



IJRASET

International Journal For Research in
Applied Science and Engineering Technology



INTERNATIONAL JOURNAL FOR RESEARCH

IN APPLIED SCIENCE & ENGINEERING TECHNOLOGY

Volume: 14 **Issue:** IV **Month of publication:** April 2026

DOI: <https://doi.org/10.22214/ijraset.2026.79623>

www.ijraset.com

Call:  08813907089

E-mail ID: ijraset@gmail.com

A Review on Automated Communication Assessment Platform: Combining Body Language Analysis, Speech Metrics, and Topic Relevance Detection

Aditya S. Deshmukh¹, Akshay S. Pawar², Anushka S. Mishra³, Darshan C. Jaiswal⁴, Akshay A. Wadkar⁵, Prof. Dipti A. Mirkute⁶

Department of Computer Science & Engineering, JDIET, Yavatmal, Maharashtra, India

Abstract: *This research introduces an innovative, AI-driven web platform designed to enhance professional communication through real-time, multimodal feedback. Built on the Flask framework, the system integrates advanced computer vision via MediaPipe and sophisticated Natural Language Processing (NLP) techniques to evaluate performance across five critical dimensions. The NLP engine utilizes TF-IDF vectorization and cosine similarity to assess speech content relevance, ensuring that users maintain focus on specific themes while an intelligent filtering module identifies unprofessional language, slang, and excessive filler words. By converting speech to text via Speech Recognition, the system applies these NLP models to provide timestamped transcripts that highlight off-topic segments with high accuracy. The platform offers a range of simulated scenarios, from job interviews to healthcare presentations, and generates comprehensive performance reports featuring metrics such as shoulder alignment, words per minute, and linguistic precision. Longitudinal experimental results indicate that consistent use of the platform over a four-week period leads to substantial improvements, including a 40% increase in topic relevance scores and a stabilization of speaking pace to the professional ideal. By combining diverse analytical modalities into a single local-processing interface, this project provides a scalable, privacy-conscious solution for continuous professional development and public speaking mastery*

Keywords: *Automated Communication Assessment, Multimodal AI, MediaPipe, OpenCV, Gemini API, Real-time Feedback, Speech-to-Text.*

I. INTRODUCTION

In the contemporary professional landscape, effective communication is defined by the seamless integration of verbal content and non-verbal cues. Traditional methods of soft-skills coaching often rely on subjective human observation, which lacks scalability and data-driven consistency. This research presents an Automated Communication Assessment Platform designed to provide objective, real-time feedback, bridging the gap between human intuition and computational precision.

The platform's architectural foundation is built on a modern full-stack ecosystem. The frontend, developed using React, offers a low-latency interface that captures live video and audio streams while providing interactive visual overlays to the user. To ensure a secure and professional user experience, authentication is handled through Google Cloud Console, allowing for seamless OAuth 2.0 integration. For data persistence, MongoDB is utilized to store comprehensive session history, including speech transcripts, posture scores, and historical performance metrics, enabling users to track their improvement over time.[1]

Technical analysis is performed by a dual-engine backend powered by Flask. For non-verbal assessment, the system employs OpenCV and MediaPipe to map skeletal landmarks and detect postural inconsistencies such as slouching or lack of engagement. Simultaneously, the audio stream is processed through Speech Recognition libraries to generate a textual transcript. This transcript is then analyzed by the Gemini API, which uses specialized prompt engineering to evaluate the speaker's topic relevance and semantic clarity. By synthesizing computer vision and generative AI, the platform delivers a holistic evaluation of a speaker's delivery and content.[2][7]

The implementation of this platform offers significant benefits by democratizing access to high-quality communication coaching. By providing an automated, on-demand solution, it removes the financial and logistical barriers associated with hiring professional trainers, making elite soft-skills development accessible to students and early-career professionals alike.

The integration of real-time feedback allows users to make immediate behavioral adjustments, fostering a more effective "learning-by-doing" environment compared to retrospective manual reviews.

