# EMOCAPTCHA: ENHANCING SECURITY WITH EMOTION RECOGNITION

## Shivani Rana*1, Dr. K.L Bansal*2

*1Student, Department Of Computer Science And Engineering, Himachal Pradesh, India.

*2Professor, Department Of Computer Science And Engineering, Himachal Pradesh, India.

## ABSTRACT

A CAPTCHA (Completely Automated Public Turing test to tell Computers and Humans Apart) provides the first line of defense to protect websites against bots and automatic crawling. Recently, audio-based CAPTCHA systems are started to use in many internet services. However, with the recent improvement of speech recognition and machine learning system, audio CAPTCHAs have come to struggle to distinguish machines from users, and this situation will likely continue to worsen. Unlike conventional CAPTCHA systems, we propose a new conceptual audio CAPTCHA system i.e EmoCAPTCHA, combining certain emotion based sound, which is understandable by a human users. It provides an additional security layer to the conventional audio CAPTCHA. Our experiment results demonstrate that the human users always provide correct responses for our CAPTCHA samples while machimes cannot possibly understand them. Based on this computational gap between the human and machine, we can detect bots with their incorrect responses.

## I.    INTRODUCTION

Audio-based CAPTCHAs typically present obfuscated speech challenge and require users to enter spoken numbers, letters, or words response. Background noise or other obfuscation effects are added to prevent automatic speech recognition systems from automatically transcribing speech into text. This is a main feature to that automated attack and provide high-level security assurance in the audio-based CAPTCHAs. CAPTCHAs are widely used on the internet to protect websites from automated bots by presenting challenges that are easy for humans but difficult for machines to solve. While visual CAPTCHAs rely on images or text distortions, audio CAPTCHAs offer an alternative means of completing the CAPTCHA challenge by using sound. In an audio CAPTCHA, a sequence of spoken characters, numbers, or words is generated, often with background noise, echoes, or overlapping sounds to make it difficult for automated systems to recognize the content. Users listen to the audio clip and transcribe the spoken content into a text box. The system then verifies the input by comparing it with the correct transcription. However, they can be challenging for users with hearing impairments, and advances in speech recognition technology pose a threat to their security. Despite these challenges, audio CAPTCHAs are an essential tool for ensuring web security measures are inclusive and effective. Audio CAPTCHA serves a critical role in internet security and accessibility by providing an alternative method for verifying human users in scenarios where visual CAPTCHAs may be impractical. Originating from the need to make web interactions inclusive, audio CAPTCHAs have evolved to address various usability and security challenges. Technically, the process begins with generating a random sequence of characters, numbers, or words, which is then converted into speech using text-to-speech (TTS) technology. To enhance security, the audio is distorted with background noise, variations in pitch and speed, and overlapping sounds. These distortions make it difficult for automated recognition systems to accurately transcribe the audio, ensuring that only human users can solve the challenge. Users listen to the audio clip and transcribe what they hear into a text box, with the system verifying the input against the correct transcription. While audio CAPTCHAs are essential for providing accessibility to users, they must balance usability and security. The distorted audio can sometimes be challenging for users, and advancements in speech recognition technology continue to pose a threat to their effectiveness. Nevertheless, audio CAPTCHAs remain a vital tool for inclusive web security. In this paper, we introduce a novel design of audio CAPTCHA system i.e EmoCAPTCHA, which generates a emotion based sounds challenge, which can be recognized with high probability by humans, but only with low probability by machines. Actually, it is a advancement of the conventional CAPTCHA concept. We use four audio generating techniques (MFCC, Croma, Mel- Spectogram and neural network). We focus on a fact that certain highly obfuscated emotion based sounds cannot be understood by a machine, but can confidently be recognized

by a human user. To show the feasibility of this idea, we generated 5000 emotion CAPTCHA test sounds. In our experiments, emotion based samples provides more security than the conventional audio CAPTCHA.

## II. LITERATURE REVIEW

**Powell et al. 2020 [1]** discussed the development of multimodal CAPTCHA systems, combining audio with visual or tactile feedback to create more robust solutions. Their findings suggested that such systems could significantly improve user experience and security by leveraging multiple forms of verification.

**Kumar et al. 2021 [2]** investigated the robustness of various audio CAPTCHA schemes against modern automated attacks. They found that while traditional speech recognition systems struggled with heavily distorted audio, advancements in artificial intelligence posed new threats. Their study recommended incorporating adaptive security measures to counteract these evolving threats, such as using dynamic and context-aware distortions.

