
Datasheet - 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

1 Motivation

2 The main purpose of this data set is to facilitate research into audio DeepFakes. These generated
3 media files have been increasingly used to commit impersonation attempts [2], influencing opposition
4 movements [8] to justify military actions [3] or online harassment [1]. We hope that this work helps
5 in finding new detection methods to prevent such attempts. The creation of this data set was supported
6 by the Deutsche Forschungsgemeinschaft (DFG, German Research Foundation) under Germany's
7 Excellence Strategy–EXC-2092 CASA–390781972.

8 The data set is distributed through zenodo² with a CC-BY-SA 4.0 license.

9 Composition

10 The data set consists of 88,600 generated audio clips (16-bit PCM wav) in total. We examine
11 multiple networks trained on two reference data sets. First, the LJSPEECH [5] data set consisting
12 of 13,100 short audio clips (on average 6 seconds each; roughly 24 hours total) read by a female
13 speaker. It features passages from 7 non-fiction books and the audio was recorded on a MacBook Pro
14 microphone. Second, we include samples based on the JSUT [13] data set, specifically, basic5000
15 corpus. This corpus consists of 5,000 sentences covering all basic kanji of the Japanese language
16 (4.8 seconds on average; roughly 6.7 hours total). The recordings were performed by a female native
17 Japanese speaker in an anechoic room. Thus, our data set consists of approximately 157 hours of
18 generated audio files in total. Note that we do not redistribute the reference data. They are freely
19 available online [5, 13].

20 We included a range of architectures in our data set:

- 21 • **MelGAN**: We include MelGAN [7], which is one of the first GAN-based generative models
22 for audio data. It uses fully convolutional feed-forward network as generator and operates on
23 Mel spectrograms. Their discriminator is a combination of three different discriminators that
24 operates on the original, and two downsampled versions of the raw audio input. Additionally,
25 they use an auxiliary loss over the feature space of the three discriminators.
- 26 • **Parallel WaveGAN (PWG)**: WaveNet [9] is one of the earliest and most common archi-
27 tectures, We include samples from one of its variants, Parallel WaveGAN [14]. It uses
28 GAN training paradigm, with a non-autoregressive version of WaveNet as its generator. In a
29 similar vein to MelGAN, it uses an auxiliary loss, but in contrast, matches the *Short-Time*
30 *Fourier Transform* (STFT) of the original training sample and the generated waveform over
31 multiple resolutions.
- 32 • **Multi-band MelGAN (MB-MelGAN)**: Incorporating more fine-grained frequency analysis,
33 might lead to more convincing samples. We include MB-MelGAN, which computes its
34 auxiliary (frequency-based; inspired by PWG) loss in different sub-bands. Its generator is

*Corresponding author joel.frank@rub.de.

² zenodo.org/record/4904579 - DOI: 10.5281/zenodo.4904579

35 based on a bigger version of the MelGAN generator, but instead of predicting the entire
36 audio directly, the generator produces multiple sub-bands, which are then summed up to the
37 complete audio signal.

- 38 • **Full-band MelGAN (FB-MelGAN)**: We include a variant of MB-MelGAN which gener-
39 ates the complete audio directly and computes its auxiliary loss (the same as PWG) over the
40 full audio instead of its sub-bands.
- 41 • **WaveGlow**: The training procedure might also influence the detectability of fake samples.
42 Therefore, we include samples from WaveGlow to investigate maximum-likelihood-based
43 methods. It is a flow-based generative model based on Glow [6], whose architecture is
44 heavily inspired by WaveNet.

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

50 Collection Process

51 For WaveGlow, we utilize the official implementation [11] (commit 8afb643) in conjunction with
52 the official pre-trained network on PyTorch Hub [10]. We use a popular implementation available
53 on GitHub [4] (commit 12c677e) for the remaining networks. The repository also offers pre-trained
54 models. We used the pre-trained networks to generate samples that are similar to their respective
55 training distributions, LJSPEECH [5] and JSUT [13]. When sampling the data set, we first extract
56 Mel spectrograms from the original audio files, using the pre-processing scripts of the corresponding
57 repositories. We then feed these Mel spectrograms to the respective models to obtain the data set.

58 Uses & Ethical Considerations

59 The intended use of this data set is to facilitate research into detecting audio DeepFakes. Our data
60 set consists of phrases from non-fiction books (LJSPEECH) and everyday conversational Japanese
61 (JSUT), which are already available online. The same is true for all models used to generate this data
62 set. Thus, we cannot think of an immediate way to misuse our data. On the contrary, we hope it can
63 accelerate research into malicious usage of generative models that already cause damage to society.

