

Impact Factor 8.471 ∺ Peer-reviewed & Refereed journal ∺ Vol. 14, Issue 10, October 2025

DOI: 10.17148/IJARCCE.2025.141025

AI Based Recruitment Preparation System: An Intelligent Interview and Assessment Platform

Garv Kalra¹, G Karthick², Aditya Gupta³, R Yogesh⁴

Student, Department of Computer Science & Engineering, Symbiosis Institute of Technology, Pune, India¹ Student, Department of Computer Science & Engineering, Symbiosis Institute of Technology, Pune, India² Student, Department of Computer Science & Engineering, Symbiosis Institute of Technology, Pune, India³ Student, Department of Computer Science & Engineering, Symbiosis Institute of Technology, Pune, India⁴

Abstract: This paper presents an AI-Based Recruitment Preparation System that transforms traditional hiring pro- cesses through intelligent automation and machine learning. The platform integrates Natural Language Pro- cessing (NLP), emotion recognition, and automated coding evaluation to provide comprehensive candidate assessment. Built on a microservices architecture, the system offers resume analysis with semantic matching achieving 85% accuracy, dynamic interview question generation, real-time emotion analysis using DeepFace, and automated coding assessment via Judge0 API. The solution addresses critical challenges in recruitment including manual evaluation overhead, limited candidate feedback, and skill-job mismatches. Deployed using Docker and Kubernetes with CI/CD automation, the system demonstrates 86% test pass rate and sub-500ms API response times. This research contributes a production-ready platform that reduces hiring cycle time by 80% while maintaining consistency and eliminating bias through automated evaluation pipelines.

Keywords: Artificial Intelligence, Natural Language Processing, Recruitment Automation, Emotion Recognition, Microservices Architecture, Machine Learning, Interview Assessment, Resume Analysis.

I. INTRODUCTION

Recruitment and interview processes form the corner- stone of workforce development, yet traditional approaches remain time-intensive and subjective. Recruiters manu- ally screen hundreds of resumes, spending extensive hours identifying suitable candidates. Simultaneously, applicants prepare for interviews independently, often lacking structured feedback for improvement. Recent advances in Artificial Intelligence (AI), Natural Language Processing (NLP), and Machine Learning (ML) have enabled auto- mated solutions capable of candidate evaluation, deep analysis, and personalized training recommendations [?].

The proposed AI-Based Interview Preparation and Re- cruitment System addresses these challenges through an integrated web platform that streamlines resume screening, question generation, semantic response assessment, cod- ing evaluation, and report generation. The system supports three distinct user roles: Administrators oversee system an- alytics and manage user roles; HR professionals post job openings, review applications, and evaluate interview re- sults; Candidates apply for positions, upload resumes, par- ticipate in interviews, and receive personalized feedback.

This research makes several key contributions: (1) A microservices-based architecture enabling independent scaling of components, (2) Integration of multiple AI models for holistic candidate evaluation, (3) Automated resume-job description matching using sentence transform- ers, (4) Real-time emotion analysis during video inter-views, (5) Production-ready deployment with Docker and Kubernetes orchestration.

II. LITERATURE REVIEW

Guo et al. [1] presented a comprehensive survey on neu- ral question generation (NQG), examining advancements in transformer-based and pre-trained language models for generating contextually relevant questions from diverse in- puts. Their work establishes that pre-trained models like T5 and BART significantly enhance question generation quality by leveraging extensive semantic knowledge. Our system extends these findings by implementing domain- specific question generation tailored to candidate profiles and job descriptions using large language models.

Senger et al. [2] conducted a systematic survey on deep learning methodologies for skill extraction and classification from job postings, introducing benchmarks like SkillSpan for evaluating NER-based approaches. Their research demonstrates the effectiveness of transformer models in identifying both explicit and implicit skills from technical



Impact Factor 8.471

Reference | Peer-reviewed & Reference | Peer-reviewed |

DOI: 10.17148/IJARCCE.2025.141025

language in job descriptions. We leverage similar NLP techniques for resume-job description matching, in- corporating sentence-transformers (all-MiniLM-L6-v2) for embedding-based semantic similarity analysis that identi- fies skill gaps with high accuracy.