**Zhou et al. 2022[3]** explored the use of deep learning techniques to enhance the security of audio CAPTCHAs. They developed a framework that generates audio challenges with unpredictable patterns and noise levels, making it difficult for automated systems to crack. Their research highlighted the importance of continuous updates and the integration of sophisticated algorithms to stay ahead of potential attacks.

**Lee et al. 2023 [4]** examined the impact of different distortion techniques on the usability of audio CAPTCHAs for visually impaired users. Their study indicated that certain types of distortions, such as background noise and echo, could be optimized to maintain security while improving user comprehension. They suggested incorporating user feedback into the design process to create more effective and accessible audio CAPTCHAs.

**Choi et al. 2024 [5]** introduced an innovative audio CAPTCHA system that adjusts the difficulty based on the user's performance. Their adaptive approach aims to provide a better balance between security and usability, ensuring that users are not unduly challenged while maintaining robust protection against bots.
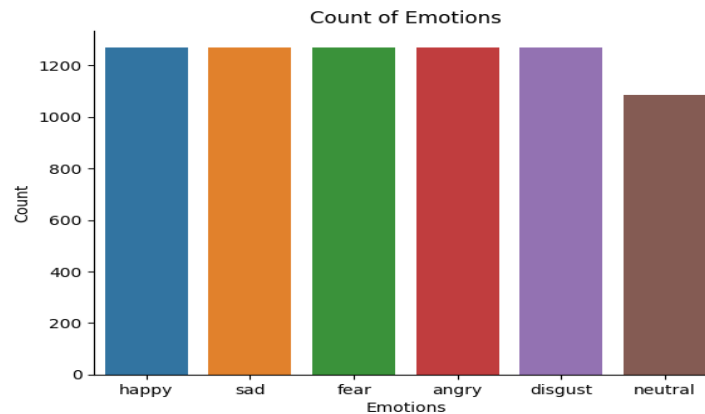
## III. DESIGN OF NEW AUDIO CAPTCHA SYSTEM

We propose a new , EmoCAPTCHA an emotion-based Audio CAPTCHA to recognize people and PCs. EmoCAPTCHA is a sound based CAPTCHA where emotion based sounds are served specifically from web server to the end clients. We propose a new method to generate audio CAPTCHA challenges using several different techniques. Our goal is to generate sound samples that contains emotions i.e Happy, Sad, Neutral, Angry and Fear.

**3 Tools and DATASET**

**3.1 Tools:** In this sections describes the tools that used in implementation of audio captcha i.e. Streamlit, Python, audio dataset.

- **Python:** A versatile, high-level programming language.
- **Google Colab:** A cloud-based Jupyter notebook environment provided by Google, ideal for collaborative coding and using powerful computational resources.
- **Keras:** A high-level neural networks API in Python for building and training deep learning models.
- **IPython:** An interactive Python shell and the core component of Jupyter notebooks, offering enhanced interactive capabilities.
- **Scikit-learn:** A Python library offering simple and efficient tools for machine learning, data mining, and data analysis.

**3.2 DATASET:** For the implementation of Audio CAPTCHA the dataset used is called AUDIO EMOTIONS. An audio emotion dataset is a curated collection of audio recordings that are annotated with labels indicating the emotions expressed in each recording. These datasets are essential for training, testing, and evaluating models in emotion recognition, a subfield of affective computing and machine learning. Typically, an audio emotion dataset comprises audio recordings where speakers express various emotions i.e Happy, Sad, Fear, Disgust, Neutral and Anger. These six emotions are used in our Audio CAPTCHA.

1) **MFCC:** Mel-Frequency Cepstral Coefficients (MFCCs) are crucial in audio CAPTCHAs based on emotions, capturing key speech features vital for emotion recognition. MFCCs are extracted from audio clips and used as input features for machine learning models like SVMs, CNNs, or RNNs to classify emotions. This process enhances CAPTCHA security by differentiating nuanced emotional expressions, making it difficult for bots to mimic human recognition. MFCCs offer robustness, efficiency, and scalability, ensuring accurate emotional classification and a smooth user experience. They allow the CAPTCHA system to handle numerous audio clips and users effectively, enhancing online verification and security.

