
WaveFake: A Data Set to Facilitate Audio DeepFake Detection

Joel Frank*

Ruhr University Bochum
Horst Görtz Institute for IT-Security

Lea Schönherr

Ruhr University Bochum
Horst Görtz Institute for IT-Security

Abstract

1 Deep generative modeling has the potential to cause significant harm to society.
2 Recognizing this threat, a magnitude of research into detecting so-called “Deep-
3 Fakes” has emerged. This research most often focuses on the image domain,
4 while studies exploring generated audio signals have—so-far—been neglected. In
5 this paper we make three key contributions to narrow this gap. First, we provide
6 researchers with an introduction to common signal processing techniques used
7 for analyzing audio signals. Second, we present a novel data set, for which we
8 collected audio samples from five different network architectures, across two lan-
9 guages. Finally, we supply practitioners with two baseline models, adopted from
10 the signal processing community, to facilitate further research in this area.

11 1 Introduction

12 \$243,000 were lost, when criminals used a generated voice recording to impersonate the CEO of a UK
13 company [76]. This is just one of several reports where current state-of-the-art generative modeling
14 was used in harmful ways. Other examples include: impersonation attempts [20], influencing
15 opposition movements [36], being used to justify military actions [24, 46], or online harassment [9].
16 While there is a multitude of beneficial use cases, for example, enhancing data sets for medical
17 diagnostics [18, 22], medical image segmentation [87], or designing DNA to optimize protein
18 bindings [29], finding effective ways to detect fraudulent usage is of utmost importance to society.

19 Detection in the image domain has received tremendous attention [41, 45, 91, 78, 83, 43, 47, 42, 17,
20 21]. However, the audio domain is severely lacking. While there does exist prior work exploring
21 image and sound together (i.e., videos) [13], an analysis of audio in isolation is missing. This is a
22 critical gap. When examining the domains jointly, we can utilize synergies, for example, analyzing
23 how well spoken audio matches video on screen.

24 **To encourage more researchers to also explore the audio domain, we make three key contributions**
25 **in this paper:** First, we provide an overview of common signal processing techniques used for
26 analyzing audio signals. We give an introduction to spectrograms, which are commonly used as an
27 intermediate representation for generative models [35, 60, 88, 89], Additionally, we review common
28 feature representations used for automatic speech recognition [56] or speaker verification [67], and
29 provide a survey of current state-of-the-art generative models.

30 Second, our main contribution is a novel data set. We collected eight sample sets from five different
31 network architectures across two languages. In this paper, we focus on analyzing samples which
32 resemble (i.e., recreate) the training distributions. This allows for one-to-one comparisons of audio
33 clips between the different architectures, in which we find subtle differences between the generators.

*Corresponding author joe1.frank@rub.de.

34 Additionally, we expect good performing classifiers to transfer well to other contexts, since recreating
35 the training distribution should yield the most quality samples.

36 Finally, we implement two classifiers, which we adopted from best practices in the signal processing
37 community [67], to give future researchers a baseline to compare against ². Furthermore, we
38 implemented BlurIG [86] a popular attribution methods, so practitioners can inspect their predictions
39 when building on our results.

40 We summarize our main contributions as follows:

- 41 • An introduction into common signal processing techniques and a survey of the current
42 landscape of audio generative modeling.
- 43 • A novel data set consisting of samples from several state-of-the-art generative network
44 architectures.
- 45 • An implementation of two baseline models for future researchers to compare against.

46 2 Background

47 In this section we provide an introduction into common techniques used for analyzing speech audio
48 signals. The list is far from exhaustive, but it provides a starting point for researchers to explore
49 the field. The interested reader is referred to the excellent books by Rabiner et al. [63] or Quatieri
50 [62]. Additionally, we provide a survey on current state-of-the-art generative models and summarize
51 related work.

52 2.1 Analyzing speech signals

53 We start by giving an introduction to commonly used techniques and representations used to analyze
54 audio signals.

55 **(Mel) spectrograms:** A spectrogram is a visual representation of the frequency information of
56 a signal over time (cf. Section 3, Figure 2 for an example). To calculate a spectrogram for an
57 audio signal, we proceed as follows: First we divide the waveform into *frames* (e.g., 20 ms) with
58 an overlap (e.g., 10 ms) between two adjacent frames. We then apply a window function $w(n)$ to
59 avoid spectral leakage ³. These functions (e.g., Hamming, Hann, Blackman window) are a trade-off
60 between frequency resolution and spectral leakage and their choice depends on the task and the signal
61 properties, cf. Prabhu [57] for a detailed overview. We multiply each individual frame from our audio
62 signal with the windowing function:

$$x_w(t, n) = x(t, n) \cdot w(n) \quad \forall n = 0, \dots, N - 1, \quad (1)$$

63 where N is the frame length and $t = 1, \dots, T$ the frame index of the signal sample $x(t, n)$. The
64 frames are then transformed individually using the *Discrete Fourier Transform* (DFT) to obtain a
65 representation in the frequency domain:

$$X(t, k) = \sum_{n=0}^{N-1} x_w(t, n) e^{-i2\pi \frac{kn}{N}} \quad \forall k = 0, \dots, K - 1, \quad (2)$$

66 with K DFT coefficients. This procedure of dividing the input signal, applying the window func-
67 tion and computing the DFT is referred to as the *Short-Time Fourier Transform* (STFT). Finally,
68 we calculate the squared magnitude $|X(t, k)|^2$ of the complex-valued signal to obtain our final
69 representation—the spectrogram.

70 A commonly used variant is the so-called Mel spectrogram. It is motivated by studies which have
71 shown that humans do not perceive frequencies on a linear scale. In particular, they can detect
72 differences in lower frequencies on a more fine grade scale when compared to higher frequencies [97].
73 The Mel scale is an empirically determined non-linear transformation which approximates this
74 relationship:

$$f_{\text{mel}} = 2595 \cdot \log_{10} \left(1 + \frac{f}{700} \right), \quad (3)$$

²Our code can be found at github.com/RUB-SysSec/WaveFake

³Energies from one frequency leak into other frequency bins.

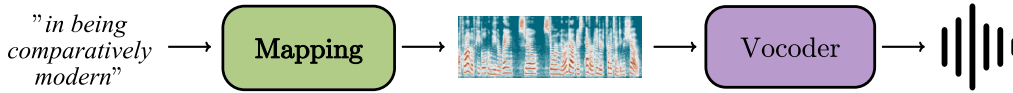


Figure 1: A typical TTS pipeline. One model takes a textual prompt with the desired audio transcription (we call it the “mapping” model) and outputs an intermediate representation, for example Mel spectrograms. This intermediate representation is then fed to a second model (in the literature often referred to as “vocoder”) to obtain the final raw audio.

75 where f is the frequency in Hz and f_{mel} the Mel-scaled frequency. To obtain Mel spectrogram, we
 76 apply an ensemble of S triangular filters H_{mel} (we provide a visual representation in the supplement-
 77 ary material). These filters have a linear distance between the triangle mid frequencies in the Mel
 78 scale, which results in a logarithmic increasing distance of the frequencies in the frequency domain

$$X_{\text{mel}}(t, s) = \sum_{k=0}^{K-1} |X(t, k)| H_{\text{mel}}(s, k) \quad \forall s = 1, \dots, S. \quad (4)$$

79 Which gives us the final Mel spectrogram. Based on it, we can compute a common feature represen-
 80 tation for audio analysis:

81 **Mel Frequency Cepstral Coefficients:** *Mel Frequency Cepstral Coefficients* (MFCC) are derived
 82 from a Mel-scaled spectrogram by applying a *Discrete Cosine Transform* (DCT) to the logarithm of
 83 the Mel-filtered signal

$$c(t, r) = \sum_{s=0}^{S-1} \log [X_{\text{mel}}(t, s)] \cdot \cos \left[\frac{\pi \cdot r \cdot (s + 0.5)}{S} \right] \quad \forall r = 0, \dots, R - 1, \quad (5)$$

84 where R is the number of DCT coefficients.

85 **Linear Frequency Cepstral Coefficients:** We can also calculate *Linear Frequency Cepstral Coef-*
 86 *ficients* (LFCC). As the name suggest these coefficients are derived by applying a linear filterbank
 87 (instead of a Mel filterbank) to the spectrogram of the signal. This results in retaining more high
 88 frequency information. Except for the replacement of the filter bank, all other step remain the same
 89 as for MFCC features.

90 **(Double) delta:** MFCCs and LFCCs are often augmented by their first and second derivatives to
 91 represent temporal structure of the input. These are referred to as delta and double delta features,
 92 respectively. In practice these are often calculated by central difference approximation via

$$d(t) = \frac{\sum_{n=1}^N n \cdot [c(t+n) - c(t-n)]}{2 \cdot \sum_{n=1}^N n^2} \quad \forall t = 0, \dots, T - 1, \quad (6)$$

93 where $d(t)$ is the delta at time t and N is a user-defined window length for computing the delta, and c
 94 is either the MFCCs/LFCCs or the delta features (when calculating the double delta features).

