

Design and implementation of covert communication based on intelligent voice system

Song Lipeng¹⁾ Zhang Ruichen¹⁾ Liu Wenxiang¹⁾ Wang Shuo¹⁾ Zhang Ye¹⁾ Ren Guanyu¹⁾

¹⁾(School of Mechanical, Electrical and Information Engineering, Shandong University (Weihai) , Weihai 2 64200)

Abstract The development of deep learning has led to the emergence of adversarial sample attacks, but traditional communication encryption technology and physical layer security technology cannot solve privacy issues well. The current results of many teams have confirmed that selective adversarial samples have the characteristics of information hiding and can be well used to construct covert communications. Moreover, this covert communication model is different from the traditional covert communication model based on rules and cryptography. A new inspiration for covert communication construction. Based on this feature, our team applies adversarial sample attack technology to the construction of covert channels. Based on the DeepSpeech deep learning open source speech-to-text engine, gradient descent is used to generate speech-selective adversarial samples. And by combining the covert model and the surface model, covert communication is achieved with good results. This covert communication model is different from the traditional covert communication model based on rules and cryptography. It is a new inspiration for the construction of covert communication and has high research value.

Keywords communication encryption technology, selective confrontation, sample information hiding, covert communication

Design and Implementation of Covert Communication based on Intelligent Voice System

Lipeng Song¹⁾ Ruichen Zhang¹⁾ Wenxiang Liu¹⁾ Shuo Wang¹⁾ Ye Zhang¹⁾ Guanyu Ren¹⁾

¹⁾(School of Mechanical, Electrical and Information Engineering, Shandong University, Weihai 264209)

Abstract The development of deep learning promotes the emergence of anti-sample attacks, while the traditional communication encryption technology and physical layer security technology can not solve the privacy problem. The results of many teams at present confirm that selective adversarial samples have the characteristics of hidden information, and can be used in the construction of covert communication. This covert communication mode is different from the traditional covert communication mode based on rules and cryptography, which is a new inspiration for the construction of covert communication. Based on this feature, the counter sample attack technique is applied to the construction of covert channel. An open source speech-to-text engine based on DeepSpeech deep learning uses gradient descent to generate speech selective adversarial samples. The covert communication is realized by the combination of covert model and surface model , and the effect is good. This kind of covert communication mode is different from the traditional covert communication mode based on rules and cryptography, which is a new inspiration for the construction of covert communication and has high research value.

Key words communication encryption technology; selective confrontation samples; information hiding; covert

Song Lipeng (corresponding author), male, born in 1975, is a doctoral supervisor and professor. His main research fields are artificial intelligence security and software vulnerability analysis. E-mail : slp880@sdu.edu.cn . Zhang Ruichen, male, born in 2002 , undergraduate, E-mail: 3033824368@qq.com. Liu Wenxiang, male, born in 2002 , undergraduate, E-mail: 2425281352@qq.com. Wang Shuo, male, born in 2002 , undergraduate student, E-mail: 2404266246@qq.com. Zhang Ye, male, born in 2002 , undergraduate student, E-mail: 1498792816@qq.com. Ren Guanyu, male, born in 2002 , undergraduate student, E-mail : guanyuren@mail.sdu.edu.cn.

1 Introduction

Recently, deep neural network (DNN) [1] has provided better performance in image recognition [2], speech recognition [3], pattern analysis [4] and intrusion detection [5]. Although DNN works well, it is easily affected by adversarial examples [6]. Adding some imperceptible perturbation information to the sample and indicating the impact of the perturbation on different models, that is, selective adversarial samples, can cause the neural network model under specific parameter conditions to make wrong predictions.

Adversarial examples have been widely studied in the image field. However, in recent years, research on adversarial examples has expanded to the field of speech. Many scenarios in the audio domain demonstrate the potential threat of adversarial examples. For example, Alexa [7] developed by Amazon can provide services that support voice interaction after the user is authenticated, such as ordering products, obtaining information, and controlling multiple smart devices. However, there may be some problems with such voice interaction. For example, they might reveal a user's personal privacy to other users, or mistakenly order an item after hearing a sound on the TV or radio. In order to exploit these weaknesses of speech recognition systems, many studies propose adversarial examples and generate adversarial samples by adding a small amount of noise to the original audio samples.

