



IJEAST

INTERNATIONAL JOURNAL
OF ENGINEERING APPLIED SCIENCE
AND TECHNOLOGY



VOLUME : 10 ISSUE : 01 Print / Issue Publication Date: 30-Jun-2025



ISSN : 2455-2143



DOI : 10.33564/IJEAST.2025.v10i01.014

Indexed In



WWW.IJEAST.COM

editor@ijeast.com



ADVANCED RESUME PARSER USING NLP

Vedansh Vinod Singh, Yukti Gupta
Department of CSE (AI & ML)
Inderprastha Engineering College, U.P

Abstract— This project focuses on developing an advanced resume parser using Natural Language Processing (NLP) techniques to automate the extraction of structured data from unstructured resumes. By applying methods such as Named Entity Recognition and keyword extraction, the system accurately identifies key sections like personal details, skills, education, and experience. The goal is to streamline recruitment processes and enhance hiring efficiency through intelligent data processing.

Keywords-- Resume Parsing, Natural Language Processing (NLP), Named Entity Recognition (NER), Information Extraction, Recruitment Automation, Structured Data

I. INTRODUCTION

The recruitment landscape has undergone a significant transformation with the advent of digital technologies and data-driven decision-making [1]. Organizations now receive hundreds or even thousands of resumes for a single job posting, making manual screening inefficient, inconsistent, and prone to bias [2]. As a result, automating the resume screening process has become a critical need for modern human resource (HR) departments.

A resume parser is a tool that converts unstructured resume data into a structured format, enabling easier and faster candidate evaluation [3]. Traditional parsers often rely on keyword matching or rule-based techniques, which lack the flexibility and accuracy required to handle diverse resume formats and writing styles [4]. To address these limitations, this research explores the use of Natural Language Processing (NLP) techniques to build a more advanced and intelligent resume parsing system [5].

NLP enables machines to understand and interpret human language [6]. By leveraging techniques such as Named Entity Recognition (NER), part-of-speech tagging, and dependency parsing, the proposed system aims to extract relevant information such as personal details, educational background, work experience, skills, and certifications with high precision [7]. The integration of machine learning models further enhances the parser's ability to generalize across varied resume structures and terminology [8].

This project aims to design, develop, and evaluate an NLP-based resume parser that improves the efficiency and effectiveness of recruitment processes. By automating the extraction of meaningful data from resumes, the system can

significantly reduce the time and effort required by HR professionals, enabling them to focus more on strategic decision-making [9].

II. LITERATURE REVIEW

The need for automated resume processing has gained momentum due to the exponential increase in digital job applications and the limitations of manual screening [10]. Traditional applicant tracking systems (ATS) have long relied on rule-based or keyword-matching methods for resume parsing. However, these systems often fail to accurately interpret context, leading to high false positive or negative rates in candidate filtering [11].

Recent advances in Natural Language Processing (NLP) have significantly improved the ability to extract structured data from unstructured text, making NLP a core technology in modern resume parsers [12]. Named Entity Recognition (NER), a key NLP technique, has been widely adopted to identify and classify essential components of resumes, such as names, locations, organizations, skills, and dates [13]. Researchers have demonstrated that combining NER with part-of-speech tagging and dependency parsing enhances information extraction accuracy, particularly in varied and semi-structured document formats [14].

Several studies have proposed machine learning-based resume parsers, which train models on annotated datasets to identify resume sections and extract relevant details [15]. Supervised learning models such as Support Vector Machines (SVM), Conditional Random Fields (CRF), and, more recently, deep learning architectures like BiLSTM and BERT have shown promising results in segmenting and understanding resume content [16]. These models adapt better to different writing styles, layouts, and industry-specific jargon compared to rule-based methods [17].

Another growing area of research focuses on integrating resume parsing tools with job recommendation systems. By matching extracted candidate data with job descriptions using similarity scoring or semantic analysis, systems can suggest suitable roles to candidates and rank applicants for recruiters more effectively [18]. This fusion of NLP and recommendation algorithms has proven beneficial in streamlining both hiring and application processes.

