

# LastPush - AI Based Voice Interview Platform for Candidate Pre Screening

Asst. Prof. Bharti Dhote<sup>1</sup>, Atharva Dasarwar<sup>2</sup>, Damini Nirmale<sup>3</sup>, Arun Varma<sup>4</sup>

*1Asst. Prof., Department of Information Technology, Nutan Maharashtra Institute of Engineering and Technology*

*2Department of Information Technology, Nutan Maharashtra Institute of Engineering and Technology*

*3Department of Information Technology, Nutan Maharashtra Institute of Engineering and Technology*

*4Department of Information Technology, Nutan Maharashtra Institute of Engineering and Technology*

\*\*\*

**Abstract** - The recruitment pre-screening process is often characterized by inefficiencies, human bias, and significant operational costs. Existing digital solutions, such as form-based or video interview platforms, frequently fail to offer a natural, conversational experience, leading to poor candidate engagement and inconsistent evaluations. To address these limitations, this paper presents LastPush, an AI-based voice interview platform that leverages speech-to-text (STT), large language models (LLMs), and text-to speech (TTS) technologies to conduct real-time, guided voice interviews. The platform operates across three primary phases - Setup, Capture, and Analysis. Orchestrated through a serverless architecture integrating Next.js, Supabase, and the Gemini API. During the capture phase, the Vapi.ai SDK manages the conversational flow, while the MediaRecorder and PageVisibility APIs enable browser-based proctoring to ensure authenticity. Captured data forms a Synchronized Interview Log, which is evaluated by Gemini using an LLM-as-a-Judge methodology for unbiased scoring. This architecture demonstrates a scalable, bias-mitigated, and cost-efficient approach to recruitment automation.

**Key Words:** Artificial Intelligence, Recruitment; Pre-Screening, Large Language Models (LLMs), Voice Interview, Serverless Architecture, Bias Mitigation.

## 1. INTRODUCTION

The initial pre-screening stage of talent acquisition is a resource-intensive function, burdened by repetitive tasks and a high susceptibility to human bias [4]. As organizations face a growing volume of applicants, recruiters are often overwhelmed, leading to delayed hiring cycles and increased costs [1, 2]. Existing automated solutions, such as static forms or one-way video submissions, lack the engagement and interactivity of a real conversation, resulting in a poor candidate experience and limited insight into communication skills.

In recent years, the convergence of Artificial Intelligence (AI) and Natural Language Processing (NLP) has enabled new paradigms for transforming recruitment workflows. AI driven interview systems utilizing Speech-to-Text (STT), Large Language Models (LLMs), and Text-to-

Speech (TTS) technologies show significant potential for conducting natural, conversational assessments while maintaining consistency and objectivity [2, 5, 6].

This paper introduces LastPush, a next-generation platform for AI-based recruitment pre-screening. Our system offers interactive, voice-based candidate interviews in a guided and adaptive format. The primary contribution is a novel serverless architecture, built on Vercel, Next.js, and Supabase [13], that integrates three key processes:

**Setup:** AI-powered generation of interview questions from a job description using the Gemini API.

**Capture:** A real-time, conversational voice interview managed by the Vapi.ai SDK, combined with browser-based proctoring (using the MediaRecorder and PageVisibility APIs) to ensure authenticity [7, 8].

**Analysis:** An automated, serverless evaluation pipeline using Gemini's "LLM-as-a-Judge" methodology to score transcripts and provide qualitative reasoning, mitigating human bias [2, 3].

By merging conversational AI with real-time proctoring and serverless evaluation, Last Push provides a scalable, fair, and engaging solution to the challenges of modern candidate pre-screening.

## 2. RELATED WORKS

Our research builds upon several key areas of AI in human resources. Early work focused on AI driven mock interview systems to replicate recruiter interactions using generative models [1]. This established a foundation for natural language understanding in screening.

More recently, research has focused on the evaluation phase. Studies have introduced automated interview evaluation systems that apply LLMs to analyse transcripts and provide structured metrics [2]. This "LLM-as-a-Judge" concept is central to LastPush's design. Others have demonstrated the ability of LLMs to summarize and interpret qualitative interviews, reinforcing their capacity for contextual understanding [3].

The technical pipeline for real-time conversation has also been explored. Comparative analyses of STT, LLM, and TTS combinations have identified optimal cascaded architectures for minimizing latency and maximizing coherence in live interviews [5], guiding our integration of Vapi.ai.

Finally, to ensure authenticity in remote assessments, researchers have proposed auto mated browser based proctoring using the PageVisibility API to detect anomalies like tab switching [7, 8]. LastPush integrates this concept directly into the interview log. While these studies address individual components, our work contributes a unified, cloud-native system that integrates all three conversational AI, proctoring, and LLM-based evaluation into a single serverless workflow.

