
Attacking Speaker Recognition with Deep Generative Models

Anonymous Author(s)

Affiliation

Address

email

Abstract

In this paper we investigate targeted and untargeted attacks on speaker recognition systems. We first investigate the efficiency of SampleRNN and Wavenet in fooling GMM-UBM and CNN speaker recognizers. We design untargeted and targeted attacks based on the GAN framework. We propose a modification of the WGAN objective function to make use of data that is real but not from the class being learned. Our method is efficient in performing targeted and untargeted attacks, thus raising attention to issues related with security.

1 Introduction

Speaker recognition and authentication systems are being deployed for security critical applications such as banking, forensics, home automation, etc. Like other domains, these systems benefit from recent advancements in deep learning that lead to improved accuracy and trainability of such systems. Despite the improvement in the efficiency of these systems, evidence shows that they can be susceptible to adversarial attacks[20], thus motivating a current trend focused understanding adversarial attacks ([17], [6]) and finding countermeasures to detect or deflect them.

Parallel to these advancements, neural speech *generation* (the process of using deep neural networks to generate human-sounding speech) has also seen huge progress ([19], [2]). Generative Adversarial Networks (GANs) have recently been found to produce incredibly authentic samples in a variety of fields. The core idea of GANs - namely the minimax game played between a generator network and a discriminator network - extends very naturally to the field of speaker authentication and spoofing. The combination of these advancements begs a natural question that has, to the best of our knowledge, not been answered:

Are recent and old state-of-the-art speech recognition systems robust to adversarial attacks by speech generative models?

More specifically, we contemplate this question and offer in this paper the following contributions:

- We evaluate SampleRNN and WaveNet in their ability to fool text-independent state-of-the-art speaker recognizers.
- We propose strategies for untargeted attacks using Generative Adversarial Networks.
- We propose strategies for targeted attacks using a new objective function based on the improved Wasserstein GAN.

2 Related work

Generative models for speech can produce fake¹ samples that look very similar to real samples, leading humans to believe that the fake samples are real. WaveNet [18] is a generative neural network trained end-to-end to model quantized audio waveforms that has produced impressive results for conditional, speaker and text, generation of speech audio. The model is fully probabilistic and autoregressive, using a stack of causal convolutional layers to condition the predictive distribution for each audio sample on all previous ones;

SampleRNN [12], is another autoregressive architecture that has been successfully used to generate both speech and music samples. SampleRNN uses a hierarchical structure of deep RNNs to model dependencies in the sample sequence. Each deep RNN operates at a different temporal resolution so as to model both long term and short term dependencies. Another impressive model is Adobe's VoCo [1]. Although VoCo's research is unpublished, Adobe its hability to textitedit a recorded audio clip via text and make a audio clip of someone saying sentence that he never said!. Training the model requires as litle as 20 minutes of speech.

Another interesting framework for generative models are Generative Adversarial Networks (GAN) framework proposed by [5], in which a *generator* network is trained to learn a function from noise to samples that approximate the real data distribution. Simultaneously, a *discriminator* network is trained to identify whether a sample came from the real distribution or not - i.e., it is trained to try to output 1 if a sample is real, and 0 if a sample is fake. The generator and discriminator can be arbitrary networks.

In the context of deep learning architectures, adversarial examples use small perturbations to the original inputs, normally imperceptible to humans, to obtain an incorrect, target or untargeted, output from the neural network. In their brilliant papers, [17] and [6] analyze the origin of adversarial attacks and describe simple and very efficient techniques for adversarial attacks, such as the fast gradient sign method.

In the vision domain, [16] describe a technique for attacking facial recognition systems. Their attacks are physically realizable and inconspicuous, allowing an attacker to impersonate another individual. In the speech domain, [4] describe attacks on speech-recognition systems that use sounds that are hard to recognize by humans but interpreted as specific commands by speech-recognition systems.

3 Method

In this section we will describe the datasets used and the pipeline, including pre-processing and feature extraction,

3.1 Datasets

In our experiments we use three datasets, each assigned to a model as described in Table1. The datasets used are public and provide audio clips of different lengths and quality.

