

IFIP WG Meeting

 The First University in Korea
Soongsil University

 AISRC
Artificial Intelligence Security Research Center

DeepVoice Detection: A Practical Approach

Souhwan Jung

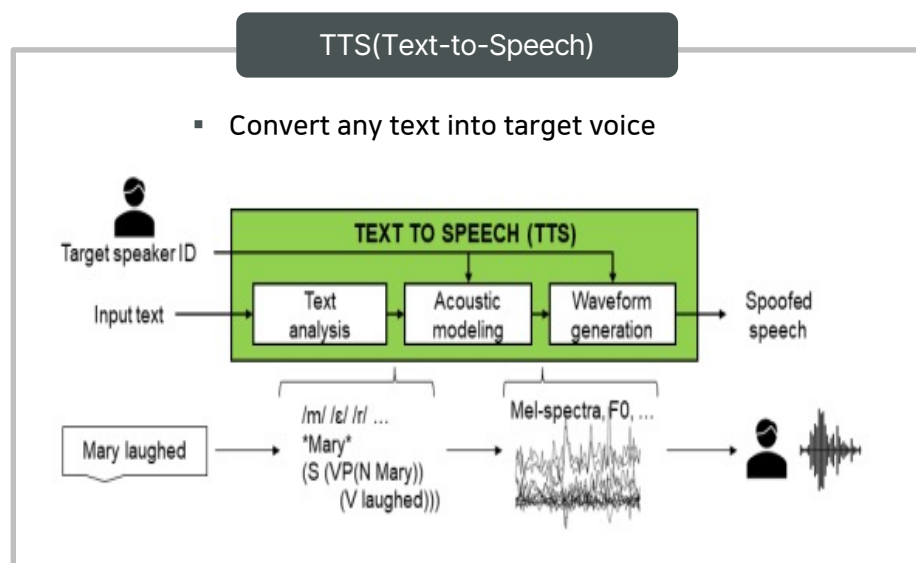
2024.06.30

QUIZ : Which of the three is an AI-generated voice?



How is a DeepVoice generated?

- DeepVoice (Audio Deepfake) : **DeepLearning** + Voice
 - The creation of a DeepVoice is achieved by using AI-powered TTS & VC technology
 - DeepVoice technology now enables the generation of deepfake audio in seconds
 - OpenAI Voice Engine requires only 15-seconds target voice sample for generating voice clones



Diffusion

- **Tortois-TTS**
 - Two main technologies : Autoregressive Diffusion Decoder
 - Focuses on creating diverse, natural-sounding voices

Transformer

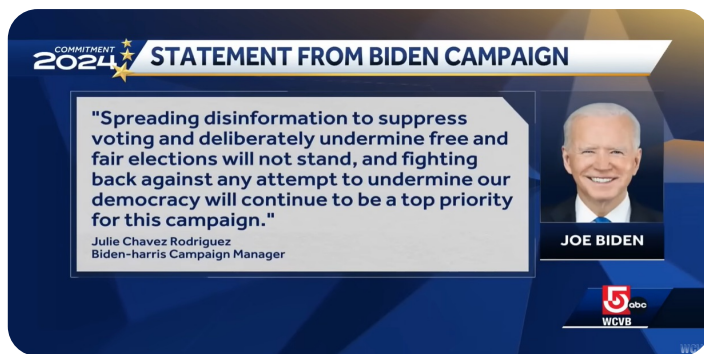
- **Transformer-TTS**
 - Parallel processing
 - Using Attention technology to learn array relationships

GAN

- **GAN-TTS**
 - Uses two competing networks to generate realistic synthetic speech.