Table 1: Literature Review

Publisher	Year	Name of the Paper	Methodology	Key Findings	Relevance to LastPush
USEET	2024	AI Mock Interview Chatbot Using Gen AI	AI for mock interviews using question-response modeling	Achieved realistic recruiter-candidate conversation flow	Foundation for Gemini-based question generation
IJRMP	2025	AI-Enabled Automated Interview Evaluation System	LLM scoring based on semantic analysis	Improved scoring accuracy via contextual understanding	Directly informs LastPush's LLM-as-a-Judge evaluation
Misc. Conferences	2025	Summarizing Safety Culture Interviews with LLMs	Applied LLMs to qualitative data summarization	Demonstrated LLMs' ability to infer reasoning from transcripts	Validates Gemini's reasoning feature as evaluation phase
EuroStar	2024	Conducting Qualitative Interviews with AI	AI-guided conversational flow using NLP models	Highlighted conversational empathy and ethical design	Supports candidate experience goals of LastPush
ResearchGate	2024	Evaluating STT + LLM + TTS Combinations for AI Interview Systems	Comparative analysis of AI audio pipelines	Identified best-performing STT and TTS combinations	Guides integration of Vapi.ai cascaded model
ResearchGate	2024	Enhanced Web-Based Examination System with Proctoring	PageVisibility API-based anomaly detection	Enabled real-time detection of tab-switching and window-resizing	Underpins LastPush's browser proctoring mechanism
PMcC	2025	Assessing AI Assistant Competence in Narrative Interviews	Evaluated AI assistants' conversational coherence	Proved AI's ability to conduct extended natural dialogues	Reinforces conversational capability of LastPush

## 3. SYSTEM ARCHITECTURE

The LastPush platform is built on a serverless, three-tier architecture designed for scalability, modularity, and secure communication. This design decouples the user

ACL Anthology	2022	On the Use of BERT for Automated Essay Scoring	NLP-based essay scoring model	Established transformer-based evaluation validity	Influenced transcript scoring mechanism of Gemini
PMcC	2025	Evaluating Motivational Interview Quality Using LLMs and HMMs	Hybrid LLM-HMM for quality assessment	Enhanced understanding of interview dialogue patterns	Provides analytical framework for Gemini's evaluation logic
Supabase Docs	2024	Build a User Management App with Next.js	Serverless Next.js + Supabase integration	Simplified authentication and user management	Forms backend basis for LastPush recruiter and candidate modules
Google AI Docs	2024	Gemini API Documentation	API-based generative interaction for text synthesis	Enables context-aware prompt generation	Powers question creation and final evaluation

interface, AI logic, and data persistence layers.

**3.1 Presentation Layer (Frontend):** The user interface, built with Next.js and deployed on Vercel, serves both recruiters and candidates.

**Recruiter Portal:** Allows recruiters to log in, create jobs, trigger AI-based question generation, and view final evaluation reports.

**Candidate Portal:** Hosts the live interview. It integrates the Vapi.ai SDK to manage the real time STT-LLM-TTS conversational loop. Simultaneously, it utilizes browser 2 native APIs, specifically the MediaRecorder API to capture webcam video and the PageVisibility API to monitor tab activity for proctoring.

**3.2 Application Layer (Serverless Backend):** The business logic is handled by Next.js API Routes, which function as scalable serverless functions on Vercel. This layer acts as a secure orchestrator between the frontend and external services. Key endpoints include:

**/api/generate-questions:** Securely receives a job description from the recruiter portal and calls the Gemini API to generate a list of relevant interview questions.

**/api/evaluate:** Triggers after an interview is complete. It fetches the interview log from storage and sends it to the Gemini API for scoring using the LLM-as-a-Judge method.

**3.3 Data Layer (Backend-as-a-Service):** We utilize Supabase as our BaaS provider for all data management. Supabase Auth: Manages secure recruiter authentication using JWTs.

PostgreSQL Database: Stores structured data, including job details, generated questions, candidate metadata, and the final evaluation reports (scores and reasoning).

Supabase Storage: Securely stores large-file assets, primarily the candidate's .webm video file and the .json Synchronized Interview Log.