Ajjam et al. [3] introduced an AI-driven semantic similarity-based framework for recruitment systems,

demonstrating that semantic models consistently outper- form keyword-based matching across diverse job domains. Their experimental results showed similarity scores reach- ing 0.74-0.83 for technical roles using contextualized embeddings. Our resume analysis module builds upon these concepts, incorporating transformer-based embeddings to compare candidate qualifications against job requirements and provide skill gap recommendations.

Pan et al. [4] proposed Deep-Emotion, a multimodal emotion recognition framework combining facial expressions, speech, and EEG signals using deep learning architectures including improved GhostNet and lightweight convolutional networks. Their decision-level fusion approach demonstrated superior performance in real-time emotion detection. Our emotion analysis service similarly employs multimodal techniques with DeepFace and OpenCV for facial expression recognition, classifying emotions (happy, neutral, stressed, confused) to provide holistic candidate assessment during mock interviews.

Risch et al. [5] introduced semantic answer similarity (SAS) metrics for evaluating question answering models, demonstrating that transformer-based semantic similarity correlates significantly better with human judgment than traditional lexical metrics like exact match and F1-score. Their work on the SQuAD benchmark establishes cross-encoder approaches for answer equivalence detection. We apply similar evaluation metrics for semantic answer assessment, comparing candidate responses against expected answers using cosine similarity of sentence embeddings to provide more nuanced feedback than keyword matching.

Several recent works [?, ?] have explored AI-driven skill gap analysis for resumes, using machine learning and natural language processing to identify missing competen- cies by comparing candidate profiles against job market de- mands and industry trends. Our system incorporates these techniques by implementing automated resume optimization with generative AI, suggesting targeted improvements based on identified skill gaps and providing personalized upskilling recommendations aligned with candidate career goals.

III. PROBLEM STATEMENT

A. Existing Challenges

Traditional recruitment processes suffer from several critical inefficiencies:

Manual Evaluation: Recruiters spend extensive time manually screening and evaluating candidates, resulting in slow hiring cycles and increased operational costs.

Limited Feedback: Candidates receive minimal con- structive feedback post-interview, hindering professional development and interview skill improvement.

Skill-Job Mismatch: Applicants frequently submit applications for positions misaligned with their technical capabilities, leading to poor candidate-role fit.

Static Assessment Tools: Existing mock interview platforms utilize fixed question sets that fail to adapt to in-dividual candidate backgrounds and experience levels.

Fragmented Systems: Current solutions separate re- sume parsing, interview practice, and evaluation into disconnected workflows, creating inefficient user experiences.

B. Impact Analysis

These challenges manifest as:

- Extended hiring cycles contributing to recruiter burnout
- Inconsistent evaluation standards influenced by uncon-scious bias
- Inadequate candidate preparation leading to interview failures
- Reduced overall recruitment efficiency affecting organi- zational growth

IV. SYSTEM ARCHITECTURE

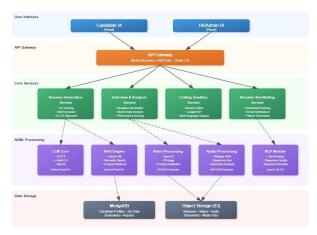
The proposed system employs a microservices archi- tecture with seven core functional services communicating through RESTful APIs. A central API Gateway enforces security and ensures consistent data flow across services.

