

Problem Statement

- Identify the synthesis algorithm used to generate the given audio
- Extract a fixed-length "synthesizer signature" from varying-length audio and classify into an algorithm.

Dataset

- Clean data
- Five different synthesis algorithms + Unknown class
- 1000 utterances synthesized using each class are given.
- Noisy data: Noise augmented(AWGN, Reverberation, MP3 compression) version of clean
- Natural samples from VCTK and CMU datasets

Speech Production

- Speech signal is an outcome of time-varying vocal-tract system driven by time-varying excitation
- Synthetic speech algorithms seek to mimic this natural system
- May fail to incorporate subtle features of speech production
- At system-level: coarticulation of phoneme sequences
- At source-level: prosodic variations, jitter, shimmer etc.,

Vocal-Tract System Features

- VTS has concatenation of acoustic tubes
- Magnitude spectral envelope of speech represents resonances of VTS
- Mel-wrapped magnitude spectrum is used as features

Vocal Tract & Voice Source Features for Synthetic Speech Detection

Tadipatri Uday Kiran Reddy¹, Sahukari Chaitanya Varun¹, Pranav Kumar Kota¹, Muhammed Fayis², Sankala Sreekanth¹, K. Sri Rama Murty¹

1 Electrical Engineering, IIT Hyderabad, 2 Engineering Science, IIT Hyderabad



Voice Source Features

- Speech signal is inverse filtered through estimated VTS response V(z) (Linear prediction filter)
- The output of the inverse filter is referred to as LP residual
- Natural signal: Higher uncertainty around the glottal closure instant
- Lesser predictability across the instants higher long-term residual



Feature Extractor

- Strided convolutional filter-bank for feature extraction
- Filter sizes and strides decide the effective frame size and shift
- Empirical studies were conducted with different filter sizes and strides

• Filter weights are estimated to minimize the classification loss

Results

	Clean		Noisy	
Model (Parameters)	Train %	Eval%	Train%	Eval%
VGGish($\sim 72M$)	99	89	-	70
$YAMNet(\sim 3.7M)$	99	92	-	76
X-vector($\sim 0.3M$)	98	98	98	93.3
LP Residual X-Vectors ($\sim 1.3M$)	98	95.4	98	96.6



- response.
- below 800 Hz.
- regions
- signatures
- ICASSP-2018

Effective Receptive Field

Receptive Field	Accuracy		
35ms	94.7%		
131 ms	96.67%		
310ms	92.9%		

Frequency Filter Maps



• Frequency response of the first layer of convolutional filter bank.

• Sorted according to peak in the magnitude

• This follows Mel-structure: around 50 filters

Future Course

• Analyze the attention weights to locate unnatural

Combine evidence from source and system level

References

• Snyder et al, X-Vectors: Robust DNN Embeddings for Speaker Recognition, • Albadawy, Ehab et. al., *Detecting Al-Synthesized* Speech Using Bispectral Analysis