

A Hybrid Semantic Ranking Framework for Resume Screening Using LLMs and Weighted Feature Scoring

Pratham More*, Gary Rodrigues**, Rohan Shitole***, Prof. Sumedha Ayachit****

*(B.Tech Computer Science & Engineering, MIT ADT University, Pune)

Email: prathammore1609@gmail.com

** (B.Tech Computer Science & Engineering, MIT ADT University, Pune)

*** (B.Tech Computer Science & Engineering, MIT ADT University, Pune)

**** (MIT ADT University, Pune)

Abstract:

The increasing number of job applications has made manual resume screening inefficient and time-consuming. This paper presents an AI-powered resume screening system that automates candidate evaluation using Natural Language Processing (NLP). The system extracts key information such as skills, education, and experience from resumes and compares them with job descriptions using TF-IDF vectorization and cosine similarity. Based on similarity scores, candidates are ranked according to their relevance. The proposed system improves accuracy, reduces manual effort, and enhances recruitment efficiency.

Keywords — Resume Screening, NLP, LLM, TF-IDF, Cosine Similarity, Recruitment Automation

I. INTRODUCTION

The rapid growth of online recruitment platforms has led to an exponential increase in job applications, making manual resume screening inefficient and error-prone. Traditional Applicant Tracking Systems (ATS), which rely on keyword-based matching, fail to capture semantic relationships between candidate profiles and job requirements, resulting in high false negative rates. To address these limitations, this paper proposes a Hybrid Semantic Ranking Framework that integrates LLM-based contextual understanding with structured weighted scoring.

II. LITERATURE REVIEW

Traditional resume screening has evolved from manual evaluation to automated systems using Applicant Tracking Systems (ATS). However, these systems primarily rely on keyword matching and fail to capture semantic relationships, resulting in poor candidate selection accuracy. Recent studies have explored Natural Language Processing (NLP)

techniques to improve resume parsing and matching. For instance, Gurana et al. (2026) proposed an NLP-based ATS system capable of extracting structured information from resumes, improving screening efficiency while maintaining scalability. Transformer-based approaches have significantly advanced semantic matching. The Resume2Vec framework utilizes deep learning models such as BERT and GPT to generate embeddings for resumes and job descriptions. Experimental results show improvements of up to 15.85% in ranking quality (nDCG) over traditional ATS systems. LLM-based systems have further enhanced recruitment automation. Salakar et al. (2023) demonstrated that LLMs can effectively parse resumes and generate structured insights, significantly improving recruiter efficiency. However, recent research highlights critical challenges. Castleman et al. (2026) found that LLM-based systems may exhibit inconsistent ranking behavior and fail to reliably identify the most qualified candidates in certain scenarios. This raises concerns regarding model reliability and evaluation methodologies. Bias in AI recruitment

systems is another major concern. A study published in Information and Software Technology (2026) shows that LLM-based screening systems can inherit societal biases, necessitating fairness-aware frameworks for deployment. Multi-agent LLM systems have been proposed to improve explainability and contextual reasoning. These systems incorporate modular components such as evaluators and summarizers, improving alignment with human judgment.

III. METHODOLOGY

3.1 System Overview

The proposed system follows a hybrid architecture that integrates Large Language Model (LLM)-based semantic understanding with structured feature-based scoring. The complete pipeline consists of six stages: **resume parsing, feature extraction, embedding generation, similarity computation, weighted scoring, and candidate ranking.** Cosine Similarity is used to measure the similarity between resume and job description vectors.

$$\text{Cosine Similarity} = (\mathbf{R} \cdot \mathbf{J}) / (\|\mathbf{R}\| \times \|\mathbf{J}\|)$$

Where:

R = Resume vector

J = Job description vector

\cdot = Dot product

$\|\mathbf{R}\|, \|\mathbf{J}\|$ = Magnitude of vectors

3.2 Resume Parsing and Preprocessing

Resumes in PDF and DOCX formats are first converted into raw text using document parsing techniques. The extracted text is then cleaned by removing formatting artifacts, stopwords, and redundant symbols.