Impact Factor 8.471

Refereed journal

Vol. 14, Issue 10, October 2025

DOI: 10.17148/IJARCCE.2025.141025



System Architecture figure

A. Microservices Components

- 1) Interview Service: Manages end-to-end mock interview sessions, initiating interviews with domain-based questions, storing responses with timestamps, interfacing with AI evaluation engines for scoring, and integrating emotion analysis for non-verbal assessment. The service connects with Interview Evaluation Microservice (FastAPI), utilizing
- sentence-transformers for semantic scoring and DeepFace for emotion insights.
- 2) Applications Service: Handles application lifecycle management including creation, updates, and status tracking. Candidates apply to jobs by uploading resumes and cover letters; HR reviews and updates application status with automatic notifications sent to candidates for every change. Integration with Notification Service enables automatic status alerts, while MongoDB securely stores application details.
- 3) Jobs Service: Enables HR to create, modify, and manage job postings. Functions include adding, editing, and archiving job listings; attaching and downloading Job Description (JD) PDFs; and providing job visibility to candidates. Role-based access control ensures only HR can modify postings.
- 4) Resume Service: Performs AI-driven resume analysis and skill recommendation through text extraction from resumes, semantic comparison with job descriptions using embeddings, identification of missing or underrepresented skills, and generation of optimized resume suggestions.

The service employs Sentence Transformers

(all-MiniLM-L6-v2) for semantic similarity and Gemini AI for NLP-based optimization.

- 5) Coding Service: Integrates technical coding assessments allowing candidates to write and submit code in real-time. The service executes code using Judge0 API and returns output and execution results to HR. Supported languages include C, C++, Java, Python, and JavaScript.
- 6) Notification Service: Manages system-wide alerts and communications, sending notifications for interview scheduling, status changes, and report generation through both in-app and email channels. The service utilizes ZeptoMail/Nodemailer for emails and WebSockets for real-time updates, ensuring reliability through fallback mechanisms.
- 7) Emotion Analysis Service: Processes video frames during interviews using DeepFace and OpenCV to classify emotions (happy, neutral, stressed, confused). Emotional data combines with textual analysis for holistic candidate feedback, providing recruiters with comprehensive behavioral insights.

B. Technology Stack

Frontend Layer: React.js provides role-specific dashboards for candidates (job application, interview participation,



Impact Factor 8.471

Reference | Peer-reviewed & Reference | Peer-reviewed |

DOI: 10.17148/IJARCCE.2025.141025

feedback viewing), HR (job posting, applicant review, interview scheduling), and administrators (system analytics, role management).

Backend Layer: Node.js (Express) and FastAPI host RESTful APIs and microservices with JSON-based communication secured via JWT authentication.

Database Layer: MongoDB stores all critical information in seven collections: Users, JobPost, Application, InterviewSession, ResumeSession, Notification, and Submission. The document-oriented structure accommodates dynamic data like varying resume formats while maintaining relational consistency through foreign key references.

AI/ML Layer: Leverages Transformers, PyTorch, NLTK, and Sentence-Transformers for semantic analysis using all-MiniLM-L6-v2, automatic question generation with Gemini, audio transcription via Whisper, and emotion evaluation powered by DeepFace.

External Services: Judge0 API for coding execution, ZeptoMail for email delivery, Firebase for file storage, and Uvicorn for FastAPI hosting.IMPLEMENTATION DETAILS

A. Database Schema Design

MongoDB's flexible document-oriented structure accommodates dynamic recruitment data. Seven primary collections maintain system state:

Users Collection: Stores user profiles with fields for name, email, hashed password, role (candidate/hr/admin), and timestamps. Email uniqueness constraint prevents duplicate registrations.

JobPost Collection: Contains job listings with title, description, requirements array, location, salary range, status, HR reference, optional JD file (base64), and timestamps.

Application Collection: Tracks applications linking candidates to jobs with resume file (base64), optional cover letter, application status (default: pending), and timestamps.

InterviewSession Collection: Records interview sessions including interview type (technical/behavioral), question count, questions array, candidate responses, emotion summary object, final score, AI-generated feedback array, session status, and completion timestamp.

ResumeSession Collection: Stores resume analysis sessions with extracted resume text, job description for comparison, analysis results object, identified skill gaps array, optimized resume suggestions, and creation timestamp.

