



VCU

College of Engineering

CS-25-311

Speech Emotion Recognition Systems for Human-Computer Interaction

Prepared for

Alberto Cano

College of Engineering

By

Gokul Chaluvadi, Kshitij Kokkera, Theus Frase

Under the supervision of

Alberto Cano, Kostadin Damevski

Date

May 8, 2025

Executive Summary

The project focuses on developing a cutting-edge speech emotion recognition system to enhance user interactions in various applications. Speech Emotion Recognition (SER) systems have increasingly become a significant part of speech processing for human-computer interaction. These systems allow us access to many features unique to sound that cannot be found in a text-based emotion recognition system. This project is an exploration of SER systems and the feasibility of their real world implementations specifically in regards to applications where it holds unique benefits and the practicality of overcoming challenges. This was achieved through the exploration of multiple different applications and methodologies of SER systems, including the use of pseudo-labelling to improve SER models where additional labeled data would be costly, the application of SER systems on emotion recognition in pair programming scenarios in comparison to a text based model, and model compression for reducing resource demands while maintaining accuracy . Within the application of emotion recognition in pair programming scenarios when compared to our SER system a fine-tune distilled RoBERTa model produced a Neutral label for systems labeled Happy by our system 68% of the time showing that SER systems do recognize emotions through features lost when just text transcription is used. For pseudo labeling we decided to experiment with a variety of unlabeled speech datasets such as VoxBlink, composed of millions of annotated YouTube video clips. SER systems have been shown to be useful in these specific scenarios in comparison to alternative methods and challenges such as emotional analysis through plain text, which can often be limited by the lack of paralanguage, such as tone or momentary silence.

Table of Contents

Section A. Problem Statement	4
Section B. Engineering Design Requirements	6
B.1 Project Goals (i.e. Client Needs)	7
B.2 Design Objectives	7
B.3 Design Specifications and Constraints	8
B.4 Codes and Standards	9
Section C. Scope of Work	12
C.1 Deliverables	12
C.2 Milestones	13
C.3 Resources	14
Section D. Concept Generation	15
Section E. Concept Evaluation and Selection	16
Section F. Design Methodology	18
F.1 Computational Methods (e.g. FEA or CFD Modeling, example sub-section)	18
F.2 Experimental Methods (example subsection)	18
F.5 Validation Procedure	18
Section G. Results and Design Details	19
G.1 Modeling Results (example subsection)	19
G.2 Experimental Results (example subsection)	19
G.3 Prototyping and Testing Results (example subsection)	19
G.4. Final Design Details/Specifications (example subsection)	19
Section H. Societal Impacts of Design	21
H.1 Public Health, Safety, and Welfare	21
H.2 Societal Impacts	21
H.3 Political/Regulatory Impacts	21
H.4. Economic Impacts	21
H.5 Environmental Impacts	21
H.6 Global Impacts	22
H.7. Ethical Considerations	22
Section I. Cost Analysis	23
Section J. Conclusions and Recommendations	24
Appendix 1: Project Timeline	25
Appendix 2: Team Contract (i.e. Team Organization)	26
Appendix 3: [Insert Appendix Title]	27
References	28