95 2.2 Text-to-speech (TTS)

96 In this Section we want to give a broad overview over different research direction for *Text-To-*
 97 *Speech* (TTS) models. Due to the rapid developments of the field, this is a non-exhaustive list.
 98 However, it serves as a starting point for interested researchers.

99 While there has been some research into end-to-end models [16, 77], typical TTS models consist of a
 100 two-stage approach, represented in Figure 1. First, we enter the text sequence which we want to gener-
 101 ate. This sequence gets mapped by some model (or feature extraction method) to a low-dimensional
 102 intermediate representation, often linguistic features [7] or Mel spectrograms [49]. Second, we use
 103 an additional model (often referred to as vocoder), to map this intermediate representation to raw
 104 audio. We focus on the literature on vocoders, since it directly connects to our work.

105 Today, vocoders are typically implemented by Deep Neural Networks (DNNs). The first DNN [93, 19]
106 approaches adopted the parametric vocoders of earlier HMM-based models [94, 80, 90]. Here the
107 DNN was used to predict the statistics of a given time frame, which are then used in traditional
108 speech parameter generation algorithms [80]. Later variants replaced each component in traditional
109 pipelines with neural equivalents [7, 6, 64, 65, 84, 4]. The first architectures who started using
110 DNNs exclusively as the vocoder were auto-regressive generative models, such as WaveNet [49],
111 WaveRNN [27], SampleRNN [44], Char2Wav [75] or Tacotron 2 [72].

112 Due to their auto-regressive nature, these models do not leverage the parallel structure of modern
113 hardware. There have been several attempts to circumvent this problem: One direction is to distill
114 trained auto-regressive decoders into flow-based [32] convolutional student networks, as done by
115 Parallel WaveNet [49] and Clarinet [54]. Another method is to utilize direct maximum likelihood
116 training as done by several flow-based models, for example, WaveGlow [60], FloWaveNet [30] or
117 WaveFlow [55]. Other probabilistic approaches include those based on variational auto-encoders [50,
118 53] or diffusion probabilistic models [34, 12]. Another family of methods is based on Generative
119 Adversarial Networks (GANs) [23], examples include, MelGAN [35], GAN-TTS [8], WaveGAN [15],
120 Parallel WaveGAN [88] or Multi-Band MelGAN [89].

121 2.3 Related Work

122 There have been several previous proposals which collected DeepFake data: The FaceForensics++
123 dataset [66] curated 1.8 million manipulated images and provides a benchmark for automated facial
124 manipulation detection. Celeb-DF [40] contains high-quality face-swapping DeepFake videos of
125 celebrities with more than 5,000 fake videos. Dolhansky et al. [14] released the DeepFake detection
126 challenge that contains more than 100,000 videos, generated with different methods.

127 There exists a multitude of research into identifying GAN-generated images: Several approaches use
128 CNNs in the image domain [41, 45, 91, 78, 83], others use statistics in the image domain [43, 47].
129 Another group of systems employs handcrafted features from the frequency domain: steganalysis-
130 based features [42], spectral centroids [82] or frequency analysis [96, 17, 21, 61]. Li and Lyu
131 [39] proposed a CNN-based DeepFake video detection framework which utilizes artefacts that are
132 consequences of the generation process. Another strain of research combines image analysis with
133 audio analysis. Chintha et al. [13] combined a DeepFake detection with an audio spoofing detection
134 to identify fake videos. At the time of writing and to the best of our knowledge no work has analyzed
135 DeepFake audio in isolation.

136 A related line of research is undertaken by the signal processing community. The biyearly ASVspoof
137 challenges [85, 79, 48] promotes countermeasure against spoofing attacks that aim to fool speaker
138 verification systems via different kinds of attacks. Their benchmarking data sets include replay
139 attacks, voice conversion and synthesized audio files. Note that the 2021 edition of the challenge
140 features an audio DeepFake track, but does not provide specific training data for it. **We imagine our
141 data set to be used complementary with the training data of the challenge. At the time of writing the
142 2021 edition is still on-going, but evaluating the best performing models in conjunction with our data
143 set is an interesting direction for future work. In the mean time, we adopt one of the baseline models
144 of the ASVspoof challenge to enable a direct comparison.** These efforts have lead to several proposed
145 models for detecting spoofing attacks, for example, CNN-based methods [81, 38, 37], ensemble
146 methods on different feature representations [52, 28, 69] or methods which detect unusual pauses
147 in human speech [95, 3]. Additionally, another data set is proposed by Kinnunen et al. [33]. They
148 released a re-recorded version of the RedDots database for replay attack detection text-dependent
149 speaker verification.

150 3 The data set

151 In this Section we provide an overview of our data set. It consists of 88,600 generated audio clips
152 (16-bit PCM wav) and can be found on zenodo ⁴. **We mostly base our work on the LJSPEECH [26]
153 data set. While TTS models often get trained on private data sets, LJSPEECH is the most common
154 public data set among the publication listed in Section 2.2. Additionally, we consider the JSUT [74]
155 data set, a Japanese speech corpus.**

⁴ zenodo.org/record/4904579 - DOI: 10.5281/zenodo.4904579

156 **Reference data:** We examine multiple networks trained on two reference data sets. First, the
157 LJSPEECH [26] data set consisting of 13,100 short audio clips (on average 6 seconds each; roughly
158 24 hours total) read by a female speaker. It features passages from 7 non-fiction books and the audio
159 was recorded on a MacBook Pro microphone. Second, we include samples based on the JSUT [74]
160 data set, specifically, basic5000 corpus. This corpus consists of 5,000 sentences covering all basic
161 kanji of the Japanese language (4.8 seconds on average; roughly 6.7 hours total). The recordings
162 were performed by a female native Japanese speaker in an anechoic room. Thus, our data set consists
163 of approximately 157 hours of generated audio files in total. Note that we do not redistribute the
164 reference data. They are freely available online [26, 74].

165 **Architectures:** We included a range of architectures in our data set:

- 166 • **MelGAN:** We include MelGAN [35], which is one of the first GAN-based generative models
167 for audio data. It uses fully convolutional feed-forward network as generator and operates on
168 Mel spectrograms. Their discriminator is a combination of three different discriminators that
169 operates on the original, and two downsampled versions of the raw audio input. Additionally,
170 they use an auxiliary loss over the feature space of the three discriminators.
- 171 • **Parallel WaveGAN (PWG):** WaveNet [49] is one of the earliest and most common archi-
172 tectures, We include samples from one of its variants, Parallel WaveGAN [88]. It uses
173 GAN training paradigm, with a non-autoregressive version of WaveNet as its generator. In a
174 similar vein to MelGAN, it uses an auxiliary loss, but in contrast, matches the STFT of the
175 original training sample and the generated waveform over multiple resolutions.
- 176 • **Multi-band MelGAN (MB-MelGAN):** Incorporating more fine-grained frequency analysis,
177 might lead to more convincing samples. We include MB-MelGAN, which computes its
178 auxiliary (frequency-based; inspired by PWG) loss in different sub-bands. Its generator is
179 based on a bigger version of the MelGAN generator, but instead of predicting the entire
180 audio directly, the generator produces multiple sub-bands, which are then summed up to the
181 complete audio signal.
- 182 • **Full-band MelGAN (FB-MelGAN):** We include a variant of MB-MelGAN which gener-
183 ates the complete audio directly and computes its auxiliary loss (the same as PWG) over the
184 full audio instead of its sub-bands.
- 185 • **WaveGlow:** The training procedure might also influence the detectability of fake samples.
186 Therefore, we include samples from WaveGlow to investigate maximum-likelihood-based
187 methods. It is a flow-based generative model based on Glow [31], whose architecture is
188 heavily inspired by WaveNet.

189 Additionally, we examine MelGAN both in a version similar to the original publication, which we
190 denote as MelGAN, and in a larger version with a bigger receptive field, MelGAN (L)arge. This
191 version is similar to the one used by FB-MelGAN, allowing for a one-to-one comparison. In total, we
192 sample eight different data sets, six based on LJSPEECH (MelGAN, MelGAN (L), FB-MelGAN,
193 WaveGlow, PWG) and two based on JSUT (MB-MelGAN, PWG).

194 **Sampling procedure:** For WaveGlow we utilize the official implementation [59] (commit 8afb643)
195 in conjunction with the official pre-trained network on PyTorch Hub [58]. We use a popular imple-
196 mentation available on GitHub [25] (commit 12c677e) for the remaining networks. The repository
197 also offers pre-trained models. When sampling the data set, we first extract Mel spectrograms from
198 the original audio files, using the pre-processing scripts of the corresponding repositories. We then
199 feed these Mel spectrograms to the respective models to obtain the data set. Intuitively, the networks
200 are asked to "recreate" the original data sets.

