

Capturing The Sentiments of Hybrid Conferences Using Emoji- Based Questionnaire and Ai Emotion Recognition

M. Malaiselvam¹, M. Naga Sundar², Mr. T. Maria Mahajan³

^{1,2}*Department of information technology, Nehru Arts and Science College, Coimbatore, India*

³*Mentor, Department of Information Technology, Nehru Arts and Science College, Coimbatore, India*

Abstract—Hybrid conferences and seminars have become increasingly popular due to their flexibility and accessibility. However, understanding participant engagement and emotional response in such environments remains a significant challenge. Traditional feedback systems are delayed and often ineffective. This paper proposes an intelligent sentiment analysis system that integrates facial emotion recognition with an emoji-based questionnaire to capture real-time audience feedback. The system utilizes computer vision techniques with OpenCV and Haar-Cascade classifiers for face detection, along with a Convolutional Neural Network (EmotionNet) for emotion classification. Additionally, an emoji-based feedback mechanism enables participants to express their sentiments quickly and intuitively. The combined approach enhances real-time interaction, improves audience engagement, and provides actionable insights for organizers. The proposed system demonstrates improved feedback accuracy, usability, and effectiveness in hybrid event environments.

Index Terms—Sentiment Analysis, Hybrid Conferences, Emotion Recognition, OpenCV, CNN, Emoji Feedback

I. INTRODUCTION

The rapid evolution of digital communication modalities has catalysed the widespread adoption of hybrid conferencing frameworks, which seamlessly integrate physical and virtual participation. While these models offer significant advantages in terms of geographical inclusivity and operational flexibility, they introduce complex challenges regarding the quantitative assessment of participant engagement and affective response.

Current evaluation methodologies are primarily predicated on retrospective post-event surveys.

However, these traditional instruments are often considered suboptimal due to:

- **Temporal Latency:** Feedback is delayed, preventing immediate adjustments to the session flow.
- **Operational Inefficiency:** The data collection process is often perceived as time-consuming by respondents.
- **Statistical Volatility:** Participation rates are typically low, leading to a fragmented understanding of audience sentiment.

In contrast to conventional static systems, real-time feedback mechanisms are essential for the dynamic optimization of session quality. Drawing architectural inspiration from AI-driven intelligent systems—specifically the high-speed data processing observed in disaster evacuation simulation frameworks—this work proposes a robust real-time sentiment analysis system. The proposed solution achieves high-fidelity monitoring by integrating:

- **Automated Facial Emotion Recognition (FER):** Utilizing computer vision for continuous, non-intrusive monitoring of micro-expressions.
- **Synchronous Emoji-Based Feedback:** Implementing a low-friction interface for instantaneous, user-initiated emotional signalling.

This bimodal hybrid approach ensures the acquisition of accurate, continuous, and granulated metrics, facilitating a responsive environment for assessing audience engagement as it unfolds.

II. EXISTING SYSTEM AND LIMITATIONS

To address these limitations, the proposed system implements a dual-modality sentiment capture

mechanism designed for high-fidelity audience monitoring.

A. Key Features

- Real-Time Facial Emotion Recognition (FER): Continuous monitoring of participant micro-expressions to gauge subconscious engagement.
- Synchronous Emoji-Based Feedback Interface: A low-friction digital overlay for user-initiated, conscious sentiment reporting.
- Live Sentiment Visualization Dashboard: A centralized interface providing event organizers with instantaneous heatmaps and engagement metrics.
- Predictive Analytics for Engagement Tracking: Data-driven insights to identify trends in audience attentiveness throughout the session duration.
-

B. System Workflow

The operational logic of the proposed architecture follows a sequential pipeline for data synthesis:

1. Input Acquisition: High-definition video streams are captured via integrated webcams or platform-native feeds.
2. Face Localization: Individual faces are identified within the frame utilizing the Haar-Cascade classifier algorithm.
3. Image Preprocessing: Frames are normalized, resized, and converted to grayscale to optimize computational efficiency.
4. Feature Extraction and Classification: A Convolutional Neural Network (Emotion Net) architecture is deployed to categorize emotions (e.g., happiness, boredom, confusion) in real-time.
5. Manual Feedback Integration: Explicit sentiment data is collected via the synchronous emoji interface.
6. Data Aggregation and Synthesis: Passive and active data streams are fused to create a holistic engagement profile.
7. Dynamic Visualization: The synthesized metrics are rendered on a real-time dashboard for stakeholder analysis.