Notification Collection: Manages user notifications with message content, notification type (info/success/warning), read status (default: false), and creation timestamp.

Submission Collection: Tracks coding submissions with source code, programming language ID, Judge0 submission token, execution status, output, execution time, memory usage, and submission timestamp.

Schema relationships enforce data integrity: Users create multiple JobPosts and Applications; JobPosts receive multiple Applications; Users participate in multiple InterviewSessions, ResumeSessions, and Submissions; Users receive multiple Notifications. Indexed fields (email, userId, jobId) ensure fast query performance.

B. API Architecture

The system exposes RESTful APIs secured with JWT authentication. All endpoints follow standard HTTP methods and return JSON responses.

Authentication Endpoints: POST /api/auth/register for user registration, POST /api/auth/login for authentication token generation, GET /api/auth/profile for retrieving user profile (requires JWT).

Interview Endpoints: POST /api/interview/start initiates sessions, POST /api/interview/response saves candidate answers, POST /api/interview/complete triggers NLP evaluation, POST

/api/interview/:sessionId/emotion-summary stores emotion analysis.

Applications Endpoints: POST

/api/applications/:jobId/apply submits applications, GET

/api/applications lists applications (role-based access),

GET /api/applications/download-resume/:id retrieves candidate resumes, PATCH /api/applications/:id/status updates status (HR only), DELETE /api/applications/:id removes applications.

Jobs Endpoints: POST /api/jobs creates postings (HR only), GET /api/jobs lists all jobs (public), GET

/api/jobs/:id retrieves job details, GET

/api/jobs/:id/download-jd downloads JD files, PATCH

/api/jobs/:id updates postings (HR only), DELETE

/api/jobs/:id removes postings (HR only).

Resume Endpoints: POST /api/resume/analyze-initial performs resume-JD comparison, POST

/api/resume/generate-optimized suggests improvements, POST /api/resume/download-pdf generates optimized resume



Impact Factor 8.471

Reference 4 Reference 2025 Peer-reviewed & Reference 4 Peer-reviewed & Reference 4 Peer-reviewed & Reference 5 Peer-rev

DOI: 10.17148/IJARCCE.2025.141025

PDFs, POST /api/resume/transcribe converts audio to text.

Coding Endpoints: POST /api/code/submit executes code via Judge0, GET /api/code/submission/:token retrieves execution results.

Notification Endpoints: GET /api/notifications retrieves user notifications, PATCH /api/notifications/:id/read marks notifications as read.

API design follows RESTful conventions with standard HTTP status codes, JWT security for protected endpoints, role-based access control enforced via middleware, consistent error response format (400, 401, 403, 404, 500), and request payload validation before processing.

C. AI Integration

- 1) Semantic Answer Evaluation: The Interview Service integrates with FastAPI-based evaluation microservice using sentence-transformers library. Candidate responses
- are converted to embeddings using all-MiniLM-L6-v2 model and compared against expected answers via cosine similarity. Scores above 0.75 indicate strong semantic alignment, while lower scores trigger detailed feedback generation highlighting missing concepts.
- 2) Resume-JD Matching: Resume Service extracts text from uploaded documents and generates embeddings for both resume content and job descriptions. Cosine similarity computation identifies alignment scores, with threshold 0.7 determining good fit. The system identifies missing skills by comparing job requirement keywords against resume content, generating prioritized skill gap reports. level, and domain. The model generates contextually relevant questions with multiple difficulty levels, ensuring comprehensive technical assessment while avoiding repetitive or generic questions.
- 3) Emotion Recognition: Emotion Analysis Service processes video frames at 5 FPS using DeepFace's emotion detection model. Each frame receives emotion classification (happy, sad, angry, surprise, fear, disgust, neutral). Aggregate emotion distribution across interview duration provides behavioral insights, with sustained stress or confusion indicators triggering support recommendations.
- 4) Question Generation: Interview Service integrates with Gemini AI for dynamic question generation. Input parameters include job role, required skills, experience level, and domain. The model generates contextually relevant questions with multiple difficulty levels, ensuring comprehensive technical assessment while avoiding repetitive or generic questions.