Fraud Cases using DeepVoice

- Recent fraud case using deepfake



DeepVoice Robocaller

- Someone using President Biden's voice
- Says "your vote makes a difference in November, not this Tuesday"



Fake Information/News

- Someone created a fake advertisement using Taylor Swift's image and voice
- This fake ad led to payments and caused harm to the fans

A DeepVoice Fraud Case in Korea

- The first voice phishing case using DeepVoice in Korea

- Impersonating a daughter to her mother
- Claiming kidnapped and asking for money to be sent

- Classification model and daughter's speech with a deep learning model

- The daughter's speech was false with an overall probability of 60% (individually, typically over 90%)
- The mother's speech was classified to be less than 1% false

File Name	Fake (%)
1 딸_엄마.wav	87.98%
2 딸.wav	97.33%
3 딸.wav	18%
4 딸.wav	93%
5 딸.wav	93.90%
6 딸.wav	0.04%
7 엄마.wav	0.04%
8 딸.wav	0.04%
9 딸.wav	0.04%
10 딸.wav	1.16%
11 딸.wav	53%
12 딸_피싱범.wav	18.98%
13 딸_엄마.wav	39.02%
14 딸.wav	99.78%
15 딸_엄마.wav	0.41%
16 딸.wav	99.91%
17 딸.wav	0.01%
18 딸_피싱범.wav	0.05%
딸_목소리_전체모음.wav	60.94%

How can we Solve these Problems?



Approach 1 : Regulation of Watermarking on DeepFake

- Policies for regulating deepfakes in various countries



- IT Companies' Response to Generative AI Content

Company	Response
NAVER	Blocking Harmful Deepfakes in Real-time Using AI Filtering Technology ' GreenEye '
Meta	Plans to Identify AI-generated Content on Instagram and Facebook
OpenAI	Inserting Invisible Watermarks in Images Created by the AI Image Generator 'DALL-E'

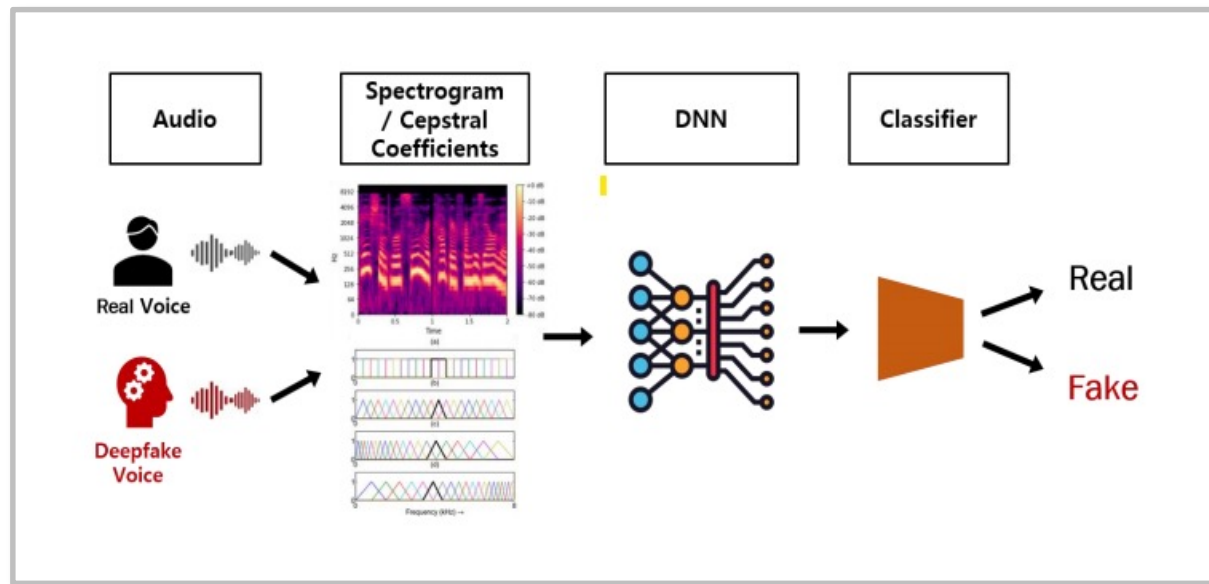
Approach 2 : DeepVoice Detection Technology