201 **Differences between the architectures:** We analyze differences between the architectures by
202 plotting the spectrograms of an audio file in Figure 2 (LJSPEECH 008-0217). Larger plots can be
203 found in the supplementary material. Generally, all architectures produce spectrograms different to
204 the original. The networks seem to generally struggle with the absent of information (solid circles in
205 Figure 2a). They also seem to consistently produce differing results in the higher frequency, especially
206 for vocals (dashed circle). Additionally, MelGAN and WaveGlow seem to cause a repeating horizontal
207 pattern. The remaining networks (all using an auxiliary loss in the frequency domain) do not seem to

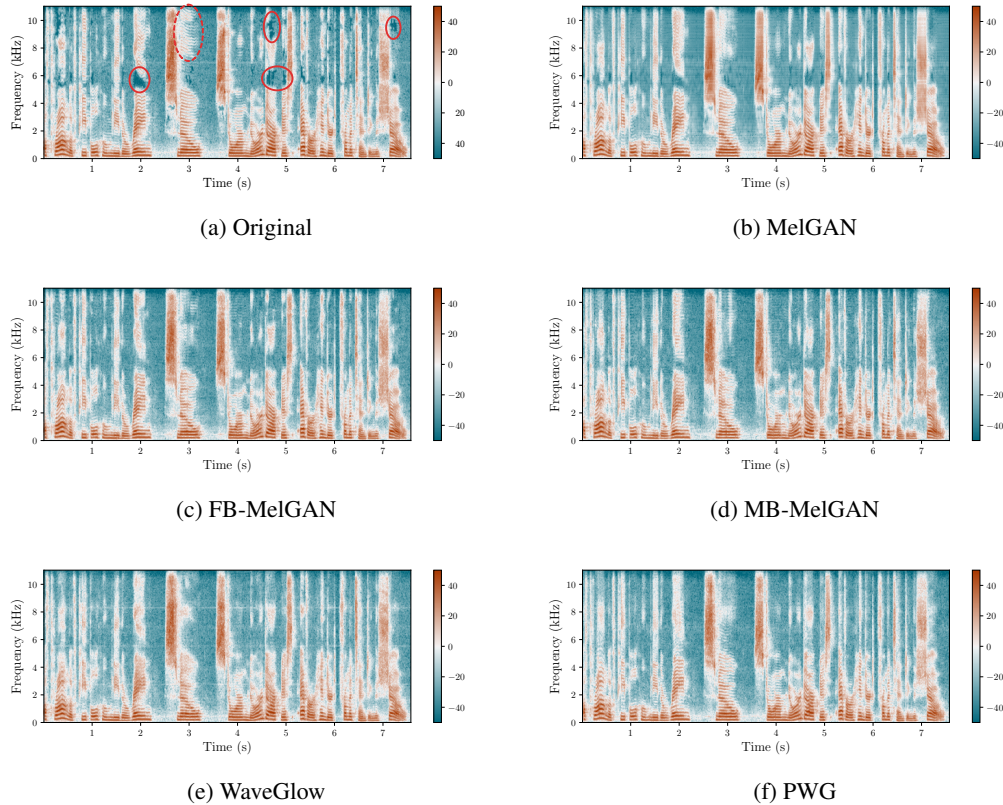


Figure 2: **Spectrograms for the same sample, for different generating models.** They show the frequencies of a signal, plotted over the time of a signal. Lower frequencies at the bottom, higher at the top. Best viewed in color.

208 exhibit this behaviour. However, they still produce clear differences. **Note that these differences are**
 209 **visible when plotting the audio but generally inaudible when listening to the samples.**

210 **A note on licensing:** During the collection of our data set, we ran into an interesting questions
 211 which we could not find a satisfying answer to. We generated samples which are intrinsically designed
 212 to be as close as possible to the original data set. So, when distributing these samples (or the models
 213 that generated them), it is not clear whether the original license does still apply. The data is obviously
 214 not the original data. Yet, it sounds remarkably like it. To the best of our knowledge this question has
 215 not been addressed by the machine learning or legal community.

216 For our sake, the LJSPEECHdata set is in the public domain. The JSUTcorpus is licensed by CC-BY-
 217 SA 4.0, with a note that redistribution is only permitted in certain cases. We contacted the author,
 218 who saw no conflict in distributing our fake samples, as long as its for research purposes.

219 To comply with JSUT we license our data set under the CC-BY-SA 4.0 license.

220 **Ethical considerations:** Our data set consists of phrases from non-fiction books (LJSPEECH) and
 221 every-day conversational Japanese (JSUT), which are already available online. The same is true for
 222 all models used to generate this data set.

223 **One might wonder if releasing research into detecting DeepFakes might contribute negatively towards**
 224 **the detection ”arms race“.** This is a long standing debate in the security community and the overall
 225 **consensus it that ”security through obscurity“ does not work.** Intuitively, **withholding information**
 226 **from the research community is in-fact more harmful, since attackers will eventually adapt to any**
 227 **defense one deploys.** We have provided a more thorough discussion of this topic in the supplementary
 228 **material and we hope that this examination contributes to the overall dialogue on security analysis of**
 229 **machine learning systems.**

230 4 Providing a baseline

231 To provide a point of reference for future researchers, we adopt the baseline model of the ASVspoof
232 challenge [79]. A bi-yearly challenge on detecting spoofed audio samples.

233 4.1 Experimental setup

234 We start by training six different classifiers, one for each vocoder in our data set (MelGAN, MelGAN
235 (L), FB-MelGAN, MB-MelGAN, PWG and WaveGlow). For training our classifiers, we exclusively
236 use the data sets based on LJSPEECH. Additionally, we use the JSUT data as a hold-out set for
237 accessing the classifiers ability to generalize to an unknown setting (different speaker, language, and
238 recording setup). While we do not explicitly asses completely novel phrases, the JSUT experiments
239 give us a good approximation. We follow Sahidullah et al. [67] and train two Gaussian Mixture
240 Models (GMMs), one fitting the real distribution (the original LJSPEECH data set) and one fitting the
241 generated audio samples (the respective vocoder-samples from our data set). In addition to the LFCC
242 features used by Sahidullah et al. [67], we evaluate MFCC features, since they are a commonly used
243 feature representation for audio tasks. We calculate the likelihood $\Lambda(\mathbf{X})$ of a test sample via

$$\Lambda(\mathbf{X}) = \log p(\mathbf{X}|\theta_n) - \log p(\mathbf{X}|\theta_s), \quad (7)$$

244 where \mathbf{X} are the input features (namely MFCC or LFCC) and θ_n and θ_s are the GMM model
245 parameter of the real and the generated audio distributions, respectively.

246 For each classifier we evaluate the performance on all vocoders over a hold-out set of 20% of the
247 data. We use the *Equal Error Rate* (EER) as our evaluation metric. This metric is also been used by
248 the ASVspoof challenge. It is defined as the point on the ROC curve, where false acceptance rate and
249 false rejection rate are equal and is commonly used to assess the performance of binary classifications
250 tasks like biometric security systems [68]. The best possible value is 0.0 (no wrong predictions),
251 worst 1.0 (everything wrong), guessing is 0.5. The lower the EER the better the system performs.
252 Additionally, we compute average EER over all test sets.

253 Finally, we train six additional models in a leave-one-out setting to access if the models picked up on
254 vocoder-specific characteristics when trained on data produced by only one model. These models are
255 exclusively trained on LFCC features.

256 **Training details:** We train the GMMs using the *Expectation Maximization* (EM) algorithm on
257 1,000 samples for a maximum of 100 iterations (the models generally converge after approximately
258 60 iterations), we use 128 mixture components and learn the diagonal covariance matrix of each
259 distribution. To ensure we do not get stuck in a local minima, we randomly reinitialized the EM
260 algorithm 10 times, picking the model with the highest log likelihood on the training data. We also
261 trained GMMs using gradient descent on a larger training corpora ($\sim 10,000$ audio samples), to
262 control for the size of our training set. The EM version obtained strictly better results. **Training EM-**
263 **based models for the leave-on-out experiments proved difficult due to numerical instability. Thus, we**
264 **exclusively rely on gradient descent based models. We doubled the amount of mixture components**
265 **(256) and epochs (20) to compensate for the more difficult task of fitting a more diverse training set.**

266 We resample all audio files to 16kHz and remove silence parts which are longer than two sec-
267 onds. When converting the audio files to MFCC/LFCC features, we use the parameters proposed
268 by Sahidullah et al. [67]. We extract 20 LFCC/MFCC features and compute delta-/double-delta-
269 features, cf. Section 2.1.