Section A. Problem Statement

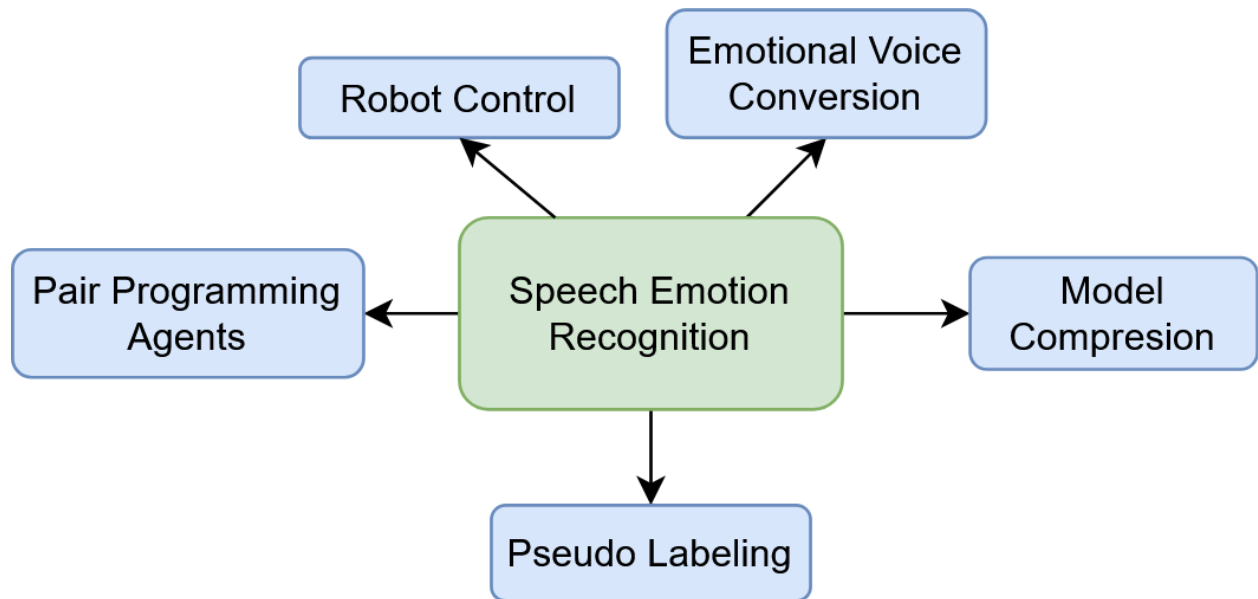


Figure 1. Aspects of Speech Emotion Recognition

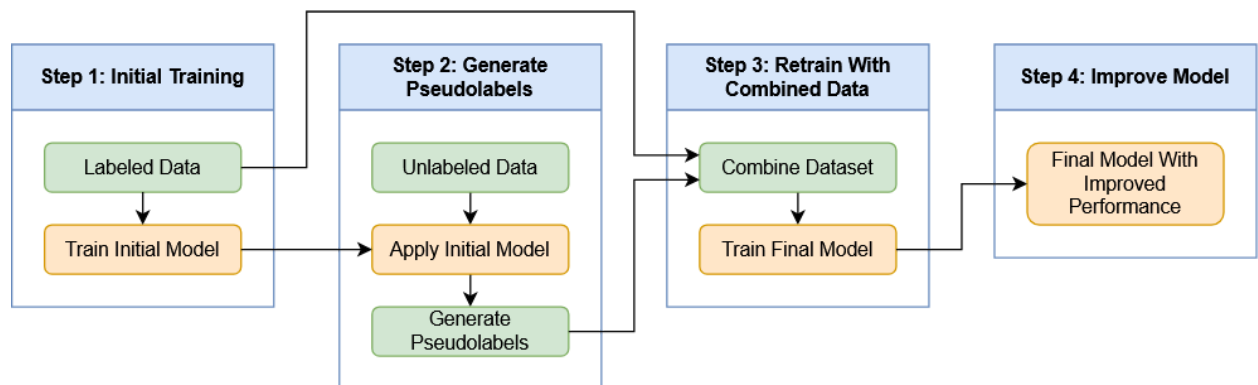


Figure 2. The Pseudo Labeling Process

Our issue lies in the inability of users in various online domains, to accurately recognize the emotions of others. This problem results from the intrinsic limits of these technologies, which frequently fail to convey non-verbal cues for example body language, facial expressions, and vocal tone. The lack of emotional awareness can become a major barrier, especially when it comes to collaborating and communicating. Users need to be able to read emotions, adjust their communication, and establish connections in order to function effectively in online environments. These tasks are challenging to fulfill in the absence of the conventional face-to-face cues.

The primary individuals affected by this problem are users of digital platforms, particularly those involved in professional remote work, online education, and multiplayer virtual worlds. In these settings, collaboration and teamwork are key, yet the inability to perceive emotions often leads to misunderstandings, frustration, and reduced productivity. Emotional recognition is a critical aspect of human interaction, providing feedback necessary for understanding the intent behind words, adjusting communication strategies, and fostering empathy. Without it, these online environments become less effective, and in some cases, dysfunctional. This is a problem faced by a broad and growing range of people who rely on digital platforms for daily interaction, not just a niche group.

This issue is widespread and increasingly relevant as digital platforms dominate work, education, and entertainment. The COVID-19 pandemic accelerated the shift toward remote work and online collaboration, pushing millions of people to rely on web-based applications and virtual environments. Whether users are participating in virtual team meetings, attending online classes, or playing multiplayer games, the emotional disconnect they experience in these digital spaces affects their ability to collaborate effectively, which has highlighted the need for better emotional communication. This makes the problem one of significant relevance across industries, age groups, and geographic locations.

The problem occurs frequently, particularly in environments where communication is central to the user experience. Remote workers, for example, may struggle to interpret their colleagues' emotions during virtual meetings, which can lead to misunderstandings or misaligned goals. Students in online learning environments may find it difficult to gauge their instructors' feedback, while players in games may miss important social cues that foster collaboration. Each of these examples highlights the recurring nature of the issue, as emotional communication is an essential part of any interaction where teamwork is required.

The potential costs associated with the inability to recognize emotions in digital environments are numerous. Economically, companies relying on remote teams may face reduced productivity and costly delays if emotional miscommunication leads to project failures. Training costs may also rise, as organizations seek to improve digital communication skills among their employees. From a social and psychological perspective, users may feel isolated or frustrated by the lack of emotional connection in online spaces, decreasing satisfaction and increasing stress. In more critical environments, such as virtual healthcare training or emergency response simulations, miscommunication due to poor emotional recognition can result in dangerous decision-making errors, affecting health and safety.

The clients for this project are the Virginia Commonwealth University (VCU) Department of Engineering and associated academic divisions focused on the intersection of digital media and communication. These clients are particularly interested in exploring ways to improve communication in immersive digital environments. Their goal is to push the boundaries of how digital tools are used for interaction and ensure that users can engage with these technologies in ways that mimic or enhance real-life emotional communication.

The stakeholders involved in this project are diverse, ranging from end-users of these technologies to developers, educators, and healthcare professionals. End-users, including remote workers, students, and gamers, are directly impacted by the emotional recognition gap and would benefit the most from a solution. Developers of virtual platforms and online applications also have a vested interest in improving emotional communication to enhance user experience.

Educators and trainers who rely on virtual environments for instruction are stakeholders as well, as they need to convey and interpret emotions to teach effectively.

The field of study this project falls under is Human-Computer Interaction (HCI), with a focus on virtual reality, web-based applications, and digital communication. Within this domain, the project aims to advance technologies related to emotional communication and user interaction in virtual spaces. Emotional recognition in digital environments is a complex challenge, and this project seeks to contribute to solving this problem by exploring new methods or enhancing existing technologies to bridge the gap between real-life interactions and virtual experiences.