➤ Framework for DeepVoice Detection

- General Framework

The diagram shows a process for detecting deepvoice

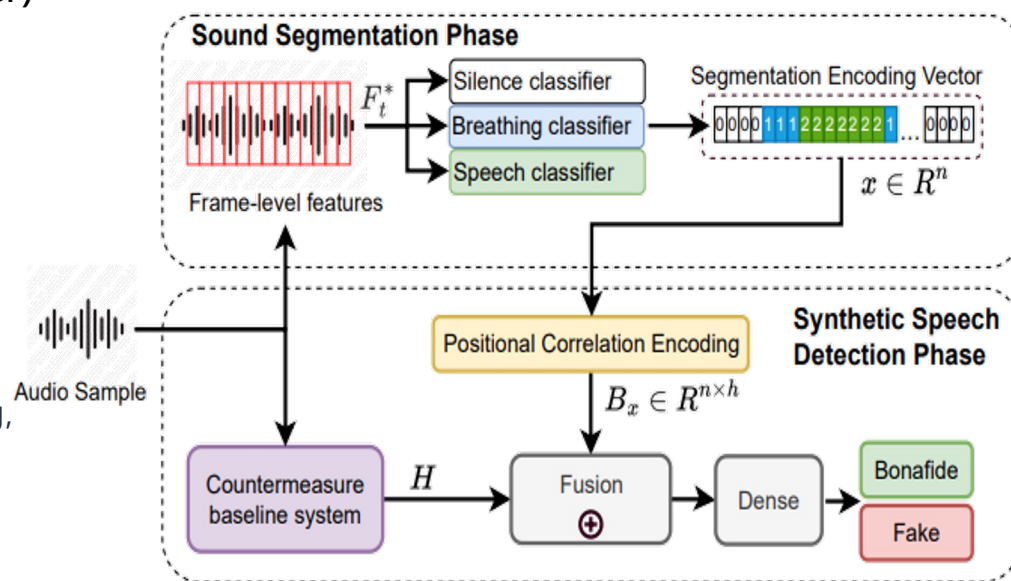
- Handcrafted Feature-based Models, DNN-based Models, End-to-End Models



Our works for DeepVoice Detection(1)

- **BTS-E (Breathing-Talking-Silence Encoder)**

- TTS focuses on linguistics content
- While training, all non-speech segments are removed
- **No biological signal** in the synthesized speech waveform
- Apply VAD to encode the positional of breathing, talking and silence signal at frame-level

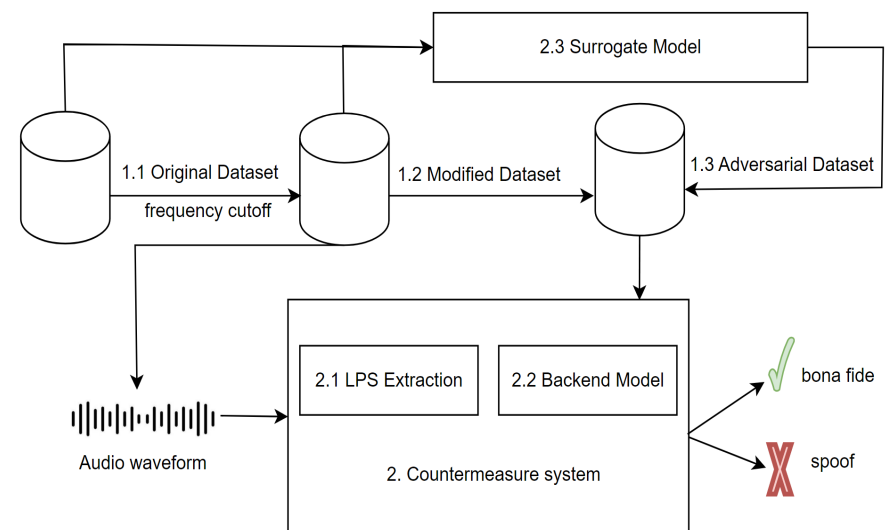


Our works for DeepVoice Detection(2)