Table 1: Description of datasets used in our experiments. Book narr. refers to book narratives. Newspaper ++ refers to newspapers and other documents. Conversational tel. refers to conversational telephone speech.

	Speakers	Language	Duration	Context	Model
2013 Blizzard	1	English	73 h	Book narr.	SampleRNN
CSTR VCTK	109	English	400 Sentences	Newspaper ++	WaveNet
2004 NIST	100	Multiple	5 min / speaker	Conversational tel.	WGAN

3.2 Pre-processing

Data pre-processing is dependent on the model being trained. For SampleRNN and WaveNet, the raw audio is reduced to 16kHz and quantized using the μ - law companding transformation as

¹We use the term fake to refer to computer generated samples

68 referenced in SampleRNN [12] and WaveNet [18]. For the model based on the Wasserstein GAN, we
 69 pre-process the data by converting it to 16kHz and removing silences by using the WebRTC Voice
 70 Activity Detector (VAD) as referenced in [21].

71 3.3 Feature extraction

72 SampleRNN and WaveNet operate at the sample level, i.e. waveform, thus requiring no feature
 73 extraction. The features used for the neural speaker recognition system is based on Mel-Spectrograms
 74 with dynamic range compression. The Mel-Spectrogram is obtained by projecting a spectrogram onto
 75 a mel scale. We use the python library librosa [11] to project the spectrogram onto 64 mel bands, with
 76 window size equal to 1024 samples and hop size equal to 160 samples, i.e. frames of 100ms long.
 77 Dynamic range compression is computed as described in [10], with $\log(1 + C * M)$, where C is the
 78 a compression constant scalar set to 1000 and M is the matrix representing the Mel-Spectrogram.

79 4 Classification Model

80 4.1 Gaussian Mixture Model - Universal Background Model

81 4.2 Neural speaker recognition system

82 The speaker recognition system used in our experiments is based on the state-of-the-art framework
 83 by [10] and is described in Figure 1. The first module at the bottom is a pre-processing step that
 84 extracts Mel-Spectrograms features from the waveform as described in section 3.2. The second
 85 module is a convolutional neural network (CNN) that performs multi-speaker classification using the
 86 Mel-Spectrograms. The CNN is a modified version of Alexnet [8].

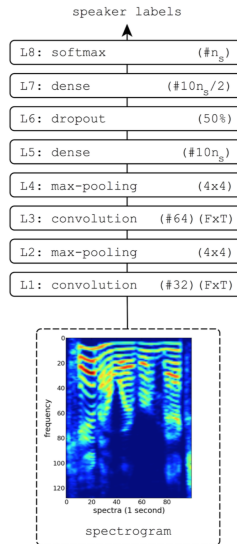


Figure 1: Architecture for CNN speaker verifier

87 We train the CNN on our training set using 64*64 Mel-Spectrograms ² consisting of balanced samples
 88 from 101 speakers from the NIST 2004 and Blizzard datasets. Our model achieves 85% test set
 89 accuracy.

²64 mel bands and 64 frames, 100 ms each

90 5 Generative Model

91 5.1 WaveNet

92 WaveNet is a generative neural network trained end-to-end to model quantized audio waveforms.
 93 It has produced impressive results for conditional, speaker and text, generation of speech audio.
 94 The model is fully probabilistic and autoregressive, using a stack of causal convolutional layers to
 95 condition the predictive distribution for each audio sample on all previous ones;

96 5.2 SampleRNN

97 SampleRNN [12], another autoregressive architecture that has been successfully used to generate
 98 both speech and music samples. SampleRNN uses a hierarchical structure of deep RNNs to model
 99 dependencies in the sample sequence. Each deep RNN operates at a different temporal resolution so
 100 as to model both long term and short term dependencies.