Despite these advancements, challenges remain in handling multilingual resumes, varying document formats (PDF, DOCX, scanned images), and distinguishing between similar entities (e.g., distinguishing between a university and a company with similar names) [19]. Addressing these



issues requires a combination of advanced NLP techniques, data normalization strategies, and large, diverse training datasets.

S.NO.	Authors/Year	Methodology	Key Findings	Relevance to Current Research
1.	Smith et al. (2020)	Review of ATS systems	Traditional ATS systems rely on keyword matching, leading to low accuracy and high bias.	Highlights limitations of manual resume screening and the need for automation.
2.	Johnson & Lee (2019)	Comparative study of resume parsers	Rule-based parsers have limited scalability and struggle with context interpretation.	Shows need for more advanced parsing systems using NLP.
3.	Patel et al. (2021)	NLP-based resume parser using NER	NER improves the extraction of candidate information (e.g., skills, experience).	Supports the use of NER for structured data extraction in resume parsing.
4.	Kumar & Singh (2018)	NLP and machine learning combination	Combining NER with POS tagging improves data extraction accuracy in semi-structured resumes.	Supports the use of NLP techniques in enhancing resume parsing performance.
5.	Zhang et al. (2020)	Machine learning with CRF models	CRFs outperform rule-based models in resume segmentation.	Relevant for integrating machine learning into the resume parser.
6.	Wu et al. (2022)	BiLSTM and BERT for text classification	Deep learning models such as BiLSTM and BERT offer better generalization across varying resume formats.	Directly supports the use of deep learning models in resume parsing.
7.	Nguyen & Tran (2021)	Deep learning for resume section extraction	BiLSTM and BERT provide superior performance in extracting detailed resume sections.	Demonstrates the effectiveness of deep learning for detailed data extraction.
8.	Gupta et al. (2019)	Study of machine learning models in recruitment	SVM, CRF, and BiLSTM models significantly enhance the recruitment process by improving data extraction.	Justifies the use of machine learning techniques to improve recruitment efficiency.
9.	Fernandez & Gupta (2020)	Job recommendation system integration	Resume parsers combined with recommendation systems improve the ranking and matching of candidates to job roles.	Supports the integration of resume parsing with job recommendation systems.



10.	Lee et al. (2019)	Multilingual and multi-format resume parsing	Challenges exist in parsing multilingual and non-standard resume formats (e.g., PDFs, scanned images).	Highlights the challenge of handling diverse resume formats in real-world applications.
-----	-------------------	--	--	---

Table 1: Literature Review

III. METHODOLOGY

This research project was conducted through a structured methodology, encompassing iterative development phases and a user-centered design approach. The following outlines the key stages involved in the design and implementation of the advanced resume parsing system.

1. Project Phases

1.1 Planning and Requirement Analysis

The project began with a comprehensive planning phase, which involved identifying core objectives and collecting system requirements through stakeholder engagement. Feedback from prospective users, academic mentors, and industry experts was synthesized to define functional and non-functional needs. The output of this phase included a detailed project roadmap, feature list, and a high-level system blueprint.

1.2 System Design

In this phase, wireframes and low-fidelity prototypes were developed to visualize the system's interface. The architecture was designed with modularity in mind, separating front-end, back-end, and data processing components. Key considerations included user experience (UX), scalability, and integration with external services (e.g., email APIs, browser extensions). Architectural models and database schema were also finalized during this phase.

1.3 Development

The system was built using a full-stack development approach. The front-end was developed using React.js to ensure a dynamic and responsive user experience. The back-end was implemented using Spring Boot and Java, providing RESTful API services for core functionalities such as resume parsing, user management, and communication tracking. A MySQL database was used for structured data storage. The parsing logic used Python-based NLP libraries (e.g., spaCy or NLTK) integrated through microservices.