- Frequency cutoff technique for robust CM against Adversarial Attacks

- Extract the subband only (remove high/low frequencies) since full bandwidth audio has noise
- Train these with adversarial model (surrogate model) to add generated data to the dataset
- Train the detections system with this dataset

- The model trained on subband is much more robust against adversarial attacks (up to 56% reduction in EER)



Our works for DeepVoice Detection(3)

- **Balance, Multiple Augmentation, and Copy-synthesis Method**

Utilize **Supervised Contrastive Learning**, but in a better training strategy, which proactively set the number of samples in each training mini-batch as follows:

(1) Balances samples between real/fake

(2) Utilizes multiple augmentation methods

(3) Copy-synthesis to generate fake samples

Our works for DeepVoice Detection(3)

- Why it's better?

01

Deepfake Dataset is unbalanced

So keep it balance for every training mini-batch improve the optimization progress

02

Using Multi-augmentation methods to improve the model generalization

03

Copy-synthesis

To generate more hard negative samples, which have the same linguistic content to real sample but have AI artifact, **useful to find better decision boundary.**

Table 5: Comparison of our best system with SOTA models in ASVSpooof 2021 DF track. The result is shown in EER (%).

	Methods	EER
[22]	Wav2Vec + lightDART	7.86
[23]	Wav2Vec + FeedForward (FF) + Attn.Pool	4.98
[24]	Wav2Vec + biLSTM	4.75
[25]	Wav2Vec + ViT-based + FF	3.18
[21]	Wav2Vec + AASIST	2.84
[7]	Wav2Vec + Conformer	2.58
	ours SCL conf-3 (Wav2Vec + Linear layers)	2.17

Challenges to DeepVoice Detection

3 Core Challenges



Data
Generalization &
Model
Robustness

Partially
Fake

Adversarial
Attack

Practical DeepVoice Detection Systems

■ Our DeepVoice Detection System

[Web-Based Detection System]

AISRC - Audio Deepfake Detection System

Model Version:
E_202401

Upload inspecting sample

Drop audio files here or
[Select Files](#)

Audio	File Name	Fake (%)
▶ 0:00 / 0:19	Fake_박영선.wav	E_202401: 99.86%
▶ 0:00 / 0:20	Real_7.wav	E_202401: 0.11%

- Supervised Contrastive Learning (SCL)
- Breathing-Talking-Silence Encoding (BTS-E)
- User-Friendly Interface

[Smartphone-Based real-time Detection]

4:28 60%

FAKE NEWS

FAKE 99.35%

GEOMFILMS
A.I. VOICE

Donald Trump Joe Biden Interview AI Voice

조회수 8.8만회 1년 전 #donaldtrump 더보기

- Efficient Lightweight Model
- Real-time Monitoring
- Enhance Voice Activity Detection(VAD)

[PC-Based real-time Detection]