270 We trained all our models on a machine running Ubuntu 18.04.5 LTS, with a AMD Ryzen 7 3700X
271 8-Core Processor and 64GB of RAM. The implementation of our models was performed in PyTorch
272 1.8.1, using the torchaudio extension in version 0.8.1 [51]. The EM version of the GMM models can
273 be trained exclusively on the CPU, taking roughly two and a half hours to train a single model (100
274 iterations; 10 reruns). When training the gradient descent version, we used a GeForce RTX 2080Ti.
275 Training a model for 10 epochs on 10,000 audio samples, takes roughly half an hour.

276 4.2 Results

277 In a first experiment, we evaluate the performance on MFCC features. The results are presented in
278 Table 1. The rows show the respective training sets and the columns the different test set. Gray values

Table 1: **Equal Error Rate (EER) of the baseline classifier on different subset (MFCC)**. We train a new GMM model for each training set and use the log-likelihood ratio to score every sample. For each data set we compute the EER, best possible result is 0.0, worst is 1.0, the lower the better. Additionally, we compute the average EER (aEER) over all sets.

Training Set	LJSPEECH						JSUT		aEER
	MelGAN	MelGAN (L)	FB-MelGAN	MB-MelGAN	WaveGlow	PWG	MB-MelGAN	PWG	
MelGAN	0.254	0.218	0.389	0.378	0.362	0.480	0.686	0.717	0.436
MelGAN (L)	0.286	0.126	0.402	0.347	0.345	0.478	0.456	0.492	0.364
FB-MelGAN	0.413	0.379	0.177	0.196	0.225	0.286	0.430	0.450	0.320
MB-MelGAN	0.460	0.430	0.321	0.007	0.110	0.060	0.251	0.315	0.244
WaveGlow	0.405	0.379	0.294	0.074	0.026	0.083	0.237	0.259	0.220
PWG	0.499	0.493	0.395	0.055	0.147	0.006	0.190	0.229	0.252

We highlight in-distribution results in gray and the best out-distribution results per column in **bold**. (L) denotes Large.

Table 2: **Equal Error Rate (EER) of the baseline classifier on different subset (LFCC)**. Again, we train a new GMM model for each data set and compute the EER as well as the aEER.

Training Set	LJSPEECH						JSUT		aEER
	MelGAN	MelGAN (L)	FB-MelGAN	MB-MelGAN	WaveGlow	PWG	MB-MelGAN	PWG	
MelGAN	0.087	0.056	0.120	0.112	0.095	0.177	0.112	0.262	0.128
MelGAN (L)	0.082	0.024	0.089	0.092	0.079	0.162	0.142	0.370	0.130
FB-MelGAN	0.178	0.103	0.007	0.015	0.013	0.024	0.053	0.153	0.068
MB-MelGAN	0.332	0.278	0.099	0.000	0.011	0.003	0.011	0.043	0.097
WaveGlow	0.257	0.204	0.047	0.011	0.001	0.006	0.023	0.064	0.077
PWG	0.379	0.349	0.005	0.148	0.018	0.000	0.005	0.026	0.116

We highlight in-distribution results in gray and the best out-distribution results per column in **bold**. (L) denotes Large.

279 indicate that the same generative model is used for the training of the GMM classifier as for the test
280 set.

281 **MFCC:** When comparing the overall performance, i.e., the lowest average EER (aEER), we can
282 observe that PWG (0.252), MB-MelGAN (0.244), and, WaveGlow (0.220) serve as the best priors for
283 the entire data set. However, they all perform significantly worse on the MelGAN, the MelGAN (L)
284 and (to a lesser extend) the FB-MelGAN data sets. This trend is reversed for MelGAN and MelGAN
285 (L), where they achieve the best results on each other (0.218 and 0.286, respectively) and dropping
286 performance on other data sets (~ 0.400 ; up to 0.717 on JSUT). FB-MelGAN does not perform
287 particularly well on any data set.

288 The similarities between PWG and WaveGlow are intuitive. The WaveGlow architecture is heavily
289 inspired by WaveNet (the generator network of PWG). Yet, the best results for both PWG (0.060)
290 and WaveGlow (0.110) are obtained by the model trained on MB-MelGAN. We hypothesize that the
291 auxiliary loss computed over sub-bands forces MB-MelGAN to generate samples more in line with
292 WaveGlow and PWG. Surprisingly FB-MelGAN, generalizes neither to the MelGAN (L) data sets
293 nor to MB-MelGAN. FB-MelGAN uses the same architecture as MelGAN (L) and a similar auxiliary
294 loss to MB-MelGAN, albeit not computing it over sub-bands.

295 When examining completely novel data (JSUT), all classifier drop in performance. However, PWG,
296 WaveGlow, and, MB-MelGAN still serve as a good prior, implying that the generating architectures
297 exhibit common patterns across different training data sets. A similar pattern was also observed in
298 the image domain [83].

299 **LFCC:** For comparison we train an additional batch of models on LFCC features. The results can
300 be found in Table 2. LFCC features seem to be a strictly better feature representation, improving
301 performance significantly across the board. Additionally, they allow the classifier trained on FB-
302 MelGAN to become the best performing classifier (0.068). It strikes a balance between generalizing
303 to PWG, WaveGlow, MB-MelGAN, while also retraining a fairly good performance on MelGAN
304 and MelGAN (L). LFCC features contain a significantly higher amount of high-frequency features.
305 Thus, we hypothesize that this fact allows FB-MelGAN to recognize its architecture similarities with
306 MelGAN and the changes caused by the auxiliary loss. Again, similar patterns were also observed in
307 the image domain [21], implying that methods might transfer between the two.

Table 3: **Equal Error Rate (EER) for the baseline classifier in an out-of-distribution setting. We train a new GMM model for each but one distribution on LFCC features.**

Left-out Set	LJSPEECH						JSUT		aEER
	MelGAN	MelGAN (L)	FB-MelGAN	MB-MelGAN	WaveGlow	PWG	MB-MelGAN	PWG	
MelGAN	0.237	0.164	0.045	0.003	0.004	0.004	0.003	0.014	0.059
MelGAN (L)	0.233	0.166	0.037	0.002	0.004	0.002	0.002	0.014	0.058
FB-MelGAN	0.194	0.122	0.056	0.004	0.005	0.004	0.003	0.007	0.049
MB-MelGAN	0.177	0.106	0.040	0.015	0.006	0.006	0.003	0.012	0.046
WaveGlow	0.182	0.110	0.040	0.003	0.012	0.006	0.005	0.027	0.048
PWG	0.176	0.106	0.033	0.004	0.005	0.017	0.003	0.015	0.045

We highlight the distribution not present in the training set in **bold**. For JSUT, we highlight the entry when the generating network architecture was not part of the training set. (L) denotes Large.

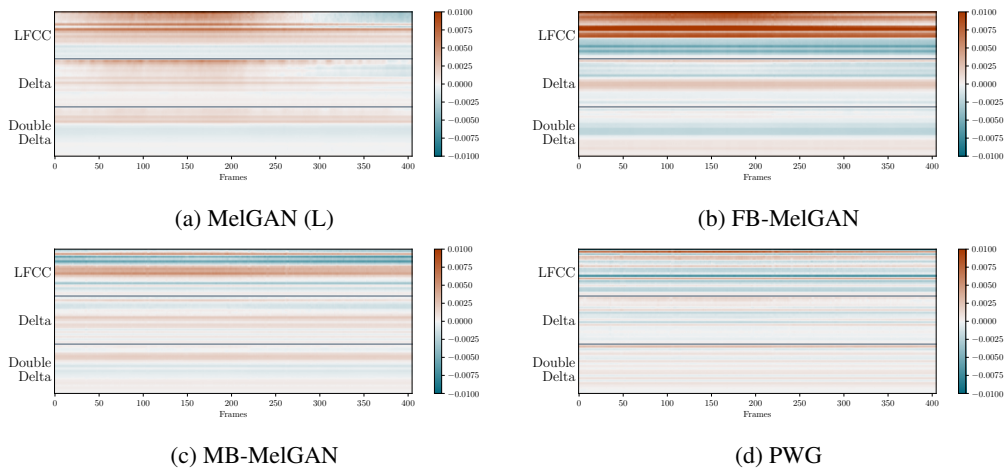


Figure 3: **Attribution of the different models on a real audio sample.** We show the LFCC, delta, and, double delta features. Since we use a linear filter bank, the plot can be read similarly to the spectrogram plots, i.e., features computed over lower frequencies are at the bottom of their respective plots, features over higher frequencies are at the top. Best viewed in color.

308 **Leave-one-out:** Finally, Table 3 present the results of the leave-one-out experiment. We highlight
 309 the distribution which was not present in the training data in bold. While we never train on JSUT, we
 310 only highlight the distribution if the generating network architecture was not part of the training set.
 311 Overall the results improve on the aEER (0.068 \rightarrow 0.045). Also, the generalization results to a novel
 312 setting (JSUT) increase significantly. However, WaveGlow seems to be a key ingredient for good
 313 performance on the JSUT-PWG data and the MelGAN and MelGAN (L) data sets still prove to be a
 314 challenge, even when included in the training set.