At present, many teams have achieved certain results in the field of intelligent voice attack and defense. After reading and discussing papers in related fields and in-depth study of the characteristics of adversarial samples, they found that the existing

selective adversarial samples have the characteristics of hiding information and are suitable for military applications. communication, automated phone eavesdropping and covert channel [8] scenarios. For example, when military communicators need to send a specific voice message, they can deliberately cause the enemy to misunderstand the message while allowing friendly forces to correctly identify the message.

Based on the fact that selective adversarial samples have the characteristics of hiding information, can be used for covert communication, and have high security characteristics, the team proposed a solution that combines audio adversarial sample attacks with covert communication ideas. This covert communication model is different from the traditional covert communication model based on rules and cryptography. It is a new inspiration for the construction of covert communication and has high research value.

The team used the audio adversarial model to hide the actual content that needs to be sent in the audio file. This method has the following characteristics:

1. High security. Only when both communicating parties master the same model and model parameters at the same time can the hidden information be correctly parsed.
2. Strong concealment. Even if it is intercepted midway, the hidden information cannot be discovered.

2 Related work

Covert communication technology has become a research hotspot in the field of information security.

At present, covert communication technology based on digital images can be said to be relatively complete and mature, while related research on voice covert communication technology has not started yet. In recent years, the attack technology of speech models has been continuously improved and improved by various teams, so voice covert communication technology based on this counterattack technology has begun to gradually develop.

The text information hiding method of modifying the text format mainly achieves information hiding by modifying the carrier text format. In recent years, most of the improvements in this type of method have been to reduce the text scaling ratio and embed secret information more evenly in the carrier text to increase the concealment and robustness of the confidential text. This type of method has relatively good visual performance. High concealment. However, if the text format is modified or the text is re-entered during the text delivery process, the secret information will be lost.

Partala[9] first tried to use blockchain as a medium to build a covert communication channel, and proposed a blockchain covert channel (BLOCCE) model, which hides information to the last digit of the transaction address and uses the corresponding cycle to generate it sequentially. The transaction address ensures the order of secret information. Since then, researchers have improved on the shortcomings of this model and proposed an improved blockchain covert communication method, and tried to reduce its communication cost [10 , 11].

Torki proposed a blockchain covert communication scheme that does not require manual changes to source data and only requires repeated execution of hidden algorithms. Guo et al. [12] and Lan Yiqin et al. [13] implemented hybridization by combining multi-layer linkable spontaneous anonymous group signatures and introduced the new elliptic curve algorithm Monero to build a covert communication channel in their blockchain application, using Monero has high security to improve the concealment of covert channels.

For deep learning-based covert communication,

Szegedy et al. [6] first proposed an adversarial sample in which the image was slightly transformed by the attacker. The main purpose of using adversarial samples is to cause DNN errors by adding a small amount of noise to the original image; however, humans cannot tell the difference between the original image and the distorted image.

Image covert communication algorithms based on deep learning generally include two information hiding algorithms based on Convolutional Neural Networks (CNN) and Generative Adversarial Networks (GAN). GAN-based algorithms usually use A variant of CNN. Based on the confrontational characteristics between information hiding and steganalysis, many scholars have carried out targeted optimization of the model structure of deep learning to achieve the embedding and extraction of secret data. The current research status shows that image information hiding methods based on deep learning have better hiding effects and performance indicators than traditional methods.

In the field of audio, the "Devil Music" attack refers to the 2018 research team of Chen Kai of the Institute of Information Engineering, Chinese Academy of Sciences, who used music as a carrier to generate audio adversarial samples [15]. This malicious audio sample inserts implicit instructions into normal music to achieve the purpose of attacking the intelligent voice system.

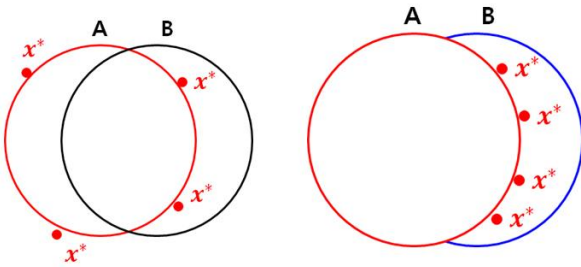
Speech-text adversarial sample attack refers to the speech adversarial sample generation method studied by Google Brain artificial intelligence scientist Carlini and his team. This attack method is to guide the speech recognition model to incorrectly identify speech samples as arbitrary specified text during the process of recognizing audio files by the intelligent speech recognition system [16] .