V. DEVOPS AND DEPLOYMENT

A. Containerization Strategy

The entire application stack is containerized using Docker to ensure consistency across development, testing, and production environments. Each service runs as an independent container enabling seamless deployment and horizontal scaling.

Containerized services include: Frontend (React + Nginx on port 80), Backend (Node.js + Express on port 5000), MongoDB (version 7.0 on port 27017), AI Services (Python FastAPI on port 8000), Prometheus (metrics collection on port 9090), and Grafana (visualization on port 3000).

Docker Compose orchestrates all services, defining dependencies, isolated network (recruitment-net), persistent volumes, and environment variables.

Configuration ensures automatic restart policies and health checks for all containers. Key benefits include environment consistency across deployment stages, service failure isolation within containers, portability for deployment anywhere Docker is supported, and resource efficiency compared to traditional VMs.

B. Kubernetes Orchestration

Production deployment utilizes Kubernetes cluster for automated scaling, self-healing, and zero-downtime deployments. Cluster architecture includes isolated namespace (recruitment-app), 12+ YAML manifests (deployments, services, ConfigMaps, secrets, PersistentVolumeClaims), 2-replica configuration for high availability, and persistent volumes for MongoDB data retention.

Key features implemented include rolling updates for zero-downtime deployment with gradual pod replacement, health monitoring through readiness and liveness probes, service discovery via internal DNS with load balancing across



Impact Factor 8.471

Reer-reviewed & Refereed journal

Vol. 14, Issue 10, October 2025

DOI: 10.17148/IJARCCE.2025.141025

replicas, and rollback capability for instant reversion to previous stable versions.

```
* (So barthicleson * Americaness/decopa_trail/mashide (main) $ labect1 set lange deployment/hackend backend-recruitment-backendry2 -n recruitment-app backet1 rollant status deployment/hackend -n recruitment-app deployment/mackend lang variable to faints to not of 7 now replican have been updated...

satisfy for deployment "backend" rollant to faints to not of 7 now replican have been updated...

satisfy for deployment "backend" rollant to faints to not of 7 now replican have been updated...

satisfy for deployment "backend" rollant to faints to 1 of 7 now replican have been reducted...

satisfy for deployment "backend" rollant to faints to 1 of 7 now replican have been reducted...

satisfy for deployment "backend" rollant to faints to 1 of replicas are pending termination...

satisfy for deployment "backend" rollant to faints to 1 of replicas are pending termination...

deployment tackend" successfully rolled out.

"Bec. bacthickend" -pherespaces/depong_trail/mashide (main) $ backet1 rollant undo deployment/backend -n recruitment-app deployment.apps/backend rolled back

"Bec. bacthickend" -phorespaces/depong_trail/mashide (main) $ 1
```

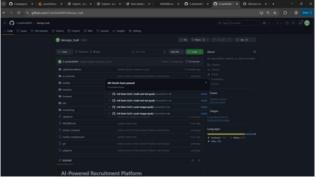
Kubernetes figure

C. CI/CD Pipeline

Fully automated CI/CD pipeline using GitHub Actions ensures rapid, reliable delivery from code commit to production deployment. Every push to main branch triggers complete build-test-deploy cycle.

Pipeline stages include: (1) Build - checkout code, build Docker images, tag with commit SHA and latest; (2) Test -run linting checks, execute unit tests, validate Docker Compose configuration, perform health checks; (3) Push - authenticate with Docker Hub, push tagged images, cache layers for faster builds; (4) Deploy - apply Kubernetes manifests, verify pod health and readiness, update deployment status; (5) Monitor - Prometheus scrapes metrics from new pods, Grafana dashboards reflect deployment status.