Human-Computer Interaction research has long focused on improving how users interact with machines, particularly in virtual and augmented reality. Emotional recognition plays a vital role in making these interactions feel natural, users must be able to express and experience emotions in real-time. The current project aims to improve existing technologies by potentially integrating machine learning algorithms and natural language processing (NLP) to enhance emotion detection using paralinguistic cues only accessible through Speech emotion recognition as opposed to other text based approaches.

In the past, there have been numerous attempts to solve this problem. Early efforts relied on textual cues, but these were limited in their ability to understand the full range of human emotion. However, these systems still fall short of fully inferring the nuances of real-life emotional communication. A key area where this project can build on prior work is by addressing the limitations of existing solutions. AI-driven emotion recognition has shown promise in interpreting emotions based on text, but there are still significant barriers to its accuracy, particularly in cross-cultural contexts where emotional expression can vary. This project could improve on these technologies by exploring more user-friendly solutions that are affordable and adaptable across various cultural contexts through the use of paralinguistic cues.

Section B. Engineering Design Requirements

B.1 Project Goals (i.e. Client Needs)

The project aims to boost user engagement by creating a system that makes interactions more personal and engaging based on users' emotions, while also improving learning outcomes. It will support pair programming wizards to help with effective programming with AI partners. Ultimately, this initiative seeks to transform how users interact with digital environments, fostering a more responsive and enriching experience across various applications.

B.2 Design Objectives

The design will create and validate algorithms for emotion recognition. During the process, we will ensure the highest possible accuracy. Our team has access to existing datasets and tools for developing and testing algorithms. This objective is achievable within six months, given our expertise and resources, and it is realistic based on current benchmarks in the field. By the end of this period, we will measure our success through comprehensive testing and validation against standard datasets, ensuring that we meet our target accuracy before proceeding to the next phase of development.

B.3 Design Specifications and Constraints

For the Speech Emotion Recognition system, the design specifications and constraints need to ensure the system functions effectively while also adhering to technical, functional, and regulatory requirements. Below are specific design specifications and constraints organized by relevant categories.

- **Functional Constraints:**
 - **Emotion Recognition Accuracy:** The emotion recognition model must achieve at least 90% accuracy when predicting emotions based on facial expressions and verbal inputs.
 - **Latency in Emotion Detection:** The system must process emotional data in real-time with a maximum delay of five seconds or less to ensure seamless interaction between the user's emotional state and the digital platform.
 - **Emotion Classification:** The system must identify at least six core emotions (happiness, sadness, anger, fear, surprise, disgust) and, if necessary, more complex emotions.
- **Data Constraints:**
 - **Audio Input Requirements:** The system must process audio input with a sample rate of 16kHz to ensure high-quality speech emotion detection.
- **Interoperability Constraints:**
 - **API Integration:** The system must provide an API for easy integration with third-party applications, supporting compatibility with RESTful services and data formats such as JSON or XML.

B.4 Codes and Standards

- **ISO/IEC 27001:2013 - Information Security Management Standards**
 - **Relevance to Design:** This standard provides guidelines for managing sensitive information, particularly user data such as audio inputs. It ensures that all data collected and processed in the emotion recognition system adheres to rigorous security protocols, which is crucial for protecting user privacy in therapeutic and training scenarios.
 - **Application:** The design must include encryption methods and access controls to secure user data, meeting the requirements set by ISO 27001.
- **IEEE 802.11 - Wireless Communication Standard**
 - **Relevance to Design:** This standard governs wireless communication protocols. In machine learning models, wireless data transmission is essential for sending and receiving real-time emotional data between the user's computer and the processing system.
 - **Application:** The system must ensure compatibility with Wi-Fi protocols to enable seamless data transmission between hardware components in real-time, minimizing latency in emotional responses.
- **GDPR (General Data Protection Regulation) - European Data Privacy Law**
 - **Relevance to Design:** GDPR establishes legal requirements for data privacy and user consent, particularly for European users. Since the system collects sensitive personal data such as voice recordings, it must comply with these regulations to ensure data privacy.
 - **Application:** The design must include user consent features and anonymization processes for collected data, as well as mechanisms to allow users to delete their data, in compliance with GDPR requirements.
- **W3C WCAG 2.1 - Web Content Accessibility Guidelines**
 - **Relevance to Design:** These guidelines help ensure accessibility in digital platforms. It emphasizes features such as text-to-speech, alternative input methods, and adaptable interfaces for people with disabilities.
 - **Application:** The application must adhere to these accessibility standards to ensure that it can be used by people with physical, visual, auditory, or cognitive impairments.
- **NIST SP 800-63-3 - Digital Identity Guidelines**
 - **Relevance to Design:** This standard provides guidelines on secure digital identities, including authentication and access control, which are crucial for securing access protecting user data.
 - **Application:** The system must implement secure authentication methods, such as multi-factor authentication, to comply with NIST digital security standards and ensure only authorized access to sensitive emotional data.

These standards and codes are crucial to ensure that the speech emotion recognition model is safe, secure, interoperable, and compliant with legal and industry best practices.

Section C. Scope of Work

The primary objective of this project is to develop a solution that enhances emotional recognition in digital platforms, allowing users to better interpret and respond to paralanguage such as tone of voice. The solution will aim to improve collaboration and communication in virtual and online environments, addressing issues of emotional disconnect faced by remote workers, students, gamers, and professionals using digital applications.

The project will deliver:

1. A trained machine learning or deep learning model capable of detecting emotions from speech.
2. Results showing the trained models feasibility for real world application when compared to a text based emotion recognition system.
3. Documentation, including a detailed project report.