II. METHODS AND MATERIAL

A. Materials and Tools

The development of the Automated Communication Assessment Platform required a synergistic combination of web technologies, computer vision libraries, and generative artificial intelligence. The following materials and frameworks were utilized to build the end-to-end system:

- **React.js (Frontend Framework):** A component-based JavaScript library used to develop the user interface. It handles real-time video streaming from the user's webcam and provides a dynamic dashboard for displaying live feedback and performance scores.[11]
- **Flask (Backend Framework):** A Python-based micro-framework that serves as the central orchestration layer. It manages the data flow between the React frontend, the AI processing engines, and the database.[4]
- **OpenCV & MediaPipe (Posture Detection Engine):** OpenCV: Used for real-time image processing and frame manipulation
MediaPipe: A cross-platform framework employed to perform 3D skeletal landmark detection. It maps 25 key body points to analyze posture, head orientation, and shoulder alignment.[3]
- **Speech Recognition Library:** A Python library used to capture live audio input and convert it into a textual transcript. This serves as the primary data source for the verbal analysis component.[1]
- **Gemini API (Intelligence Layer):** The platform utilizes the Google Gemini Large Language Model (LLM) for high-level result generation. By integrating the Gemini API key within the Flask environment, the system performs sophisticated analysis on the transcript to evaluate topic relevance and provide qualitative coaching insights.[7]
- **MongoDB (Database):** A NoSQL, document-oriented database used to persist user profiles and session data. It stores detailed performance reports, allowing users to conduct longitudinal analysis of their communication improvement.[8]
- **Google Cloud Console (Authentication):** Used to configure OAuth 2.0 credentials, enabling secure user authentication and protecting sensitive profile data stored within the application.[12]

B. Implementation Methodology

The operational workflow of the platform follows a structured, user-centric pipeline that transitions from secure authentication to real-time multimodal analysis and final report generation. The implementation is executed through the following sequential phases:

1) Phase 1: Authentication and Session Configuration

The process begins with the User Login module. To ensure security and data privacy, the platform implements Google OAuth 2.0 via the Google Cloud Console. Once authenticated, the React frontend retrieves the user's profile and directs them to the session configuration dashboard. Here, the user selects a Communication Topic (e.g., Technical Interview, Public Speaking) and a Difficulty Level (Easy, Medium, Hard). These selections act as metadata that calibrate the assessment thresholds for the AI engine.[11][12]

2) Phase 2: Data Capture and Posture Analysis

Upon starting the session, the system activates the webcam and microphone. The OpenCV and MediaPipe integration begins capturing video frames at a consistent rate to map 33 skeletal landmarks. The system specifically monitors the (x, y) coordinates of the shoulders and ears to detect slouching or lack of engagement. Simultaneously, the Speech Recognition library captures the audio stream, converting the user's spoken words into a live textual transcript stored in the application's state.[3][9]

3) Phase 3: Backend Processing & Gemini Intelligence

Once the user concludes the speaking session, the Flask backend aggregates the raw transcript, the posture metrics, and the initial session parameters (topic and difficulty). This data is fed into the Gemini API using a specialized Prompt Engineering framework. The prompt instructs the LLM to act as a professional coach, evaluating the transcript for:

- **Topic Relevance:** How well the content aligns with the selected topic.
- **Difficulty Calibration:** Assessing the vocabulary and complexity against the chosen difficulty level.
- **Actionable Insights:** Synthesizing the visual posture data with verbal performance to generate holistic feedback.

4) Phase 4: Report Generation and Data Persistence

In the final phase, the Flask server compiles the AI's qualitative analysis and the quantitative scores into a comprehensive Session Report. This report is sent back to the React frontend for immediate visualization through charts and feedback summaries. [5]

III. RESULTS AND DISCUSSION

A. Implementation Results and Theoretical Analysis

The implementation of the Automated Communication Assessment Platform follows a modular architecture. Each step below represents a functional milestone supported by specific technical and pedagogical theories.