After focused analysis of the above articles and related work, the team decided to combine audio adversarial sample attacks with covert communication ideas to achieve a targeted attack designed on the audio-to-text engine to hide the content sent by the user into an audio file. . Only when both communicating parties master the same model and

model parameters at the same time can the hidden information be correctly parsed. Even if it is intercepted by an eavesdropper on purpose, the hidden information cannot be discovered and deciphered .

3 background

Adversarial examples deliberately add imperceptible perturbations to input samples in the data set, causing the model to give an incorrect output with high confidence. That is, adversarial audio samples only need to make small perturbations on a piece of audio, and the classifier will mistranscribe the audio with a high degree of confidence, or even transcribe it into a specified text (not a text that the audio is correctly transcribed). The reason this happens is that neural networks are easily "spoofed." In Figure 1(a), model A is the target model with a neural network. The corresponding line is the decision boundary of target model A. If the sample is within the boundaries of target model A, then the sample will be correctly classified by target model A. Figure 1(b) is an example of a selective adversarial example that is correctly classified by target model B but misclassified by target model A. In Figure 1(b), the selective adversarial example x^* is within the decision boundary of target model B, but deviates from the decision boundary of target model A.



(a) Examples of adversaries. (b) Example of safe confrontation between friends.

Figure 1: Transferability example: single enemy target model A and friendly target model B.

The transferability of adversarial examples was first proposed by Szegedy et al. [6] in the literature. The transferability of adversarial examples means that the adversarial examples are misclassified by model A

and can also be misclassified by model B. The migration attributes of adversarial examples This means that an attacker can choose to attack a machine learning model to cause samples to be misclassified without direct contact with the basic model. Szegedy et al. [6] studied the transferability of different models on the same data set. In addition, they also trained the same or different models on disjoint subsets of data and studied the transfer issues between them. However, the shortcoming is that their experimental results All are implemented on the MNIST data set. Goodfellow et al. [17] proposed that the generalization ability of adversarial samples between different models is due to the high consistency of the vectors of the adversarial interference and the model, so when training the same task, the opponent can learn similar functions on different models. This generalization feature means that if the enemy wants to attack the model, it does not need to access the target model. It can only be achieved by sending the adversarial samples generated by its own model training to the target model.

the team found that the transferability of adversarial samples can be well used in the construction of covert communications. Therefore, a new architecture is needed, which has a target model that does not require access and a self-model for attack, and uses the self-model to generate adversarial samples.

4 models

In this section, the selective audio adversarial examples and information hiding model architecture are introduced in detail.

4.1 Selective Audio Adversarial Examples

In Figure 2(a), model A is a hidden model based on neural networks. The curve in the figure is the decision boundary of the hidden model. If the sample is within the decision boundary of the hidden model, it means that the sample is correctly classified by the hidden model. Otherwise, the classification is wrong. In Figure 2(b), model B is a protection model based on neural networks, which has the same classification

decision as model A. The situation where sample x^* is between the decision boundaries of model A and model B means that sample x^* is classified incorrectly by model A and correctly classified by model B. This leads to the concept of selective adversarial examples.

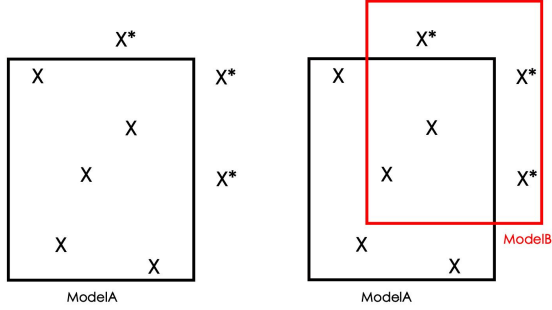


Figure 2 Single selective adversarial sample x , concealment model A and protection model B

Figure 2(b) shows that the selective adversarial example x^* is outside the decision boundary of model A, while at the same time, x^* is within the decision boundary of model B. But it is worth noting that x^* should be as close as possible to the decision boundary of model A, which means that the distance metric of the original sample x and the selective adversarial sample x^* is closer. At this time, if it further converges to the minimum distance on the decision boundary, the selective adversarial sample obtained will have the smallest noise and be selectively adversarial.