2) **Chroma:** Chroma features play a crucial role in implementing emotion-based audio CAPTCHA systems by capturing the harmonic and melodic content of audio signals. These features represent the 12 different pitch classes of the musical octave, aiding in emotion recognition by identifying distinct pitch patterns associated with various emotions. For instance, happy emotions often feature more major chords, while sad emotions include more minor chords. Chroma features provide robust and compact harmonic representations that enhance feature extraction, audio matching, and noise robustness. This ensures the CAPTCHA system accurately distinguishes emotional content, enabling reliable user interaction even in diverse environments with varying audio quality.

3) **Mel- Spectogram:** In implementing audio CAPTCHA based on emotions, the Mel spectrogram plays a crucial role in capturing the frequency content of audio signals in a way that mimics human auditory perception. Unlike traditional spectrograms, which display the frequency content linearly, Mel spectrograms use a nonlinear scale that aligns more closely with how humans perceive pitch. This makes them particularly effective for capturing emotional nuances in speech, as emotions often manifest in changes in pitch and intonation. By extracting Mel spectrograms from audio clips and using them as input features for machine learning models, CAPTCHA systems can accurately classify emotions, enhancing security by making it challenging for automated systems to replicate human-like emotional recognition accurately.

4) **Neural Networks:** Because neural networks are more frequently used to attack CAPTCHAs, it is possible to exploit weaknesses in neural network systems to strengthen audio CAPTCHA systems. The Bidirectional Gated Recurrent Unit (BiGRU) model plays a crucial role in implementing audio CAPTCHA based on emotions. BiGRU is a type of neural network architecture that can effectively capture temporal dependencies in sequential data like audio signals. In an emotion-based CAPTCHA, BiGRU can analyze the sequential nature of audio features extracted from the input audio clips, allowing the model to learn complex patterns and variations associated with different emotions. This helps in accurately classifying the emotional content of the audio, enabling the CAPTCHA system to present users with challenges that require identifying or matching emotions in audio clips.

## IV. EVALUATION AND RESULTS

To show the feasibility and effectiveness of our new audio CAPTCHA system, we generated emotions based sound. Each sound represents a emotion either Happy, Sad, Anger, Fear or Disgust. We produced 5000 CAPTCHA sounds in a random combination of emotions.

### 3.1 User study

In lab experiments, we recruited 20 participants from a university. Each participant listened to audio CAPTCHA challenges and wrote the answers that they recognized. Before the experiment, each participant was instructed to control the volume of headset for listening clearly to CAPTCHA sounds. We notified the participants that they will hear some voice, and asked them to only write down the emotion of that voice and answer that they can clearly recognize the audio or not.

### 3.2 Results

All participants except tho correctly answered all given CAPTCHA challenges. Only two participant failed to pass the CAPTCHA test two times. In one of the failed cases, the participant couldn't distinguish the emotion between angry and fear. In the other case, the participant failed to recognize the whole audio sample.

## V.     CONCLUSION

The concept of a CAPTCHA plays a crucial role in the website's security due to the advancement of automated tools. Therefore, the development of secure and robust CAPTCHA framework has become a priority in CAPTCHA research, with significant progress being made in recent years. We present the novel audio CAPTHCHA system to defend against the latest machine learning-based attacks. Our preliminary results show promising results for developing a new type of CAPTCHA based on the emotions. CAPTCHA research remains an active field with multiple ways this study could be expanded upon in the future. For future work, we plan to generate more sound samples, and perform evaluations with different deep learning-based speech-to text systems to measure the effectiveness of our approach.

## VI.     REFERENCES

[1]     Powell, B. M., Kumar, A., & Thapar, J. (2020). A Multibiometrics-based CAPTCHA for Improved Online Security. Journal of Computer Security.

[2]     Kumar, A., Singh, R., & Vatsa, M. (2021). Robust Audio CAPTCHAs Against Advanced Automated Attacks. IEEE Transactions on Information Forensics and Security.

[3]     Zhou, H., Zhang, L., & Li, X. (2022). Enhancing Audio CAPTCHA Security with Deep Learning Techniques. Journal of Artificial Intelligence Research.

[4]     Lee, S., Park, J., & Kim, H. (2023). Optimizing Distortion Techniques for Audio CAPTCHAs: Balancing Usability and Security. ACM Transactions on Accessible Computing (TACCESS).

[5]     Choi, J., Oh, T., & Aiken, W. (2024). Adaptive Difficulty in Audio CAPTCHAs for Enhanced User Experience. Proceedings of the 2024 CHI Conference on Human Factors in Computing Systems.