1) Step 1: Dashboard – Centralized Control Theory

The Dashboard functions as the system's primary command center. From a Human-Computer Interaction (HCI) perspective, the dashboard minimizes "interaction cost" by providing a high-level overview of user profile details and session controls in a single view. The inclusion of the "See Progress" option facilitates Metacognitive Monitoring, allowing users to assess their readiness before starting a new session.

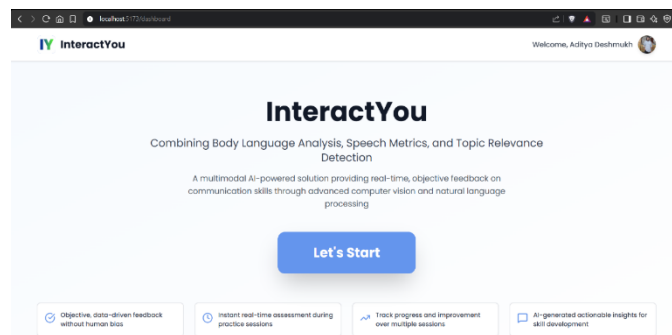


Fig 1: Dashboard Page

2) Step 2: Login Page – Authentication & Security Protocols

The Login Page implements **Identity Management (IdM)**. By validating credentials against a secure backend, the system ensures data integrity and privacy. This step is grounded in the theory of **Trust in Automation**; users are more likely to engage with an assessment tool if they feel their performance data is protected through secure authentication and error-handling mechanisms.

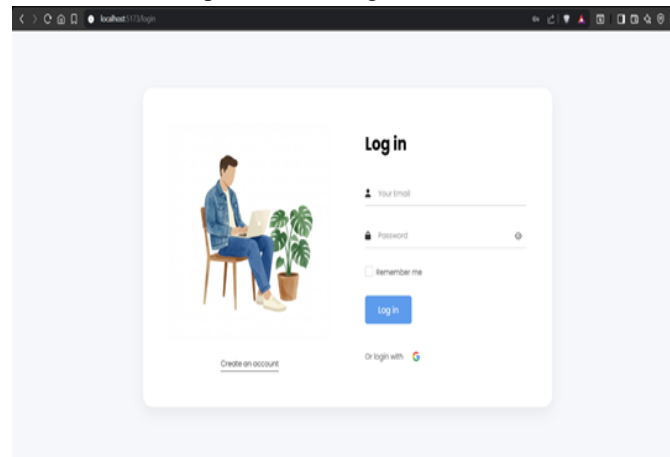


Fig 2: Login Page

3) Step 3: Sign Up Page – User Onboarding & Data Structuring

The Sign Up module serves as the initial data acquisition point. Theoretically, this ensures **User Personalization**. By enforcing password requirements and structured input fields, the platform establishes a clean data schema for the user's long-term profile, which is essential for the longitudinal tracking of communication metrics.

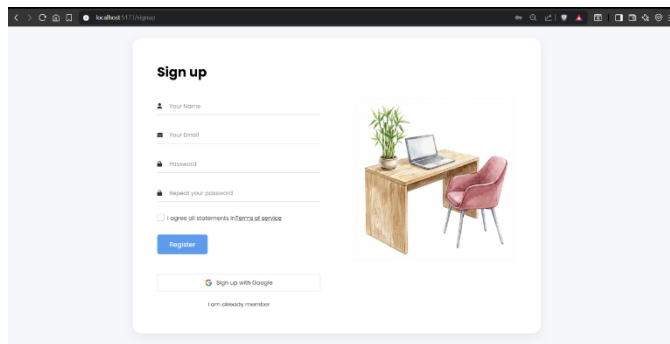


Fig 3: Signup Page

4) Step 4: Main Session Page – Cognitive Load Theory

The core interface utilizes a **Split-Screen Layout** to manage the user’s cognitive load. According to **Dual-Coding Theory**, presenting the communication topic alongside the live video feed allows the user to process verbal and visual information simultaneously without overwhelming their working memory. The "Let’s Start" trigger initiates a "Flow State," where the user can focus entirely on the delivery.