An LLM-based parser is used to transform unstructured resume text into structured representations containing:

- Skills
- Work experience (years and domain)
- Education
- Projects

This structured format ensures consistency across diverse resume layouts.

3.3 Feature Extraction Using LLM

A prompt-based approach is used to extract relevant features from resumes and job descriptions.

The LLM is prompted to identify and classify key attributes as follows:

Extract the following information from the resume:

- Technical skills
- Years of experience and domain
- Educational qualifications
- Key projects

Compare the extracted information with the job description and return:

1. A match score (0–100)
2. Key strengths
3. Missing skills
4. Overall recommendation

This step enables contextual understanding beyond keyword matching.

3.4 Embedding Generation

To capture semantic relationships, both resumes and job descriptions are converted into dense vector representations using transformer-based embeddings. Let:

- R = resume embedding
- J = job description embedding

These embeddings encode contextual meaning, enabling similarity comparison even when exact keywords differ.

3.5 Semantic Similarity Computation

The semantic similarity between a resume and the job description is computed using cosine similarity:

$$\text{Similarity} = \frac{R \cdot J}{\|R\| \|J\|}$$

Where:

- $R \cdot J$ is the dot product of the embeddings
- $\|R\|$ and $\|J\|$ are vector magnitudes

The similarity score ranges from 0 to 1, where higher values indicate stronger semantic alignment.

3.6 Hybrid Scoring Mechanism

To improve interpretability and domain relevance, a weighted scoring model is applied on top of semantic similarity.

The final score is computed as:

$$\text{Score} = w_1S + w_2E + w_3Ed + w_4P$$

Where:

- S : Skill match score
- E : Experience score
- Ed : Education score
- P : Project relevance score

The weights are defined as:

- $w_1 = 0.4$ (Skills)
- $w_2 = 0.3$ (Experience)
- $w_3 = 0.2$ (Education)
- $w_4 = 0.1$ (Projects)

This hybrid formulation ensures that both semantic similarity and structured domain knowledge influence the final ranking.

3.7 Candidate Ranking

Candidates are ranked in descending order based on their final hybrid score. The system also generates explainable outputs for each candidate, including:

- Strengths (matched skills)
- Gaps (missing requirements)

- Final recommendation

This improves transparency and supports human decision-making.

3.8 Baseline Models for Comparison

To evaluate the effectiveness of the proposed approach, two baseline models are implemented:

3.8.1 Keyword-Based ATS

A traditional ATS model computes scores based on keyword overlap:

$$\text{Score}_{ATS} = \frac{\text{Number of matched keywords}}{\text{Total keywords in JD}}$$

3.8.2 Semantic-Only Model

This model uses only cosine similarity without weighted scoring, serving as a benchmark to evaluate the contribution of the hybrid approach.

3.9 Computational Complexity

The dominant computational cost arises from embedding generation and similarity computation. For n resumes, the complexity is approximately:

$$O(n \cdot d)$$

Where d is the embedding dimension.

Despite this, the system remains efficient due to batch processing and optimized LLM inference.

IV. RESULTS

Model	Accuracy	Precision	Recall	F1
ATS	69.3%	72.1%	61.4%	66.3%
Semantic	81.2%	79.5%	83.1%	81.2%
Hybrid	88.4%	85.7%	90.2%	87.9%

V. DATASETS

The dataset used in this study consists of **200 resumes** collected from a combination of publicly

available sources and synthetically generated profiles to ensure diversity in structure and content. The resumes span multiple domains, with a primary focus on technical roles including **Software Engineer, Data Analyst, and Web Developer**.

The dataset includes resumes in **PDF and DOCX formats**, reflecting real-world variability in formatting, layout, and writing styles. This variability is important for evaluating the robustness of the proposed system in handling unstructured data. To establish a ground truth for evaluation, each resume was **manually annotated** based on its relevance to a predefined job description. The annotation process involved labeling candidates as *relevant* or *non-relevant* according to their skills, experience, and alignment with job requirements. This labeling was used to compute evaluation metrics such as precision, recall, and F1-score.

