

---

# Supplementary Material - 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 In this supplementary material we provide full size spectrogram and attribution plots. All the plots  
2 are in reference to the same audio sample (LJSPEECH 008-0217). For attribution, we applied BlurIG  
3 directly on the feature vectors. Additionally, we provide a visual representations of the filterbanks  
4 used and a discussion on releasing security relevant research.

## 5 **A note on releasing security research**

6 One might wonder if releasing research into detecting DeepFakes might contribute negatively towards  
7 the detection "arms race". This is a long standing debate in the security community and the overall  
8 consensus is that "security through obscurity" does not work. This is often echoed in best security  
9 practices, for example, published by the National Institute of Standards and Technology (NIST) [8].  
10 Intuitively, withholding information from the research community is in-fact more harmful, since  
11 attackers will eventually adapt to any defense one deploys anyway. Thus, contributing to the invention  
12 of new systems is more helpful in an ever changing environment [6].

13 The debate dates back to at least the 19th century where the cryptographer Auguste Kerckhoffs  
14 introduced Kerckhoffs's principle [3]. The principle states that an encryption scheme should still  
15 work if an adversary knows everything about the system but a secret passphrase. Similar thought  
16 would later be formulated by Claude Shannon [9].

17 A typical example is the advanced encryption standard (AES). The algorithm's entire specification  
18 and inner workings can be found in the standardization [7]. Yet, it is considered unbreakable as  
19 long as the password used for the encryption is not revealed. AES is also the only algorithm used to  
20 encrypt US government documents [1]. The principle also found adoption in the machine learning  
21 community, where adversarial defense papers are now advised to evaluate against so-called white box  
22 attackers [2], i.e., attackers which know the inner workings of the system and actively try to avoid it.

23 While complete openness is obviously not possible, the greater security community has adapted  
24 similar practices. For example, so-called attack papers are regularly published at security venues.  
25 The underlying motivation being, that before one can protect systems, one has to understand how to  
26 attack them. Prominent examples are the Meltdown [5] and Spectre [4] vulnerabilities which showed  
27 that certain instructions in CPUs could be used for unauthorized access.

28 Similar patterns are also used in the industry. Google's project zero team regularly analyses and  
29 finds critical vulnerabilities in commonly used software. Their standard practice is to inform the  
30 vendor and work with them to help fix the vulnerability. However, after a hard deadline of 90 days,  
31 the details of the vulnerability will be released to the public [10]. The effects are two-fold. First, the  
32 deadline encourages faster patch development by the vendor. Second, the techniques used can be  
33 studied to prevent similar vulnerabilities in the future.

---

\*Corresponding author [joel.frank@rub.de](mailto:joel.frank@rub.de).

## 34 Spectrograms

35 Here we plot the spectrograms of an audio file (LJSPEECH 008-0217) for the training data and the  
36 different generative networks. Notice the differences especially in the higher frequencies and the  
37 horizontal artifacts produced by MelGAN and WaveGlow.