4.2 Model architecture

The prerequisite of the model architecture is that all parameters of the concealment model and the protection model are known. The information hiding model architecture is shown in the figure below:

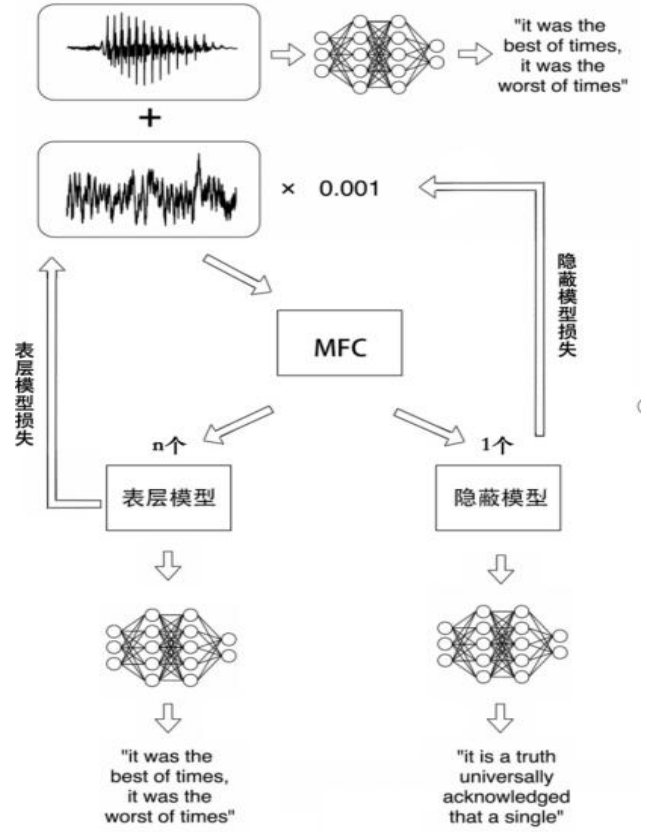


Figure 3 Model architecture diagram for generating selective adversarial examples

As can be seen from Figure (3), the architecture for generating selective adversarial samples consists of n protection models (surface models) M_{p_i} ($1 \leq i \leq n$) and a hidden model M_h . The hidden model can also be expanded to n , but considering the practicality of the product, a proprietary concealment model will suffice.

The inputs to the overall model architecture are the original sample x , the initial noise δ , and the correct classification phrase label org_phrase and the hidden classification phrase label $target_phrase$. After input, the architecture randomly generates initial noise, adds original samples, generates initial adversarial samples, and then performs mel cepstrum and forward propagation in the protection model and covert model. Mel cepstrum (MFC) is used as a preprocessing step to reduce the dimensionality of the input x . MFC splits the waveform at 50 frames per second and maps each frame to each frequency domain.

M_{p_i} and M_h keep the model unchanged in the process of generating selective adversarial samples (only model pre-training parameters and model

structure are needed in the generation process). After Mp_i and Mh get the audio array input, Mp_i outputs the classification probability of all protector models, and Mh outputs the classification probability of the covert model. Secondly, the CTC loss function is used to align the classification probability matrix and the label to obtain the loss value of each Mp_i model and the loss value of Mh . Finally, the weighted loss value of Mp_i and the loss value of Mh are weighted and fed back to the added noise to form a closed loop. The noise is further adjusted until a predetermined level of selection resistance is met.

the distance between the original sample x and the selected adversarial sample x^* on the premise of ensuring that each protection model Mp_i is correctly classified and the concealment model Mp is classified incorrectly. Thus, the objective function can be set, let $f^p(x^*)$ be the recognition behavior of the protection model, $f^h(x^*)$ be the recognition behavior of the covert model, org_phrase is the original phrase, and $target_phrase$ is the target hidden phrase. The objective function is as follows:

$$\begin{aligned} x^* : \operatorname{argmin}(x, x^*) \text{ st } f^p(x^*) &= org_phrase \text{ and } f^h(x^*) \\ &= target_phrase \\ x^* &= x + \delta \end{aligned}$$

In order to make the protection model and concealment model correctly classify the given target phrase in the objective function, and the noise delta needs to be as small as possible, the added noise adjustment requires the feedback of the loss function. The total feedback function of the loss function is as follows:

$$\begin{aligned} Loss_T = C_{dis} Loss_{distortion} + C_{Mp} \sum_{i=0}^n w_i Loss(Mp_i) \\ + C_{Mh} Loss(Mh) \end{aligned}$$

Among them C_{dis} , C_{Mp} , C_{Mh} are the weight distribution of noise, protection model and hidden model losses in the total loss respectively, which can make the generated selective adversarial samples pay more attention to a certain aspect. w_i is the loss weight of different protection models, and the distribution of weights can be adjusted according to the difficulty of

attacks by different protection models.