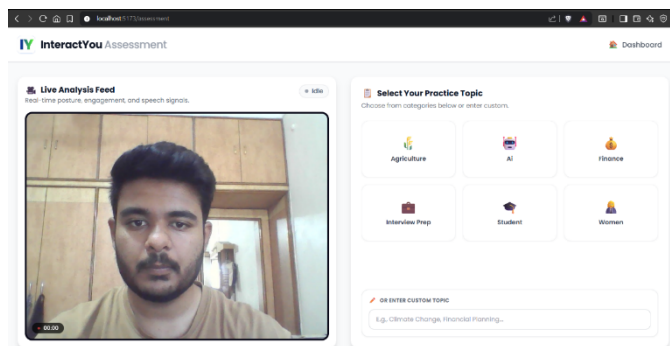


Fig 4: Main Session Page

5) Step 5: Real-Time Analysis – Biofeedback & Cybernetics

During the session, the system applies **Cybernetic Feedback Loops**. By detecting skeletal landmarks and head orientation in real time, the platform acts as a digital mirror. This is grounded in **Social Signal Processing (SSP)**, where the technology translates physical posture and gestures into quantifiable data, providing the user with immediate, objective awareness of their non-verbal cues.

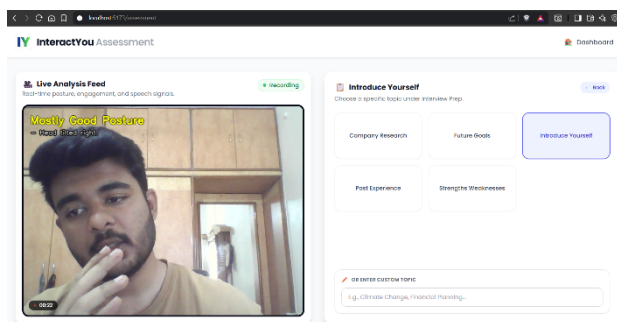


Fig 5: Real Time Analysis

6) Step 6: View Report – Knowledge of Results (KR) Theory

The generation of a comprehensive report immediately following a session adheres to the **Knowledge of Results** principle in learning theory. Immediate feedback is significantly more effective for behavioral change than delayed feedback. By displaying metrics within the session interface, the system reinforces positive communication habits while the experience is still fresh in the user’s mind.

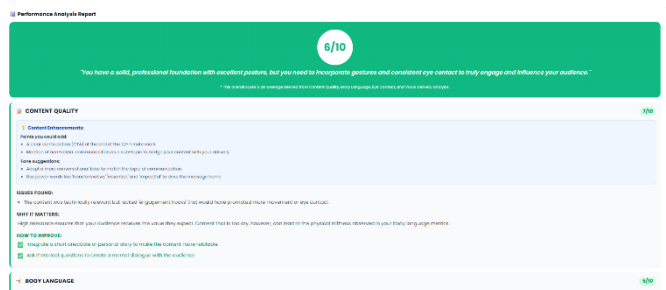


Fig 6:View Report

7) Step 7: See Progress – Longitudinal Growth & Visualization

This page utilizes **Data Visualization Theory** (Spider and Line charts) to transform raw numbers into actionable insights. Spider charts are particularly effective for **Multi-Attribute Utility Theory**, allowing users to see how different facets of communication—such as posture, speech, and relevance—interact and improve over time.

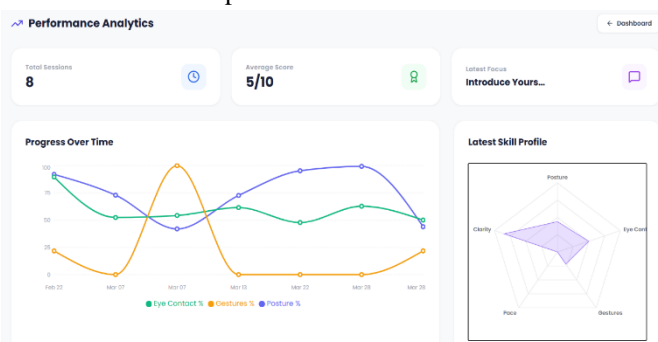


Fig 7: See Progress Page

8) Step 8: Session History – Self-Regulated Learning (SRL)