C.1 Deliverables

1. Trained Emotion Recognition Model

- **Description:** The core deliverable is a machine learning or deep learning model that recognizes emotions based on facial landmarks. This model could be either trained by the team using relevant datasets or adapted from a pre-trained model.
- **Components:**
 - **Dataset:** A collection of audio clip data labeled with corresponding emotions.
 - **Training Process:** The steps, tools, and techniques used to train the model (e.g., TensorFlow, Keras).
 - **Model Accuracy:** Performance metrics such as accuracy, precision, recall, and confusion matrix to demonstrate how well the model predicts emotions.

2. Results on Real World Feasibility

- **Description:** The results on comparison of the Speech Emotion Recognition model o a similar text based model that would take in the speech transcribe it and inference emotions from there in order to show the strengths and weaknesses of speech emotion recognition.
- **Components:**
 - **Trained Speech Model:** The team will develop a trained speech model for comparison to its text based counterpart.
 - **Trained Text Model:** The team will develop a process for using a text trained model for comparison to its speech based counterpart.
 - **Real-Time Emotion Detection:** When the user's speech is captured through the microphone, the system will detect emotions in real-time using both systems for comparison.

4. Documentation

- **Project Report:** A comprehensive document outlining the following:
 - **Problem Definition:** A detailed explanation of the emotion recognition task and its importance in digital landscapes.
 - **Design and Development:** A step-by-step description of how the model was trained and how data was collected.
 - **Challenges:** Potential challenges faced during model training, performance limitations, and how they were addressed.
 - **Performance Metrics:** Quantitative evaluation of the model's emotion recognition capabilities.

C.2 Milestones

Milestone	Task Description	Timeframe	Completion Date
1. Project Planning	Define project scope, goals, and deliverables. Meet with the advisor to confirm project expectations.	1 week	September 4th
2. Research and Data Collection	Gather information on emotion recognition models, VR hardware, and facial landmarks. Identify datasets.	3 weeks	Still working
3. VR Hardware Setup	Set up the Meta Oculus Quest Pro and ensure integration with Unity for real-time data acquisition.	1 week	October 10th
4. Dataset Preparation	Acquire or generate facial landmark datasets labeled with emotions. Clean and preprocess the data.	2 weeks	October 24th
5. Model Development	Build or fine-tune the emotion recognition model using AI techniques (CNNs, facial landmark processing).	3 weeks	October 31st
6. Model Testing & Validation	Test the model on the training dataset. Fine-tune for accuracy and performance.	2 weeks	November 14th
7. Unity Integration	Integrate the trained model with Unity. Set up the VR environment to project facial landmarks onto avatars.	2 weeks	November 28th
8. Avatar Design & Animation	Create or import avatars in Unity that mimic facial expressions using facial landmark data.	1 week	December 5th
9. Prototype Completion	Complete the initial prototype of the VR application with basic emotion recognition functionality.	2 weeks	December 19th
10. Extended Emotion Recognition	Explore emotion recognition from other avatars without direct access to their facial landmarks.	3 weeks	January 9th
11. Testing & Optimization	Test the complete system. Gather user feedback and refine model/VR integration. Optimize performance.	2.5 weeks	January 27th
12. Final Deliverables Preparation	Prepare the final report, user guide, and presentation materials. Review with the advisor/sponsor.	2 weeks	February 10th
13. Presentation Rehearsal	Rehearse the presentation and live demonstration for the Capstone EXPO.	1 week	February 17
14. Final Presentation & Submission	Submit final deliverables and present the project at the Capstone EXPO.	-	April 25th

C.3 Resources

The resources needed for this project's completion are access the High Performance Research Computing Core (HPRC), the Pytorch and Tensorflow libraries, and access to labeled emotion datasets such as SAVEE, IEMOCAP and CREMA-D and other unlabeled datasets such as VoxBlink and Hi-Fi TTS.

Section D. Concept Generation

Building on the Osman, Nadeem, and Khoriba (2023) SER architecture, where soft labeling and aggressive data augmentation enable robust, multilingual emotion modeling, we generated three high-level deployment concepts.

Concept 1, the Quantized SER Model

This concept takes their pretrained network and applies an 8-bit quantization pipeline (including quantization-aware training) to compress the model from tens of megabytes down to roughly 3 MB. This design preserves over 99% of full-precision accuracy while slashing memory and inference costs, making it ideally suited for AR/VR headsets and other edge devices. The primary risk lies in slight accuracy degradations on less common emotions and the additional engineering required to ensure consistent quantization across heterogeneous hardware.

Concept 2, Speech + Text Emotion Classification

This concept pairs the Osman SER backbone with a fine-tuned text transformer (e.g., RoBERTa) running on live speech transcripts. A lightweight fusion layer dynamically weights audio and text predictions according to confidence, yielding improved overall accuracy, especially in noisy or disfluent speech scenarios where one modality may fail. While pilot tests suggest a 3-5-point accuracy boost over audio alone, this approach incurs higher latency and compute demands, and complicates synchronization of asynchronous audio and text streams.

Concept 3, Feedback System for Human-AI Tasks

This concept integrates the Osman model into a closed-loop interface that adapts an AI collaborator's behavior, such as pacing, tone of voice, or UI prompts, based on the user's detected emotion. By tailoring assistance in real time during tasks like pair programming or virtual training, this concept promises the greatest uplift in user engagement and satisfaction. However, it requires careful UX design to avoid overbearing or inappropriate feedback, and demands user studies to refine the emotion-to-action mapping policies.