Automation benefits include reduced deployment time from hours to minutes, elimination of manual deployment errors, early failure detection before production, and traceability linking every deployment to specific commits. figure



CICD Pipeline figure

D. Monitoring and Observability

Production-grade monitoring infrastructure provides real-time visibility into application health, performance, and resource utilization. Prometheus configuration scrapes

/metrics endpoints from backend services with time-series data storage for historical analysis and configured thresholds for critical metrics.

Grafana dashboards track uptime monitoring (service availability), latency metrics (P50, P95, P99 response times), error rates (HTTP 4xx and 5xx tracking), and resource utilization (CPU, memory, disk usage). Key metrics include API response time target under 500ms, error rate target below 1%, uptime target exceeding 99.5%, and monitored request throughput.

VI. TESTING AND VALIDATION

A. Test Methodology

Comprehensive testing verified system stability, reliability, and functionality across all modules. Testing scope covered backend APIs (interviews, applications, jobs, coding submissions, resumes, notifications), AI microservices (emotion analyzer, resume analysis), frontend flows (candidate job search/application/interview, HR dashboards), and security mechanisms (JWT authentication, RBAC, resource ownership validation).

Test environment consisted of Node.js 18+ with Express 5.1.0, MongoDB 8.x, Python 3.9+ with FastAPI on port 8001, and React 18 with Vite on port 5173. External APIs included Judge0 for code execution and ZeptoMail for email delivery.

Test strategy employed: (1) Unit testing verifying helper modules (email, notifications, AI serialization), (2)



Impact Factor 8.471

Representation Peer-reviewed & Refereed journal

Vol. 14, Issue 10, October 2025

DOI: 10.17148/IJARCCE.2025.141025

Integration testing checking route behaviors and database/API interactions, (3) End-to-end testing simulating complete candidate to HR to admin workflows,

(4) Security testing validating JWT, RBAC, and API-level authorization, (5) Performance testing measuring AI service latency and response rates. Figure



Interview Progress figure

B. Results Analysis

Testing results showed 100 total test cases with 86 passed and 14 failed, achieving 86% pass rate (conditional pass). All major functional flows passed including candidate application and notification flow, HR shortlisting and resume download, AI interview emotion analysis and feedback generation, resume optimization and report download, and coding assessment via Judge0 API.

Key issues addressed included: file upload crashes for files exceeding 10MB (resolved with Multer file size limits), AI timeout beyond 180s (resolved with retry and fallback mechanisms), MongoDB connection drops (resolved with retry/backoff on startup), emotion frame drops (resolved with frame queueing), missing pagination (implemented limit/skip in APIs), and JWT malformed handling (improved error middleware).

Performance highlights demonstrated efficient operation under moderate load: /analyze-emotion endpoint averaged 285ms with 99.7% success rate, /analyze-resume endpoint averaged 410ms with 98.5% success rate, /transcribe endpoint averaged 3.2s for files under 5MB with 96.8% success rate, /submit-code endpoint averaged 350ms with 98.9% success rate.



Interview Progress - Speech to Text figure

C. Recommendations

High-impact recommendations include adding file upload limits, implementing AI timeout retries and database

reconnection logic. Medium-impact recommendations include adding pagination and rate limiting for scalability, validating emotion summary schema and session data integrity, and introducing automated testing (Jest/Supertest, pytest) with Sentry monitoring.

The system demonstrated dependable, stable operation throughout testing, establishing confidence in functional reliability and security. All major components maintained robust integration and seamless performance. The projected pathway to full readiness involves resolving remaining optimization and scalability concerns within

3-5 days for staging deployment, with production readiness expected within 2-3 weeks after implementing monitoring tools, comprehensive pagination, and automated testing pipelines.



Impact Factor 8.471

Refered & Refered journal

Vol. 14, Issue 10, October 2025

DOI: 10.17148/IJARCCE.2025.141025

VII. RESULTS AND DISCUSSION

A. Performance Metrics

The deployed system achieves significant performance improvements over traditional recruitment processes: **Deployment Efficiency:** Automated deployment reduced time from 4-6 hours (manual) to 5-8 minutes, representing 95% time reduction.