Choice of three loss functions: For noise loss deviation $Loss_{distortion}$, the team uses the L2 norm to measure the distance between audio arrays. for $Loss(Mp_i)$ and $Loss(Mh)$, the team uses the CTC loss function with automatic alignment characteristics.

5 experiments

Test the hidden information "this is a test" on a piece of English audio on the DEEPSPEECH model, and repeat it multiple times to observe the hiding effect and time-consuming situation. The demonstration audio uses the audio sample in carlini/audio_adversarial_examples of N. Carlini on github. The sample is a 102.2kb wav file with the English sentence "without the dataset the article is useless". The experiment was performed with Python 3.6 and tensorflow-gpu installed. ==1.15.4 was conducted under UBUNTU 20.04 LTS. In the time-consuming test, CUDA11.4 was used to call the display adapter NVIDIA RTX2070 SUPER for hardware acceleration, and statistics were collected for successful recognition in each test (the model correctly recognized the hidden information for the first time) There are three sets of data: time-consuming, time-consuming for complete convergence (complete convergence is considered when $CTCloss < 0.005$), and time-consuming for 60,000 iterations.

Table 1 Recognition time consumption unit seconds/S under different requirements

serial number	Time required for successful identification	Time consuming for complete convergence	6 0000th iteration takes time
1	3 8 . 33	3 57 . 39	4 23 . 62
2	3 7 . 12	3 45 . 46	3 68 . 00 _
3	3 6 . 51	344.41 _ _	3 75 . 94
4	3 2 . 38	3 52 . 22	3 69 . 32
5	3 5 . 25	3 70 . 01	3 98 . 83

As can be seen from the table above, the successful recognition time of the model is shorter. The average successful recognition time of the above

five groups is only 3 5. 92 s , while the average complete convergence time is 3 53. 90 s. The big difference between the two is mainly due to The process of minimizing noise takes a long time. It can be seen that audio with hidden information can be quickly produced without considering noise. However, the audio at this time has an audible "tingling" sound similar to electric current. If The large number of iterations required to remove this type of noise would take too long, and there is still much room for improvement in quickly reducing noise.

5.1 The impact of adding noise on the original audio

Next, we will further analyze and evaluate the noise in the fully converged audio. The waveforms of the audio samples before and after hiding the information are as follows:

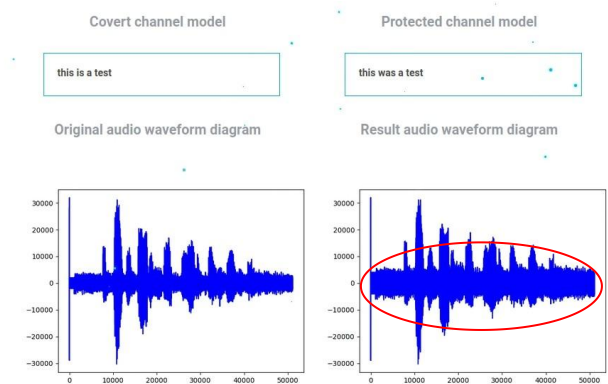


Figure 4 Audio waveform changes before and after interference

As shown in Figure 4 , the sample waveforms before and after hiding are basically the same. In particular, the high-amplitude part is difficult to see with the naked eye. However, the low-amplitude part appears slightly increased and blurred. The impact of this phenomenon on human hearing is that the audio feels blurry. You can slightly hear the noise like electric current. In order to further explore the impact on human hearing, 10 volunteers were selected for a double-blind experimental evaluation. Ten groups of short English audio, long English audio, and white noise were selected to hide phrases in them. " This is a test" , each group runs for about 3 to 50 seconds before fully converging.

5.2 Evaluation of information hiding effect

短音频组	长音频组	白噪声
1.this is a test	1.this is a test	1.this is a test
2.this is a test	2.this is a test	2.this is a test
3.this is a test	3.this is a test	3.this is a test
4.this was sist	4.this is a test	4.this is a test
5.this is a test	5.this is a test	5.this is a test

Figure 5 Audio recognition results

As shown in Figure 5, except for one sample in the short audio group that had a recognition error, all other samples were successfully identified by the model with the hidden information "this is a test", and the information hiding effect was good.

The questionnaire contains 20 sets of audio before interference, and volunteers are asked to judge whether the audio has been interfered. The questionnaire provides three options: interference, no interference, and no difference .

The results show that:

9 volunteers were unable to accurately distinguish the audio differences before and after interference (more than 20% of the audio before and after interference were correctly selected) .