101 5.3 Wasserstein GANs

102 In the original generative adversarial network (GAN) framework proposed by [5], a *generator*
 103 network is trained to learn a function from noise to samples that approximate the real data distribution.
 104 Simultaneously, a *discriminator* network is trained to identify whether a sample came from the real
 105 distribution or not - i.e., it is trained to try to output 1 if a sample is real, and 0 if a sample is fake.
 106 The generator and discriminator can be arbitrary networks.
 107 The GAN framework has been shown to be able to produce very realistic samples with low training
 108 overhead. However, since the generator is trained to minimize the Kullback-Leibler (KL) divergence
 109 between its constructed distribution and the real one, it suffers from an exploding loss term when the
 110 real distribution’s support isn’t contained in the constructed one. To counter this, the *Wasserstein GAN*
 111 [3] (WGAN) framework instead uses the Wasserstein (Earth-Mover) distance between distributions
 112 instead, which in many cases does not suffer from the same explosion of loss and gradient. Based on
 113 this, the loss functions of the generator and *critic* (which no longer emits a simple probability, but
 114 rather an approximation of the Wasserstein distance between the fake distribution and real) become:

$$L_G = - \mathbb{E}_{\tilde{x} \sim \mathbb{P}_g} [D(\tilde{x})] \quad (1)$$

$$L_C = \mathbb{E}_{\tilde{x} \sim \mathbb{P}_g} [D(\tilde{x})] - \mathbb{E}_{x \sim \mathbb{P}_r} [D(x)] \quad (2)$$

115 where P_r is the real distribution, and P_g the learnt distribution of the generator.
 116 The original WGAN framework uses weight clipping to ensure that the critic satisfies a Lipschitz
 117 condition. As pointed by [7], however, this clipping can lead to problems with gradient stability.
 118 Instead, [7] suggest adding a gradient penalty to the critic’s loss function, which indirectly tries to
 119 constrain the original critic’s gradient to have norm close to 1. Equation (2) thus becomes (taken
 120 from [7]):

$$L_C = \underbrace{\mathbb{E}_{\tilde{x} \sim \mathbb{P}_g} [D(\tilde{x})] - \mathbb{E}_{x \sim \mathbb{P}_r} [D(x)]}_{\text{Original critic loss}} + \lambda \underbrace{\mathbb{E}_{\hat{x} \sim \mathbb{P}_{\hat{x}}} [(\|\nabla_{\hat{x}} D(\hat{x})\|_2 - 1)^2]}_{\text{Gradient Penalty}} \quad (3)$$

121 5.4 Adversarial attacks

122 Attacks to classification systems can be targeted or untargeted. In targeted attacks, an adversary is
 123 interested in producing an adversarial input that makes the classification system predict a target class.
 124 In untargeted attacks, the adversary is interest in a confident prediction, regardless of the class.

125 6 Results on existing methods

126 We used the samples produced with WaveNet and SampleRNN, as well as Mel-Spectrograms
 127 generated with improved WGAN, to perform adversarial attacks to the neural speaker recognition
 128 system.

129 6.1 WaveNet

130 We attempted to replicate the model described in [18] but, unfortunately, we do not have access
 131 to Google’s computing power nor to the North American Speech dataset Google used to train
 132 the WaveNet model that produced the samples referenced in [18]. Nonetheless, we used the data
 133 from CSTR VCTK to train speaker dependent WaveNet speech synthesis model that converged to
 134 producing speech that resembled the speakers voice but sounded like babbling.

135 6.2 sampleRNN

136 We also tested samples from sampleRNN [12], another autoregressive architecture that has been
 137 successfully used to generate both speech and music samples. SampleRNN uses a hierarchical
 138 structure of deep RNNs to model dependencies in the sample sequence. Each deep RNN operates
 139 at a different temporal resolution so as to model both long term and short term dependencies. We
 140 generated samples from the three-tiered variant, trained on the Blizzard 2013 dataset [14], a 300
 141 hour corpus of a single female speaker’s narration. We also downloaded 10 second samples from
 142 the original paper’s online repository at <https://soundcloud.com/samplerenn/sets>, which we
 143 qualitatively found to have less noise than our generated ones.

