. Remarks

In summary, the following aspects distinguish SIRENATTACK from existing adversarial attacks on acoustic systems.

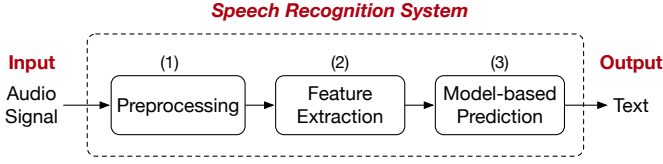


Fig. 1. A typical end-to-end speech recognition system.

First, previous work usually focuses on one acoustic system and attacks only one or two models under white-box settings. In contrast, we systematically study adversarial audios against state-of-the-art acoustic models in three kinds of popular acoustic scenarios under both white-box settings and black-box settings. To our best knowledge, this is the first large-scale evaluation on the robustness of state-of-the-art acoustic models. Second, SIRENATTACK is computationally efficient and can generate an adversarial audio within minutes. Finally, our adversarial audios can also be misinterpreted by many popular ASR platforms, while previous studies seldom evaluate their attacks' performance on real ASR platforms. This implies that SIRENATTACK is more general and robust.

III. BACKGROUND

A. Recognition-oriented Acoustic Systems

An end-to-end speech recognition model can directly map the raw audio into the output words as shown in Fig. 1. It consists of the following three steps: (1) *Pre-processing*. This step eliminates the time periods whose signal energy falls below a particular threshold. One of the most popular technique used in this step is Voice Activity Detection (VAD), which usually consists of a noise reduction stage, a block-feature calculation stage and a classification stage. (2) *Feature Extraction*. This step splits the pre-processed audio into short frames and extracts features from each frame. The most commonly used feature extraction method in speech recognition systems is MFCC [20]. (3) *Model-based Prediction*. This step takes the extracted features as input, and matches them with an existing model to generate prediction results. Modern systems usually use Recurrent Neural Networks (RNNs) with a CTC loss function [3], which only requires one input sequence and one output sequence.

B. Classification-oriented Acoustic Systems

Generally, the intent of a classification-oriented acoustic system is to categorize the sample points in a clip of audio into one of the given classes. As shown in Fig. 2, a classification-oriented acoustic system consists of the following three steps: (1) *Pre-processing*. This step is the same as that in recognition-oriented systems. (2) *Feature Extraction*. This step can extract audio-level features and frame-level features. Specifically, the audio-level features are extracted from the whole audio waveforms [21], while the frame-level features are extracted from the segmented waveform frames [22]. (3) *Model-based Classification*. This step matches the extracted features with an existing model to generate classification results. The technique used in this step can vary widely.

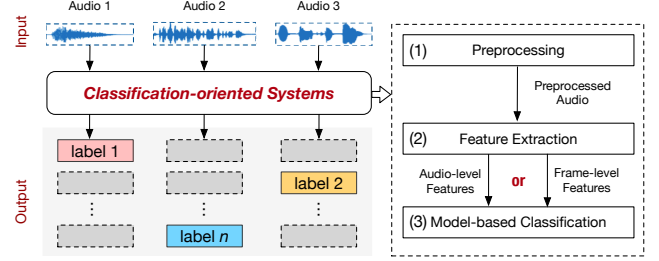


Fig. 2. A typical classification-oriented system.

Nevertheless, modern systems usually use CNNs [23] due to its outstanding performance in the computer vision domain.

IV. ATTACK DESIGN

A. Problem Formulation

Given a target classification/recognition model $f : \mathcal{X} \rightarrow \mathcal{Y}$ from a feature space \mathcal{X} to a set of prediction results \mathcal{Y} , an adversary aims to generate an adversarial audio x_{adv} from a legitimate audio $x \in \mathcal{X}$ with its ground truth label $y \in \mathcal{Y}$, so that $x_{adv} \approx x$, i.e., it is difficult for human to distinguish x_{adv} from x , while the classifier predicts $f(x_{adv}) = t$ where t is the targeted phrase or class and $t \neq y$.

B. Threat Model

Under white-box settings, attackers are assumed to have the complete knowledge of all the details including model architecture and model parameters about the victim model and can interact with it while conducting the attack. This is a common threat model adopted in most prior work [24], [12] which assumes an adversary with the most power.

Under black-box settings, attackers are assumed to know nothing about the architecture, parameters or training data of the victim model. Therefore, the query function of the victim model can be characterized as an oracle $\mathcal{O}(x)$ which returns the confidence value of the candidate classes. This assumption is practical, since many Machine-Learning-as-a-Service (MLaaS) platforms usually do not release their detailed algorithms or training data but provide the confidence value of each candidate class.

C. Preparation

SIRENATTACK is based on the Particle Swarm Optimization algorithm [25] and the fooling gradient method [10]. We begin by briefly introducing these techniques.

Particle Swarm Optimization (PSO). PSO is a heuristic and stochastic algorithm to find solutions for optimization problems by imitating the behavior of a swarm of birds [25]. It can search a very large space of candidate solutions while does not require the gradient information. At a high level, it solves a problem by iteratively making a population of candidate solutions (which we referred to as *particles*) move around in the search-space according to their fitness values. The fitness value of a particle is the evaluation result of the objective function on that particle's position in the solution space. In each iteration, each particle's movement is influenced

by its local best position P_{best} , and meanwhile is guided toward the global best position G_{best} in the search-space. This iteration process is expected to move the swarm toward the best solution. Once a termination criterion is met, G_{best} should hold the solution for a global minimum.

Fooling Gradient Method. This method is first used to generating adversarial images [10], and is later leveraged by researchers to conduct simple adversarial attacks on audios [12], [14]. In this method, the gradient is computed with respect to the input data rather than the model parameters. Then gradient descent technique is applied to iteratively modify the input data. In a nutshell, the key differences between the standard setup for training a NN and the fooling gradient method are (1) the gradients are applied only to the input data, and (2) the loss is computed between the network's predictions and the target labels rather than the ground truth labels.

D. Design of the Attack

Logically, a classification-oriented task can be regarded as a one-frame instance of the recognition-oriented task. Hence, we introduce SIRENATTACK from the aspect of white-box and black-box settings instead of application scenarios.

1) *White-box Attack:* At a high level, the white-box attack contains two phases. The goal of the first phase is to find a coarse-grained noise δ' that is close to the exact adversarial noise δ , while the goal of the second phase is to find the exact adversarial noise δ by slightly revising δ' . Such procedure is designed under the consideration of effectiveness and efficiency. The detailed white-box attack is shown in Algorithm 1. The first phase contains the steps in line 2-13 and the second phase contains the steps in line 15-19.

First, we initialize the *epoch* to zero and generate $n_{particle}$ randomized sequences from a uniform distribution (line 1). The randomized sequences are collectively referred to as *seeds*. Then we run the PSO subroutine (line 3) with the target output t and *seeds*. If any particle p_i produces the target output t when being added to the original audio x , then the attack succeeds (line 4-5), and the particle p_i is the expected noise δ . Otherwise, we will preserve the best particle that has the minimum fitness value in the current PSO run as one of the *seeds* in the next PSO run (line 7-8). From the n_{step} epoch, we calculate the standard deviation *std* of the global best fitness value from the last n_{step} PSO runs (line 10-12). Once the *std* is below the threshold ϵ , it is not efficient for continuously running the PSO subroutine to find the exact noise δ since the global best fitness value now changes slowly. Therefore, we only obtain a coarse-grained noise δ' after the first phase.

We would further emphasize two key aspects of our attacks: (1) We modify the PSO to globally keep track of the current saved best particle throughout all PSO iterations instead of using the standard PSO. (2) During each iteration, PSO aims to minimize an objective function defined as $g(x + p_i, t)$. Note that RNN-like models' output is a matrix containing the probability of the characters at each frame. Therefore, we choose the CTC loss [3] as $g(\cdot)$ in this attack, i.e., $g(x + p_i, t) = CTC - loss(x + p_i, t)$. The value of $g(x + p_i)$ at each particle is then used to move them in a new direction.