8 volunteers could not roughly distinguish the difference between the audio before and after interference (more than 10% of the audio before and after interference were correctly selected) .

At the same time, 9 people said that the difference in audio before and after was very small and difficult to distinguish .

To sum up, the current model can produce audio samples that are difficult for the human ear to distinguish whether there is interference. However, the production of such samples takes too long, and there is still much room for improvement in reducing noise reduction time and further reducing low-frequency noise .

6 Summary

Selective adversarial samples have the characteristics of information hiding. Based on this characteristic, the team applied adversarial sample attack technology to the construction of covert channels. This method is based on the DeepSpeech speech-to-text engine, uses gradient descent to generate speech-selective adversarial samples, and hides the actual content that needs to be sent in the audio file. The method of combining hidden model and surface model is used to achieve information hiding and achieve good results. This information hiding model is different from traditional information hiding methods based on rules and cryptography. It provides new inspiration for the construction of information hiding and has high research value.

At present, the research of this project has the problem of too long time to generate adversarial samples and high noise. In the future, we will focus on increasing the speed of adversarial sample generation, shortening noise reduction time, and further reducing low-frequency noise.

references

- [1] J. Schmidhuber, "Deep learning in neural networks: An overview," *Neural networks*, vol. 61, pp. 85–117, 2015.
- [2] K. Simonyan and A. Zisserman, "Very deep convolutional networks for large-scale image recognition," in *International Conference on Learning Representations*, 2015.
- [3] G. Hinton, L. Deng, D. Yu, GE Dahl, A.-r. Mohamed, N. Jaitly, A. Senior, V. Vanhoucke, P. Nguyen, TN Sainath, et al., "Deep "neural networks for acoustic modeling in speech recognition: The shared views of four research groups," *IEEE Signal Processing Magazine*, vol. 29, no. 6, pp. 82–97, 2012.
- [4] R. Collobert and J. Weston, "A unified architecture for natural language processing: Deep neural networks with multitask learning," in *Proceedings of the 25th international conference on Machine learning*, pp. 160–167, ACM, 2008.
- [5] S. Potluri and C. Diedrich, "Accelerated deep neural networks for enhanced intrusion detection system," in *Emerging Technologies and Factory Automation (ETFA)*, 2016 IEEE 21st International Conference on, pp. 1–8, IEEE, 2016.
- [6] C. Szegedy, W. Zaremba, I. Sutskever, J. Bruna, D. Erhan, I. Goodfellow, and R. Fergus, "Intriguing properties of neural networks," in *International Conference on Learning Representations*, 2014.
- [7] "<https://developer.amazon.com/alexa>"
- [8] M. Smeets and M. Koot, "Covert channels," research report for RPI university of Amsterdam MSc in System and Network Engineering, 2006
- [9] PARTALA J. Provably secure covert communication on blockchain[J]. *Cryptography*, 2018, 2(3): 18.
- [10] In the Song Dynasty, Peng Wei. BLOCCE+: An improved covert communication method based on blockchain [J], *Journal of Chongqing University of Technology (Natural Science)*, 2020, 34(9): 238-244.
- [11] ZHANG LJ, ZHANG ZJ, WANG WZ, et al. A covert communication method using special bitcoin addresses generated by vanitygen[J]. *Computers, Materials & Continua*, 2020, 65(1): 597-616.
- [12] GUO ZZ, SHI LC, XU MZ, et al. MRCC: a practical covert channel over monero with provable security[J]. *IEEE Access*, 2021, 9:31816-31825.
- [13] Lan Yiqin, Zhang Fangguo, Tian Haibo. Using Monero to achieve covert communication[J]. *Journal of Xi'an University of Electronic Science and Technology*, 2020, 47(5): 19-27.
- [15] XuejingYuan, Yuxuan Chen, Yue Zhao, Yunhui Long, Xiaokang Liu, Kai Chen, ShengzhiZhang, Heqing Huang, Xiaofeng Wang, and Carl A. Gunter, "CommanderSong: ASystematic Approach for Practical Adversarial Voice Recognition," in *USENIX Security Symposium*, pages 49–64, 2018.
- [16] N. Carlini and D. Wagner, "Audio Adversarial Examples: Targeted Attacks on Speech-to-Text," 2018 IEEE Security and Privacy Workshops (SPW), 2018, pp. 1-7, doi: 10.1109/SPW.2018.00009.
- [17] Tramèr F, Papernot N, Goodfellow I, et al. The space of transferable adversarial examples. arXiv preprint arXiv: 1704.03453, 2017