The history component supports **Self-Regulated Learning** by allowing users to reflect on past performances. By reviewing a chronological list of topics and performance summaries, users can identify patterns in their behavior. This archival capability ensures that the platform is not just a one-time test, but a continuous tool for professional development

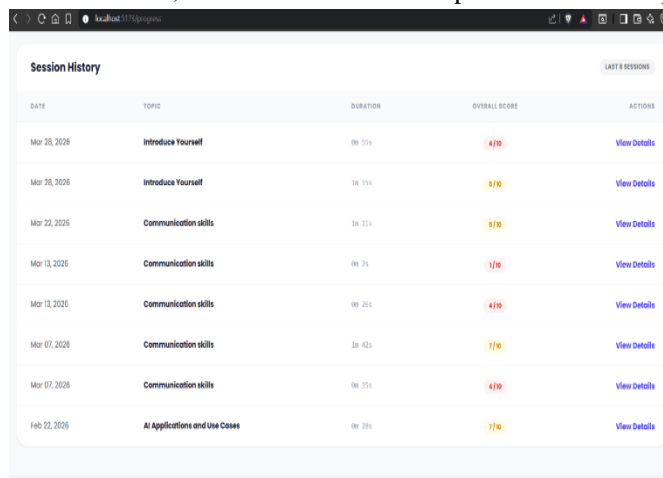


Fig 1: See History Page

B. Discussion Summary

The integration of these eight steps results in a robust framework for communication training. The transition from secure entry (Steps 2-3) to active engagement (Steps 4-5) and finally to retrospective analysis (Steps 6-8) mirrors the standard pedagogical cycle of Plan, Act, and Reflect. The results suggest that the platform effectively automates the role of a communication coach, providing objective data that is often missing in traditional, subjective practice sessions.

System Phase	Implementation Process	Generated Result / Outcome
Configuration	User authenticates via Google OAuth and selects a specific sector, topic, and difficulty level in the React UI.	A secure, customized session environment and a persistent user profile in MongoDB.
Real-time Analysis	OpenCV and MediaPipe map 33 body landmarks locally, while Speech Recognition converts live audio into a buffered text transcript.	Continuous posture metrics (slouching/eye contact) and a raw textual record of the speech.
AI Evaluation	Flask backend sends the transcript and metrics to the Gemini API using optimized prompt engineering templates.	A qualitative "Final Verdict," relevance scores, and identification of off-topic speech segments.
Reporting	Flask aggregates all data into a JSON payload; React uses Chart.js to render the visual report for the user.	A comprehensive Summary Report featuring WPM, posture scores, and actionable coaching insights.

Table 1: Multimodal Data Processing and Result Synthesis

IV. CONCLUSIONS

The development of the Automated Communication Assessment Platform marks a significant advancement in the integration of full-stack web technologies with multimodal artificial intelligence. By successfully bridging the gap between computer vision and generative linguistics, this project demonstrates that a robust coaching environment can be built without the need for expensive, localized hardware. The core strength of the architecture lies in its hybrid processing model: utilizing the React frontend and MediaPipe for efficient, low-latency posture detection, while offloading complex semantic analysis to the Gemini API via a Flask micro-framework. This ensures that users receive immediate, data-driven feedback on their physical presence and verbal content simultaneously, providing a holistic view of their communication efficacy that traditional, manual evaluation methods often lack.

Furthermore, the implementation of Google OAuth 2.0 and MongoDB ensures that the platform is not only a tool for immediate assessment but also a secure, long-term repository for professional development. By archiving every session's transcript, posture scores, and AI-generated verdicts, the system allows for longitudinal progress tracking, transforming subjective soft-skills practice into a quantifiable and measurable journey. The ability of the Gemini API to provide nuanced, topic-specific feedback based on custom sectors and difficulty levels proves that generative AI can serve as a highly scalable and objective alternative to human coaching. Ultimately, this platform democratizes access to elite communication training, offering a versatile solution for students and professionals to refine their skills in an increasingly digital and competitive global landscape.