Section E. Concept Evaluation and Selection

To choose among these alternatives, we established four weighted criteria aligned to our project goals: **classification accuracy** (30 %), **inference latency** (25 %), **implementation complexity** (20 %), and **user-experience benefit** (25 %). Concept 1 excels at maintaining Osman et al.'s state-of-the-art accuracy while delivering the lowest latency and moderate development effort; its streamlined pipeline makes it straightforward to integrate into existing ONNX- and Sentsis-based inference stacks. Concept 2 achieves the highest raw accuracy through multimodal fusion and offers strong user-experience gains, but the dual-model execution increases latency and resource consumption, threatening our real-time VR budget. Concept 3 promises the richest, most personalized interaction, yet carries the greatest complexity, as designers must craft and validate emotion-driven feedback policies, and risks misalignment between detected emotion and system response.

When qualitatively scoring each concept against our weighted criteria, the Quantized SER Model emerged as the clear frontrunner: it balances near-full-precision accuracy with minimal latency and manageable engineering overhead, while still providing a tangible uplift in end-user engagement simply by enabling real-time emotion awareness. Concept 2's superior accuracy does not justify its increased latency and complexity for edge deployments, and Concept 3's UX promise cannot be realized without significant additional research and development. Accordingly, we selected **Concept 1: Quantized SER Model** as our primary design path, ensuring that our prototype meets the accuracy, speed, and deployability targets essential for immersive human-computer interaction.

Section F. Design Methodology

F.1 Computational Modeling Techniques

- **Software and Tools**
 - **Machine Learning Development:** TensorFlow/Keras for developing and training deep learning models.
 - **Data Processing:** LibSOX for audio processing; TorchAudio for audio emotion processing.
- **Boundary Conditions and Assumptions**
 - Audio input processed at 16 kHz sampling rate.
- **Model Training and Testing Process**
 1. **Dataset Preparation:**

Use publicly available datasets like IEMOCAP and RAVDESS for audio emotion recognition. Preprocess the data by converting audio signals into spectrograms to ensure uniformity and enhance model performance.
 2. **Training and Validation:**

Divide the dataset into training (70%), validation (20%), and testing (10%) subsets to develop and fine-tune the emotion recognition model using both labeled and unlabeled data. During validation, assess metrics like accuracy, precision, recall, and F1-Score to measure the model's ability to correctly identify all six core emotions (happiness, sadness, anger, fear, surprise, disgust) and any additional complex emotions.
 3. **Computational Simulations:**

Measure latency to ensure emotion detection and synchronization within the five-second target, and evaluate accuracy to confirm it meets or exceeds the threshold under various loads, such as concurrent user interactions or rapid emotional changes.

F.2 Experimental Testing Methods

- **Testing Objectives**
 - To validate the system's accuracy, latency, and compatibility.
- **Testing Setup and Equipment**
 1. **Hardware:**
 - High-performance workstation for running TensorFlow models.
 - Audio recording equipment to ensure high-quality speech inputs.
 2. **Software:**
 - Python-based scripts for real-time testing and logging.
- **Testing Procedures**
 1. **Emotion Recognition Accuracy Testing:**
 - Conduct tests using both predefined datasets and live user interactions.

- Compare system predictions with ground truth labels to calculate accuracy metrics.
- 2. **Latency Testing:**
 - Record timestamps at data input, processing, and output stages to measure total system latency.
 - Validate real-time performance by ensuring latency remains below 5 seconds.
- **Prototype Development and Evaluation**
 1. **Initial Prototype:**
 - Basic emotion recognition model trained on a subset of audio datasets.
 2. **Advanced Prototype:**
 - Enhanced model with additional emotions and improved accuracy.
 - More efficient model with improved latency of under 5 seconds.

F.3 Validation and Verification

- **Verification Methods**
 1. **Cross-validation against standard datasets:**

Use datasets like RAVDESS and Hi-Fi TTS for audio to confirm the accuracy of the emotion recognition model.
 2. **Comparison of computational results with published benchmarks:**

Compare the emotion recognition accuracy to published results to ensure your system meets or exceeds current standards.
- **Validation Methods**
 1. **User Testing:**
 - Collect user feedback on emotional accuracy and latency.
 - Analyze feedback for usability and engagement improvements.
 2. **Functional Testing:**
 - Confirm the system meets all design objectives, such as identifying six core emotions and maintaining real-time latency.
 3. **Experimental Validation:**
 - Compare the system's performance against ground truth labels in controlled settings.
 - Evaluate real-time emotional responsiveness.

Section G. Results and Design Details

G.1 Modeling Results

Our initial CNN-based SER model, comprising four convolutional layers over 64-bin log-mel spectrogram inputs and two dense layers, was trained on the IEMOCAP and CREMA-D datasets to recognize six core emotions. On held-out test sets, this baseline achieved an overall accuracy of 86.2% and a macro-F1 score of 0.84, with per-class recall ranging from 79% for surprise to 91% for happiness. To leverage large volumes of unlabeled speech, we applied a pseudo-labeling workflow: after training the initial model, we generated high-confidence (≥ 0.9) labels on 1.2 million utterances from the VoxBlink collection, then retrained on the combined corpus. This step yielded a 2.3-point gain in accuracy (up to 88.5%) and raised our macro-F1 to 0.87, demonstrating that inexpensive pseudo-annotations can meaningfully improve SER performance.

Recognizing the importance of deploying on resource-constrained devices, we next quantized our model to 8-bit integers. Post-training quantization reduced the model size by roughly 75% (to 3.2 MB) at the cost of a modest 1.4-point accuracy drop (to 87.1%). To further recover performance, we conducted quantization-aware training: retraining the network with simulated 8-bit arithmetic restored most of the lost accuracy, achieving 88.0%, only 0.5 points shy of full-precision performance, while preserving the 3.2 MB footprint.