Algorithm 1 Generation of targeted adversarial audios under white-box settings

Input: Original audio x , target output t , $n_{particles}$, $epoch_{max}$, n_{step} , ϵ
Output: An targeted adversarial audio x_{adv}

- 1: Initialize $epoch = 0$ and *seeds* and set CTC loss as the objective function;
- 2: **while** $std \leq \epsilon$ **do**
- 3: Run PSO subroutine with t and *seeds*;
- 4: **if** any particle produce target output t during PSO **then**
- 5: Solution is found. Exit.
- 6: **else**
- 7: Clear *seeds*;
- 8: *seeds* \supseteq *best particle* that produce the minimum CTC loss value from the current PSO run;
- 9: **end if**
- 10: **if** $epoch \geq n_{step}$ **then**
- 11: Calculate *std*;
- 12: **end if**
- 13: **end while**
- 14: Obtain coarse-grained noise δ' from current *seeds*;
- 15: **while** $epoch \leq epoch_{max}$ or $\mathcal{O}(x + \delta') \neq t$ **do**
- 16: Calculate loss function according to Eq. (1),
- 17: Update δ' according to the gradient information;
- 18: $epoch = epoch + 1$;
- 19: **end while**
- 20: Get adversarial audio x_{adv} with target label t .

In the second phase (line 15-19), we use the SGD optimizer to adjust δ' until $\mathcal{O}(x + \delta') = t$ or $epoch$ reaches $epoch_{max}$. The second-stage loss function is defined as:

$$\text{minimize } \mathcal{L}(x + \delta', t) + \lambda \|\delta'\|_2^2 \quad (1)$$

where \mathcal{L} is also the CTC loss and $\lambda \|\delta'\|_2^2$ is the regularization term. This loss function can be revised to $-\mathcal{L}(x + \delta', y)$ to conduct untarget attacks, where y is the ground truth label.

2) *Black-box Attack:* The detailed black-box attack is shown in Algorithm 2. The basic procedure (line 2-11) of the black-box attack is similar to the white-box attack's first phase except for the following two things: (i) the objective function is different due to lacking of gradient information, and (ii) the termination condition is different since we should obtain the exact noise δ in this process. We experimented with several definitions of $g(\cdot)$ and found the following to be the most effective:

$$g(x + p_i, t) = \max_{j \neq t} (\max_j (\mathcal{O}(x + p_i)_j) - \mathcal{O}(x + p_i)_t, \kappa) \quad (2)$$

where $\mathcal{O}(x + p_i)_j$ is the confidence value of label j for input $x + p_i$. This hinge loss function is inspired by the state-of-the-art model evasion method – ZOO attack [26]. This function can move the particles to the position that maximizes the probability of the target label t . In addition, we can control the confidence of misprediction with the parameter κ , and a smaller κ means that the found adversarial audio will be predicted as t with higher confidence. We set $\kappa = 0$ for SIRENATTACK but we note here that a side benefit of

this formulation is that it allows one to control the desired confidence. The algorithm iterates on this process (line 2-11) till the attack succeeds or it reaches $epoch_{max}$. If succeed, we would obtain an adversarial audio x_{adv} that can be predicted as t by the victim model. Furthermore, this function can be used to conduct untargate attacks with trivial modifications.

Compared with the white-box attack, the black-box attack is less efficient and introduces more noise in the generated adversarial audios. This is because the black-box attack lacks of loss information and gradient information. Therefore, some performance decrease of the black-box attack is reasonable.

Algorithm 2 Generation of targeted adversarial audios under black-box settings

Input: Original audio x , target output t , $n_particles$ and $epoch_{max}$
Output: An targeted adversarial audio x_{adv}

- 1: Initialize $epoch = 0$ and $seeds$ and set Eq. (2) as the objective function;
- 2: **while** epoch reaches $epoch_{max}$ **do**
- 3: Run PSO subroutine with t and $seeds$;
- 4: **if** any particle produce target output t during PSO **then**
- 5: Solution is found. Exit.
- 6: **else**
- 7: Clear $seeds$;
- 8: $seeds \supseteq best\ particle$ that produce the minimum value of Eq. (2) from the current PSO run;
- 9: **end if**
- 10: $epoch = epoch + 1$;
- 11: **end while**
- 12: Get adversarial audio x_{adv} with target label t .

V. WHITE-BOX ATTACK EVALUATION

A. Datasets

In this experiment, we take the audios from the Common Voice dataset [27] and the VCTK Corpus [28] as the original samples. The Common Voice dataset is a corpus of speech data read by users based upon the text from a number of public domain sources like user submitted blog posts, old books, movies, and other public speech corpora. It has 500 hours of samples, comprising 400,000 recordings made by 20,000 people. The VCTK Corpus includes speech data from 109 native speakers of English with various accents.

B. Target Model

In this evaluation, we examine the security of DeepSpeech [16], a state-of-the-art RNN-based ASR model proposed by Baidu, which is trained on a dataset consisting of 100,000 hours of noisy speech data and can achieve around 81% accuracy in noisy environments like restaurants. Although there are other speech recognition models proposed in [29], we choose DeepSpeech as our target model due to the following reasons: (i) it is hard to reproduce the results in those papers due to lacking of sufficient implementation details, and (ii) the input data format of some available Speech-To-Text engines (e.g., WaveNet) are MFCC features instead of raw audio

waveforms, and therefore we need to rebuild the adversarial audio from the inverse MFCC process which will greatly reduce the quality of the audios. On the other hand, the DeepSpeech model implemented by the Mozilla group, which has more than 6,000 stars in the Github repository, is a proper choice for evaluating SIRENATTACK. Its input data format is raw audio waveform. Though it is a research project now, the developers of DeepSpeech claimed that Baidu would integrate DeepSpeech into automatic car, CoolBox and wearable devices in the future. Thus, it is more practical than other models.

C. Evaluation Metric

There are two objective audio quality assessment techniques [30], i.e., Signal-to-Noise Ratio (SNR) and Objective Difference Grade (ODG). As previous works usually use SNR to evaluate the quality of generated adversarial audios [15], [31], we also use SNR to evaluate the audio quality for consistency and comparison with previous works.

SNR is a metric extensively used to quantify the level of signal power to noise power, which is calculated as follows:

$$SNR(dB) = 10 \log_{10} \left(\frac{P_x}{P_\delta} \right) \quad (3)$$

where x is the original audio, δ is the added noise, and P_x and P_δ are the power of the original signal and the noise signal, respectively. A large SNR value indicates a small noise scale. For our purpose, we use it to measure the distortion of the adversarial audio relative to the original audio.

According to the International Federation of the Phonographic Industry (IFPI), the imperceptible noise requires at least 20 dB SNR value between the noise signal and the original signal. However, this is unnecessary for SIRENATTACK. SIRENATTACK tolerates the noise to some extent as long as it does not impact human perception. Therefore, the SNR of the generated adversarial audio is acceptable even though they do not reach the 20 dB threshold. To further demonstrate this, a user study was conducted in Section VII-C.

D. Implementation

We conducted the experiments on a server with two Intel Xeon E5-2640 v4 CPUs running at 2.40GHz, 64 GB memory, 4TB HDD and a GeForce GTX 1080 Ti GPU card. We set $epoch_{max} = 300$, $n_step = 5$, $\epsilon = 2$ and the iteration limit of PSO to 30 in all experiments. For the PSO subroutine, we set $n_particles = 25$, $c_1 = c_2 = 1.4961$. Specially, r_1 and r_2 are random values uniformly sampled from $[0, 1]$ to avoid consistency. In addition, we adopted the adaptive method on inertia weight w , i.e., we initially set $w = 0.9$, which makes the PSO has strong global optimization ability; with the increasing of the iteration, w is decremented, so that the PSO has strong local optimization ability; when the iteration ends, $w = 0.1$. The specific meaning of these hyper-parameters can be found in [25]. For the gradient-based phase, we did some search over hyper-parameters such as learning rate to find a trade-off between effectiveness and efficiency. In particular, we set the learning rate as 1.

TABLE I
RESULTS OF THE WHITE-BOX ATTACK ON DEEPSPEECH.