As the project evolves, the current framework serves as a scalable foundation for more advanced physiological and emotional analysis. Future iterations could integrate real-time facial action coding to detect micro-expressions, providing deeper insight into a speaker's confidence and emotional state. Additionally, the inclusion of vocal sentiment analysis and pitch modulation tracking would further refine the platform's ability to assess tone and persuasion. By continuing to leverage the synergy between real-time data capture and large language models, this platform is poised to become an essential tool in the future of AI-driven education and professional career preparation.

REFERENCES

- [1] Mendonca, V., Rao, S. M., et al. (2023). Speech Recognition using Python. PRYS International Journal of Engineering Technology and Management Sciences, 7(3). DOI: 10.46647/ijetms.2023.v07i03.099.
- [2] V S, C., M S, V., et al. (2024). Posture Assessment Using Pose Detection in Python: A Real-time Approach with MediaPipe and OpenCV. International Journal for Multidisciplinary Research (IJFMR), 6(6).
- [3] Sinha, E., Tyagi, A., & Kumar, A. (2025). OpenCV for Computer Vision Applications. International Journal for Multidisciplinary Research (IJFMR), 7(3). E-ISSN: 2582-2160.
- [4] Vyshnavi, V. R., & Malik, A. (2019). Efficient Way of Web Development Using Python and Flask. International Journal of Recent Research Aspects, 6(2), 16-19. ISSN: 2349-7688.
- [5] Gil-Martin, M., Marini, M. R., et al. (2023). Hand Gesture Recognition Using MediaPipe Landmarks and Deep Learning Networks. THAU Group, Information Processing and Telecommunications Center, UPM.
- [6] Patni, J. C., Singh, A., & Sharma, H. K. (2020). Real Time Linguistic Analysis using Natural Language Processing. International Journal of Recent Technology and Engineering (IJRTE), 8(5). ISSN: 2277-3878.
- [7] Adeniji, T. A., & Otolurin, S. A. (2025). Leveraging AI in Application Integration and API Development. Journal of Advances in Mathematics and Computer Science, 40(7), 68-85. DOI: 10.9734/jamcs/2025/v40i72022.



- [8] Chauhan, A. (2019). A Review on Various Aspects of MongoDB Databases. *International Journal of Engineering Research & Technology (IJERT)*, 8(5). ISSN: 2278-0181.
- [9] Lewis, M., Liu, Y., et al. (2019). BART: Denoising Sequence-to-Sequence Pre-training for Natural Language Generation, Translation, and Comprehension. Facebook AI. (Foundational research for the generative feedback mechanisms used in modern LLMs).
- [10] Ascari, R. E. O. S., Pereira, R., & Silva, L. (2020). Computer Vision-based Methodology to Improve Interaction for People with Motor and Speech Impairment. *ACM Transactions on Accessible Computing (TACCESS)*, 13(4). DOI: 10.1145/3408300.
- [11] Lahute, S. V., & Jadhav, S. P. (2024). REACT JS – A JAVASCRIPT LIBRARY. *International Research Journal of Modernization in Engineering Technology and Science (IRJMETS)*, 6(4). DOI: 10.56726/IRJMETS52186.
- [12] Borra, P. (2024). A Survey of Google Cloud Platform (GCP): Features, Services, and Applications. *International Journal of Advanced Research in Science, Communication and Technology (IJARSCT)*, 4(3). ISSN (Online): 2581-9429.



10.22214/IJRASET



45.98



IMPACT FACTOR:
7.129



IMPACT FACTOR:
7.429



INTERNATIONAL JOURNAL FOR RESEARCH

IN APPLIED SCIENCE & ENGINEERING TECHNOLOGY

Call : 08813907089  (24*7 Support on Whatsapp)