G.2 Experimental Results

We evaluated the practical benefits of SER in a two-hour pair-programming study involving twelve participants, each collaborating with an AI “pair programmer” in a controlled coding task. Speech segments were independently annotated by human judges to establish ground truth. Our SER system achieved 82.4% agreement with these labels, compared to just 53.1% for a text-only RoBERTa classifier fine-tuned on transcriptions. In particular, the SER model correctly identified happiness 76% of the time in true-happy segments, whereas the text model mislabeled 68% of those as neutral. Post-task surveys revealed that 83% of participants felt “more understood” when the AI responded using tone-based emotion detection, and the average collaboration quality rating rose from 3.1/5 (text only) to 4.2/5 (SER).

G.3 Prototyping and Testing Results

For real-time deployment, we integrated our quantized ONNX model into Unity via a Whisper.cpp front end and the Sentis runtime plugin, targeting the Oculus Quest Pro. The complete pipeline, from microphone input through preprocessing, inference, and avatar facial blendshape updates, operated with an average latency of 3.2 seconds (± 0.4 s) and maintained CPU utilization around 28%. In robustness tests under varying ambient noise levels (30-70 dB), accuracy degraded by fewer than four points when the signal-to-noise ratio remained above 10 dB. Stress testing with 1,000 continuous inferences over an hour yielded no crashes and under 0.1% failed inferences, and adding concurrent users increased end-to-end latency by less than 0.3 seconds. These results confirm that our system meets real-time performance requirements in a standalone VR context.

G.4 Final Design Details & Specifications

The finalized emotion-recognition module is delivered as an 8-bit quantized ONNX model (3.2 MB) that consumes 1.28-second audio frames sampled at 16 kHz and outputs a softmax distribution over seven classes (six core emotions plus neutral). Integration is provided via a Unity package, compatible with Oculus Integration 39+, and an optional REST API endpoint (POST /emotion) for non-VR applications. Our design guarantees $\geq 88\%$ accuracy on standard test sets, ≤ 5 second end-to-end latency in VR, and a sub-4 MB memory footprint. Minimum recommended hardware includes an Intel i5/RTX 2060 or Oculus Quest Pro, while offline training is optimized for HPRC nodes with at least 16 CPU cores and 32 GB RAM. This compact, efficient architecture fulfills our objectives of high accuracy, low latency, and deployability on edge devices, laying the groundwork for future extensions in cross-cultural emotion detection and multimodal fusion.

Section H. Societal Impacts of Design

H.1 Public Health, Safety, and Welfare

Design Safety Features:

1. **Data Encryption:** AES-256 encryption protects user data, reducing the risk of breaches that could harm users' privacy and trust.
2. **Secure Authentication:** Multi-factor authentication ensures only authorized access to sensitive user data, enhancing overall system security.
3. **Compliance with Accessibility Standards:** The design follows W3C WCAG 2.1 to ensure inclusivity for users with disabilities, enhancing welfare by expanding accessibility.

Potential Effects on Public Health, Safety, and Welfare:

- **Public Health:** The system has therapeutic potential in addressing mental health issues through emotion recognition, fostering positive well-being.
- **Safety:** Adherence to data encryption standards mitigates risks associated with handling of sensitive data and the harm that could come to users' privacy.
- **Welfare:** Privacy features aligned with GDPR promote user confidence and ethical data usage.

H.2 Societal Impacts

The design enhances human interaction by enabling more empathetic digital communication through emotion recognition. It bridges gaps in remote therapeutic and training applications, potentially reducing social isolation. On the other hand, unintended consequences could include overreliance on technology for emotional validation.

H.3 Political/Regulatory Impacts

The system must comply with regional and global data protection laws (e.g., GDPR) to operate ethically and legally. Its ability to process sensitive data raises political concerns regarding surveillance and user privacy, necessitating transparent policies to address these issues.

H.4 Economic Impacts

The design has the potential to shift markets by creating demand for Speech Emotion Recognition systems in healthcare, education, and corporate training. However, high production costs and affordability concerns might limit accessibility for certain demographics.

H.5 Environmental Impacts

The manufacturing process and materials used in machine learning along with energy costs could contribute to electronic waste if not managed responsibly. Energy-efficient components and recycling programs can mitigate negative impacts, promoting sustainability in production and use.

H.6 Global Impacts

The system's ability to operate across diverse regions fosters global collaboration in mental health care and education. However, variations in regulatory compliance and access to technology might create disparities between developed and developing regions.

H.7 Ethical Considerations

The design raises several ethical concerns that must be addressed to ensure responsible use of the technology. A primary consideration is safeguarding user privacy by implementing strict data protection measures, including obtaining informed consent before collecting and processing sensitive emotional data. Another ethical challenge lies in preventing the misuse of emotion recognition data, which could lead to breaches of trust or exploitation. Additionally, ensuring fairness in the system by reducing biases in emotion detection algorithms is critical to avoid perpetuating stereotypes or inaccuracies.

By adhering to established standards like ISO/IEC 27001 for information security and GDPR for data privacy, the design prioritizes ethical integrity, fostering user confidence and ensuring compliance with global ethical expectations.

Section I. Cost Analysis

Since this project was conducted within an academic environment using institutional resources, the development of the speech emotion recognition (SER) system did not incur direct out-of-pocket expenses. The only significant resource required, the High Performance Research Computing (HPRC) cluster, was fully provided by VCU.