1.4 Testing

A robust testing strategy ensured system reliability and correctness. The testing process included:

- Unit Testing: Performed on individual components using JUnit and automated tools.

- Integration Testing: Validated interactions between system modules.
- Usability Testing: Feedback from a small group of users helped improve the user interface and workflows.
- Performance Testing: Conducted using JMeter to evaluate load handling capacity.
- Security Testing: Ensured data privacy and compliance with regulations like GDPR using penetration testing and data protection audits.

1.5 Deployment

The final system was deployed on a cloud platform (AWS) to ensure scalability and public accessibility. Docker containers and Kubernetes were used for managing deployment and orchestration. Additionally, user documentation and onboarding materials were created to support adoption and training.

2. Key Methodological Approaches

2.1 Agile Development

The project followed Agile principles using Scrum methodology, allowing iterative development in sprints. This approach facilitated continuous integration of feedback, faster adaptation to changing requirements, and early detection of issues.

2.2 User-Centered Design (UCD)

A UCD approach was employed to ensure the system aligned with real user needs. Wireframes, prototypes, and iterative feedback loops with test users helped refine both UI and UX. Special attention was given to simplifying repetitive tasks and improving accessibility.

3. Tools and Technologies Used

- Frontend: React.js, HTML5, CSS3, JavaScript
- Backend: Java, Spring Boot, REST APIs
- NLP & Resume Parsing: Python (spaCy, NLTK), Named Entity Recognition (NER)
- Database: MySQL
- Testing: JUnit, Selenium (automation), JMeter (performance)
- Deployment: Docker, Kubernetes, AWS EC2, S3
- Version Control: Git, GitHub
- Browser Extension: JavaScript and Chrome API for automated form filling



4. Project Management

- **Documentation:** Continuous documentation of system architecture, APIs, and user flows.
- **Communication:** Weekly team meetings and progress tracking through project management tools.
- **Risk Management:** Identified potential risks such as data breaches, parsing inaccuracies, and integration failures. Contingency plans were established accordingly.

IV. CHALLENGES AND SOLUTIONS

1. Diverse Resume Formats

Challenge: Resumes come in a wide variety of formats—some are structured with clear headings, while others use creative layouts, multiple columns, or even images. This inconsistency makes it difficult to design a parser that performs well across all formats.

Solution: A hybrid approach combining rule-based techniques with machine learning was adopted. Regular expressions and pattern matching handle standard structures, while NLP models (e.g., spaCy or BERT) are used to identify sections and entities based on context. This combination improves robustness against layout variations.

2. Ambiguity in Text Content

Challenge: Free-form text in resumes often contains ambiguous or overlapping information. For instance, job titles and roles may not be explicitly stated, or dates may not follow a consistent format.

Solution: Context-aware NLP techniques, such as dependency parsing and named entity recognition, are employed to better understand the grammatical structure of sentences. Custom entity labels and post-processing rules help to disambiguate entities like job roles, dates, and company names.

3. Inconsistent Section Labels

Challenge: Section headers are not standardized. For example, “Professional Background,” “Work Experience,” and “Career Summary” might all refer to the same concept, making detection difficult.

Solution: A curated dictionary of synonyms and variations for common section headers was developed. During parsing, fuzzy string matching is applied to identify equivalent terms, ensuring that different headings are correctly interpreted.

4. Skill Extraction and Relevance

Challenge: Resumes often feature a mix of technical and soft skills, some of which may be embedded within job descriptions rather than listed clearly. It’s also difficult to determine whether a skill is claimed or demonstrated.

Solution: A predefined skill ontology is used alongside keyword extraction. NLP models analyze the surrounding

context to distinguish between listed skills and those inferred from responsibilities or achievements, improving precision.

5. Date Normalization and Validation

Challenge: Dates are written in various formats (e.g., “Jan 2021,” “01/2021,” “2021 – Present”), and may be incomplete or ambiguous.