Dataset	Original Length	Target Length	Performance (without VAD)			Performance with VAD		
			Success Rate	SNR(dB)	Time(s)	Success Rate	SNR(dB)	Time(s)
Common Voice	11.87 words	4.74 words	100%	18.72	1560.14	100%	20.01	1201.66
VCTK Corpus	12.93 words	4.74 words	100%	16.54	1897.49	100%	18.34	1613.87

TABLE II
ADVERSARIAL AUDIOS AGAINST DEEPSPEECH.

Number	Original Audio (Recognized result of DeepSpeech)	Adversarial Audio (Recognized result of DeepSpeech)	SNR(dB)	Time(s)
1	Follow the instructions here	Read last sms from boss	18.07	685.14
2	One can imagine these two covered with sand running up the little street in the bright sunlight	Ask capital one to make a credit card payment	17.48	2205.63
3	Nature knows me as the wisest being in creation the sun said	Please restart the phone	19.07	2079.11
4	The boy reminded the old man that he had said something about hidden treasure	Clear SMS history from my phone	21.45	1879.29
5	It was dropping off in flakes and raining down on the sand	Remove all photos in my phone	20.71	1177.88

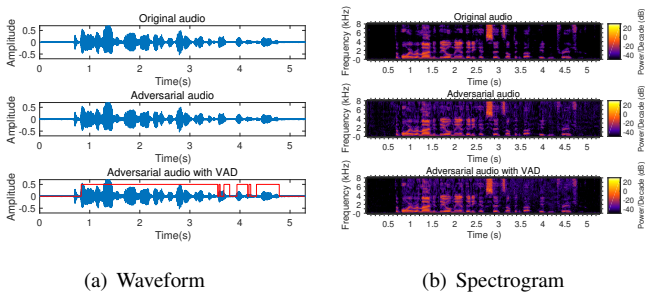


Fig. 3. Comparison of the waveform and spectrogram among the original audio, adversarial audio without VAD and adversarial audio with VAD. The original transcription is “the boy reminded the old man that he had said something about hidden treasure” while the adversarial transcription is “clear SMS history from my phone”.

E. Results and Analysis

Effectiveness and Efficiency. The main experimental results are shown in Table I, which summarizes the performance of SIRENATTACK on the two datasets. We randomly chose 200 instances from the Common Voice dataset and the VCTK Corpus as the original audios. The target commands were randomly chosen from a list of all the Google Now voice commands¹. We evaluate the average time of generating an adversarial audio, since it is important for an adversary to mount the attacks in realistic settings. From Table I, we can see that SIRENATTACK is very effective and efficient. It takes less than 1,600 seconds and 1,900 seconds on average to generate a successful adversarial audio (100% success rate) on the Common Voice dataset and the VCTK Corpus, respectively. Therefore, attackers may create plenty of adversarial audios in a short time. Furthermore, the adversarial audios have small distortion as shown in Table I. For instance, the average SNR of the generated adversarial audios on the Common Voice dataset is 18.72 dB, which means less than 2% distortion compared with the original audios.

To visualize the distortion, we plot the waveform and spectrogram of an example original audio and the corresponding adversarial audio in Fig. 3. The spectrogram of the original

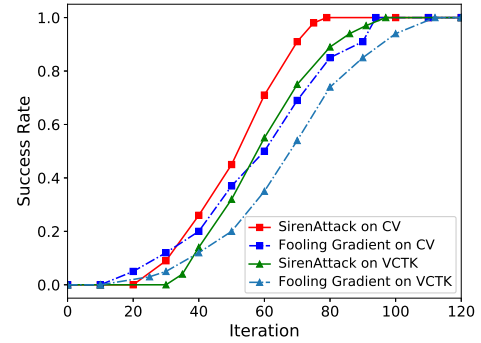


Fig. 4. The correlation between the iteration and the success rate.

audio and the corresponding adversarial audio in Fig. 3(b) are obtained from Short-Time Fourier Transform (STFT) of the waveform, where the horizontal axis represents time, the vertical axis represents frequency, and the color indicates the strength of energy. In fact, after the sound enters the human ear, the cochlea will also process the sound similar to STFT. Therefore, the sounds that people can distinguish often show specific patterns on the spectrogram. From Fig. 3(b), we can see that although the noise covers a broad spectrum, its energy is much lower than the vocal part. Hence, the noise in the adversarial audios is ignorable to humans and such attack is very stealthy.

Examples. Table II shows five examples in which the prediction results of adversarial audios are completely changed. For instance, the case of converting “follow the instructions here” to “read last sms from boss” can be used to steal users’ privacy information through their speech assistants. Therefore, this kind of attack can be leveraged by attackers to conduct malicious attacks on speech recognition systems. In addition, we observe a positive correlation between the length of the utterance and the time required to generate adversarial audios, which means generating longer adversarial audios may suffer from scaling issues to some extent.

Performance Comparison. We compare SIRENATTACK with Carlini’s attack [12] (which we refer to as “fooling gradient method” in the following) by showing the correlation between the iteration and the success rate in Fig. 4.

¹<https://www.greenbot.com/article/2359684/android/a-list-of-all-the-ok-google-voice-commands.html>

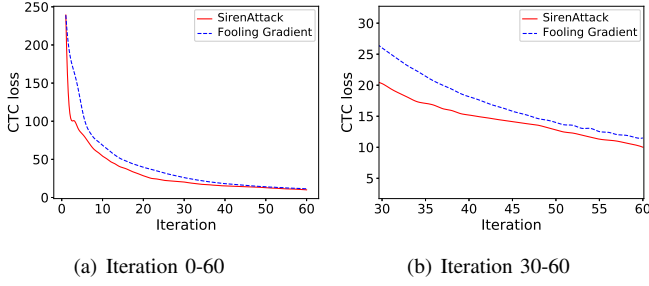


Fig. 5. CTC loss of the fooling gradient method and SIRENATTACK when converting the original audio to an adversarial audio.

Observe that SIRENATTACK on the Common Voice dataset reaches 100% success rate at iteration 79 while the fooling gradient method reaches 100% success rate at iteration 94. This suggests that SIRENATTACK is more efficient. Though the fooling gradient method finds the first adversarial audio faster than SIRENATTACK, its success rate increases slower. Specifically, taking converting “*the boy reminded the old man that he had said something about hidden treasure*” to “*clear SMS history from my phone*” as an example, we show the CTC loss of the fooling gradient method (blue dotted line) and SIRENATTACK (red solid line) in Fig. 5. We can see that the CTC loss is decreasing faster in SIRENATTACK than that in the fooling gradient method. This implies that SIRENATTACK chooses a direction that can find an adversarial audio faster than the fooling gradient method.

Improved Attack. To further improve the performance of SIRENATTACK, we use the Voice Activity Detection (VAD) Toolkit [32] to find the active part of audios and only add noise to this region. The results are also shown in Table I, from which we can see that VAD does increase the SNR of the generated adversarial audios and improve the efficiency of the generation process. For instance, the adversarial audios on the Common Voice dataset have an average SNR of 20.01 dB (18.72 dB without VAD) and an average generation time of 1201.66 seconds (1560.14 seconds without VAD) when applying VAD. Further, we compare the waveform and spectrogram of an example adversarial audio with and without VAD in Fig. 3, where the inactive voice parts occupy nearly one third of the original audio. Therefore, adding noise to the active parts of audios does increase the SNR of adversarial audios, i.e., generate better adversarial audios.

VI. BLACK-BOX ATTACK EVALUATION

A. Target Applications

1) *Speech Command Recognition*: In this scenario, we generated adversarial commands that can be recognized as target phrases for speech command recognition systems. For instance, we may start with an audio saying “yes”, which can be correctly recognized by the system. After applying the attack, the system will recognize the input as “no” while a human still clearly hears “yes”.

We used two datasets in this experiment: (i) *Speech Commands Dataset* [33]. This dataset consists of 65,000 audio files of 30 short words. Each file is a one-second audio of

TABLE III
SYNTHESIZED COMMANDS.

Number	Commands
1	Okay Google
2	Restart the phone
3	Flashlight on
4	Read email
5	Clear notification
6	Airplane mode on
7	Turn on wireless hot spot
8	Read last sms from boss
9	Open the front door
10	Turn off the light
11	Ask capital one to make a credit card payment

a single word like: “yes”, “no”, digits, and directions. (ii) *Synthesized Commands*. As shown in Table III, we synthesized 33,000 audio files of 11 long speech commands with 3,000 clips per label at different speeds and tones through several famous Text-to-Speech engine including Baidu, Google, Bing and IBM. The 11 commands are commonly used in daily life, therefore representing a variety of potential attacks against personal speech assistants.

