

VOICESHIELD: AUDIO DEEPFAKE & DURESS DETECTION EOR BECU

Objective

 BECU Contact Center needs a solution that performs voice signature match analysis of a member's voice (on the call) to identify if it's an AI deepfake, smart assistant (e.g., Google Smart Assistant), a caller under duress (e.g., pressured to withdraw money against their will) or the legitimate member.

Requirements

• The solution should provide an analysis of the voice signature match, along with a confidence level. This information would be displayed to the Contact Center Representative. The solution would prompt the Representative to ask further questions to authenticate the user when the confidence level is low.



Front-end & Back-end Design

Front-end:

 Users can drag and drop or manually upload an audio file into the upload area.



Back-end:

- The API supports a two-stage workflow. It first calls the deepfake detection service and if the prediction is real, the audio is then forwarded to the duress detection service.
- The system is built with FastAPI and deployed on AWS EC2. It uses asynchronous I/O and a thread pool to handle concurrent requests efficiently and safely.

MainAP Receive audio at analyze endpoint Deepfake API Fake or Real? Duress API Fake Return prediction and confidence score

Frontend: User uploads audio

ELECTRICAL & COMPUTER ENGINEERING

UNIVERSITY of WASHINGTON

Deepfake Detection Technical Design

• We evaluated two recent top-performing speech deepfake detection models [2] [3]. AASIST2 showed better accuracy and accent robustness. Hence, it was selected as the base model for our work.

Models	CVoiceFake+ DeepVoice Datasets	Self-Collected Dataset	2s Audio Performance	4s Audio Performance	8s Audio Performance	Accent Variation Effect		
RawNet	0.55	0.84	0.67	0.68	0.68	Yes		
AASIST2(base)	0.74	0.94	0.58	0.81	0.88	Νο		
Metric: AUC								

Correct Example Wrong Exam

• During our experiments, we found that different deepfake speech generation techniques have distinct characteristics, and strong performance on one does not ensure good results on others. To improve robustness, we expanded the dataset with diverse techniques and retrained the model.

Models	CVoiceFake+ DeepVoice Datasets Self-Collected Dataset			
AASIST2(base)	0.74	0.94		
AASIST2(finetune)	0.88	0.81		
AASIST2(retrain)	0.99	0.84		
	Metric: AUC			
Datasets	Deepfal	Deepfake Technology		
CVoiceFake+ DeepVoice	Retrieva Conversi WORLD, Pa DiffWa	• Work		
BECU	BECU Yourtts, xtts, VALL-E X			
Self-Collected	DupDub Al ByteDance' Text-t	DupDub Al Voice Generator, ByteDance's Doubao, Google Text-to-Speech Al		



FACULTY MENTOR: Jai Jaisimha INDUSTRY MENTOR: Alan Wilson & Dan Gibbons **SPONSOR: BECU**



flow Overview: User telephone is preprocessed (segmentation esampling) and then passed to the T2 model to predict the probability of being real or spoofed.

Under Duress Detection Technical Design

of multiple technical paths:

- Er fea of

Method (Key Tools)		Key Notes / Focus				
Acoustic Feature Analysis + Machine Learr Classification (OpenSMILE, ML)	ning F	Focuses on comprehensive sound characteristics; uses ML to identify duress patterns.				
X Voice Features to Natural Language + LLM Ana (Custom Extractor, LLM)	alysis LLN	LLM analyzes text-converted voice features; challenges with specific emotion detection & false positives.				
X Automatic Speech Recognition (ASR) + Sensi Word Detection (Qwen2_Audio, ASR)	itive	Relies on real-time transcription and predefined sensitive words/security codes for monitoring.				
Keyword Detection + BERT Semantic Analys (Qwen2_Audio, spaCy, BERT)	sis	ldentifies understa	keywords a nding; still r	and uses BERT [.] eliant on keywo	for contextual ord presence.	
nerging from data analysis, critical i to vocal duress markers directly info ature engineering and structural de	nsights ormed t evelopm	the nent			Audio Input	
the subsequent model.	·	(Acoustic	: Feature		
Vocal Cue Duress Indicat	tion		Extra	action	Preprocessing	
Extreme Low Pitch Drops Signals emotional c	listress		SHAP	+ P+RFE	Treprocessing	
Distinct Pitch Patterns Differentiates dures	ss types		Feature	Selection		
Steeper Loudness Slopes Indicates panic/	′fear					
Variable Loudness Increase Links to deception/	'anxiety	(Cross-\	/alidation		
Reduced Speech Variation Shows physical te	ension				Duress Detection	
This duress detection model proces audio to provide binary duress class	sses raw		LST Model	M+RF Fusion	Model	
and an associated confidence score guidance.	tor age	2	0.7 Duress	0.3 Not Duress	Prediction Output	

Future Work & References

- advanced AI model.
- real-world scenarios and populations.
- Reduce the latency of the system.
- Integrate with the company's real call scenarios.
- ACM, Florence, Italy, ISBN 978-1-60558-933-6, pp. 1459-1462, 25.-29.10.2010. (2022).

[3] C. Sun, S. Jia, S. Hou and S. Lyu, "AI-Synthesized Voice Detection Using Neural Vocoder Artifacts," 2023 IEEE/CVF Conference on Computer Vision and Pattern *Recognition Workshops (CVPRW)*, Vancouver, BC, Canada, 2023, pp. 904–912, doi: 10.1109/CVPRW59228.2023.00097.





• The selection of the current duress detection system was informed by an evaluation

• Improve the model to deal with more challenging situations such as a more

• Explore the dynamic interplay of the features and their applicability across diverse

[1] Florian Eyben, Martin Wöllmer, Björn Schuller: "openSMILE – The Munich Versatile and Fast Open-Source Audio Feature Extractor", Proc. ACM Multimedia (MM), [2] Tak, Hemlata, et al. "Automatic speaker verification spoofing and deepfake detection using wav2vec 2.0 and data augmentation." arXiv preprint arXiv:2202.12233