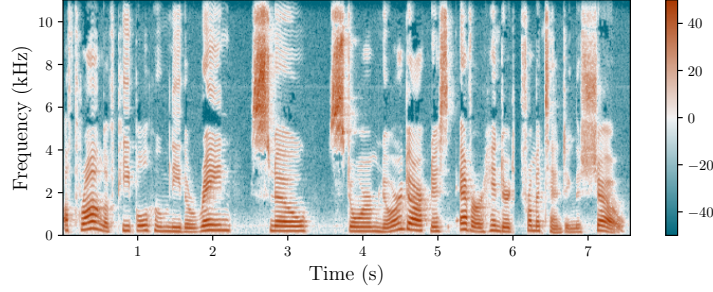


Figure 1: Original

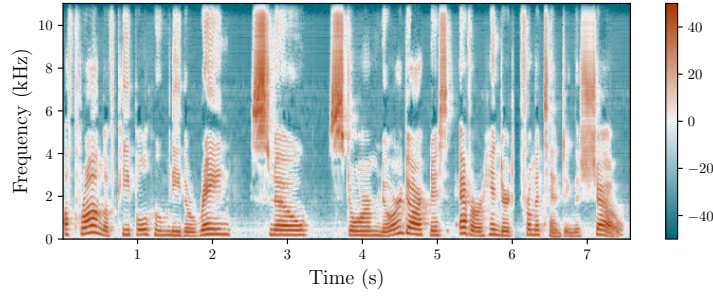


Figure 2: MelGAN

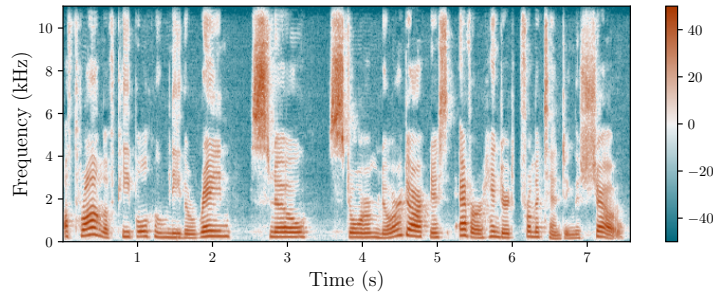


Figure 3: FB-MelGAN

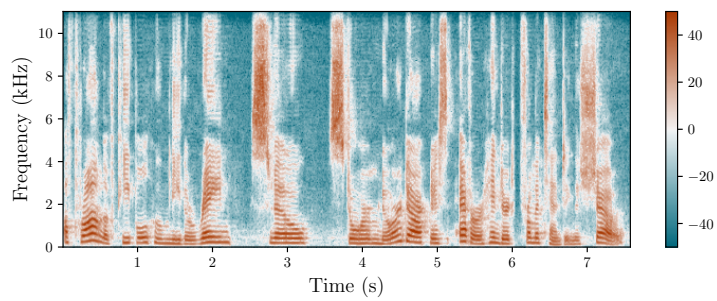


Figure 4: MB-MelGAN

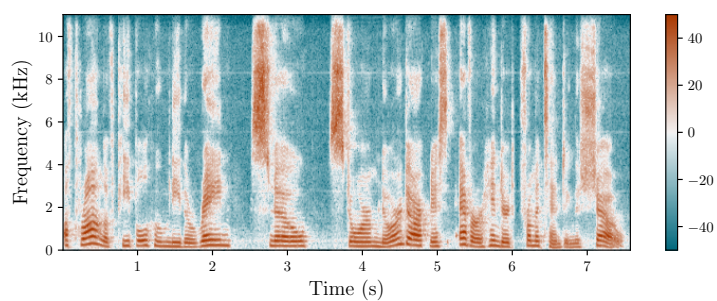


Figure 5: WaveGlow

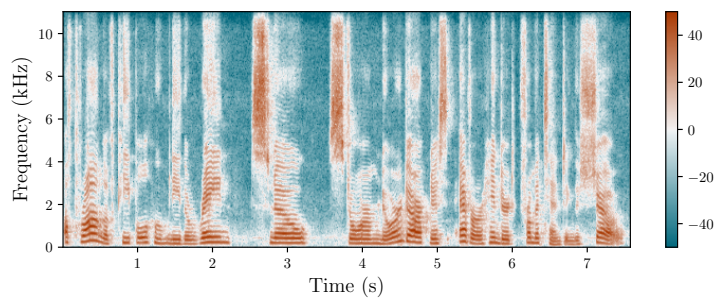
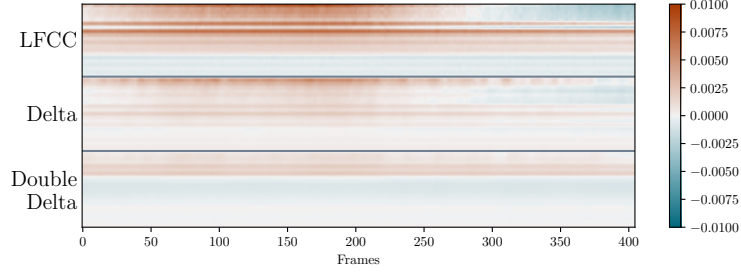