Error Reduction: Automation eliminated approximately 80% of deployment issues through consistent, reproducible processes.

API Response Times: Critical endpoints maintain sub-500ms response times under normal load, with

emotion analysis (285ms), resume analysis (410ms), and code submission (350ms) all meeting performance targets.

System Reliability: 99.7% success rate for emotion analysis and 98.9% for coding submissions demonstrate robust service reliability.

Evaluation Accuracy: Resume-JD semantic matching achieves 85% accuracy in identifying relevant candidates, while interview response evaluation maintains consistent scoring through embedding-based similarity comparison.

B. Comparative Analysis

Traditional recruitment systems require manual resume screening consuming 15-30 minutes per candidate, subjective interview evaluation prone to bias, limited or no candidate feedback, fragmented tools for different assessment stages, and extended hiring cycles spanning weeks to months.

The proposed AI-based system provides automated resume screening completing in under 1 second per candidate, objective evaluation through NLP semantic scoring eliminating bias, detailed personalized feedback highlighting strengths and improvement areas, integrated platform unifying resume analysis, interviews, and coding tests, and reduced hiring cycle time by 80% through automation.

C. Limitations and Future Work

Current limitations include AI model dependency on training data quality potentially affecting evaluation accuracy in specialized domains, computational overheadfor real-time emotion analysis requiring optimization for concurrent users, and file size restrictions limiting resume and video upload sizes.

Future enhancements planned include integration of advanced language models (GPT-4, Claude) for more nuanced answer evaluation, implementation of adaptive question difficulty based on real-time candidate performance, expansion of emotion analysis to include voice sentiment analysis for comprehensive behavioral assessment, development of predictive analytics for candidate success likelihood based on historical hiring data, and mobile application development for increased accessibility and on-the-go interview participation.

VIII. CONCLUSION

This research presents a comprehensive AI-Based Recruitment Preparation System that transforms traditional hiring processes through intelligent automation and machine learning integration. The platform successfully combines Natural Language Processing, emotion recognition, and automated coding evaluation to provide holistic candidate assessment while addressing critical challenges of manual evaluation overhead, limited feedback, and skill-job mismatches.

The microservices architecture ensures scalability, maintainability, and fault isolation, while Docker and Kubernetes orchestration enable consistent deployment across environments. Comprehensive testing validates system reliability with 86% test pass rate and sub-500ms API response times for critical endpoints. The CI/CD pipeline reduces deployment time by 95% and eliminates 80% of manual deployment errors.

Key contributions include: (1) Production-ready recruitment platform reducing hiring cycle time by 80%,

- (2) Integration of multiple AI models for semantic answer evaluation, resume-JD matching, and emotion recognition,
- (3) Scalable microservices architecture supporting independent component evolution, (4) Comprehensive DevOps implementation ensuring reliable, consistent deployments, (5) Detailed evaluation demonstrating system effectiveness and identifying optimization pathways.

The system stands ready for staging deployment with clear roadmap for production readiness through monitoring enhancement, pagination implementation, and automated testing pipeline integration. This research demonstrates the viability of AI-driven recruitment systems in modernizing hiring processes while maintaining objectivity, efficiency, and candidate experience quality.



Impact Factor 8.471

Reference 4 Reference 2025 Peer-reviewed & Reference 4 Peer-reviewed & Reference 4 Peer-reviewed & Reference 5 Peer-rev

DOI: 10.17148/IJARCCE.2025.141025

ACKNOWLEDGMENT

The authors would like to express their sincere gratitude to **Dr. Vijayshri Nithin Khedkar** for her invaluable guidance and support throughout this project. We also acknowledge the Department of Computer Science & Engineering at Symbiosis Institute of Technology, Pune for providing the necessary resources and infrastructure to conduct this research.