Solution: A date normalization module standardizes formats using date-parsing libraries. Additional logic validates date ranges and identifies inconsistencies, such as overlapping employment periods or future dates.

6. Handling Noisy or Low-Quality Data

Challenge: Some resumes may include OCR errors, typos, or irrelevant content such as headers, footers, or disclaimers that confuse the parser.

Solution: Preprocessing routines include noise filtering, character correction, and stopword removal. For scanned documents, OCR output is cleaned using language models to restore probable word forms and correct misrecognized terms.

7. Evaluation and Ground Truth Creation

Challenge: Lack of publicly available, labeled datasets makes it hard to evaluate the accuracy and performance of the parser.

Solution: A custom dataset was created by manually annotating a diverse set of resumes. This dataset is used to train and test the models, with performance measured using metrics like precision, recall, and F1 score.

8. Multilingual Support

Challenge: Parsing resumes in multiple languages poses difficulties due to linguistic differences and lack of language-specific NLP resources.

Solution: The system was designed to be modular, allowing language-specific NLP pipelines to be integrated as needed. Open-source multilingual models and translation tools are used to handle non-English resumes effectively.

V. EVALUATION

To assess the performance and reliability of the resume parser, a structured evaluation approach was used. The evaluation focused on accuracy, consistency, and the ability to generalize across diverse resume formats. Both quantitative metrics and qualitative assessments were employed to ensure comprehensive analysis.

1. Dataset Preparation

A custom dataset of 200 resumes was compiled, representing a variety of industries, experience levels, and formatting styles (e.g., chronological, functional, creative). Each resume was manually annotated to serve as ground



truth, including labeled sections such as personal details, education, experience, and skills.

2. Evaluation Metrics

The system's outputs were compared to the ground truth using the following metrics:

- Precision: Measures how many of the extracted entities were correct.
- Recall: Measures how many relevant entities were successfully extracted.
- F1 Score: The harmonic mean of precision and recall, providing a balanced view of performance.
- Accuracy: Specifically used for structured fields like email, phone number, and dates.

3. Component-wise Performance

Each major component was evaluated individually:

- Personal Information Extraction:
 - Precision: 96%
 - Recall: 94%
 - Notes: High accuracy due to consistent formatting of names and contact details.
- Education Section Extraction:
 - Precision: 88%
 - Recall: 85%
 - Notes: Performance varied depending on how clearly institutions and degrees were stated.
- Work Experience Extraction:
 - Precision: 86%
 - Recall: 82%
 - Notes: Challenges included parsing complex job descriptions and inconsistent date formats.
- Skill Extraction:
 - Precision: 91%
 - Recall: 87%
 - Notes: Improved through the use of a predefined skill dictionary and contextual analysis.

4. Error Analysis

A manual review of incorrectly parsed resumes revealed common sources of error:

- Resumes with non-standard or graphical layouts were harder to parse accurately.
- Misclassification of educational institutions as company names.
- Incomplete sentence structures leading to ambiguous skill or title extraction.

These findings were used to refine the parsing logic and retrain the NLP models on edge cases.

5. Comparison with Baseline

The parser was benchmarked against a rule-based baseline system. Results showed that the hybrid approach (NLP + rules) significantly outperformed the baseline in terms of F1

score, especially in handling complex or less structured resumes.

Component	Baseline F1	Proposed System F1
Personal Info	82%	95%
Education	70%	86%
Experience	65%	84%
Skills	75%	89%

6. User Testing and Feedback

A small user study was conducted with HR professionals and recruiters. Participants were asked to review parsed outputs and rate their usefulness on a 5-point scale. The average satisfaction score was 4.4, indicating that the parser provided relevant and structured insights suitable for professional screening.

7. Limitations

While the parser achieved high performance in most cases, its effectiveness can still be impacted by:

- Highly creative or graphic-heavy resumes.
- Resumes in languages with limited NLP support.
- Missing or ambiguous section headers.

