VoiceRadar: Voice Deepfake Detection using Micro-Frequency and Compositional Analysis

Kavita Kumari, Maryam Abbasihafshejani, Alessandro Pegoraro, Phillip Rieger, Kamyar Arshi, Murtuza Jadliwala, Ahmad-Reza Sadeghi

NDSS 2025





LAB

TECHNISCHE UNIVERSITÄT DARMSTADT



The University of Texas at San Antonio™



Audio Deepfakes

Fraud

OTREND Business

Unusual CEO Fraud via Deepfake Audio Steals US\$243,000 From UK Company

September 05, 2019

Audio Deepfakes

Fraud

OTREND Business

Unusual CEO Fraud via Deepfake Audio Steals US\$243,000 From UK Company

September 05, 2019

Disinformation

E Politics SCOTUS Congress Facts First 2024 Elections

A fake recording of a candidate saying he'd rigged the election went viral. Experts say it's only the beginning

Audio Deepfakes

Defamation AP

U.S. NEWS

Deepfake of principal's voice is the latest case of AI being used for harm

 \bigcirc



The Audio Generation Arms Race



Voice Engine (Voice Cloning with 15 sec of samples) [Open Al Blog]

VALLE-A (Text-to-speech voice generation using Language Modelling) [Meng et al. arXiv 2024]

Singing VC (Singing Voice Conversion) [Nachmani et al. Intespeech 2019]

Facebook/Meta

Google



AudioLM (Audio Generation using Language Modelling) [Borsos et al. arXiv 2022]













Categorization of Audio Deepfake Detectors



Categorization of Audio Deepfake Detectors



Evaluation of Deepfake Detectors



Speech-to-Speech Deepfake 100 90 80 70 60 50 40 30 20 -10 0 RawGAT-ST AASIST Raw PC-DARTS Wav2Vec 2.0 Whisper Features

 $TPR = \frac{TP}{TP + FN} \qquad \blacksquare \quad TNR = \frac{TN}{TN + FP}$

Evaluation of Deepfake Detectors



 $TPR = \frac{TP}{TP + FN} \qquad TNR = \frac{TN}{TN + FP}$

Human vs. Machine



"So, when did you first realize humans might have feelings?"

Intuition and Hypothesis

Human interactions



Affected by personal and environmental factors dependent on individuals involved



Intuition and Hypothesis

Human interactions



Affected by personal and environmental factors dependent on individuals involved

Machines interactions



Follow predefined rules and pattern recognition

6

S

R

Less dynamic than human interactions



VoiceRadar

Approximates the physical models

Approximates the physical models

To simulate the audio wave propagation

Approximates the physical models

Extracts the subtle micro-frequencies

To simulate the audio wave propagation







VoiceRadar: Micro-Frequency Analysis [Observer Frequency]

$$\begin{array}{c} \bullet \\ \bullet \\ \hline \end{array} \end{array} \xrightarrow{\bullet} \\ \begin{bmatrix} e_1, e_2, \dots, e_k \end{bmatrix}$$



Physical models approximation



















VoiceRadar: Translational Frequency

➢ Wave propagation of speech in straight lines from speaker to observer.





VoiceRadar: Rotational Frequency

Voice projection changes while speaking causing a rotational frequency shift, like when turning the head.



VoiceRadar: Vibrational Frequency

- Speech has subtle vibrational variations
- Stemming from individual and environmental factors, causing tremors.







Evaluation

Text-to-Speech (TTS) Dataset



Based on recordings from VCTK datasets



Used 8 most recent TTS approaches



40 Speakers

• Cover different distributions of: Gender, age, and regions



Dataset	Records
VALL-E-X	32 000
Speech T5	32 000
Bark	32 000
Style TTS2	32 000
Jenny	32 040
Vits	32 040
XTTS	32 040
Tortoise	32 000
Total	256 120

Speech-to-Speech (STS) Dataset



Based on recordings from VCTK datasets



Used 4 most recent STS approaches



40 Speakers

• Cover different distributions of: Gender, age, and regions

(((Generated for each speaker different
	combinations of STS

Dataset	Records
DiffHierVC	145 465
DiffVC	145 465
HierSpeech++	145 483
SpeechT5	145 483
Total	581 896

Comparison of Audio Deepfake Detectors

100

90



Speech-to-Speech Deepfake

80 70 60 50 40 30 20 10 0 RawGAT-ST AASIST Wav2Vec VoiceRadar Raw PC-Whisper 2.0 DARTS Features $TNR = \frac{TN}{TN + FP}$

Conclusion

Deepfakes are real threats to modern societies

We addressed existing detectors' limitations

➤We proposed VoiceRadar

Agnostic of Text-to-Speech or Voice-Conversion