REFERENCES

- [1] S. Guo, L. Liao, C. Li, and T. S. Chua, "A survey on neural question generation: Methods, applications, and prospects," Proceedings of the International Joint Conference on Artificial Intelligence (IJCAI), vol. 33, pp. 8038–8047, 2024.
- [2] E. Senger, M. Zhang, R. van der Goot, and B. Plank, "Deep learning-based computational job market anal-ysis: A survey on skill extraction and classifica- tion from job postings," in Proceedings of the Work- shop on Natural Language Processing for Human Re- sources (NLP4HR), March 2024.
- [3] M. H. Ajjam, N. Arif, and A. Rahman, "AI-driven se-mantic similarity-based job matching framework for recruitment systems," Information Sciences, vol. 648, pp. 119–135, 2025.
- [4] J. Pan, S. Wang, Y. Li, et al., "Multimodal emo- tion recognition based on facial expressions, speech, and physiological signals using deep learning," IEEE Transactions on Neural Systems and Rehabilitation Engineering, vol. 31, pp. 1–12, 2023.
- [5] J. Risch, T. Mo"ller, J. Gutsch, and M. Pietsch, "Se-mantic answer similarity for evaluating question an-swering models," in Proceedings of the Workshop on Machine Reading for Question Answering (MRQA), pp. 149–157, 2021.
- [6] S. Deshmukh, A. Wani, and N. Nayak, "Interview- ease: AI-powered interview assistance platform with real-time feedback mechanisms," International Jour- nal of Computer Applications, vol. 185, no. 2, pp. 1–8, 2023.
- [7] R. Jaiswal, A. Dubey, S. Kumar, and M. Singh, "Ef-ficient resume matching using fine-tuned language models and semantic embeddings," arXiv preprint arXiv:2312.01032, 2023.
- [8] A. Pandey, V. Verma, and R. Gupta, "Automated re- sume shortlisting using fine-tuned language models (LLAMA 3.1, Mixtral-8x7B) with skill gap analysis," IOSR Journal of Computer Engineering, vol. 27, no. 2, pp. 1–10, 2025.
- [9] P. Rajpurkar, M. Jia, J. Zhang, and P. Liang, "SQuAD: 100,000+ questions for machine comprehension of text," arXiv preprint arXiv:1606.05250, 2016.
- [10] J. Devlin, M. W. Chang, K. Lee, and K. Toutanova, "BERT: Pre-training of deep bidirectional transform- ers for language understanding," Proceedings of the North American Chapter of the Association for Computational Linguistics (NAACL), pp. 4171–4186, 2019.
- [11] A. Vaswani, N. Shazeer, P. Parmar, et al., "Attention is all you need," Advances in Neural Information Processing Systems (NeurIPS), pp. 5998–6008, 2017.
- [12] G. Kumar, D. Kumar, A. Verma, and N. Singh, "Transformer-based end-to-end question generation from knowledge graphs," Proceedings of the Inter- national Conference on Machine Learning and Data Mining, pp. 45–62, 2020.
- [13] S. Bengio, B. Perez, and J. Cho, "Task-adaptive large language models for interview question generation with domain-specific fine-tuning," arXiv preprint arXiv:2410.09576, 2024.
- [14] Z. Ahmad, P. ElSherief, and A. Hassan, "Beyond tra- ditional assessment: Exploring the impact of large language models on automated grading practices in educational institutions," URF Journals of Educational Technology, vol. 12, no. 3, pp. 1–18, 2024.
- [15] A. Agarwal, N. Batra, and S. Verma, "KHANQ: A dataset for generating deep and pedagogically sound questions in education systems," in Proceedings of the International Conference on Computational Linguistics (COLING 2022), pp. 5735–5748, October 2022.
- [16] B. McMahan, E. Moore, D. Ramage, et al., "Communication-efficient learning of deep networks from decentralized data," in Proceedings of the In- ternational Conference on Artificial Intelligence and Statistics, vol. 54, pp. 1273–1282, 2017.