The target victim models are: (i) *The CNN described in* [23]. This model, which is pre-trained by the TensorFlow team, is an efficient and light-weight keyword spotting model based on a CNN and achieves 96.10% classification accuracy on the Speech Commands Dataset. (ii) *Six State-of-the-art Speech Command Recognition Models*. We use VGG19 [34], DenseNet [35], ResNet18 [36], ResNeXt [37], WideResNet18 [38] and DPN-92 [39] as the target victim models. These models are well known for their good classification performance on image data. In addition, they have good performance in the TensorFlow Speech Recognition Challenge². Therefore, we modify them to adapt to the spectrogram input.

2) *Speaker Recognition*: Speaker recognition is the identification of a person from the characteristics of voices [40], which can be used to authenticate the identity of a speaker as part of a security process. We simplify the speaker recognition task in our experiment by limiting it to a ten-class classification problem, which is reasonable and common [17]. Then, we target the same kinds of models used in the speech command recognition task. Further, we conduct the adversarial attack using the IEMOCAP dataset [41], which consists of ten speakers (five female, five male) and is a commonly used dataset in speech paralinguistic research [17].

3) *Sound Event Classification*: The goal of sound event classification is to give a predefined label to the sound event (e.g., “dogbark”, “siren”) within an audio signal. It has numerous applications, including audio surveillance systems [42], [1], hearing aids [43], smart room monitoring [44], and pornographic content detection [45]. In this scenario, our goal is to fool the sound event classification systems into producing an incorrect target prediction. For instance, we may start with an audio correctly recognized as “gunshot”, a dangerous event that may cause attention from monitors. However, the system will classify the corresponding adversarial audio as a normal

²<https://www.kaggle.com/c/tensorflow-speech-recognition-challenge>

TABLE IV
PERFORMANCE OF THE BLACK-BOX ATTACK ON SPEECH COMMAND RECOGNITION AND SPEAKER RECOGNITION.

Model/Dataset	Speech Commands				Synthesized Commands				IEMOCAP			
	Accuracy	Success Rate	SNR(dB)	Time(s)	Accuracy	Success Rate	SNR(dB)	Time(s)	Accuracy	Success Rate	SNR(dB)	Time(s)
CNN	96.10%	95.25%	22.36	100.69	-	-	-	-	-	-	-	-
VGG19	91.39%	88.10%	18.22	332.26	93.12%	93.75%	17.04	420.89	85.01%	91.65%	16.33	376.40
DenseNet	94.93%	86.90%	15.34	458.13	93.34%	89.25%	15.13	602.81	84.28%	94.65%	15.28	572.23
ResNet18	92.06%	87.35%	15.87	340.31	93.90%	90.15%	17.28	381.27	86.37%	99.45%	23.12	386.15
ResNeXt	94.28%	90.05%	17.03	317.92	94.80%	92.60%	18.49	458.44	87.66%	95.60%	20.87	420.37
WideResNet18	90.80%	89.25%	17.57	368.29	92.68%	91.05%	17.23	403.76	92.41%	93.95%	17.06	393.56
DPN92	95.20%	83.60%	14.04	462.58	96.81%	90.55%	14.77	587.80	86.93%	92.80%	15.51	564.98

event (e.g., “dogbark”), while human beings can still hear “gunshot”.

We employ three large-scale sound event datasets to evaluate SIRENATTACK: (i) *AudioSet* [46]. It has 632 sound event classes covering a wide range of everyday environmental sounds. (ii) *ESC-50* [47]. It consists of 2,000 5-second-long environmental audio recordings organized into 50 classes with 40 audios per class. (iii) *UrbanSound8K* [48]. It contains 8,732 labeled sound excerpts (no more than four seconds for each) of urban sounds from ten classes.

As for the victim models, we use the YouTube-8M starter code³ to train three victim models including the *Logistic Model (LM)*, the *Mixture of Experts (MoE) model* and the *Frame-Level Logistic Model (FLLM)*, according to the instruction of *AudioSet*⁴.

4) *Music Genre Classification*: The goal of music genre classification is to classify music into various genres like “classical”, “jazz”, “rock”, etc. If content-based music recommendation is polluted with adversarial audios, users may receive recommendations that is not in line with their taste or even contain terrorism and pornographic content. This can be maliciously leveraged by competitors of the music recommendation system. In this scenario, we evaluate SIRENATTACK on *GTZAN* [49], which consists of 1,000 30-second music recording excerpts of ten genres, and is the most-used public dataset in the Music Information Retrieval (MIR) research. The target state-of-the-art models are *ConvNet* [50] and *ConvRNN* [51].

B. Implementation

The implementation details of our black-box attack are almost the same as that in the white-box attack. One difference is that we need to train some targeted models due to the lack of pre-trained models. Therefore, except for the CNN model, all models were trained in a hold-out test strategy, i.e., 80%, 10%, 10% of the data was used for training, validation and testing, respectively. Hyper-parameters were tuned only on the validation set, and the audios used to conduct attacks were chosen from the testing set. We emphasize again that we do not know the training data about the black-box models when conducting attacks, and we train the victim models ourselves only because there are few public available victim models.

C. Evaluation Results

Attacks on Speech Command Recognition Systems. We selected 2,000 audio clips from the Speech Commands Dataset

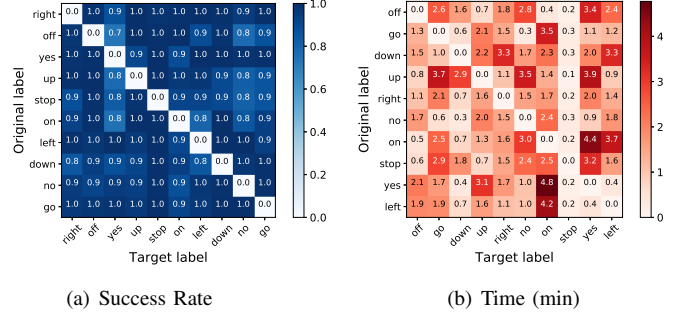


Fig. 6. Performance of SIRENATTACK for every $\{source, target\}$ pair on the Speech Commands Dataset against the CNN model.

with 200 clips per label and generated nine targeted adversarial audios for each audio file. Notice that the CNN model was pre-trained on the Speech Commands Dataset. Hence, we did not evaluate it on the Synthesized Commands dataset. The attack results are shown in Table IV with $\delta = 800$ and $epoch_{max} = 300$, including the models’ accuracy on the original dataset, the success rate of SIRENATTACK, the SNR of the generated adversarial audios and the average time to generate an adversarial audio. Figs. 6(a) and 6(b) show the pair-to-pair success rate and the average time to generate an adversarial audio of SIRENATTACK on the Speech Commands Dataset (resp., the Synthesized Commands dataset). We can observe the following from Table IV and Fig. 6.

- From Table IV and Fig. 6(a), we can see that SIRENATTACK is effective when against all the target models, even when the models have high performance on the legitimate datasets. For instance, SIRENATTACK has 95.25% success rate on the Speech Commands Dataset when against the CNN model and has 93.75% success rate on the Synthesized Commands dataset when against the VGG19 model. Therefore, SIRENATTACK is sufficiently effective to be used by an adversary. In addition, we notice that certain transformations seem to be easier than others. For instance, the conversion from “yes” to “stop” can be done in 10 iterations, while the conversion from “stop” to “yes” takes 160 iterations. We conjecture that this might result from the victim model’s different prediction robustness among different categories. Another interesting observation is that some intermediate adversarial audios appear in the attacking process. For example, when we convert “restart the phone” to “flashlight on”, the transcription result first changes to “clear notification” and then changes to “flashlight on”.
- From Table IV, the average generation time of an adversarial audio is very short. For instance, the average generation time

³<https://github.com/google/youtube-8m>

⁴<https://research.google.com/audioset>

TABLE V
PERFORMANCE OF THE BLACK-BOX ATTACK ON EVENT SOUND CLASSIFICATION.