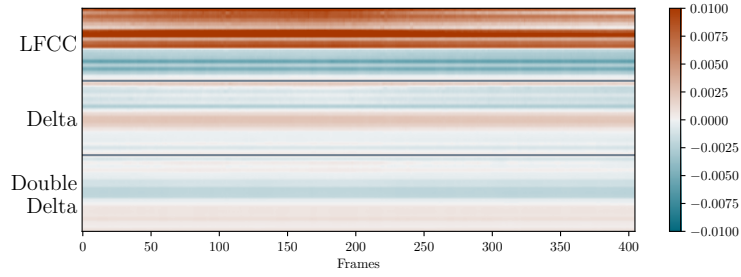
Figure 6: PWG

## 38 Attribution

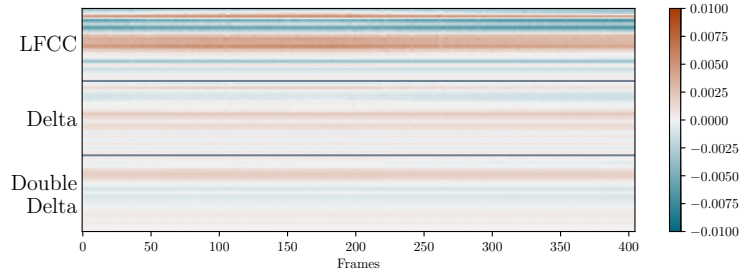
39 These are the full-size version of the attribution plots used in Section 4.3. Note the spread out  
40 attention of the MelGAN classifier, the transition to narrow band attribution and the balance of the  
41 classifier trained on FB-MelGAN.



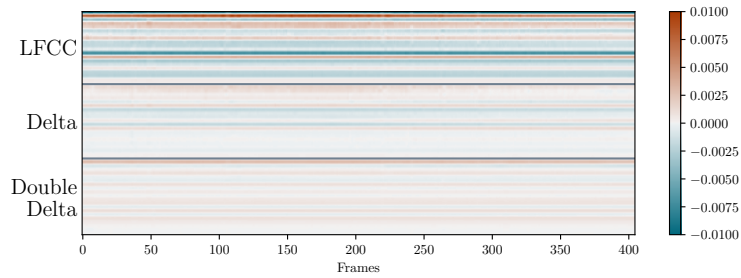
(a) MelGAN (L)



(b) FB-MelGAN



(c) MB-MelGAN



(d) PWG

## 42 Filterbanks

43 Here we show a visual representation of the triangular filterbanks used to compute the *Mel Frequency*  
 44 *Cepstral Coefficients* (MFCC) and *Linear Frequency Cepstral Coefficients* (LFCC) features.