315 While these first results are encouraging, there is still much room for improvement. Even the best
 316 performing classifier trained on multiple network architectures has a false acceptance/false rejection
 317 rate of roughly 4.5%.

318 4.3 Attribution

319 Finally, we want to investigate which parts of the audio signal influence the prediction. To this
 320 end, we implemented BlurIG [86], a popular attribution method. We inspect the attribution of four
 321 classifier (MelGAN (L), FB-MelGAN, MB-MelGAN and PWG) for the audio clip used in Section 3.
 322 The results are displayed in Figure 3, full-sized version are available in the supplementary material.
 323 We show the attribution over the LFCC, delta, and, double delta features.

324 Overall, we can see a shift from very broad attention, spread somewhat evenly across all three
 325 feature representations (MelGAN (L)), to a more narrow focused attention across very specific filters
 326 (PWG). MelGAN (L) and FB-MelGAN classifiers operate (mostly) on the higher frequencies, while
 327 MB-MelGAN and PWG also rely on low frequencies for the detection. These observation confirm
 328 our suspicion about the MFCC features. They mask higher frequencies, needed for classifying

329 MelGAN (L) and FB-MelGAN, while over representing lower frequencies, which still leads to a
330 good performance on the MB-MelGAN and PWG data sets. This also explains the significantly
331 better performance of FB-MelGAN on LFCC features, which strikes a balance between all necessary
332 features. The spread out attribution might also explain the poor in-distribution performance of the
333 classifiers trained on the MelGAN variants, since the classifier needs to focus on a broader range of
334 features.

335 All in all we can conclude that high frequencies do provide an overall advantage, but lower frequencies
336 cannot be neglected. Thus, we advice that future classifier operate on the entire spectrum.

337 5 Discussion

338 In this paper we took the first step towards research into audio DeepFakes. While we hope our data
339 set proves useful for future practitioners, there are several limitations to our work:

340 **Evaluating on realistic data:** The difficulties of obtaining realistic data set has been a long standing
341 problem in the security community [73]. Often benign data is readily available, but data actually
342 used in malicious contexts is hard to come by. This leaves us with estimating real-world performance
343 on proxy data. We argue that in our case, we might have good odds that results transfer. As of right
344 know, images generated by off-the-shelf neural networks are used in malicious attempts [9]. We
345 expect the amount of audio DeepFakes to increase as well.

346 **We also abstain from evaluating a complete TTS pipeline. Completely novel audio is not only**
347 **influenced by the vocoder but also by the model generating the intermediate representation. While**
348 **this is an interesting direction for future work, a full evaluation would probably be on the scale of an**
349 **entire new data set.**

350 An additional line of research is automatic speaker verification, which has been studied in the signal
351 processing community [48, 79, 67]. Due to the similarity of the two domain, we expect that results
352 might transfer between the two. Thus, evaluating models on data sets from both domains, might be
353 beneficial.

354 **Adversarial examples and perturbations:** DeepFake-image detectors have already been shown
355 to be vulnerable against adversarial examples [10]. There also exists a myriad of adversarial attacks
356 against automatic speech recognition [11, 70, 92, 70, 5, 71, 2] (Abdullah et al. [1] provide a survey).
357 Thus, classifiers should report their robustness against these attacks and common perturbations (noise,
358 room responses, over-the-air settings, etc.) as part of their evaluation. In this work we focused on
359 providing first steps towards audio DeepFake detection. We leave this questions as an interesting
360 direction for future work.

361 **Variety of the data set:** Our data set presents a first step towards automatic detection of audio
362 DeepFakes. We specifically choose to focus on the LJSPEECH corpus, since it is commonly used for
363 training generative audio models. This allows a one-to-one comparison. However, it only contains
364 recordings by one speaker. While we can make some observation about generalization by comparing
365 against the JSUT data set, a broader analysis focusing on different scenarios would be ideal. We
366 image our corpus being used to study multiple potential classifier designs, evaluating them in a
367 contained environment, before exploring more elaborate settings.

368 6 Conclusion

369 This paper presents a starting point for researchers who want to investigate generated audio signals.
370 We started by presenting a broad overview of signal processing techniques and common feature
371 representations. Then, we introduced a novel data set, with samples from five different state-of-the-art
372 architectures across two languages. In a first analysis, we already discovered subtle differences
373 between the different models, especially among the higher frequencies. To provide a baseline for
374 future practitioners, we trained several baseline models and evaluated their performance across the
375 different data sets. Finally, we inspected the different classifiers by using an attribution method and
376 found that, while high frequency information proved indispensable, lower frequencies cannot be
377 neglected.

References

- 378
- 379 [1] Hadi Abdullah, Kevin Warren, Vincent Bindschaedler, Nicolas Papernot, and Patrick Traynor.
380 SoK: The Faults in our ASRs: An Overview of Attacks against Automatic Speech Recognition
381 and Speaker Identification Systems. In *IEEE Symposium on Security and Privacy (S&P)*, 2020.
- 382 [2] Hojjat Aghakhani, Thorsten Eisenhofer, Lea Schönherr, Dorothea Kolossa, Thorsten Holz,
383 Christopher Kruegel, and Giovanni Vigna. VENOMAVE: Clean-Label Poisoning Against
384 Speech Recognition. *Computing Research Repository (CoRR)*, abs/2010.10682, 2020.
- 385 [3] Muhammad Ejaz Ahmed, Il-Youp Kwak, Jun Ho Huh, Iljoo Kim, Taekkyung Oh, and Hy-
386 ounghshick Kim. Void: A fast and light voice liveness detection system. In *USENIX Security*
387 *Symposium*, 2020.
- 388 [4] Yang Ai and Zhen-Hua Ling. A Neural Vocoder With Hierarchical Generation of Amplitude
389 and Phase Spectra for Statistical Parametric Speech Synthesis. *IEEE/ACM Transactions on*
390 *Audio, Speech, and Language Processing*, 2020.
- 391 [5] Moustafa Alzantot, Bharathan Balaji, and Mani Srivastava. Did you hear that? Adversarial Ex-
392 amples Against Automatic Speech Recognition. In *Advances in Neural Information Processing*
393 *Systems (NeurIPS)*, 2017.
- 394 [6] Sercan Arik, Gregory Diamos, Andrew Gibiansky, John Miller, Kainan Peng, Wei Ping, Jonathan
395 Raiman, and Yanqi Zhou. Deep Voice 2: Multi-Speaker Neural Text-to-Speech. In *Advances in*
396 *Neural Information Processing Systems (NeurIPS)*, 2017.
- 397 [7] Sercan Ö Arik, Mike Chrzanowski, Adam Coates, Gregory Diamos, Andrew Gibiansky, Yong-
398 guo Kang, Xian Li, John Miller, Andrew Ng, Jonathan Raiman, et al. Deep Voice: Real-Time
399 Neural Text-to-Speech. In *International Conference on Machine Learning (ICML)*, 2017.
- 400 [8] Mikołaj Bińkowski, Jeff Donahue, Sander Dieleman, Aidan Clark, Erich Elsen, Norman
401 Casagrande, Luis C Cobo, and Karen Simonyan. High Fidelity Speech Synthesis with Adver-
402 sarial Networks. In *International Conference on Learning Representations (ICLR)*, 2020.
- 403 [9] Matt Burgess. Telegram Still Hasn't Removed an AI Bot That's Abusing Women. *Wired*, 2020.
- 404 [10] Nicholas Carlini and Hany Farid. Evading DeepFake-Image Detectors with White-and Black-
405 Box Attacks. In *IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*,
406 2020.
- 407 [11] Nicholas Carlini and David Wagner. Audio Adversarial Examples: Targeted Attacks on Speech-
408 to-Text. In *IEEE Deep Learning and Security Workshop (DLS)*, 2018.
- 409 [12] Nanxin Chen, Yu Zhang, Heiga Zen, Ron J Weiss, Mohammad Norouzi, and William Chan.
410 WaveGrad: Estimating Gradients for Waveform Generation. In *International Conference on*
411 *Learning Representations (ICLR)*, 2020.
- 412 [13] Akash Chintha, Bao Thai, Saniat Javid Sohrawardi, Kartavya Bhatt, Andrea Hickerson, Matthew
413 Wright, and Raymond Ptucha. Recurrent Convolutional Structures for Audio Spoof and Video
414 DeepFake Detection. *IEEE Journal of Selected Topics in Signal Processing*, 2020.
- 415 [14] Brian Dolhansky, Joanna Bitton, Ben Pflaum, Jikuo Lu, Russ Howes, Menglin Wang, and
416 Cristian Canton Ferrer. The DeepFake Detection Challenge (DFDC) Dataset, 2020.
- 417 [15] Chris Donahue, Julian McAuley, and Miller Puckette. Adversarial Audio Synthesis. In
418 *International Conference on Learning Representations (ICLR)*, 2019.
- 419 [16] Jeff Donahue, Sander Dieleman, Mikołaj Bińkowski, Erich Elsen, and Karen Simonyan. End-
420 to-End Adversarial Text-to-Speech. In *International Conference on Learning Representations*
421 *(ICLR)*, 2021.
- 422 [17] Ricard Durall, Margret Keuper, and Janis Keuper. Watch your Up-Convolution: CNN Based
423 Generative Deep Neural Networks are Failing to Reproduce Spectral Distributions. In *IEEE*
424 *Conference on Computer Vision and Pattern Recognition (CVPR)*, 2020.