144 7 Results on adversarial methods

145 7.1 GAN Mel-Spectrogram Results

146 We trained this form of the GAN to construct 64x64 Mel-Spectrogram images from noise. We used
 147 two popular models for ConvNet image generation for the generator/discriminator architecture:

- 148 • *DCGAN* [15] models the generator as a series of deconvolutional layers with ReLU activa-
 149 tions, and the discriminator as a series of convolutional ones with leaky ReLU activations.
 150 Both architectures use batch normalization after each layer.
- 151 • *ResNet* [9] GAN models the generator and discriminator each as very deep convnets (30
 152 layers in our experiments) with upsampling/downsampling respectively. Residual (skip)
 153 connections are added every few layers to make training easier.

154 Visual results are demonstrated below. We generally used the same parameters as the [7] paper,
 155 namely 5 critic iterations per generator iteration, a gradient penalty weight of 10, and batch size of 64.
 156 We saw recognizable spectrogram-like features in the data after only 1000 generator iterations, and
 157 after 5000 the generated samples were indistinguishable from real ones. Training took around 10
 hours for 20000 iterations on a single 4 GB Nvidia GK104GL GPU.

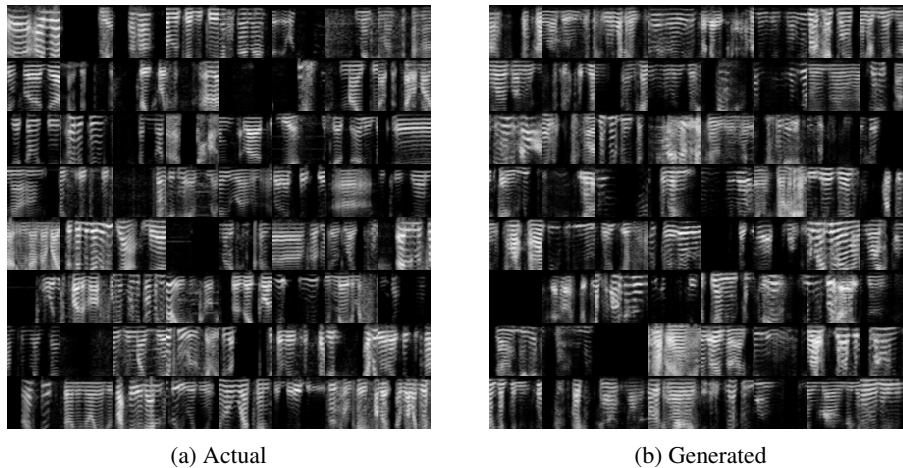


Figure 2: Comparison of real and generated (~ 5000 generator iterations) spectrogram samples. Each grid contains 64 samples.

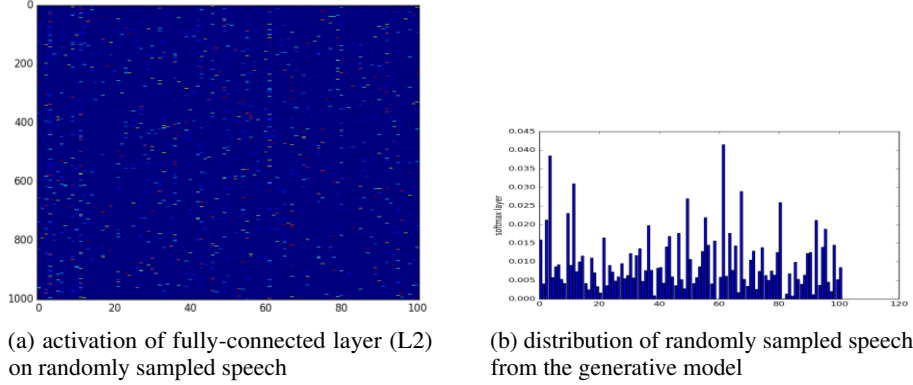


Figure 3: Summary of Untargeted Attacks

7.2 Adversarial attacks

Within the GAN framework, untargeted models are trained by using all data available, irrespective of class label. We show that an untargeted model able to generate data from the real distribution of the data and with enough variety can be used to perform adversarial attacks by classifying the samples produced by the generator. We provide details in the Untargeted attacks subsection. The models for targeted attacks can be trained in two manners. The first is based on conditioning the model on additional information, e.g. class labels, as described in [13]. The second is based on using only data from the label of interest. A drawback of the second approach is that the discriminator, and by consequence the generator, does not have access to universal³ properties of speech. To circumvent this problem, we propose a new objective function that allows using the data from all classes, although the generator is learning how to generate one class.