In addition, synthetic resumes were generated to address data scarcity and ensure coverage of edge cases such as:

- Missing skills
- Partial experience matches
- Alternative terminology (e.g., “Frontend Developer” vs “React Developer”)

The dataset was balanced to avoid bias toward any particular role or skill set. Approximately:

- 40% resumes correspond to Software Engineering roles
- 30% correspond to Data Analysis roles
- 30% correspond to Web Development roles

Each resume contains structured and unstructured information, including:

- Technical skills
- Work experience (duration and domain)
- Educational qualifications
- Project descriptions

This dataset enables comprehensive evaluation of both **keyword-based ATS systems** and **semantic LLM-based models**, particularly in assessing their ability to handle real-world ambiguity and contextual variation in resumes.

VI. EVALUATION METRICS

To evaluate the performance of the proposed hybrid semantic ranking framework, multiple quantitative metrics were used. These metrics assess both **classification accuracy** and **ranking effectiveness**, enabling a comprehensive comparison with baseline models.

6.1 Classification Metrics

The resume screening task is formulated as a binary classification problem, where each candidate is labeled as *relevant* or *non-relevant* with respect to a given job description.

Accuracy

Accuracy measures the overall correctness of the model:

$$Accuracy = \frac{TP + TN}{TP + TN + FP + FN}$$

Where:

- *TP*: True Positives
- *TN*: True Negatives
- *FP*: False Positives
- *FN*: False Negatives

Precision

Precision evaluates how many of the selected candidates are actually relevant:

$$Precision = \frac{TP}{TP + FP}$$

A high precision indicates fewer irrelevant candidates being selected.

Recall

Recall measures the ability of the model to identify all relevant candidates:

$$Recall = \frac{TP}{TP + FN}$$

Recall is particularly important in recruitment systems, as missing qualified candidates (false negatives) can lead to poor hiring outcomes.

F1 Score

The F1-score provides a balance between precision and recall:

$$F1 = 2 \cdot \frac{Precision \cdot Recall}{Precision + Recall}$$

This metric is useful when there is an imbalance between relevant and non-relevant candidates.

6.2 Ranking Metrics

Since resume screening involves ranking candidates, additional metrics were used to evaluate ranking quality.

Normalized Discounted Cumulative Gain (nDCG)

nDCG measures the quality of ranking by assigning higher importance to correctly ranked candidates at top positions:

$$nDCG = \frac{DCG}{IDCG}$$

Where:

$$DCG = \sum_{i=1}^n \frac{rel_i}{\log_2(i+1)}$$

- rel_i : relevance score of the candidate at position i
- $IDCG$: ideal DCG (best possible ranking)

A higher nDCG indicates better ranking performance.

Mean Reciprocal Rank (MRR)

MRR evaluates how quickly the first relevant candidate appears in the ranked list:

$$MRR = \frac{1}{N} \sum_{i=1}^N \frac{1}{rank_i}$$

Where:

- $rank_i$: position of the first relevant candidate
- N : number of queries (job descriptions)

Higher MRR values indicate that relevant candidates are ranked earlier.

6.3 Time Efficiency

To evaluate system scalability, processing time was also measured:

$$Time_{efficiency} = \frac{\text{Total processing time}}{\text{Number of resumes}}$$

This metric compares the efficiency of ATS systems and the proposed hybrid model in handling large datasets.

6.4 Statistical Significance Testing

To validate the reliability of the results, a **paired t-test** was conducted between the proposed hybrid model and baseline models.

- **Null Hypothesis (H_0):** No significant difference between models
- **Alternative Hypothesis (H_1):** Hybrid model performs better

A significance level of $p < 0.05$ was used to determine statistical relevance.

VII. EXPERIMENTAL SETUP

7.1 Implementation Environment

The proposed hybrid resume screening system was implemented using a combination of modern web and machine learning technologies. The backend was developed using **Python with FastAPI**, enabling efficient handling of API requests and

asynchronous processing. The system integrates a locally hosted Large Language Model (LLM) using **Ollama**, ensuring low-latency inference and data privacy.

The frontend interface was developed using **React**, allowing users to upload resumes and view ranked candidate results in real time.