- 425 [18] Cristóbal Esteban, Stephanie L Hyland, and Gunnar Rätsch. Real-Valued (Medical) Time
426 Series Generation with Recurrent Conditional GANs. In *International Conference on Learning*
427 *Representations (ICLR)*, 2018.
- 428 [19] Yuchen Fan, Yao Qian, Feng-Long Xie, and Frank K Soong. TTS Synthesis with Bidirectional
429 LSTM Based Recurrent Neural Networks. In *International Conference on Acoustics, Speech*
430 *and Signal Processing (ICASSP)*, 2014.
- 431 [20] Lorenzo Franceschi-Bicchierai. Listen to This Deepfake Audio Impersonating a CEO in Brazen
432 Fraud Attempt. *Motherboard*, 2020.
- 433 [21] Joel Frank, Thorsten Eisenhofer, Lea Schönherr, Asja Fischer, Dorothea Kolossa, and Thorsten
434 Holz. Leveraging Frequency Analysis for Deep Fake Image Recognition. In *International*
435 *Conference on Machine Learning (ICML)*, 2020.
- 436 [22] Maayan Frid-Adar, Idit Diamant, Eyal Klang, Michal Amitai, Jacob Goldberger, and Hayit
437 Greenspan. GAN-Based synthetic Medical Image Augmentation for Increased CNN Perform-
438 ance in Liver Lesion Classification. *Neurocomputing*, 2018.
- 439 [23] Ian J Goodfellow, Jean Pouget-Abadie, Mehdi Mirza, Bing Xu, David Warde-Farley, Sherjil
440 Ozair, Aaron Courville, and Yoshua Bengio. Generative Adversarial Networks. In *Advances in*
441 *Neural Information Processing Systems (NeurIPS)*, 2014.
- 442 [24] Karen Hao. The Biggest Threat of Deepfakes isn't the Deepfakes Themselves. *MIT Technology*
443 *Review*, 2019.
- 444 [25] Tomoki Hayashi. Parallel WaveGAN (+ MelGAN & Multi-band MelGAN) implementation
445 with Pytorch. <https://github.com/kan-bayashi/ParallelWaveGAN>, 2020.
- 446 [26] Keith Ito and Linda Johnson. The LJ Speech Dataset. [https://keithito.com/](https://keithito.com/LJ-Speech-Dataset/)
447 [LJ-Speech-Dataset/](https://keithito.com/LJ-Speech-Dataset/), 2017.
- 448 [27] Nal Kalchbrenner, Erich Elsen, Karen Simonyan, Seb Noury, Norman Casagrande, Edward
449 Lockhart, Florian Stimberg, Aaron Oord, Sander Dieleman, and Koray Kavukcuoglu. Efficient
450 Neural Audio Synthesis. In *International Conference on Machine Learning (ICML)*, 2018.
- 451 [28] Madhu R Kamble, Hemlata Tak, and Hemant A Patil. Effectiveness of Speech Demodulation-
452 Based Features for Replay Detection. In *Proceedings of Interspeech (INTERSPEECH)*, 2018.
- 453 [29] Nathan Killoran, Leo J Lee, Andrew DeLong, David Duvenaud, and Brendan J Frey. Generating
454 and designing DNA with deep generative models. *arXiv preprint arXiv:1712.06148*, 2017.
- 455 [30] Sungwon Kim, Sang-Gil Lee, Jongyoon Song, Jaehyeon Kim, and Sungroh Yoon. FloWaveNet:
456 A Generative Flow for Raw Audio. In *International Conference on Machine Learning (ICML)*,
457 2019.
- 458 [31] Diederik P Kingma and Prafulla Dhariwal. Glow: Generative Flow with Invertible 1x1 Convo-
459 lutions. In *Advances in Neural Information Processing Systems (NeurIPS)*, 2019.
- 460 [32] Diederik P Kingma, Tim Salimans, Rafal Jozefowicz, Xi Chen, Ilya Sutskever, and Max Welling.
461 Improving Variational Inference with Inverse Autoregressive Flow. In *International Conference*
462 *on Learning Representations (ICLR) - Workshop track*, 2016.
- 463 [33] Tomi Kinnunen, Md Sahidullah, Mauro Falcone, Luca Costantini, Rosa González Hautamäki,
464 Dennis Thomsen, Achintya Sarkar, Zheng-Hua Tan, Héctor Delgado, Massimiliano Todisco,
465 Nicholas Evans, Ville Hautamäki, and Kong Aik Lee. RedDots replayed: A new replay spoofing
466 attack corpus for text-dependent speaker verification research. In *International Conference on*
467 *Acoustics, Speech and Signal Processing (ICASSP)*, 2017.
- 468 [34] Zhifeng Kong, Wei Ping, Jiaji Huang, Kexin Zhao, and Bryan Catanzaro. DiffWave: A Versatile
469 Diffusion Model for Audio Synthesis. In *International Conference on Learning Representations*
470 *(ICLR)*, 2021.

- 471 [35] Kundan Kumar, Rithesh Kumar, Thibault de Boissiere, Lucas Gestin, Wei Zhen Teoh, Jose
472 Sotelo, Alexandre de Brébisson, Yoshua Bengio, and Aaron Courville. MelGAN: Generative
473 Adversarial Networks for Conditional Waveform Synthesis. In *Advances in Neural Information
474 Processing Systems (NeurIPS)*, 2019.
- 475 [36] The Atlantic Council’s Digital Forensic Research Lab. Inauthentic Instagram accounts with
476 synthetic faces target Navalny protests. *Medium*, 2021.
- 477 [37] Cheng-I Lai, Alberto Abad, Korin Richmond, Junichi Yamagishi, Najim Dehak, and Simon
478 King. Attentive Filtering Networks for Audio Replay Attack Detection. In *International
479 Conference on Acoustics, Speech and Signal Processing (ICASSP)*, 2019.
- 480 [38] Galina Lavrentyeva, Sergey Novoselov, Egor Malykh, Alexander Kozlov, Oleg Kudashev, and
481 Vadim Shchemelinin. Audio Replay Attack Detection with Deep Learning Frameworks. In
482 *Proceedings of Interspeech (INTERSPEECH)*, 2017.
- 483 [39] Yuezun Li and Siwei Lyu. Exposing DeepFake Videos by Detecting Face Warping Artifacts.
484 *arXiv preprint arXiv:1811.00656*, 2018.
- 485 [40] Yuezun Li, Xin Yang, Pu Sun, Honggang Qi, and Siwei Lyu. Celeb-DF: A New Dataset for
486 DeepFake Forensics. In *Proceedings of the IEEE/CVF Conference on Computer Vision and
487 Pattern Recognition*, 2020.
- 488 [41] Francesco Marra, Diego Gagnaniello, Davide Cozzolino, and Luisa Verdoliva. Detection
489 of GAN-Generated Fake Images over Social Networks. In *IEEE Conference on Multimedia
490 Information Processing and Retrieval (MIPR)*, 2018.
- 491 [42] Francesco Marra, Diego Gagnaniello, Luisa Verdoliva, and Giovanni Poggi. Do GANs Leave
492 Artificial Fingerprints? In *IEEE Conference on Multimedia Information Processing and
493 Retrieval (MIPR)*, 2019.
- 494 [43] Scott McCloskey and Michael Albright. Detecting GAN-Generated Imagery Using Color Cues.
495 *arXiv preprint arXiv:1812.08247*, 2018.
- 496 [44] Soroush Mehri, Kundan Kumar, Ishaan Gulrajani, Rithesh Kumar, Shubham Jain, Jose Sotelo,
497 Aaron Courville, and Yoshua Bengio. SampleRNN: An Unconditional End-to-End Neural
498 Audio Generation Model. *arXiv preprint arXiv:1612.07837*, 2016.
- 499 [45] Huaxiao Mo, Bolin Chen, and Weiqi Luo. Fake Faces Identification via Convolutional Neural
500 Network. In *ACM Workshop on Information Hiding and Multimedia Security*, 2018.
- 501 [46] Peter Mwai. Tigray conflict: The fake UN diplomat and other misleading stories. *BBC Reality
502 Check*, 2021.
- 503 [47] Lakshmanan Nataraj, Tajuddin Manhar Mohammed, BS Manjunath, Shivkumar Chandrasekaran,
504 Arjuna Flenner, Jawadul H Bappy, and Amit K Roy-Chowdhury. Detecting GAN Generated
505 Fake Images Using Co-Occurrence Matrices. *Electronic Imaging*, 2019.
- 506 [48] Andreas Nautsch, Xin Wang, Nicholas Evans, Tomi H. Kinnunen, Ville Vestman, Massimiliano
507 Todisco, Héctor Delgado, Md Sahidullah, Junichi Yamagishi, and Kong Aik Lee. ASVspoof
508 2019: Spoofing Countermeasures for the Detection of Synthesized, Converted and Replayed
509 Speech. *IEEE Transactions on Biometrics, Behavior, and Identity Science*, 2021.
- 510 [49] Aaron van den Oord, Sander Dieleman, Heiga Zen, Karen Simonyan, Oriol Vinyals, Alex
511 Graves, Nal Kalchbrenner, Andrew Senior, and Koray Kavukcuoglu. WaveNet: A Generative
512 Model for Raw Audio. *arXiv preprint arXiv:1609.03499*, 2016.
- 513 [50] Aaron van den Oord, Oriol Vinyals, and Koray Kavukcuoglu. Neural Discrete Representation
514 Learning. In *Advances in Neural Information Processing Systems (NeurIPS)*, 2017.
- 515 [51] Adam Paszke, Sam Gross, Francisco Massa, Adam Lerer, James Bradbury, Gregory Chanan,
516 Trevor Killeen, Zeming Lin, Natalia Gimelshein, Luca Antiga, Alban Desmaison, Andreas
517 Kopf, Edward Yang, Zachary DeVito, Martin Raison, Alykhan Tejani, Sasank Chilamkurthy,
518 Benoit Steiner, Lu Fang, Junjie Bai, and Soumith Chintala. PyTorch: An Imperative Style,
519 High-Performance Deep Learning Library. In *Advances in Neural Information Processing
520 Systems (NeurIPS)*, 2019.