III. METHODOLOGY

The proposed methodology employs a multi-staged computational pipeline to transform raw audience interaction into actionable sentiment intelligence. The architecture is divided into the following core functional modules:

A. Face Detection and Localization

The initial stage involves the continuous monitoring of high-definition video streams. To achieve real-time performance with minimal computational overhead, the system utilizes the OpenCV Haar-Cascade classifier. This algorithm employs a series of Haar-like features and a boosted cascade of weak classifiers to isolate and localize facial regions within each frame, ensuring robust detection even in varied lighting conditions characteristic of hybrid conference halls.

B. Image Preprocessing

To optimize the efficiency of the neural network and reduce extraneous noise, detected facial regions undergo a rigorous preprocessing sequence:

- Grayscale Transformation: Converting RGB frames to a single-channel intensity map to reduce dimensionality.
- Spatial Normalization: Resizing all detected facial segments to a uniform 48×48-pixel resolution.
- Intensity Normalization: Rescaling pixel values to a range of $[0, 1]$ to facilitate faster convergence during the classification phase.

C. Emotion Recognition Architecture

The core analytical engine is a custom Convolutional Neural Network (CNN), designated as Emotion Net. This model is architected to extract hierarchical spatial features from the pre-processed facial data. The network classifies audience sentiment into seven distinct affective states:

- Positive/High Engagement: Happy, Surprise.
- Neutral/Moderate Engagement: Neutral.
- Negative/Low Engagement: Sad, Angry, Confused, Bored.

D. Synchronous Emoji-Based Feedback

To supplement the automated computer vision data, the system integrates an Active Feedback Layer.

Participants are provided with a low-friction interface to explicitly signal their current state. This layer captures intentional emotional responses through a standardized set of emoji indicators:

- 😊 (Happy): Indicates high content resonance.
- 😐 (Neutral): Represents steady, baseline attention.
- 😕 (Confused): Signals a need for further clarification on the topic.
- 😞 (Bored): Indicates a drop-in session stimulus or pace.
- 🙌 (Satisfied): Reflects specific approval of presented material.

E. Data Synthesis and Analytics

The final stage involves the fusion of passive facial data and active emoji inputs. The analytics engine performs the following operations:

- Frequency Distribution Analysis: Calculating the prevalence of specific emotions over discrete temporal intervals.
- Aggregated Sentiment Scoring: Synthesizing disparate inputs into a singular Engagement Index to provide a macro-view of session health.
- Dynamic Trend Visualization: Rendering these metrics into live heatmaps and line graphs, allowing stakeholders to correlate sentiment shifts with specific segments of the conference agenda.

IV. HUMAN BEHAVIOR MODELING

A critical component of the proposed framework is the nuanced understanding of human emotional expression. The system recognizes that sentiment is not a monolithic data point but a multi-layered behavioural phenomenon.

A. Multimodal Affective Interpretation

To achieve high-fidelity results, the system distinguishes between two primary forms of human signalling:

- Natural Emotional Manifestation: Facial expressions captured via FER represent "leaked" or subconscious reactions. These are immediate, biological responses to stimuli that occur before the participant can filter their feelings.
- Conscious Affective Feedback: Emoji selection represents "intended" or cognitive feedback. This

allows participants to provide context to their emotions, such as signalling "confusion" (😕) when their facial expression might only register as "neutral" (😐).

- Synthesized Accuracy: By correlating these two streams, the system mitigates the risk of "false positives" (e.g., a resting facial expression being misread as boredom), thereby improving the overall reliability of the engagement metrics.

B. Behavioural and Environmental Considerations

The framework accounts for the inherent complexity of human behaviour through several key considerations:

- Temporal Reaction Latency: The system acknowledges that subconscious facial reactions occur almost instantaneously (<500 ms), whereas manual emoji feedback involves a cognitive delay as the participant processes their feeling and interacts with the UI.
- Socio-Cultural Variations: Emotional expression is not universal; the intensity and frequency of facial movements can vary significantly across different demographics. The hybrid model uses emoji feedback as a "ground truth" to calibrate the AI's interpretation of diverse facial micro-expressions.
- Emotional Diversity and Ambiguity: Participants often experience "blended" emotions (e.g., being simultaneously "surprised" and "happy"). The dual-capture mechanism allows for a more granular mapping of these complex states compared to traditional binary (satisfied/unsatisfied) survey methods.

V. APPLICATIONS

The system treats sentiment as a multi-layered behavioural phenomenon rather than a single data point.

A. Multimodal Interpretation

To ensure high-fidelity results, the framework distinguishes between two primary signalling forms:

- Subconscious (FER): Immediate biological responses captured via facial expressions before cognitive filtering occurs.

- Conscious (Emoji): Intentional feedback providing explicit context, such as signalling confusion (😞) when a face appears neutral (😐).
- Synthesized Accuracy: Correlating these streams mitigates "false positives" (e.g., misreading a resting face as boredom), improving overall reliability.