7.2.1 Untargeted attacks

For each of the speech in test set, we use the Mel-Spectrogram algorithm in section 3.2 to transform them into Mel-Spectrograms. The resulting mel-spectrogram is then fed into the CNN verifier and we extract a 505-dimensional feature G from the penultimate fully-connected layer (L7) in the pre-trained CNN model (1) trained on the real speech dataset with all speaker ID. The deep feature G is then used to train a K-nearest-neighbor (KNN) classifier.

To control the generator trained by our WGAN, we feed the generated Mel-Spectrograms into the same CNN-L2 pipeline to extract their corresponding feature \hat{G} . Utilizing the pre-trained KNN, each sample is assigned to the nearest speaker in the deep feature space. Therefore, we know which speaker our generated sample belongs to when we attack our CNN verifier. We evaluate our controlled WGAN samples against the state-of-the-art CNN verifier, and the confusion matrix can be found in Figure 4a. Although not included in the figure, **neither WaveNet samples nor SampleRNN samples were able to attack the recognition model in the same way.**⁴

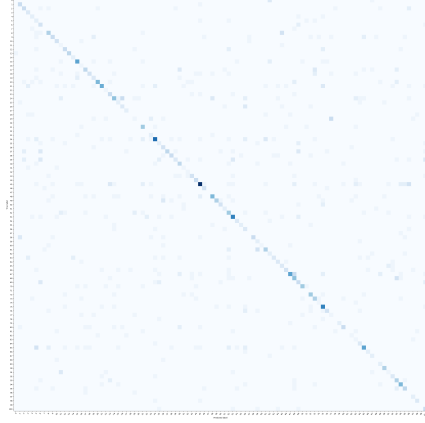
Figure 3b depicts that our GAN-trained generator successfully learns all speakers across the dataset, without mode collapsing.

7.2.2 Targeted attacks

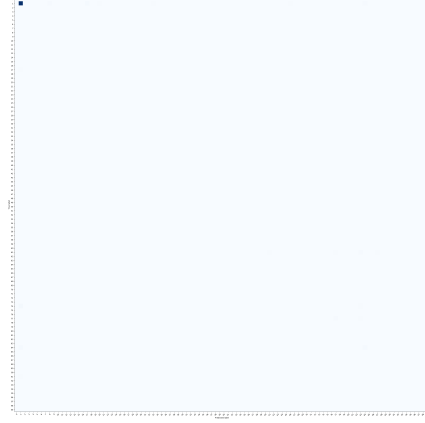
However, the most natural attack is one in which we train a GAN to directly fool a speaker-authentication system, i.e., to produce samples that the system classifies as matching a target speaker with reasonable confidence. To train our WGAN in this way, we propose a modification to the critic that allows it to learn differentiate between not only real samples and generated samples, but also between real speech samples from a target speaker and real speech samples from other speakers. We do this by adding a term to the critic’s loss that encourages its discriminator to classify real speech

³We draw a parallel with Universal Background Models in speech

⁴We also use the old state-of-the-art GMM-UBM speaker verifier.



(a) Confusion matrix of targeted attacks using untargeted model



(b) Confusion matrix of targeted attacks using targeted model

Figure 4: Confusion matrix of targeted attacks

192 samples from untargeted speakers as fake. The critic’s loss L_C becomes:

$$\underbrace{\mathbb{E}_{\tilde{x} \sim \mathbb{P}_g} [D(\tilde{x})]}_{\text{Generated Samples}} + \underbrace{\alpha * \mathbb{E}_{\hat{x} \sim \mathbb{P}_{\hat{x}}} [D(\hat{x})]}_{\text{Different Speakers}} - \underbrace{\mathbb{E}_{x \sim \mathbb{P}_r} [D(x)]}_{\text{Real Speaker}} + \underbrace{\lambda \mathbb{E}_{\hat{x} \sim \mathbb{P}_{\hat{x}}} [(\|\nabla_{\hat{x}} D(\hat{x})\|_2 - 1)^2]}_{\text{Gradient Penalty}} \quad (4)$$