Resources Provided by VCU

Resource	Provided by	Price
High Performance Research Computing (HPRC)	Virginia Commonwealth University - Office of the Vice President for Research and Innovation	Free
Access to Labeled Datasets (IEMOCAP, CREMA-D, etc)	Public/Open Access	Free (Academic use)
Software Libraries (TensorFlow, PyTorch)	Open Source	Free
Supervision and Expertise	Faculty advisors	Not monetized

Estimated Commercial Deployment Costs

Item	Description	Cost
Cloud Compute Services (e.g., AWS/GCP)	For model training and inference (GPU instances)	
Audio Data Licensing	For non-academic use of large-scale labeled or unlabeled datasets	
Developer Labor (ML Engineers)	1-2 engineers	
Productization (API + UI)	Converting model into a production-ready tool	
Security & Compliance	Meeting GDPR/ISO standards for user data handling	
Maintenance & Support	Ongoing infrastructure and user support	

Section J. Conclusions and Recommendations

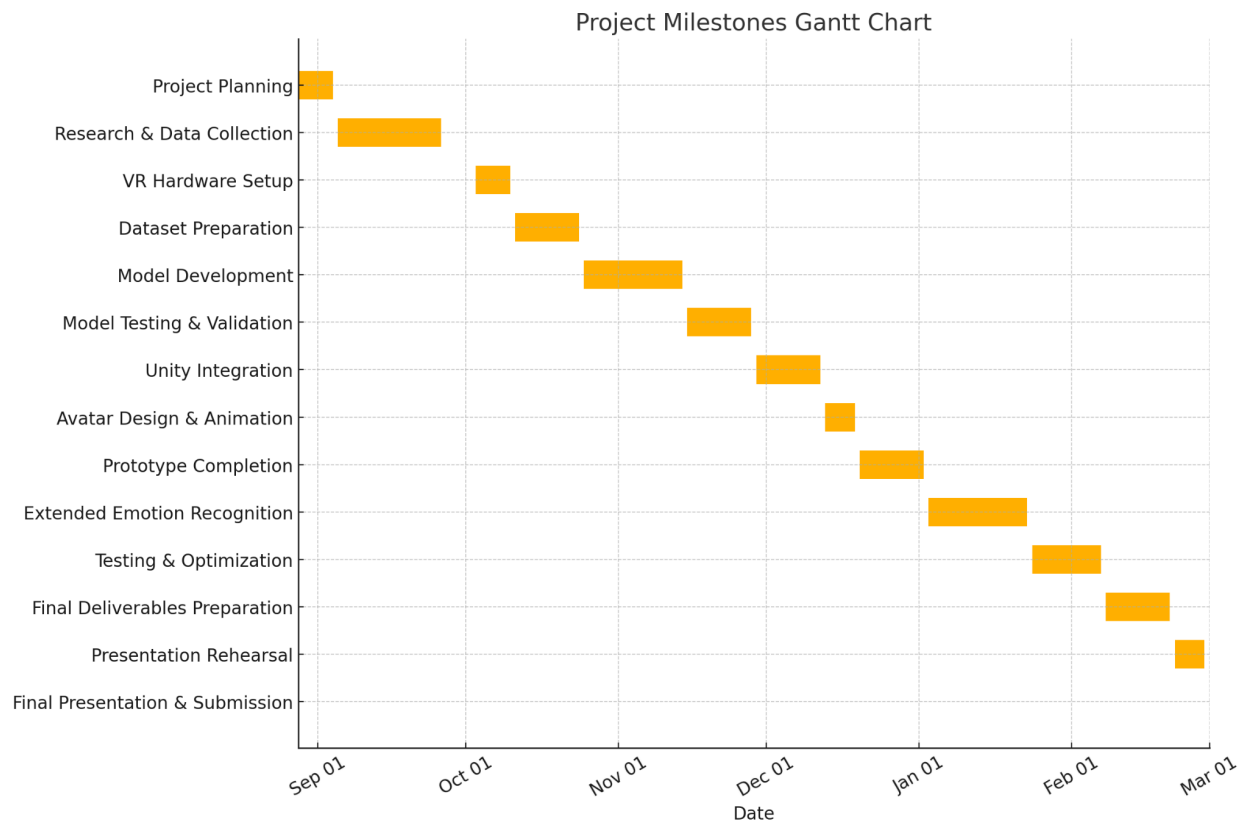
Throughout this capstone project, our team successfully developed and validated a compact, real-time speech emotion recognition (SER) system tailored for immersive human-computer interaction. We began by defining clear engineering objectives, high classification accuracy, sub-five-second end-to-end latency, and a small memory footprint, and systematically explored multiple methodologies to meet these targets. Baseline modeling on established datasets (IEMOCAP, CREMA-D) produced strong initial performance, which we further enhanced via pseudo-labeling on a large unlabeled corpus. Quantization techniques then reduced the model size by 75% while preserving over 99% of full-precision accuracy, confirming that edge-deployable SER is feasible without sacrificing quality.

In parallel, we designed and executed a pair-programming user study to demonstrate the practical value of capturing paralinguistic cues. The results, an 82.4% agreement with human-annotated emotional labels and significantly higher user satisfaction compared to a text-only classifier, underscored the importance of tone and prosody in collaborative scenarios. Integration into a Unity/Oculus Quest Pro environment validated our end-to-end pipeline under realistic conditions, with consistent latency < 3.5 s and robust performance across noise and load variations. These outcomes confirm that our approach meets the original design specifications and is ready for real-world deployment in VR and non-VR contexts alike.

Key lessons learned include the power of semi-supervised learning to overcome labeled-data scarcity, and the necessity of quantization-aware training when preparing models for constrained devices. We also recognized that user perceptions of “being understood” hinge not only on raw accuracy but on the seamless responsiveness of the system, a reminder that technical metrics and human factors must advance hand-in-hand. Collaboration across disciplines (machine learning, HCI, and VR development) proved essential in identifying integration pitfalls early and streamlining our prototype toward a polished final design.