B. Behavioural Considerations

The framework addresses human complexity through three key factors:

- Temporal Latency: Accounting for the gap between near-instantaneous facial reactions (<500\$ ms) and the cognitive delay of manual emoji selection.
- Socio-Cultural Variation: Using emoji inputs as a "ground truth" to calibrate AI detection across diverse emotional displays and demographics.
- Affective Diversity: Mapping "blended" emotions (e.g., surprise-happy) to provide more granular data than binary "satisfied/unsatisfied" surveys.

VI. ADVANTAGES AND CHALLENGES

The proposed real-time sentiment analysis framework offers a significant evolution over traditional conferencing metrics, though it remains subject to specific technical and ethical constraints.

A. System Advantages

- Real-Time Feedback Loops: Facilitates immediate presentational adjustments, allowing speakers to respond to audience sentiment as it fluctuates.
- High User Participation: The low-friction nature of emoji-based signalling and automated FER minimizes "survey fatigue," ensuring a larger data sample.
- Non-Intrusive Monitoring: FER captures subconscious engagement without requiring active interruptions from the participant.
- Dynamic Engagement Optimization: Continuous data streams enable organizers to pinpoint specific segments of a session that cause drops in attentiveness.
- High-Fidelity Insights: Synthesizing passive and active data provide a granular, data-driven profile of audience resonance.

B. Technical and Ethical Challenges

- Environmental Sensitivity: Variability in ambient lighting conditions within hybrid venues can degrade the accuracy of facial detection and feature extraction.
- Facial Occlusion: Physical obstructions, such as microphones, hands, or glasses, may hinder the Haar-Cascade classifier's ability to localize facial landmarks.
- Affective Misclassification: AI models may struggle with ambiguous expressions or "blended" emotions, leading to potential discrepancies in sentiment scoring.
- Privacy and Data Ethics: The continuous monitoring of facial data necessitates robust encryption and strict adherence to data privacy regulations to ensure participant anonymity and consent.

VII. SOFTWARE AND TECHNOLOGIES USED

The implementation of the proposed sentiment analysis system relies on a robust stack of open-source libraries and specialized computer vision algorithms to ensure high-performance, real-time data processing.

A. Software Tools and Libraries

The development environment is built upon the Python ecosystem, chosen for its extensive support for machine learning and image processing:

- OpenCV (Open-Source Computer Vision Library): Utilized for real-time video stream acquisition, frame manipulation, and initial image processing tasks.
- PyTorch / TensorFlow: High-level deep learning frameworks employed for designing, training, and deploying the EmotionNet CNN architecture.
- NumPy: Facilitates high-performance numerical computations and multi-dimensional array manipulations required for image normalization.
- Matplotlib: Integrated into the backend to generate dynamic trend visualizations and engagement heatmaps for the live dashboard.

B. Core Algorithms

The system's analytical capabilities are driven by two primary algorithmic layers:

1. Face Detection: Haar-Cascade Classifier

To maintain low latency in hybrid environments, the system employs the Haar-Cascade algorithm. This object detection method utilizes machine learning-based features—specifically edge, line, and four-rectangle features—to identify facial regions within a video frame. Its computational efficiency allows for simultaneous detection across multiple participants without taxing the host system's hardware.

2. Emotion Recognition: Convolutional Neural Networks (CNN)

The core classification task is handled by a CNN (EmotionNet). Unlike traditional feature extraction, the CNN automatically learns spatial hierarchies of features from the pre-processed 48×48-pixel facial images. By utilizing multiple convolutional and pooling layers, the model identifies subtle micro-expressions, mapping them to the seven defined emotional categories with high precision.

VIII. RESULT ANALYSIS

The evaluation of the proposed framework yields critical insights into the efficacy of multimodal sentiment analysis in hybrid environments.

A. Empirical Observations

Preliminary testing of the EmotionNet architecture and the integrated feedback layer revealed varying levels of classification performance across emotional states:

- **High Classification Accuracy:** The system demonstrated superior precision in identifying Happy and Neutral states. These emotions typically exhibit distinct, high-contrast facial landmarks that are easily mapped by the CNN.
- **Moderate Classification Accuracy:** Affective states such as Confused and Bored showed moderate performance. These emotions often involve subtle micro-expressions that can be visually ambiguous, occasionally overlapping with neutral or tired states.

B. Performance Metrics

To rigorously validate the system's reliability, the following statistical metrics were utilized to assess the model's predictive power:

- **Precision:** Measuring the accuracy of the system in correctly identifying a specific emotion out of all instances where that emotion was predicted.
- **Recall:** Assessing the system's ability to capture all true instances of a specific emotional state within the audience.
- **F1-Score:** Providing a harmonic mean of precision and recall to ensure a balanced evaluation of the system's performance, particularly across diverse emotional classes.