Feature Type	Target Model	ESC-50				UrbanSound8K				AudioSet			
		Accuracy	Success Rate	SNR(dB)	Time(s)	Accuracy	Success Rate	SNR(dB)	Time(s)	Accuracy	Success Rate	SNR(dB)	Time(s)
Audio-level	LM	99.95%	85.33%	18.28	269.34	88.04%	84.00%	14.51	598.62	85.49%	92.67%	18.04	283.57
	MoE	99.82%	84.00%	17.91	283.41	92.36%	82.00%	12.80	609.63	88.32%	91.33%	13.58	287.19
Frame-level	FLLM	87.71%	80.67%	18.11	328.75	84.28%	80.00%	19.73	487.04	81.02%	90.67%	21.62	403.41

TABLE VI
EXAMPLES OF THE BLACK-BOX ATTACK ON EVENT SOUND CLASSIFICATION.

Feature Type	Target Model	ESC-50				UrbanSound8K				AudioSet			
		Original	Target	SNR(dB)	Time(s)	Original	Target	SNR(dB)	Time(s)	Original	Target	SNR(dB)	Time(s)
Audio-level	LM	Breaking	Crickets	29.51	166.02	Gunshot	Dog bark	17.55	223.76	Breaking	Clock alarm	22.19	214.93
		Breaking	Crickets	25.23	167.11	Gunshot	Dog bark	17.57	224.75	Breaking	Clock alarm	11.51	219.16
	MoE	Siren	Wind	12.09	222.51	Siren	Street music	16.26	735.41	Gunshot	Children playing	12.88	224.84
		Siren	Wind	12.21	219.08	Siren	Street music	14.03	745.96	Gunshot	Children playing	11.34	229.21
Frame-level	FLLM	Breaking	Crickets	15.85	289.61	Gunshot	Dog bark	20.05	306.03	Breaking	Frog	28.60	340.73
		Siren	Crickets	15.81	308.62	Siren	Street music	20.99	385.69	Gunshot	Dog bark	17.10	309.11

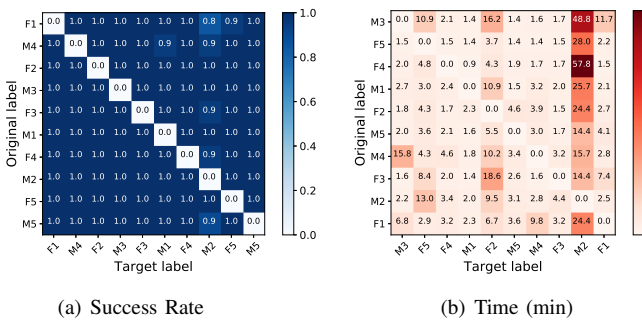


Fig. 7. Performance of SIRENATTACK for every $\{source, target\}$ pair on the IEMOCAP dataset against the ResNet18 model.

of the Speech Commands Dataset when against the CNN model is 100.69 seconds. In addition, from Fig. 6(b) we can see that all of the adversarial audios can be generated in less than 5 minutes, and some $\{source, target\}$ pairs like $\{go, stop\}$ can be done within one minute. Therefore, SIRENATTACK is very efficient in practice.

- From Table IV, we can also see the noise is slight, e.g., the SNR of the adversarial audios ranges from 14 dB to 22 dB on both datasets when against the target models. This implies the noise in the adversarial audios is less than 3%.

Attacks on Speaker Recognition Systems. In this evaluation, we used the 1,000 audio clips from the IEMOCAP dataset with 100 clips per speaker and generated nine targeted adversarial audios for each audio file. The attack results are shown in Table IV with $\delta = 800$ and $epoch_{max} = 300$. Fig. 7 further shows the pair-to-pair success rate of SIRENATTACK and the average time to generate an adversarial audio. Similar to the attack on the speech command recognition systems, SIRENATTACK is also very effective against all the target models. For instance, SIRENATTACK has 99.45% success rate when against the ResNet18 model. In addition, SIRENATTACK is also efficient in this task, e.g., the average generation time of an adversarial audio of the IEMOCAP dataset against the VGG19 model is 376.40 seconds.

Attacks on Sound Event Classification Systems. In this evaluation, we trained the LM model and the MoE model on the audio-level features, which were extracted by the method

TABLE VII
MUSIC GENRE CLASSIFICATION RESULTS.

	Accuracy	Success Rate	SNR(dB)	Time(s)
ConvNet	89.10%	89.30%	15.39	452.21
ConvRNN	91.40%	91.20%	17.24	520.61

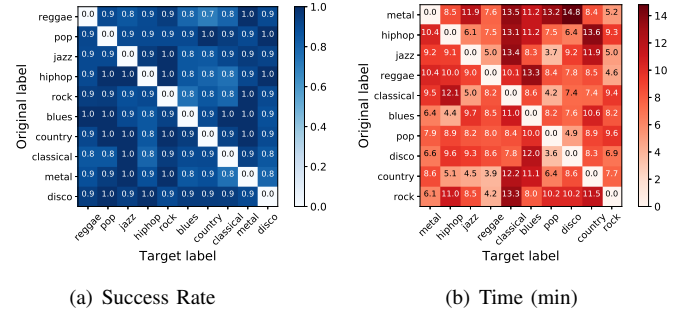


Fig. 8. Performance of SIRENATTACK for every $\{source, target\}$ pair on the GTZAN dataset against the ConvRNN model.

in [21]. In addition, we trained the FLLM model on the frame-level features, which were extracted by the VGGish model in [22]. For ESC-50 and UrbanSound8K datasets, the victim models were trained on their own datasets respectively; for AudioSet, we took the pre-trained models on UrbanSound8K as the victim models. To demonstrate that SIRENATTACK can convert the threatening events to normal events, we randomly picked 150 $\{source, target\}$ pairs matching the $\{threatening\ event, normal\ event\}$ pattern from each of the three datasets to evaluate SIRENATTACK. The results are shown in Table V, from which we can see that SIRENATTACK is also effective and efficient against the target models. For instance, SIRENATTACK has 92.67% success rate on the AudioSet dataset when against the LM model with the average generation time of 283.57 seconds. Table VI demonstrates some examples of SIRENATTACK to convert threatening events to normal events, e.g., SIRENATTACK can convert the threatening event “gunshot” to the normal event “dogbark”, which can be used as an attack on acoustic surveillance systems.

Attacks on Music Genre Classification Systems. In this evaluation, we used 1,000 music clips from the GTZAN

TABLE VIII
TRANSFERABILITY EVALUATION RESULTS.

	Sphinx	Google	Bing	Houndify	Wit.ai	IBM
Success Rate	39.60%	10.00%	14.00%	12.80%	21.20%	20.40%

TABLE IX
EXAMPLE RESULTS OF TRANSFERABILITY EVALUATION.

Number	Original Text	Advesarial Text	ASR Platforms	Results
1	stop	no	Sphinx	no
2	off	on	IBM	on
3	down	no	Wit.ai, Bing	no
4	go	no	Wit.ai	no
5	go	yes	Sphinx	yes
6	left	yes	Wit.ai, IBM	yeah
7	on	right	Wit.ai	alright
8	right	on	Google, Bing	play
9	right	down	Google, Bing	play
10	off	no	Bing	call
11	on	stop	Bing	call
12	on	up	Wit.ai	okay
13	stop	off	Wit.ai	the
14	down	up	Bing	phone
15	stop	go	Wit.ai	tell

dataset with 100 clips per label and generated nine targeted adversarial audios for each music file. The attack results are shown in Table VII. In addition, Fig. 8 shows the pair-to-pair success rate of SIRENATTACK and the average time to generate an adversarial audio. From Table VII and Fig. 8, we can see that SIRENATTACK is effective when against both two target models. For instance, SIRENATTACK has 91.20% success rate when against the ConvRNN model.