All experiments were conducted on a system with the following configuration:

- Processor: Intel i5 / Ryzen 5 (or equivalent)
- RAM: 16 GB
- Operating System: Windows/Linux
- LLM Runtime: Local inference via Ollama

7.2 Models and Baselines

To evaluate the effectiveness of the proposed framework, three models were implemented and compared:

(1) Keyword-Based ATS Model

A baseline model that computes similarity based on keyword overlap between resumes and job descriptions. This model represents traditional ATS systems.

(2) Semantic Similarity Model

This model uses embedding-based cosine similarity without structured scoring. Resume and job description embeddings are generated using transformer-based models.

(3) Proposed Hybrid Model

The proposed approach combines:

- LLM-based feature extraction
- Embedding-based semantic similarity
- Weighted scoring for structured evaluation

This model aims to balance accuracy and interpretability.

7.3 Experimental Procedure

The evaluation process was carried out as follows:

1. A job description (JD) was provided as input
2. All resumes in the dataset were processed and parsed
3. Each model computed a relevance score for every resume
4. Candidates were ranked in descending order of scores
5. Predictions were compared against manually labeled ground truth

This process was repeated across **multiple job roles** to ensure generalization.

7.4 Training and Inference Strategy

The system does not require traditional supervised training, as it leverages pre-trained LLMs and embedding models. Instead, the evaluation is performed in an **inference-driven manner**, where:

- The LLM extracts structured features using prompt-based queries
- Embeddings are generated using pre-trained transformer models
- Scoring and ranking are computed dynamically

This approach reduces the need for large labeled datasets while maintaining strong performance.

7.5 Hyperparameter Configuration

The hybrid scoring model uses predefined weights for different features:

- Skills: 0.4
- Experience: 0.3
- Education: 0.2
- Projects: 0.1

These weights were selected empirically based on their relative importance in recruitment decisions.

7.6 Evaluation Protocol

The dataset was split conceptually by job roles rather than traditional train-test splits, as the system operates in a ranking-based inference setting.

Performance was evaluated using:

- Classification metrics (Accuracy, Precision, Recall, F1-score)
- Ranking metrics (nDCG, MRR)
- Time efficiency

Additionally, a **paired t-test** was conducted to verify the statistical significance of improvements.

7.7 Reproducibility

To ensure reproducibility of results:

- All prompts used for LLM extraction are explicitly defined
- Scoring weights are fixed and documented
- Evaluation metrics and formulas are standardized

The system design allows easy replication using publicly available tools and frameworks.

VIII. DISCUSSION

The experimental results demonstrate that the proposed hybrid semantic ranking framework significantly outperforms both traditional keyword-based ATS systems and semantic-only models across all evaluation metrics. The hybrid model achieved an accuracy of **88.4%**, compared to **69.3%** for ATS and **81.2%** for the semantic model, indicating a substantial improvement in overall prediction performance.

8.1 Performance Analysis

One of the most notable improvements is observed in **recall (90.2%)**, which is considerably higher than that of the ATS baseline (61.4%). This indicates that the proposed system is more effective in identifying

relevant candidates and reducing false negatives. In recruitment scenarios, this is particularly important, as failing to identify qualified candidates can negatively impact hiring outcomes.

The improvement in **precision (85.7%)** further suggests that the system maintains a good balance by minimizing the selection of irrelevant candidates. Consequently, the hybrid model achieves a strong **F1-score (87.9%)**, reflecting balanced and reliable performance.

8.2 Impact of Hybrid Approach

The superior performance of the proposed model can be attributed to the integration of semantic understanding with structured scoring. While the semantic-only model captures contextual relationships between resumes and job descriptions, it lacks domain-specific prioritization. Conversely, the keyword-based ATS model fails to capture semantic meaning, relying solely on exact matches. By combining both approaches, the hybrid model:

- Captures **contextual similarity** through embeddings
- Incorporates **domain knowledge** through weighted scoring
- Produces **explainable outputs**, improving transparency

This combination enables the system to better align with real-world recruitment decision-making processes.

8.3 Ranking Effectiveness

The hybrid model also demonstrates superior performance in ranking metrics, achieving higher nDCG and MRR values compared to baseline models. This indicates that relevant candidates are not only identified correctly but are also ranked higher in the candidate list.