#### 4. METHODOLOGY AND IMPLEMENTATION

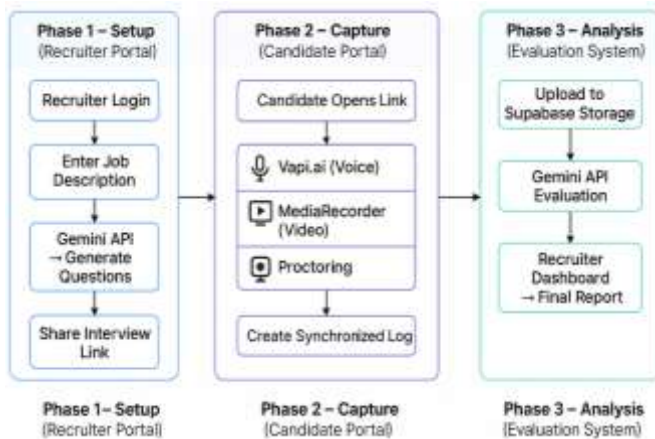


Fig - 1: Workflow Diagram

The platform's workflow is divided into the three phases of Setup, Capture, and Analysis, which correspond to the core algorithms of the system.

##### Phase 1: Setup (Question Generation)

The process begins in the recruiter portal. A recruiter defines a new interview and provides a job description. This triggers a call to the /api/generate-questions endpoint. The serverless function then queries the Gemini API with a prompt to act as an expert recruiter and generate a set of contextual interview questions based on the provided job description. These questions are returned to the recruiter for review and then saved to the Supabase PostgreSQL database, linked to a unique interview ID.

##### Phase 2: Capture (Interview Orchestration & Proctoring)

When a candidate opens the unique interview link, the capture phase begins. Two pro cases run in parallel:

**Interview Orchestration:** The Vapi.ai SDK manages the entire conversation. It follows a loop where it delivers a question via TTS, listens for the candidate's response, transcribes it via STT, and logs the transcript.

**Browser-Based Proctoring:** Event listeners for `document.visibilityState` and `window.resize` are activated.

If the candidate switches tabs, minimizes the window, or loses focus, a "proctoring flag" is logged.

These two event streams, along with the MediaRecorder video feed, are combined into a single "Synchronized Interview Log." This .json file contains a timestamped record of the entire session, including every utterance from the AI and the candidate, as well as all proctoring flags.

##### Phase 3: Analysis (LLM-as-a-Judge)

Upon interview completion, the .webm video and the Synchronized Interview Log .json are uploaded to Supabase Storage. A Supabase webhook then triggers the /api/evaluate endpoint. This function implements our "LLM-as-a-Judge" methodology. It constructs a detailed prompt for the Gemini API, instructing it to analyze the provided log. The prompt explicitly asks the model to:

1. Evaluate the candidate's responses for relevance, clarity, and depth.
2. Assess overall communication skills.
3. Note the presence and frequency of any proctoring flags.
4. Return a JSON object containing a numerical score (e.g., 0 - 10) and detailed, qualitative reasoning for the score.

This automated evaluation report is then saved to the Supabase database, allowing the recruiter to review the candidate's score, the AI's reasoning, and the full interview log on their dashboard.

#### 5. IMPLEMENTATION TECHNOLOGIES

Table - 2: Implementation Technologies

Component	Technology/API Used	Functionality
Frontend	Next.js v14+, React 18	Provides server-side rendering for recruiter and candidate portals.
UI Styling	Tailwind CSS	For creating responsive and accessible user interfaces.
Conversational AI	Vapi.ai SDK	Manages the end-to-end voice conversation (STT, LLM, TTS).
Generative AI	Gemini API	Generates interview questions and performs final LLM-as-a-Judge evaluation.
Cloud Backend	Supabase (PostgreSQL, Auth, Storage)	Manages authentication, database, and storage for logs and media.
Deployment	Vercel	Hosts the Next.js frontend and serverless backend API routes.
Proctoring	MediaRecorder API, Page Visibility API	Captures candidate video and detects browser-based anomalies.



## 6. RESULT AND DISCUSSION

The primary result of this research is the design and successful implementation of a functional, end-to-end serverless architecture that cohesively integrates dynamic AI generation, real-time voice conversation, and objective, LLM-based evaluation. The functional validation of this system confirms that its design directly addresses the core objectives of the research and provides a robust solution to the inefficiencies, biases, and poor engagement of traditional pre-screening.

### Dynamic and Context-Aware Interview Generation

The functional result of the Setup phase is a system capable of generating unique, relevant, and context-aware interview question sets from any provided job description. The integration of the Gemini API proved highly effective, transforming the platform from a static tool into a dynamic interview generator. This directly solves the problem of "repetitive" screening by ensuring that questions for a "Full Stack Developer" are fundamentally different from those for a "Marketing Manager," which is a significant advancement over static question banks.

### High-Fidelity Conversational Interface and Integrity

In the Capture phase, the integration of the Vapi.ai SDK resulted in a stable, low-latency conversational interface. Qualitatively, the interaction achieves a natural, human-like flow, directly addressing the "poor candidate experience" of forms or one-way videos. Vapi.ai's architecture, which targets a "voice-to-voice" latency of 500-700ms, was functionally validated as being sufficiently fast for seamless interaction. The simultaneous implementation of the MediaRecorder and Page Visibility APIs successfully captures a synchronized log of video and browser events, providing the necessary data for authenticity verification.