64 One might wonder if releasing research into detecting DeepFakes might contribute negatively towards
65 the detection "arms race". This is a long standing debate in the security community and the overall
66 consensus is that "security through obscurity" does not work. This is also often echoed in best security
67 practices, for example, published by the National Institute of Standards and Technology (NIST) [12].
68 Intuitively, withholding information from the research community is in-fact more harmful, since
69 attackers will eventually adapt to any defense one deploys.

70 Distribution & Licensing

71 The LJSPEECH data set is in the public domain. The JSUT corpus is licensed by CC-BY-SA 4.0, with
72 a note that redistribution is only permitted in certain cases. We contacted the author, who saw no
73 conflict in distributing our fake samples, as long as its for research purposes. To comply with JSUT
74 we license our data set under the CC-BY-SA 4.0 license.

75 We do not redistribute any models or training distributions and bear all responsibility in case of
76 violation of rights, etc.

77 References

- 78 [1] Matt Burgess. Telegram Still Hasn't Removed an AI Bot That's Abusing Women. *Wired*, 2020.
- 79 [2] Lorenzo Franceschi-Bicchierai. Listen to This Deepfake Audio Impersonating a CEO in Brazen
80 Fraud Attempt. *Motherboard*, 2020.
- 81 [3] Karen Hao. The Biggest Threat of Deepfakes isn't the Deepfakes Themselves. *MIT Technology*
82 *Review*, 2019.

- 83 [4] Tomoki Hayashi. Parallel WaveGAN (+ MelGAN & Multi-band MelGAN) implementation
84 with Pytorch. <https://github.com/kan-bayashi/ParallelWaveGAN>, 2020.
- 85 [5] Keith Ito and Linda Johnson. The LJ Speech Dataset. [https://keithito.com/
86 LJ-Speech-Dataset/](https://keithito.com/LJ-Speech-Dataset/), 2017.
- 87 [6] Diederik P Kingma and Prafulla Dhariwal. Glow: Generative Flow with Invertible 1x1 Convo-
88 lutions. In *Advances in Neural Information Processing Systems (NeurIPS)*, 2019.
- 89 [7] Kundan Kumar, Rithesh Kumar, Thibault de Boissiere, Lucas Gestin, Wei Zhen Teoh, Jose
90 Sotelo, Alexandre de Brébisson, Yoshua Bengio, and Aaron Courville. MelGAN: Generative
91 Adversarial Networks for Conditional Waveform Synthesis. In *Advances in Neural Information
92 Processing Systems (NeurIPS)*, 2019.
- 93 [8] The Atlantic Council’s Digital Forensic Research Lab. Inauthentic Instagram accounts with
94 synthetic faces target Navalny protests. *Medium*, 2021.
- 95 [9] Aaron van den Oord, Sander Dieleman, Heiga Zen, Karen Simonyan, Oriol Vinyals, Alex
96 Graves, Nal Kalchbrenner, Andrew Senior, and Koray Kavukcuoglu. WaveNet: A Generative
97 Model for Raw Audio. *arXiv preprint arXiv:1609.03499*, 2016.
- 98 [10] Ryan Prenger, Rafael Valle, and Bryan Catanzaro. WaveGlow: a Flow-based
99 Generative Network for Speech Synthesis. [https://pytorch.org/hub/nvidia_
100 deeplearningexamples_waverglow/](https://pytorch.org/hub/nvidia_100_deeplearningexamples_waverglow/), 2018.
- 101 [11] Ryan Prenger, Rafael Valle, and Bryan Catanzaro. WaveGlow: a Flow-based Generative
102 Network for Speech Synthesis. <https://github.com/NVIDIA/waverglow>, 2018.
- 103 [12] Karen Scarfone, Wayne Jansen, Miles Tracy, et al. Guide to General Server Security. *NIST
104 Special Publication*, 2008.
- 105 [13] Ryosuke Sonobe, Shinnosuke Takamichi, and Hiroshi Saruwatari. JSUT Corpus: Free
106 Large-Scale Japanese Speech Corpus for End-to-End Speech Synthesis. *arXiv preprint
107 arXiv:1711.00354*, 2017.
- 108 [14] Ryuichi Yamamoto, Eunwoo Song, and Jae-Min Kim. Parallel WaveGAN: A Fast Waveform
109 Generation Model Based on Generative Adversarial Networks with Multi-Resolution Spec-
110 trogram. In *International Conference on Acoustics, Speech and Signal Processing (ICASSP)*,
111 2020.