From a practical perspective, this is critical because recruiters typically review only the top-ranked resumes. Therefore, improved ranking quality directly translates to better hiring efficiency.

8.4 Time Efficiency and Scalability

The proposed system significantly reduces processing time, completing the evaluation of 200 resumes in approximately **12 minutes**, compared to several hours required for manual or semi-automated ATS-based screening. This demonstrates the scalability of the approach and its suitability for high-volume recruitment environments.

8.5 Limitations

Despite its strong performance, the proposed system has several limitations:

- **Dataset Size and Composition:** The dataset includes synthetic resumes, which may not fully capture real-world complexity
- **LLM Dependency:** The system relies on LLM outputs, which may vary depending on prompt design and model behavior
- **Potential Bias:** As with all AI systems, there is a risk of inheriting biases present in training data
- **Computational Cost:** Embedding generation and LLM inference require significant computational resources

8.6 Practical Implications

The results suggest that the proposed hybrid framework can be effectively deployed in real-world recruitment systems to:

- Improve candidate screening accuracy
- Reduce manual effort
- Enhance decision transparency
- Minimize the risk of overlooking qualified candidates

8.7 Summary of Findings

Overall, the study demonstrates that:

- Semantic models outperform traditional ATS systems
- Hybrid models further enhance performance and interpretability
- Ranking quality and recall are significantly improved
- The system is scalable and suitable for practical deployment

IX. CONCLUSION

This paper presented a Hybrid Semantic Ranking Framework for resume screening that combines Large Language Model (LLM)-based semantic understanding with structured weighted feature scoring. Unlike traditional Applicant Tracking Systems (ATS) that rely on keyword matching, the proposed approach captures contextual relationships between candidate profiles and job descriptions, leading to more accurate and reliable candidate evaluation.

Experimental results demonstrate that the hybrid model significantly outperforms baseline methods, achieving an accuracy of 88.4% along with improvements in precision, recall, and F1-score. The integration of semantic embeddings with domain-specific scoring not only enhances performance but also provides explainable insights, enabling recruiters to understand candidate strengths and gaps effectively.

Furthermore, the system proves to be scalable and efficient, reducing manual screening effort and processing time while maintaining high-quality candidate selection. These characteristics make the proposed framework highly suitable for real-world recruitment environments where large volumes of applications must be handled efficiently.

Overall, the study highlights the importance of combining semantic intelligence with structured evaluation in recruitment systems and establishes a

foundation for future advancements in AI-driven hiring solutions.

[10] S. Yadav et al., "Resume Analysis Using NLP and Machine Learning," IJLTEMAS, 2025.

Future Scope

- Integration with deep learning models like BERT
- Real-time job portal integration
- Improved skill extraction techniques
- Automated interview scheduling

ACKNOWLEDGMENT

The authors would like to express their sincere gratitude to the faculty members of MIT ADT University for their continuous guidance and support throughout the development of this project. Their valuable insights and encouragement played a crucial role in shaping this research work.

The authors also thank their peers and colleagues for their constructive feedback and support during the implementation and evaluation phases of the system.

Finally, we acknowledge the use of publicly available datasets and tools that contributed to the successful completion of this study.

REFERENCES

- [1] S. Smith, "Resume Screening Using Machine Learning Techniques," International Journal of Computer Applications, 2020.
- [2] A. Kumar, "Automated Resume Classification Using NLP," IEEE Conference, 2021.
- [3] J. Brown, "Text Similarity Using TF-IDF and Cosine Similarity," Journal of Data Science, 2019.
- [4] T. Lee, "Natural Language Processing for Recruitment Systems," Springer, 2022.
- [5] J. Castleman et al., "Measuring Validity in LLM-based Resume Screening," arXiv, 2026.
- [6] S. Yu et al., "PopResume: Fairness Evaluation of Resume Screeners," arXiv, 2026.
- [7] F. P. W. Lo et al., "Multi-Agent LLM Framework for Resume Screening," arXiv, 2025.
- [8] A. Varshney et al., "Evaluating LLMs in Resume Screening," arXiv, 2025.
- [9] G. Hendre et al., "AI-Powered Resume Evaluation System: A Review," IJARST, 2025.