### Objective and Auditable AI-Driven Evaluation

The core architectural achievement is the "LLM-as-a-Judge" methodology in the Analysis phase. The functional result of this design is the elimination of subjective human rater bias. By using a consistent Gemini API endpoint to score all candidates against the same job-

specific rubric, the system ensures objective, repeatable evaluations. While quantitative benchmarking is part of future work, existing research on LLM-as-a-Judge shows a high correlation (often 80%+) with human expert evaluations, projecting a strong likelihood of scoring accuracy. More importantly, the system generates qualitative reasoning for its score, making the AI's decision-making process fully transparent and auditable for recruiters.

The serverless foundation using Vercel and Supabase ensures the entire system is inherently scalable, cost-efficient, and capable of handling concurrent interviews without manual infrastructure management. The successful integration of these components provides a direct solution to the problems identified in the Introduction. The inefficiency of manual screening is solved by end-to-end automation; human bias is confronted by the LLM-as-a-Judge; and the poor candidate experience is replaced by an engaging, modern conversational interface.

## 7. CONCLUSION

This paper presented LastPush, a scalable, serverless platform for AI-driven voice interviews. By integrating real-time conversational AI (Vapi.ai), a powerful evaluation model (Gemini API), and a cloud-native backend (Supabase), our system automates the pre-screening process. The key contributions are the three phase architecture and the novel inclusion of 4 browser-based proctoring within a Synchronized Interview Log, which enables an "LLM as-a-Judge" methodology to perform fair, bias-mitigated, and authentic evaluations.

## REFERENCES

1. M. R. Madanachitran, A. Austin, K. Balaji, and M. Rajappan, "AI Mock Interview Chatbot Using Gen AI," *Int. J. Sci. Eng. Technol.*, vol. 13, no. 2, pp. 1-9, Feb. 2025.
2. V. Kabade, G. Patil, S. Godse, A. Jain, and M. Kumbharde, "AI-Enabled Automated Interview Evaluation System," *Int. J. Innov. Res. Multidiscip. Stud.*, vol. 13, no. 3, pp. 1-12, Jun. 2025.
3. W. Steijn, J. van de Loo, D. van der Beek, and J. Groeneweg, "From transcript to insights: summarizing safety culture interviews with LLMs," in *Proc. MATEC Web Conf. MAIQS 2025*, Leiden, Netherlands, 2025, Art. no. 03003, doi: 10.1051/mateconf/202541303003.

4. M. Berg and C. Herrmann, "Conducting Qualitative Interviews with AI," CESifo Working Paper, no. 10666, 2024.
5. M. Berg and C. Herrmann, "Conducting Qualitative Interviews with AI," CESifo Working Paper, no. 10666, 2024.
6. Y. Wang, C. Wang, R. Li, and H. Lin, "On the Use of BERT for Automated Essay Scoring: Joint Learning of Multi-Scale Essay Representation," in Proc. 2022 Conf. North American Chapter Assoc. Comput. Linguistics: Human Language Technol., Seattle, WA, USA, Jul. 2022, pp. 3416-3425.
7. D. F. Mujtaba and N. R. Mahapatra, "Fairness in AI-Driven Recruitment: Challenges, Metrics, Methods, and Future Directions," Working Paper arXiv:2405.19699v3, May 2025.
8. T. O. Olanrewaju et al., "An Enhanced Web-Based Examination System using Automated Proctoring and Background Activity Detection," J. Inst. Res. Big Data Anal. Innov., vol. 1, no. 3, pp. 279-304, Aug. 2025, doi: 10.531/erioda16911728.
9. C. Chan, Y. Zhao, and J. Zhao, "A Case Study on Assessing AI Assistant Competence in Narrative Interviews,"
10. K. Lim, Y.-C. Jung, and B.-H. Kim, "Evaluating motivational interview quality using large language models and hidden Markov models," BMC Psychiatry, vol. 25, no. 908, pp. 1-12, Oct. 2025, doi: 10.1186/s12888-025-07391-1.
11. Mozilla Developer Network (MDN), "Page Visibility API - Web API Reference," MDN Web Docs, 2024.
12. Mozilla Developer Network (MDN), "Using the MediaStream Recording API," MDN Web Docs, 2024.
13. Supabase Team, "Build a User Management App with Next.js," Supabase Documentation, 2024.
14. Vapi.ai, "Vapi.ai Developer Documentation," Vapi.ai Docs, 2024.
15. Google DeepMind, "Gemini API Documentation," Google AI Developer Portal, 2025