DEEFAKE VOICE GENERATOR

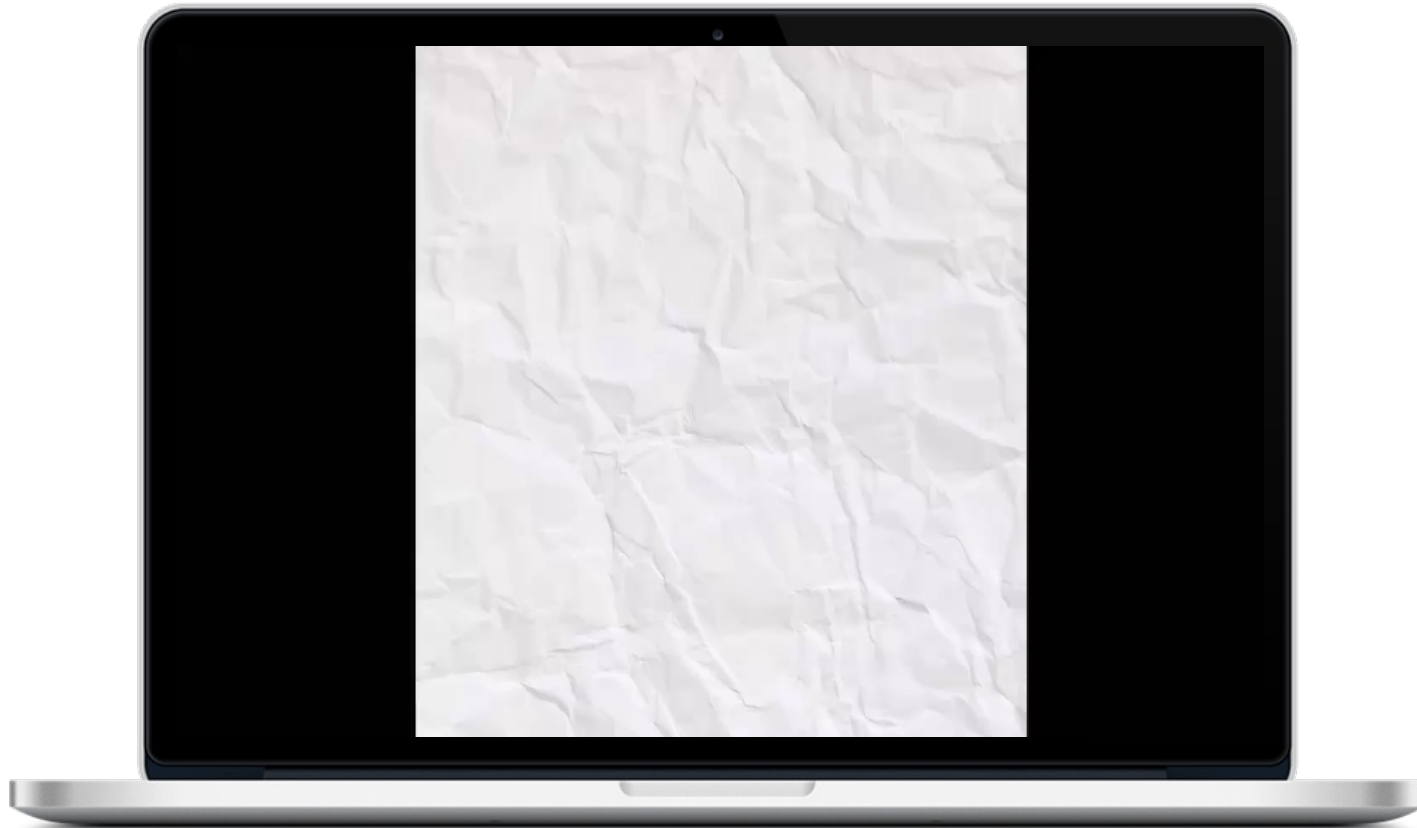
DIGITAL FORENSIC EXPERT: "I DON'T THINK IT TAKES A LONG TIME TO LOOK AT THE RISKS"

CNN

CNN reporter calls his parents using an AI deepfake voice. Watch what happens next

- Optimized Lightweight Model
- Continuous Real-Time Analysis
- Interactive Graphical Display
- Alert Notifications

Video Clip of DeepVoice Detection Systems



Take-Away

- ✓ You could be a victim anytime by your deepvoice in the near future.
- ✓ Two approaches to respond to these threats :
Regulation(watermarking) and Technology(detection system)
- ✓ Government regulations are already going on.
- ✓ Urgent to develop a robust detection system
- ✓ Still many challenges to robust detection system in real-time environment

Thank You!

Any questions can be directed to my email
souhwanj@ssu.ac.kr