193 where $P_{\hat{x}}$ is the distribution of samples from other speakers, and α is a tunable scaling factor. In
 194 our experiments, we trained this GAN with 1 target speaker and a set of 99 ‘other’ speakers. On
 195 each of the critic, we would feed the critic a batch of samples from one target speaker, and a batch of
 196 data uniformly sampled from the other speakers. We used an α of 0.1 (we found that not including
 197 this scaling factor led to serious the critic seriously overfitting and poor convergence of the GAN).
 198 Results are demonstrated in figure 4 b): The confusion matrix demonstrates all energy is concentrated
 199 in speaker id 2.

200 8 Discussion and Conclusion

201 In this paper we have investigated the use of speech generative models to perform adversarial attacks
 202 on speaker recognition systems. We show that the autoregressive models we trained, i.e. SampleRNN
 203 and WaveNet, were not able to fool the GMM-UBM and CNN speaker recognizers we built. On
 204 the other hand, we show that adversarial examples generated with GAN networks are successful in
 205 performing targeted and untargeted adversarial attacks. A pertinent argument against the validity of
 206 the GMM-UBM tests lies on the fact that GMM-UBM models have high precision and would not
 207 generalize to speech in different conditions, e.g. different room and microphone conditions. First, it
 208 is not within the scope of this paper to build a speaker recognition system that is invariant to room and
 209 microphone conditions. Second, given that the speaker recognition models, GMM-UBM and CNN,
 210 have good performance on test data and that WaveNet and SampleRNN goal is to replicate speech
 211 data that is from a speaker with similar and fixed microphone and room conditions, it is expected
 212 that the outputs of these generative models should be properly classified by the speaker recognition
 213 system.

214 On the other hand, we show that targeted and untargeted adversarial attacks with the GAN framework
 215 are efficient on the CNN speaker classifier trained by us. With this paper we hope to raise attention
 216 to issues that generative models bring to security and biometric systems. We foresee that samples
 217 produced with generative models have a signature that can be used to identify the source of the data
 218 and leave this investigation for future work.