VII. FURTHER ANALYSIS

A. Perturbation Analysis

Now, we evaluate the impact of the noise scale δ and $epoch_{max}$ on the effectiveness and efficiency of generating adversarial audios. Specifically, we generated adversarial audios on the Speech Commands Dataset with different bound values of noises as well as $epoch_{max} = 100, 200, 300$. The targeted model is CNN. The success rate and required time are shown in Fig. 9, from which we can see that the trend of success rate is generally consistent with the noise scale. For instance, when $epoch_{max} = 100$, SIRENATTACK has 82% success rate with $\delta = 100$ while having 90% success rate with $\delta = 1000$. This implies that attackers can use larger noise scale to improve the success rate of their attacks. On the other hand, a larger δ also implies lower utility, i.e., human may notice the changes of the audio. From Fig. 9, we can also see that when $\delta = 800$, all the three time curves reach the minimum value. In addition, when $epoch_{max} = 300$, the overall success rate is higher than that of $epoch_{max} = 100, 200$. These findings help us derive better parameter settings. Hence, we use $\delta = 800, epoch_{max} = 300$ in our evaluation.

B. Transferability Evaluation

Previous studies have shown that adversarial images generated for one model can be misclassified by other models, even when they have different architectures [9], i.e., adversarial images exhibit transferability, which can be used to conduct black-box attacks. Therefore, we are interested in (i)

whether the transferability also exists in adversarial audios, and (ii) whether this property can be used to conduct black-box attacks. Specifically, we used 500 adversarial audios that are generated from the Speech Commands Dataset with the target model VGG19 to conduct proof-of-concept attack on several famous ASR platforms, including Sphinx, Google Cloud Speech Recognition, Microsoft Bing Voice Recognition, Houndify, Wit.ai and IBM Speech-to-Text. Note that we do not directly conduct black-box attacks on these ASR platforms since they are all recognition-oriented models which do not give any information except for the final transcription. In this scene, it is very difficult, if possible, to directly conduct black-box attacks on these models while guaranteeing the added noise is human-imperceptible.

The evaluation results are shown in Table VIII, from which we can see that the adversarial audios generated by SIRENATTACK can also be misinterpreted as the target text by the target ASR platforms to some extent. For instance, SIRENATTACK achieves 39.60% success rate on the Sphinx platform. This implies that the adversarial audios generated by SIRENATTACK can be used to mount targeted black-box attacks to other ASR platforms. Line 1-7 in Table IX show some examples that are successfully transferred against other ASR platforms. In addition, line 8-15 in Table IX show some additional misclassification results, which imply that the adversarial audios generated by SIRENATTACK may pose threats to people's privacy when being concatenated with other words, such as "call 911", "okay Google", "restart the phone", and "tell me the phone number of Jack".

C. Human Perceptual Study

To quantify the perceptual realism of the adversarial audios generated by SIRENATTACK, we also perform a user study with human participants on Amazon Mechanical Turk (MTurk). Before the study, we consulted with the IRB office and this study was approved and we did not collect any other information of participants except for necessary result data.

In the study, we recruited 200 native English speakers whose age ranges from 18 to 40 to participate in our survey. Each participant was asked to listen to 20 legitimate audios and 20 adversarial audios generated from the Speech Command dataset with CNN as the target model in a quiet environment. During each trial, participants are given unlimited time to replay audios and make their decisions. For each audio, a series of questions need to be answered, i.e., (1) what they heard from this audio (choose one option from the given ten options, i.e., *stop*, *go*, *yes*, *no*, *left*, *right*, *off*, *on*, *up*, and *down*); (2) whether they heard anything abnormal than a regular command (the four options are *no*, *not sure*, *a little noisy*, and *noisy*); (3) if choosing *a little noisy* or *noisy* option in (2), where they believe the noise comes from (the three options are *the device (speaker, radio, etc.)*, *the sample itself*, *other*).

After examining the results, we find that 93.50% legitimate audios can be recognized correctly while 92.00% adversarial audios can be recognized as their original labels. None of the adversarial audios is classified as its adversarial label.

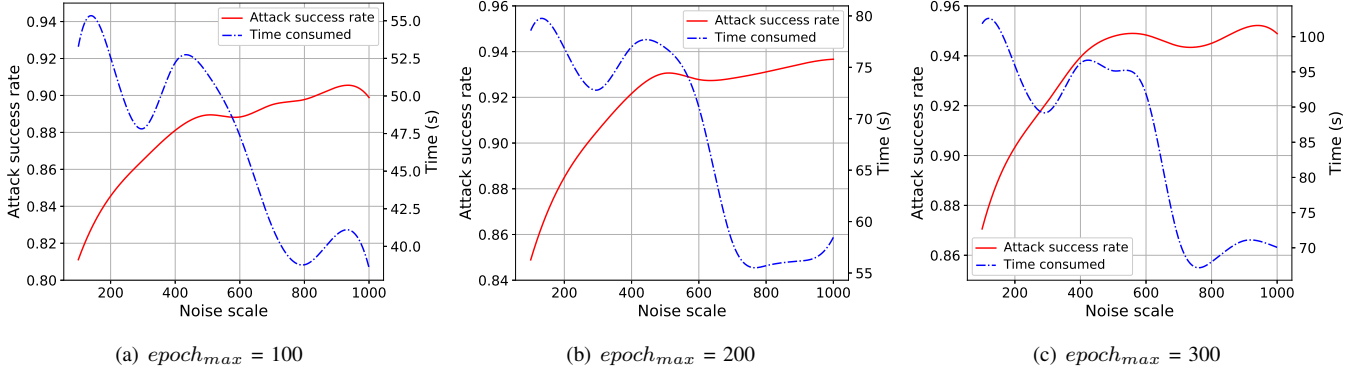


Fig. 9. The success rate and required time with different noise scale for Speech Commands Dataset.

TABLE X
ADVERSARIAL TRAINING AS A DEFENSE STRATEGY.

Dataset	# of Legitimate Audio	# of Adversarial Audio	Target Model	Accuracy	Success rate
Speech Commands	20,000	1,000	CNN	94.22%	17.90%
Synthesized Commands	33,000	1,100	VGG19	90.84%	23.30%
IEMOCAP	10,000	1,000	ResNet18	85.79%	20.50%

This indicates that the generated adversarial audios have little impact on human perception. What's more, 38.5% of participants think the adversarial audios are a little noisy and only 4.5% participants think the noise are from the samples themselves. Furthermore, 10.5% of participants think the adversarial audios are noisy and only 2.5% participants think the noise are from the samples themselves. This implies that SIRENATTACK is stealthy.

VIII. POTENTIAL DEFENSES

As there are few defense methods for adversarial audio attacks to the best of our knowledge, we conduct a preliminary exploration of potential defense schemes. By default, all the adversarial audios are generated using our black-box attack and we use the same implementation and evaluation settings as that in Section VI.

Adversarial Training. Adversarial training means training a new model with both legitimate and adversarial examples. We show the performance of this scheme along with detailed settings in Table X, where the accuracy means the prediction accuracy of the new models on the legitimate audios. From Table X, we can see that the success rate of adversarial audios decreases while the models' performance on clean samples does not change too much. However, a limitation of adversarial training is that it needs to know the details of the attack strategy and to have sufficient adversarial audios for training. In practice, however, attackers usually do not make their approaches or adversarial audios public. Further, they can change the parameters of the attack frequently (e.g., the perturbation factor [17]) to evade the defense. Therefore, adversarial training is limited in defending unknown adversarial attacks.

Audio Downsampling. The second potential defense method is to reduce the sampling rate of the input audio x . We

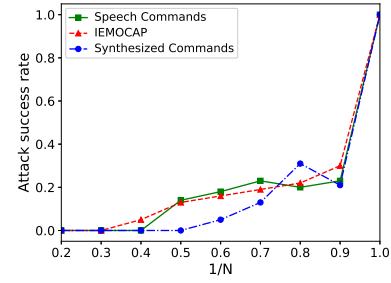


Fig. 10. Results of the audio downsampling defense on three datasets.

TABLE XI
MOVING AVERAGE FILTERING AS A DEFENSE STRATEGY.

Dataset	# of Adversarial Audios	Model	k	Success Rate
Speech Commands	1,000	CNN	5	20.60%
Synthesized Commands	1,100	VGG19	5	4.90%
IEMOCAP	1,000	ResNet18	5	28.50%

denote the downsampled audio as $D(x)$, and its recognition result is referred to as y_D . When we feed an acoustic system with an adversarial audio x_{adv} with label y_{adv} , if $y_D \neq y_{adv}$, x_{adv} will be determined as successfully defended. The results of this defense are shown in Fig. 10, where x-axis means that the original sampling rate is N times of the downsampled rate. From Fig. 10, we can see that this defense can reduce the success rate of SIRENATTACK. For instance, when $\frac{1}{N} = 0.8$, the success rate of SIRENATTACK is 20%. However, according to the Nyquist sampling theorem [52], this method would cause distortion when the sampling rate is below twice of the highest frequency of the original audios.