- 521 [52] Hemant A Patil, Madhu R Kamble, Tanvina B Patel, and Meet H Soni. Novel Variable Length
522 Teager Energy Separation Based Instantaneous Frequency Features for Replay Detection. In
523 *Proceedings of Interspeech (INTERSPEECH)*, 2017.
- 524 [53] Kainan Peng, Wei Ping, Zhao Song, and Kexin Zhao. Non-Autoregressive Neural Text-to-
525 Speech. In *International Conference on Machine Learning (ICML)*, 2020.
- 526 [54] Wei Ping, Kainan Peng, and Jitong Chen. Clarinet: Parallel Wave Generation in End-to-End
527 Text-to-Speech. In *International Conference on Learning Representations (ICLR)*, 2019.
- 528 [55] Wei Ping, Kainan Peng, Kexin Zhao, and Zhao Song. WaveFlow: A Compact Flow-based
529 Model for Raw Audio. In *International Conference on Learning Representations (ICLR)*, 2020.
- 530 [56] Daniel Povey, Arnab Ghoshal, Gilles Boulianne, Lukas Burget, Ondrej Glembek, Nagendra
531 Goel, Mirko Hannemann, Petr Motlicek, Yanmin Qian, Petr Schwarz, Jan Silovsky, Georg
532 Stemmer, and Karel Vesely. The Kaldi Speech Recognition Toolkit. In *IEEE Workshop on*
533 *Automatic Speech Recognition and Understanding (ASRU)*, 2011.
- 534 [57] KM Muraleedhara Prabhu. *Window Functions and their Applications in Signal Processing*.
535 Taylor & Francis, 2014.
- 536 [58] Ryan Prenger, Rafael Valle, and Bryan Catanzaro. WaveGlow: a Flow-based
537 Generative Network for Speech Synthesis. [https://pytorch.org/hub/nvidia_](https://pytorch.org/hub/nvidia_deeplearningexamples_waverglow/)
538 [deeplearningexamples_waverglow/](https://pytorch.org/hub/nvidia_deeplearningexamples_waverglow/), 2018.
- 539 [59] Ryan Prenger, Rafael Valle, and Bryan Catanzaro. WaveGlow: a Flow-based Generative
540 Network for Speech Synthesis. <https://github.com/NVIDIA/waverglow>, 2018.
- 541 [60] Ryan Prenger, Rafael Valle, and Bryan Catanzaro. Waveglow: A Flow-based Generative
542 Network for Speech Synthesis. In *International Conference on Acoustics, Speech and Signal*
543 *Processing (ICASSP)*, 2019.
- 544 [61] Yuyang Qian, Guojun Yin, Lu Sheng, Zixuan Chen, and Jing Shao. Thinking in Frequency: Face
545 Forgery Detection by Mining Frequency-Aware Clues. In *European Conference on Computer*
546 *Vision (ECCV)*, 2020.
- 547 [62] Thomas Quatieri. *Discrete-Time Speech Signal Processing: Principles and Practice*. Pearson
548 Education India, 2006.
- 549 [63] Lawrence Rabiner, Bernard Gold, and CK Yuen. *Theory and Application of Digital Signal*
550 *Processing*. Prentice-Hall, 2016.
- 551 [64] Yi Ren, Yangjun Ruan, Xu Tan, Tao Qin, Sheng Zhao, Zhou Zhao, and Tie-Yan Liu. FastSpeech:
552 Fast, Robust and Controllable Text to Speech. In *Advances in Neural Information Processing*
553 *Systems (NeurIPS)*, 2019.
- 554 [65] Yi Ren, Chenxu Hu, Xu Tan, Tao Qin, Sheng Zhao, Zhou Zhao, and Tie-Yan Liu. FastSpeech 2:
555 Fast and High-Quality End-to-End Text to Speech. *arXiv preprint arXiv:2006.04558*, 2020.
- 556 [66] Andreas Rossler, Davide Cozzolino, Luisa Verdoliva, Christian Riess, Justus Thies, and Matthias
557 Nießner. Faceforensics++: Learning to Detect Manipulated Facial Images. In *Proceedings of*
558 *the IEEE/CVF International Conference on Computer Vision*, 2019.
- 559 [67] Md Sahidullah, Tomi Kinnunen, and Cemal Hanilçi. A Comparison of Features for Synthetic
560 Speech Detection. In *Proceedings of Interspeech (INTERSPEECH)*, 2015.
- 561 [68] Dirk Scheuermann, Scarlet Schwiderski-Grosche, and Bruno Struif. *Usability of Biometrics*
562 *in Relation to Electronic Signatures*. GMD-Forschungszentrum Informationstechnik Sankt
563 Augustin, 2000.
- 564 [69] Lea Schönherr, Steffen Zeiler, and Dorothea Kolossa. Spoofing Detection via Simultaneous
565 Verification of Audio-Visual Synchronicity and Transcription. In *2017 IEEE Automatic Speech*
566 *Recognition and Understanding Workshop (ASRU)*, 2017.