Acknowledgments

References

References

- [1] Adobe. Adobe voco, 2017.
- [2] Sercan O Arik, Mike Chrzanowski, Adam Coates, Gregory Diamos, Andrew Gibiansky, Yongguo Kang, Xian Li, John Miller, Jonathan Raiman, Shubho Sengupta, et al. Deep voice: Real-time neural text-to-speech. *arXiv preprint arXiv:1702.07825*, 2017.
- [3] Martin Arjovsky, Soumith Chintala, and Léon Bottou. Wasserstein gan. *arXiv preprint arXiv:1701.07875*, 2017.
- [4] Nicholas Carlini, Pratyush Mishra, Tavish Vaidya, Yuankai Zhang, Micah Sherr, Clay Shields, David Wagner, and Wenchao Zhou. Hidden voice commands. In *25th USENIX Security Symposium (USENIX Security 16)*, Austin, TX, 2016.
- [5] Ian Goodfellow, Jean Pouget-Abadie, Mehdi Mirza, Bing Xu, David Warde-Farley, Sherjil Ozair, Aaron Courville, and Yoshua Bengio. Generative adversarial nets. In *Advances in neural information processing systems*, pages 2672–2680, 2014.
- [6] Ian J Goodfellow, Jonathon Shlens, and Christian Szegedy. Explaining and harnessing adversarial examples. *arXiv preprint arXiv:1412.6572*, 2014.
- [7] Ishaan Gulrajani, Faruk Ahmed, Martin Arjovsky, Vincent Dumoulin, and Aaron Courville. Improved training of wasserstein gans. *arXiv preprint arXiv:1704.00028*, 2017.
- [8] Alex Krizhevsky, Ilya Sutskever, and Geoffrey E Hinton. Imagenet classification with deep convolutional neural networks. In *Advances in neural information processing systems*, pages 1097–1105, 2012.
- [9] Christian Ledig, Lucas Theis, Ferenc Huszár, Jose Caballero, Andrew Cunningham, Alejandro Acosta, Andrew Aitken, Alykhan Tejani, Johannes Totz, Zehan Wang, et al. Photo-realistic single image super-resolution using a generative adversarial network. *arXiv preprint arXiv:1609.04802*, 2016.
- [10] Yanick Lukic, Carlo Vogt, Oliver Dürr, and Thilo Stadelmann. Speaker identification and clustering using convolutional neural networks. In *Machine Learning for Signal Processing (MLSP), 2016 IEEE 26th International Workshop on*, pages 1–6. IEEE, 2016.
- [11] Brian McFee, Colin Raffel, Dawen Liang, Daniel PW Ellis, Matt McVicar, Eric Battenberg, and Oriol Nieto. librosa: Audio and music signal analysis in python. In *Proceedings of the 14th python in science conference*, 2015.
- [12] Soroush Mehri, Kundan Kumar, Ishaan Gulrajani, Rithesh Kumar, Shubham Jain, Jose Sotelo, Aaron Courville, and Yoshua Bengio. Samplernn: An unconditional end-to-end neural audio generation model. *arXiv preprint arXiv:1612.07837*, 2016.
- [13] Mehdi Mirza and Simon Osindero. Conditional generative adversarial nets. *arXiv preprint arXiv:1411.1784*, 2014.
- [14] Kishore Prahallad, Anandaswarup Vadapalli, Naresh Elluru, G Mantena, B Pulugundla, P Bhaskararao, HA Murthy, S King, V Karaiskos, and AW Black. The blizzard challenge 2013–indian language task. In *Blizzard Challenge Workshop*, volume 2013, 2013.
- [15] Alec Radford, Luke Metz, and Soumith Chintala. Unsupervised representation learning with deep convolutional generative adversarial networks. *arXiv preprint arXiv:1511.06434*, 2015.
- [16] Mahmood Sharif, Sruti Bhagavatula, Lujo Bauer, and Michael K Reiter. Accessorize to a crime: Real and stealthy attacks on state-of-the-art face recognition. In *Proceedings of the 2016 ACM SIGSAC Conference on Computer and Communications Security*, pages 1528–1540. ACM, 2016.

- 265 [17] Christian Szegedy, Wojciech Zaremba, Ilya Sutskever, Joan Bruna, Dumitru Erhan, Ian Goodfel-
 266 low, and Rob Fergus. Intriguing properties of neural networks. *arXiv preprint arXiv:1312.6199*,
 267 2013.
- 268 [18] Aäron van den Oord, Sander Dieleman, Heiga Zen, Karen Simonyan, Oriol Vinyals, Alex
 269 Graves, Nal Kalchbrenner, Andrew Senior, and Koray Kavukcuoglu. Wavenet: A generative
 270 model for raw audio. *CoRR abs/1609.03499*, 2016.
- 271 [19] Yuxuan Wang, RJ Skerry-Ryan, Daisy Stanton, Yonghui Wu, Ron J Weiss, Navdeep Jaitly,
 272 Zongheng Yang, Ying Xiao, Zhifeng Chen, Samy Bengio, et al. Tacotron: A fully end-to-end
 273 text-to-speech synthesis model. *arXiv preprint arXiv:1703.10135*, 2017.
- 274 [20] Zhizheng Wu, Nicholas Evans, Tomi Kinnunen, Junichi Yamagishi, Federico Alegre, and
 275 Haizhou Li. Spoofing and countermeasures for speaker verification: a survey. *Speech Commu-*
 276 *nication*, 66:130–153, 2015.
- 277 [21] Adham Zeidan, Armin Lehmann, and Ulrich Trick. Webrtc enabled multimedia conferenc-
 278 ing and collaboration solution. In *WTC 2014; World Telecommunications Congress 2014;*
 279 *Proceedings of*, pages 1–6. VDE, 2014.