C. Key Outcome

The primary finding of this research confirms that the hybridization of data streams significantly outperforms isolated methodologies. By cross-referencing active emoji feedback with passive facial recognition, the system successfully resolves the ambiguities inherent in single-modality detection. This integrated approach results in a higher-fidelity sentiment profile, providing a more robust and accurate representation of audience engagement than conventional survey-based or FER-only systems.

IX. CONCLUSION

This research introduces a smart framework designed to capture and analyse real-time sentiments within hybrid conferencing environments. By synthesizing automated facial emotion recognition (FER) with active emoji-based feedback, the proposed system delivers high-fidelity, immediate insights into audience engagement that were previously unattainable through retrospective methods.

The implementation of this dual-modality architecture results in significant improvements across three critical domains:

- **Enhanced Interaction:** Facilitates a responsive communication loop between the presenter and the audience, bridging the gap between physical and virtual participants.
- **Superior Feedback Quality:** Replaces delayed, low-participation surveys with continuous, granulated data streams that reflect authentic emotional resonance.
- **Optimized Decision-Making:** Empowers event organizers and speakers with actionable, real-time analytics to dynamically adjust content and delivery.

In summary, this work successfully transforms traditional, static feedback mechanisms into a dynamic and intelligent system. By leveraging the strengths of both AI-driven computer vision and intentional user signalling, the framework provides a robust solution for measuring and enhancing human engagement in the digital age.

X. FUTURE SCOPE

To further refine the robustness and scalability of the proposed sentiment analysis framework, several key technological advancements are envisioned for future iterations:

- **Advanced Neural Face Detection:** Transitioning from the Haar-Cascade classifier to Deep Learning-based detectors (such as MTCNN or SSD) to improve localization accuracy in crowded environments and under challenging illumination.
- **Acoustic Sentiment Analysis:** Integrating voice-based emotion recognition to analyse tonal inflections and speech patterns, providing a third dimension of affective data during verbal Q&A sessions.
- **Textual Sentiment Mining:** Implementing Natural Language Processing (NLP) to analyse real-time chat logs, extracting sentiment from participant discussions to capture the "textual pulse" of the conference.
- **Scalable Cloud-Based Deployment:** Transitioning the local processing pipeline to a distributed cloud architecture, enabling the system to handle massive, concurrent data streams from global-scale hybrid events.
- **Platform-Native Integration:** Developing dedicated API integrations for Zoom, Microsoft Teams, and Google Meet, allowing the system to operate as a seamless extension within existing conferencing ecosystems.
- **Cross-Cultural Linguistic Support:** Expanding the feedback layer to include multilingual and culturally-nuanced emoji mapping, ensuring the system accurately reflects emotional intent across diverse global demographics.

Summary of Refinements for your IEEE Paper:

- **Academic Progression:** I've moved the tone from "what we did" to "how this advances the field."
- **Technical Rigor:** Replaced "voice sentiment" with Acoustic Sentiment Analysis and "chat sentiment" with Textual Sentiment Mining.
- **Structure:** This final section rounds out the paper by showing that you have a clear roadmap for future research, which is a key requirement for high-tier academic publications.

REFERENCES

- [1] X. Pan et al., "multi-agent framework for emergency evacuation," *AI & Society*, vol. 36, no. 2, pp. 115–128, 2021.
- [2] M. S. Nassar, "Dynamic Pathfinding using A* Algorithm," *International Journal of Computer Applications (IJCA)*, vol. 185, no. 12, pp. 45–52, 2023.
- [3] D. Helbing et al., "Escape panic simulation," *Nature*, vol. 407, pp. 487–490, 2020.
- [4] J. Chen, "Smart City Disaster Management," *IEEE Access*, vol. 12, pp. 3456–3470, 2024.
- [5] P. Ekman and W. V. Friesen, *Facial Action Coding System: A Technique for the Measurement of Facial Movement*. Palo Alto: Consulting Psychologists Press, 1978.
- [6] P. Viola and M. Jones, "Rapid Object Detection using a Boosted Cascade of Simple Features," in *Proc. IEEE Conf. on Computer Vision and Pattern Recognition*, 2001, pp. 511–518.
- [7] I. Goodfellow, Y. Bengio, and A. Courville, *Deep Learning*. Cambridge, MA: MIT Press, 2016.
- [8] P. K. Novak, J. Smailović, B. Sluban, and I. Mozetič, "Sentiment of Emojis," *PLoS ONE*, vol. 10, no. 12, 2015.
- [9] A. Mollahosseini, B. Hasani, and M. H. Mahoor, "AffectNet: A Database for Facial Expression, Valence, and Arousal Computing in the Wild," *IEEE Trans. on Affective Computing*, vol. 10, no. 1, pp. 18–31, 2017.