- 567 [70] Lea Schönherr, Katharina Kohls, Steffen Zeiler, Thorsten Holz, and Dorothea Kolossa. Adversarial Attacks Against Automatic Speech Recognition Systems via Psychoacoustic Hiding. In
568 *Symposium on Network and Distributed System Security (NDSS)*, 2019.
569
- 570 [71] Lea Schönherr, Thorsten Eisenhofer, Steffen Zeiler, Thorsten Holz, and Dorothea Kolossa. Imperio: Robust Over-the-Air Adversarial Examples for Automatic Speech Recognition Systems. In
571 *Annual Computer Security Applications Conference (ACSAC)*, 2020.
572
- 573 [72] Jonathan Shen, Ruoming Pang, Ron J Weiss, Mike Schuster, Navdeep Jaitly, Zongheng Yang, Zhifeng Chen, Yu Zhang, Yuxuan Wang, Rj Skerrv-Ryan, et al. Natural TTS Synthesis by
574 Conditioning WaveNet on Mel Spectrogram Predictions. In *International Conference on*
575 *Acoustics, Speech and Signal Processing (ICASSP)*, 2018.
576
- 577 [73] Robin Sommer and Vern Paxson. Outside the Closed World: On Using Machine Learning for
578 Network Intrusion Detection. In *IEEE Symposium on Security and Privacy (S&P)*, 2010.
- 579 [74] Ryosuke Sonobe, Shinnosuke Takamichi, and Hiroshi Saruwatari. JSUT Corpus: Free
580 Large-Scale Japanese Speech Corpus for End-to-End Speech Synthesis. *arXiv preprint*
581 *arXiv:1711.00354*, 2017.
- 582 [75] Jose Sotelo, Soroush Mehri, Kundan Kumar, Joao Felipe Santos, Kyle Kastner, Aaron Courville,
583 and Yoshua Bengio. Char2wav: End-to-End Speech Synthesis. In *International Conference on*
584 *Learning Representations (ICLR) Workshop Track*, 2017.
- 585 [76] Catherine Stupp. Fraudsters Used AI to Mimic CEO’s Voice in Unusual Cybercrime Case. *The*
586 *Wall Street Journal*, 2019.
- 587 [77] Yaniv Taigman, Lior Wolf, Adam Polyak, and Eliya Nachmani. VoiceLoop: Voice Fitting and
588 Synthesis via a Phonological Loop. In *International Conference on Learning Representations*
589 *(ICLR)*, 2017.
- 590 [78] Shahroz Tariq, Sangyup Lee, Hoyoung Kim, Youjin Shin, and Simon S Woo. GAN is a Friend
591 or Foe? A Framework to Detect Various Fake Face Images. In *ACM/SIGAPP Symposium on*
592 *Applied Computing*, 2019.
- 593 [79] Massimiliano Todisco, Xin Wang, Ville Vestman, Md Sahidullah, Hector Delgado, Andreas
594 Nautsch, Junichi Yamagishi, Nicholas Evans, Tomi Kinnunen, and Kong Aik Lee. ASVspoof
595 2019: Future Horizons in Spoofed and Fake Audio Detection. *Computing Research Repository*
596 *(CoRR)*, abs/1904.05441, 2019.
- 597 [80] Keiichi Tokuda, Takayoshi Yoshimura, Takashi Masuko, Takao Kobayashi, and Tadashi Ki-
598 tamura. Speech Parameter Generation Algorithms for HMM-Based Speech Synthesis. In
599 *International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, 2000.
- 600 [81] Francis Tom, Mohit Jain, and Prasenjit Dey. End-To-End Audio Replay Attack Detection Using
601 Deep Convolutional Networks with Attention. In *Proceedings of Interspeech (INTERSPEECH)*,
602 2018.
- 603 [82] Rafael Valle, Wilson Cai, and Anish Doshi. TequilaGAN: How to Easily Identify GAN Samples.
604 *arXiv preprint arXiv:1807.04919*, 2018.
- 605 [83] Sheng-Yu Wang, Oliver Wang, Richard Zhang, Andrew Owens, and Alexei A Efros. CNN-
606 generated images are surprisingly easy to spot... for now. In *IEEE Conference on Computer*
607 *Vision and Pattern Recognition (CVPR)*, 2020.
- 608 [84] Xin Wang, Shinji Takaki, and Junichi Yamagishi. Neural Source-Filter Waveform Models
609 for Statistical Parametric Speech Synthesis. *IEEE/ACM Transactions on Audio, Speech, and*
610 *Language Processing*, 2019.
- 611 [85] Zhizheng Wu, Junichi Yamagishi, Tomi Kinnunen, Cemal Hanilçi, Mohammed Sahidullah,
612 Aleksandr Sizov, Nicholas Evans, Massimiliano Todisco, and Héctor Delgado. ASVspoof: The
613 Automatic Speaker Verification Spoofing and Countermeasures Challenge. *IEEE Journal of*
614 *Selected Topics in Signal Processing*, 2017.

- 615 [86] Shawn Xu, Subhashini Venugopalan, and Mukund Sundararajan. Attribution in Scale and Space.
616 In *IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2020.
- 617 [87] Yuan Xue, Tao Xu, Han Zhang, L Rodney Long, and Xiaolei Huang. SeGAN: Adversarial
618 Network with Multi-Scale L1 Loss for Medical Image Segmentation. *Neuroinformatics*, 2018.
- 619 [88] Ryuichi Yamamoto, Eunwoo Song, and Jae-Min Kim. Parallel WaveGAN: A Fast Waveform
620 Generation Model Based on Generative Adversarial Networks with Multi-Resolution Spec-
621 trogram. In *International Conference on Acoustics, Speech and Signal Processing (ICASSP)*,
622 2020.
- 623 [89] Geng Yang, Shan Yang, Kai Liu, Peng Fang, Wei Chen, and Lei Xie. Multi-Band Melgan:
624 Faster Waveform Generation For High-Quality Text-To-Speech. In *2021 IEEE Spoken Language
625 Technology Workshop (SLT)*, 2021.
- 626 [90] Takayoshi Yoshimura, Keiichi Tokuda, Takashi Masuko, Takao Kobayashi, and Tadashi Ki-
627 tamura. Simultaneous Modeling of Spectrum, Pitch and Duration in HMM-Based Speech
628 Synthesis. In *Sixth European Conference on Speech Communication and Technology*, 1999.
- 629 [91] Ning Yu, Larry S Davis, and Mario Fritz. Attributing Fake Images to GANs: Learning and
630 Analyzing GAN Fingerprints. In *IEEE International Conference on Computer Vision (ICCV)*,
631 2019.
- 632 [92] Xuejing Yuan, Yuxuan Chen, Yue Zhao, Yunhui Long, Xiaokang Liu, Kai Chen, Shengzhi
633 Zhang, Heqing Huang, Xiaofeng Wang, and Carl A. Gunter. CommanderSong: A Systematic
634 Approach for Practical Adversarial Voice Recognition. In *USENIX Security Symposium*, 2018.
- 635 [93] Heiga Ze, Andrew Senior, and Mike Schuster. Statistical Parametric Speech Synthesis Using
636 Deep Neural Networks. In *International Conference on Acoustics, Speech and Signal Processing
637 (ICASSP)*, 2013.
- 638 [94] Heiga Zen, Keiichi Tokuda, and Alan W Black. Statistical Parametric Speech Synthesis. *Speech
639 Communication*, 2009.
- 640 [95] Linghan Zhang, Sheng Tan, Jie Yang, and Yingying Chen. Voicelive: A Phoneme Localization
641 Based Liveness Detection for Voice Authentication on Smartphones. In *ACM Conference on
642 Computer and Communications Security (CCS)*, 2016.
- 643 [96] Xu Zhang, Svebor Karaman, and Shih-Fu Chang. Detecting and Simulating Artifacts in GAN
644 Fake Images. In *2019 IEEE International Workshop on Information Forensics and Security
645 (WIFS)*, 2019.
- 646 [97] Eberhard Zwicker and Hugo Fastl. *Psychoacoustics: Facts and Models*. Springer, Heidelberg,
647 Germany, 2 edition, 2007.

648 **Checklist**

- 649 1. For all authors...
- 650 (a) Do the main claims made in the abstract and introduction accurately reflect the paper's
651 contributions and scope? [Yes]
- 652 (b) Did you describe the limitations of your work? [Yes] See Section 5
- 653 (c) Did you discuss any potential negative societal impacts of your work? [Yes] See
654 Section 3
- 655 (d) Have you read the ethics review guidelines and ensured that your paper conforms to
656 them? [Yes]
- 657 2. If you are including theoretical results...
- 658 (a) Did you state the full set of assumptions of all theoretical results? [N/A]
- 659 (b) Did you include complete proofs of all theoretical results? [N/A]
- 660 3. If you ran experiments (e.g. for benchmarks)...
- 661 (a) Did you include the code, data, and instructions needed to reproduce the main ex-
662 perimental results (either in the supplemental material or as a URL)? [Yes] For the
663 submission we will include them in the supplementary material. For the publication
664 we will move the code to a public GitHub repository. The data can be found on zenodo,
665 see Section 3
- 666 (b) Did you specify all the training details (e.g., data splits, hyperparameters, how they
667 were chosen)? [Yes] See Section 4.1
- 668 (c) Did you report error bars (e.g., with respect to the random seed after running experi-
669 ments multiple times)? [No] We did chose the results over multiple runs and controlled
670 for the data set size, see Section 4.1.
- 671 (d) Did you include the total amount of compute and the type of resources used (e.g., type
672 of GPUs, internal cluster, or cloud provider)? [Yes] See Section 4.1
- 673 4. If you are using existing assets (e.g., code, data, models) or curating/releasing new assets...
- 674 (a) If your work uses existing assets, did you cite the creators? [Yes] See Section 3
- 675 (b) Did you mention the license of the assets? [Yes] See Section 3
- 676 (c) Did you include any new assets either in the supplemental material or as a URL? [Yes]
677 See Section 3
- 678 (d) Did you discuss whether and how consent was obtained from people whose data you're
679 using/curating? [Yes] See Section 3
- 680 (e) Did you discuss whether the data you are using/curating contains personally identifiable
681 information or offensive content? [Yes] See Section 3
- 682 5. If you used crowdsourcing or conducted research with human subjects...
- 683 (a) Did you include the full text of instructions given to participants and screenshots, if
684 applicable? [N/A]
- 685 (b) Did you describe any potential participant risks, with links to Institutional Review
686 Board (IRB) approvals, if applicable? [N/A]
- 687 (c) Did you include the estimated hourly wage paid to participants and the total amount
688 spent on participant compensation? [N/A]