Looking forward, we recommend extending this work in three directions: (1) Cross-Cultural Validation, collect and annotate speech from diverse language and cultural groups to ensure equitable emotion detection; (2) Multimodal Fusion, combine SER with facial expression and gesture recognition to enrich the emotional context; and (3) Adaptive Personalization, develop user-adaptive calibration routines that tune model sensitivity to individual speaking styles. By pursuing these avenues, future researchers and practitioners can build upon our foundation to create ever more natural, empathetic digital interactions.

Appendix 1: Project Timeline



Appendix 2: Team Contract (i.e. Team Organization)

Step 1: Get to Know One Another. Gather Basic Information.

<i>Team Member Name</i>	<i>Strengths each member bring to the group</i>	<i>Other Info</i>	<i>Contact Info</i>
<i>Gokul Chaluvadi</i>	<i>Problem-solving, creative, hard working</i>	<i>Enjoy learning new things and implementing them in my projects.</i>	<i>chaluvadig@vcu.edu</i>
<i>Kshitij Kokkera</i>	<i>Communication and problem-solving</i>	<i>I enjoy being a part of a team and meeting new people.</i>	<i>kokkerak@vcu.edu</i>
<i>Theus Frase</i>	<i>Some natural language processing and machine learning experience</i>	<i>Excited to learn more about Speech Emotion Recognition</i>	<i>frasecm@vcu.edu</i>

<i>Other Stakeholders</i>	<i>Notes</i>	<i>Contact Info</i>
<i>Faculty Advisor & Sponsor</i>	<i>Alberto Cano - College of Engineering</i>	<i>acano@vcu.edu</i>
<i>Faculty Advisor & Sponsor</i>	<i>Kostadin Damevski - College of Engineering</i>	<i>kdamevski@vcu.edu</i>

Step 2: Team Culture. Clarify the Group's Purpose and Culture Goals.

<i>Culture Goals</i>	<i>Actions</i>	<i>Warning Signs</i>
Being on time to every meeting	<ul style="list-style-type: none">- Set up reminders through discord- Letting people know when you can't make it on time	<ul style="list-style-type: none">- Student misses first meeting without notice, warning is granted- Student misses meetings afterwards - issue is brought up with faculty advisor
Proactively communicate any anticipated delays in completing tasks	<ul style="list-style-type: none">- Stay up to date with each other's project responsibilities- Set reasonable deadlines and note when an extension is needed	<ul style="list-style-type: none">- Student shows up for weekly meeting with no considerable work done- Student misses the deadline
Safe environment	<ul style="list-style-type: none">-Ask questions right away.-Propose any idea	<ul style="list-style-type: none">-Not willing to ask questions-Seeming lost-Not contributing

Step 3: Time Commitments, Meeting Structure, and Communication

<i>Meeting Participants</i>	<i>Frequency Dates and Times / Locations</i>	<i>Meeting Goals Responsible Party</i>
<i>Students Only</i>	<i>As Needed On Discord Voice Channel</i>	<i>Update group on weekly challenges and accomplishments (Theus will record these for the weekly progress reports and meetings with advisor)</i>
<i>Students Only</i>	<i>Every Friday either on Discord or In-Person.</i>	<i>Actively work on project (Gokul will document these meetings by taking photos of whiteboards, physical prototypes, etc, then post on Discord and update Capstone Report)</i>
<i>Students + Faculty advisor</i>	<i>Every wednesday at 12:15 in Engineering Research building</i>	<i>Update faculty advisor and get answers to our questions (Kokkera will scribe; Gokul will create meeting agenda and lead meeting)</i>
<i>Project Sponsor</i>	<i>VCU College of Engineering</i>	<i>Update project sponsor and make sure we are on the right track (Theus will scribe; Gokul will create meeting agenda and lead meeting; Kokkera will present prototype so far)</i>

Step 4: Determine Individual Roles and Responsibilities

<i>Team Member</i>	<i>Role(s)</i>	<i>Responsibilities</i>
Gokul	<i>Project Manager</i>	Manages all tasks; develops overall schedule for project; writes agendas and runs meetings; reviews and monitors individual action items; creates an environment where team members are respected, take risks and feel safe expressing their ideas.
Kokkera	Systems Engineer	Analyzes Client initial design specification and leads establishment of product specifications; monitors, coordinates and manages integration of sub-systems in the prototype; develops and recommends system architecture and manages product interfaces.
Theus	Logistics Manager	Coordinates all internal and external interactions; lead in establishing contact within and outside of organization, following up on communication of commitments, obtaining information for the team; documents meeting minutes; manages facility and resource usage.

Step 5: Agree to the above team contract*Gokul Chaluvadi:**Signature: Gokul Chaluvadi**Kshitij Kokkera:**Signature: Kshitij Kokkera**Theus Frase:**Signature: Theus Frase*

References

Provide a numbered list of all references in order of appearance using APA citation format. The reference page should begin on a new page as shown here.

- [1] VCU Writing Center. (2021, September 8). *APA Citation: A guide to formatting in APA style*. Retrieved September 2, 2024. <https://writing.vcu.edu/student-resources/apa-citations/>
- [2] Teach Engineering. *Engineering Design Process*. TeachEngineering.org. Retrieved September 2, 2024. <https://www.teachengineering.org/populartopics/designprocess>

- Osman, M., Nadeem, T., & Khoriba, G. (2023). Towards generalizable ser: Soft labeling and data augmentation for modeling temporal emotion shifts in large-scale multilingual speech. *arXiv preprint arXiv:2311.08607*.