Moving Average Filtering (MAF). Now, we use a sliding window with a fixed length for MAF to reduce the impact of adversarial noise. Specifically, for a sampling point x_i , we consider the $k - 1$ points before and after it as local reference points and replace x_i by the average value of its reference points. The results are shown in Table XI, from which we can see that MAF can reduce the success rate of SIRENATTACK. For instance, when $k = 5$, the success rate of SIRENATTACK decreases to 20.60% on the Speech Commands Dataset. However, MAF might reduce the quality of audios, thus having a negative impact on the models' performance.

IX. DISCUSSION

Universal Adversarial Perturbations. In SIRENATTACK, we need to find special adversarial noises for each audio. In the image domain, it is possible to construct a single perturbation δ that can lead to misclassifications for various images when applied to them [53]. This kind of attack would also be incredibly powerful to audio if it is possible. We take this as a future research direction.

More Threatening Attacks. In our evaluation, we assume that an attacker can directly feed the audio files to the victim model. This is realistic since many speech content monitors can directly censor the raw audios. Therefore, SIRENATTACK would indeed pose a threat on the web environment. However, a more powerful attack scene is “over-the-air”, where an attacker calculates and plays an adversarial noise signal $\delta(t)$ according to a legitimate audio $x(t)$ in real time, so that the superimposed audio $x(t) + \delta(t)$ would be interpreted as a malicious command. In addition, SIRENATTACK can be combined with other attacks to form more dangerous ones, e.g., combine SIRENATTACK with GVS-Attacks [18] so that the malware can replay an adversarial audio when it finds an opportunity. We also plan to study this attack in the future.

Other Limitations and Future Work. Although the generated adversarial audios can be misclassified by popular ASR platforms, the success rate is not very high. Therefore, how to generate adversarial audios with better transferability deserves further research. Furthermore, developing effective and robust defense schemes is also a promising future work.

X. CONCLUSION

In this paper, we study targeted adversarial attacks against acoustic systems in both white-box and black-box settings. To the best of our knowledge, this is the first systematical study on generating adversarial audios for various acoustic systems, including speech recognition, speaker recognition, sound event classification and music genre classification. Extensive experimental results show that SIRENATTACK is effective and efficient, and has potential threats on many real applications. We also discuss three potential approaches to defend against such attacks.

REFERENCES

- [1] “Shotspotter,” <http://www.shotspotter.com/>.
- [2] P. K. Atrey, N. C. Maddage, and M. S. Kankanalli, “Audio based event detection for multimedia surveillance,” in *ICASSP*, vol. 5. IEEE, 2006, pp. 813–816.
- [3] A. Graves, S. Fernández, F. Gomez, and J. Schmidhuber, “Connectionist temporal classification: labelling unsegmented sequence data with recurrent neural networks,” in *ICML*. ACM, 2006, pp. 369–376.
- [4] D. Bahdanau, J. Chorowski, D. Serdyuk, P. Brakel, and Y. Bengio, “End-to-end attention-based large vocabulary speech recognition,” in *ICASSP*, 2016, pp. 4945–4949.
- [5] N. Carlini, P. Mishra, T. Vaidya, Y. Zhang, M. Sherr, C. Shields, D. Wagner, and W. Zhou, “Hidden voice commands,” in *USENIX Security*, 2016, pp. 513–530.
- [6] G. Zhang, C. Yan, X. Ji, T. Zhang, T. Zhang, and W. Xu, “Dolphinattack: Inaudible voice commands,” in *CCS*. ACM, 2017, pp. 103–117.
- [7] X. Lei, G.-H. Tu, A. X. Liu, C.-Y. Li, and T. Xie, “The insecurity of home digital voice assistants-amazon alexa as a case study,” *arXiv preprint arXiv:1712.03327*, 2017.
- [8] H. Feng, K. Fawaz, and K. G. Shin, “Continuous authentication for voice assistants,” in *MobiCom*, 2017, pp. 343–355.
- [9] I. J. Goodfellow, J. Shlens, and C. Szegedy, “Explaining and harnessing adversarial examples,” in *ICLR*, 2015.
- [10] C. Szegedy, “Intriguing properties of neural networks,” in *ICLR*, 2014.
- [11] A. Kurakin, I. Goodfellow, and S. Bengio, “Adversarial examples in the physical world,” in *ICLR Workshop*, 2017.
- [12] N. Carlini and D. Wagner, “Audio adversarial examples: Targeted attacks on speech-to-text,” *arXiv preprint arXiv:1801.01944*, 2018.
- [13] M. Cisse, Y. Adi, N. Neverova, and J. Keshet, “Houdini: Fooling deep structured prediction models,” *arXiv preprint arXiv:1707.05373*, 2017.
- [14] D. Iter, J. Huang, and M. Jermann, “Generating adversarial examples for speech recognition,” http://web.stanford.edu/class/cs224s/reports/Dan_Iter.pdf.
- [15] X. Yuan, Y. Chen, Y. Zhao, Y. Long, X. Liu, K. Chen, S. Zhang, H. Huang, X. Wang, and C. A. Gunter, “Commandersong: A systematic approach for practical adversarial voice recognition,” in *USENIX Security*, 2018, pp. 49–64.
- [16] 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.
- [17] Y. Gong and C. Poellabauer, “Crafting adversarial examples for speech paralinguistics applications,” *arXiv preprint arXiv:1711.03280*, 2017.
- [18] W. Diao, X. Liu, Z. Zhou, and K. Zhang, “Your voice assistant is mine: How to abuse speakers to steal information and control your phone,” in *SPSM*. ACM, 2014, pp. 63–74.
- [19] G. Petracca, Y. Sun, T. Jaeger, and A. Atamli, “Audroid: Preventing attacks on audio channels in mobile devices,” in *ACM ACSAC*, 2015, pp. 181–190.
- [20] D. O’Shaughnessy, “Automatic speech recognition: History, methods and challenges,” *Pattern Recognition*, vol. 41, no. 10, pp. 2965–2979, 2008.
- [21] A. Kumar, M. Khadkevich, and C. Fugen, “Knowledge transfer from weakly labeled audio using convolutional neural network for sound events and scenes,” *arXiv preprint arXiv:1711.01369*, 2017.
- [22] S. Hershey, S. Chaudhuri, D. P. Ellis, J. F. Gemmeke, A. Jansen, R. C. Moore, M. Plakal, D. Platt, R. A. Saurous, B. Seybold *et al.*, “Cnn architectures for large-scale audio classification,” in *ICASSP*. IEEE, 2017, pp. 131–135.
- [23] T. N. Sainath and C. Parada, “Convolutional neural networks for small-footprint keyword spotting,” in *INTERSPEECH*, 2015, pp. 1478–1482.
- [24] A. Nguyen, J. Yosinski, and J. Clune, “Deep neural networks are easily fooled: High confidence predictions for unrecognizable images,” in *CVPR*, 2015, pp. 427–436.
- [25] R. Eberhart and J. Kennedy, “A new optimizer using particle swarm theory,” in *MHS*. IEEE, 1995, pp. 39–43.
- [26] P.-Y. Chen, H. Zhang, Y. Sharma, J. Yi, and C.-J. Hsieh, “Zoo: Zeroth order optimization based black-box attacks to deep neural networks without training substitute models,” in *ACM AISec*, 2017, pp. 15–26.
- [27] “Common voice dataset,” <https://voice.mozilla.org/zh-CN/data>.
- [28] C. Veaux, J. Yamagishi, K. MacDonald *et al.*, “Cstr vetk corpus: English multi-speaker corpus for cstr voice cloning toolkit,” 2017.
- [29] E. Battenberg, J. Chen, R. Child, A. Coates, Y. Gaur, Y. Li, H. Liu, S. Satheesh, D. Seetapun, A. Sriram *et al.*, “Exploring neural transducers for end-to-end speech recognition,” *arXiv preprint arXiv:1707.07413*, 2017.
- [30] P. K. Dhar and T. Shimamura, *Advances in audio watermarking based on singular value decomposition*. Springer, 2015.
- [31] C. Kereliuk, B. L. Sturm, and J. Larsen, “Deep learning and music adversaries,” *IEEE Transactions on Multimedia*, vol. 17, no. 11, pp. 2059–2071, 2015.
- [32] J. Kim and M. Hahn, “Voice activity detection using an adaptive context attention model,” *IEEE Signal Processing Letters*, 2018.
- [33] P. Warden, “Speech commands: A public dataset for single-word speech recognition,” *Dataset available from http://download.tensorflow.org/data/speech_commands_v0.01.tar.gz*, 2017.
- [34] K. Simonyan and A. Zisserman, “Very deep convolutional networks for large-scale image recognition,” *ICLR*, 2015.
- [35] G. Huang, Z. Liu, K. Q. Weinberger, and L. van der Maaten, “Densely connected convolutional networks,” in *CVPR*, 2017.
- [36] K. He, X. Zhang, S. Ren, and J. Sun, “Deep residual learning for image recognition,” in *CVPR*, 2016, pp. 770–778.
- [37] S. Xie, R. Girshick, P. Dollár, Z. Tu, and K. He, “Aggregated residual transformations for deep neural networks,” in *CVPR*, 2017, pp. 5987–5995.
- [38] S. Zagoruyko and N. Komodakis, “Wide residual networks,” *arXiv preprint arXiv:1605.07146*, 2016.
- [39] Y. Chen, J. Li, H. Xiao, X. Jin, S. Yan, and J. Feng, “Dual path networks,” in *NIPS*, 2017, pp. 4470–4478.