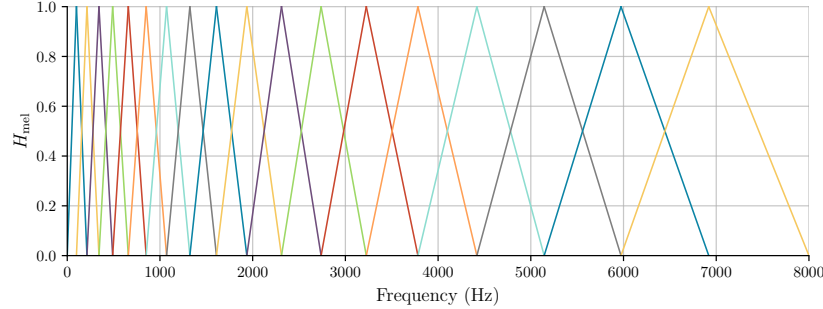


Figure 8: Mel filterbank

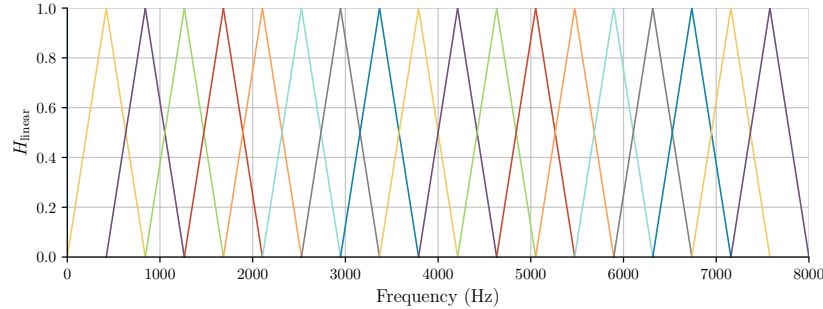


Figure 9: Linear filterbank

## 45 References

- 46 [1] Elaine Barker et al. Guideline for Using Cryptographic Standards in the Federal Government:  
 47 Cryptographic Mechanisms. *NIST special publication*, 2016.
- 48 [2] Nicholas Carlini, Anish Athalye, Nicolas Papernot, Wieland Brendel, Jonas Rauber, Dimitris  
 49 Tsipras, Ian Goodfellow, and Aleksander Madry. On Evaluating Adversarial Robustness.  
 50 *Computing Research Repository (CoRR)*, abs/1902.06705, 2019.
- 51 [3] Auguste Kerckhoffs. *La cryptographie militaire, ou, Des chiffres usités en temps de guerre:*  
 52 *avec un nouveau procédé de déchiffrement applicable aux systèmes à double clef*. Librairie  
 53 militaire de L. Baudoin, 1883.
- 54 [4] Paul Kocher, Jann Horn, Anders Fogh, , Daniel Genkin, Daniel Gruss, Werner Haas, Mike  
 55 Hamburg, Moritz Lipp, Stefan Mangard, Thomas Prescher, Michael Schwarz, and Yuval Yarom.  
 56 Spectre attacks: Exploiting speculative execution. In *40th IEEE Symposium on Security and*  
 57 *Privacy (S&P'19)*, 2019.
- 58 [5] Moritz Lipp, Michael Schwarz, Daniel Gruss, Thomas Prescher, Werner Haas, Anders Fogh,  
 59 Jann Horn, Stefan Mangard, Paul Kocher, Daniel Genkin, Yuval Yarom, and Mike Hamburg.  
 60 Meltdown: Reading kernel memory from user space. In *27th USENIX Security Symposium*  
 61 *(USENIX Security 18)*, 2018.
- 62 [6] Bill McCarty. The honeynet arms race. *IEEE Security & Privacy*, 2003.

- 63 [7] Vincent Rijmen and Joan Daemen. Advanced Encryption Standard. *Proceedings of Federal*  
64 *Information Processing Standards Publications, National Institute of Standards and Technology,*  
65 2001.
- 66 [8] Karen Scarfone, Wayne Jansen, Miles Tracy, et al. Guide to General Server Security. *NIST*  
67 *Special Publication*, 2008.
- 68 [9] Claude E Shannon. Communication Theory of Secrecy Systems. *The Bell system technical*  
69 *journal*, 1949.
- 70 [10] Tim Willis. Project Zero Policy and Disclosure: 2020 Edi-  
71 tion, 2020. [https://googleprojectzero.blogspot.com/2020/01/](https://googleprojectzero.blogspot.com/2020/01/policy-and-disclosure-2020-edition.html)  
72 [policy-and-disclosure-2020-edition.html](https://googleprojectzero.blogspot.com/2020/01/policy-and-disclosure-2020-edition.html), as of July 11, 2021.