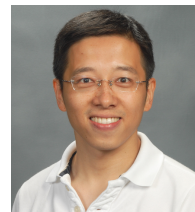
- [40] A. Poddar, M. Sahidullah, and G. Saha, "Speaker verification with short utterances: a review of challenges, trends and opportunities," *IET Biometrics*, 2017.
- [41] C. Busso, M. Bulut, C.-C. Lee, A. Kazemzadeh, E. Mower, S. Kim, J. N. Chang, S. Lee, and S. S. Narayanan, "Iemocap: Interactive emotional dyadic motion capture database," *Language resources and evaluation*, vol. 42, no. 4, p. 335, 2008.
- [42] K. Łopata, P. Zwan, and A. Czyżewski, "Dangerous sound event recognition using support vector machine classifiers," in *Advances in Multimedia and Network Information System Technologies*. Springer, 2010, pp. 49–57.
- [43] E. Alexandre, L. Cuadra, M. Rosa, and F. Lopez-Ferreras, "Feature selection for sound classification in hearing aids through restricted search driven by genetic algorithms," *IEEE Transactions on Audio, Speech, and Language Processing*, vol. 15, no. 8, pp. 2249–2256, 2007.
- [44] M. Vacher, J.-F. Serignat, and S. Chaillol, "Sound classification in a smart room environment: an approach using gmm and hmm methods," in *SpeD*, vol. 1. Publishing House of the Romanian Academy, 2007, pp. 135–146.
- [45] J.-D. Lim, J.-N. Kim, Y.-G. Jung, Y.-D. Yoon, and C.-H. Lee, "Improving performance of x-rated video classification with the optimized repeated curve-like spectrum feature and the skip-and-analysis processing," *Multimedia Tools and Applications*, vol. 71, no. 2, pp. 717–740, 2014.
- [46] J. F. Gemmeke, D. P. Ellis, D. Freedman, A. Jansen, W. Lawrence, R. C. Moore, M. Plakal, and M. Ritter, "Audio set: An ontology and human-labeled dataset for audio events," in *ICASSP*, 2017, pp. 776–780.
- [47] K. J. Piczak, "Esc: Dataset for environmental sound classification," in *ACM MM*, 2015, pp. 1015–1018.
- [48] J. Salamon, C. Jacoby, and J. P. Bello, "A dataset and taxonomy for urban sound research," in *ACM MM*, 2014, pp. 1041–1044.
- [49] G. Tzanetakis and P. Cook, "Musical genre classification of audio signals," *IEEE Transactions on Speech and Audio Processing*, vol. 10, no. 5, pp. 293–302, 2002.
- [50] K. Choi, G. Fazekas, and M. Sandler, "Automatic tagging using deep convolutional neural networks," in *ISMIR*, 2016.
- [51] K. Choi, G. Fazekas, M. Sandler, and K. Cho, "Convolutional recurrent neural networks for music classification," in *ICASSP*, 2017, pp. 2392–2396.
- [52] Wikipedia, "Wikipedia, nyquist-shannon sampling theorem," 2018. [Online]. Available: https://en.wikipedia.org/wiki/Nyquist%E2%80%9393Shannon_sampling_theorem
- [53] S.-M. Moosavi-Dezfooli, A. Fawzi, O. Fawzi, and P. Frossard, "Universal adversarial perturbations," in *CVPR*. IEEE, 2017, pp. 86–94.



Jinfeng Li is currently a postgraduate student in College of Computer Science and Technology at Zhejiang University, P.R. China, under the supervision of Prof Shouling Ji. He received his BS degree from Wuhan University of Technology in 2017. His research interests include Big Data Driven Security, Adversarial Learning, and AI Security.



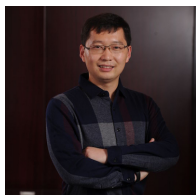
Qinchun Gu received an MS degree in electrical and computer engineering from Georgia Institute of Technology (Georgia Tech). He is currently a PhD student in the School of Electrical and Computer Engineering at Georgia Tech, and a graduate research assistant of the Communications Assurance and Performance (CAP) group. His research primarily focuses on the security for cyber-physical systems. Contact him at qgu7@gatech.edu.



Ting Wang is an assistant professor at Lehigh University. Prior to joining Lehigh, he was a Research Staff Member at IBM Thomas J. Watson Research Center. He conducts research at the intersection of machine learning, computational privacy, and cyber-security. His current work focuses on enforcing security assurance for machine learning systems. He obtained his doctoral degree from Georgia Institute of Technology. He is a member of IEEE and ACM.



Tianyu Du is currently a Ph.D. student in College of Computer Science and Technology at Zhejiang University, P.R. China, under the supervision of Prof Shouling Ji. She received her BS degree from Xiamen University in 2017. Her research interests include Big Data Driven Security, Adversarial Learning, and AI Security.



Shouling Ji is a ZJU 100-Young Professor in the College of Computer Science and Technology at Zhejiang University and a Research Faculty in the School of Electrical and Computer Engineering at Georgia Institute of Technology. He received a Ph.D. in Electrical and Computer Engineering from Georgia Institute of Technology and a Ph.D. in Computer Science from Georgia State University. His current research interests include Big Data Security and Privacy, Big Data Driven Security and Privacy, and Adversarial Learning. He also has interests in Graph

Theory and Algorithms and Wireless Networks. He is a member of IEEE and ACM and was the Membership Chair of the IEEE Student Branch at Georgia State (2012-2013).



Raheem Beyah is the Motorola Foundation Professor and Associate Chair in the School of Electrical and Computer Engineering at Georgia Tech, where he leads the Communications Assurance and Performance Group (CAP) and is a member of the Communications Systems Center (CSC). Prior to returning to Georgia Tech, Dr. Beyah was an Assistant Professor in the Department of Computer Science at Georgia State University, a research faculty member with the Georgia Tech CSC, and a consultant in Andersen Consultings (now Accenture) Network Solutions Group. He received his Bachelor of Science in Electrical Engineering from North Carolina A&T State University in 1998. He received his Masters and Ph.D. in Electrical and Computer Engineering from Georgia Tech in 1999 and 2003, respectively. His research interests include network security, wireless networks, network traffic characterization and performance, and critical infrastructure security. He received the National Science Foundation CAREER award in 2009 and was selected for DARPA's Computer Science Study Panel in 2010. He is a member of AAAS and ASEE, is a lifetime member of NSBE, and is a senior member of ACM and IEEE.