VI. FUTURE DIRECTIONS

While the current implementation of the resume parser delivers reliable and accurate results, there are several areas that offer potential for further enhancement and expansion. These improvements can help increase its adaptability, scalability, and intelligence in real-world applications.

1. Integration of Advanced Language Models

In future iterations, the parser can be enhanced by integrating more advanced transformer-based language models such as BERT, RoBERTa, or GPT variants. These models can significantly improve the system's understanding of context, making it better at handling ambiguous or unconventional language found in resumes.

2. Multilingual Support

Currently optimized for English, the system can be extended to support multiple languages. This will require training or fine-tuning models on multilingual corpora, enabling parsing of resumes from global candidates and increasing the system's applicability in international recruitment.

3. Real-Time Resume Parsing API

Developing a real-time REST API would allow seamless integration of the parser into applicant tracking systems (ATS), job portals, and HR software. This would enable instant resume analysis and structured data generation during job application submissions.



4. Continuous Learning with User Feedback

Incorporating a feedback loop where users (e.g., recruiters) can correct or validate extracted information would allow the parser to improve over time. A semi-supervised or active learning approach could be implemented to retrain models on this validated data, boosting long-term accuracy.

5. Improved Handling of Graphical and Visual Resumes

Resumes with heavy graphical elements, columns, or non-standard layouts still pose challenges. Future work could explore the use of computer vision techniques alongside NLP to better segment and interpret visual resume components, especially in scanned or image-based documents.

6. Skill Proficiency Detection

Currently, the parser identifies listed skills but does not estimate the candidate's proficiency level. Future enhancements could involve analyzing contextual indicators (e.g., number of years of experience, frequency of skill usage) to infer proficiency levels for more informed decision-making.

7. Bias and Fairness Audits

As with any automated system in recruitment, it's important to ensure fairness and reduce potential biases. Future development could include tools to audit and monitor the system for demographic or linguistic bias, ensuring equitable treatment of all candidates.

8. Cross-Validation with External Sources

To improve data accuracy, especially for education and work experience, future versions could include optional cross-validation with public profiles like LinkedIn (with user permission). This would help verify candidate information and reduce fraudulent claims.

9. Modular Customization for Industries

Different industries emphasize different qualifications. Future iterations could allow custom configuration based on the target sector (e.g., IT, healthcare, education), tailoring the extraction logic and scoring models to prioritize industry-specific details.

VII. APPLICATIONS

1. Automated Candidate Screening

Recruiters and hiring managers can use the parser to automatically extract key information such as skills, experience, and education from resumes. This allows for faster and more consistent initial screening of candidates, reducing manual effort and time-to-hire.

2. Integration with Applicant Tracking Systems (ATS)

The parser can be embedded into ATS platforms to process resumes uploaded by applicants. This enables structured storage of candidate profiles and supports advanced search and filtering capabilities, making the recruitment process more scalable and organized.

3. Job Matching and Recommendation Systems

By extracting skills, qualifications, and job histories, the parser helps power intelligent job matching engines. Candidates can be automatically matched to job openings based on their profile, improving relevance and personalization in job recommendations.

4. Talent Analytics and Workforce Planning

Organizations can use structured resume data to analyze workforce trends, skill gaps, and talent distribution. HR departments can generate insights from aggregated resume data to support strategic planning and workforce development initiatives.

5. Educational and Career Counseling

Career services in universities and training centers can use the parser to analyze student resumes, identify strengths and weaknesses, and recommend improvements. This helps guide students and job seekers in tailoring their resumes to industry expectations.

6. Resume Quality Scoring Tools

The parser can be integrated into tools that assess the quality of resumes by evaluating completeness, relevance of skills, and alignment with job descriptions. These tools can offer real-time suggestions for improvement, benefiting job seekers.

7. Fraud Detection and Verification

By comparing parsed data with public records or professional networks (with consent), the system can assist in identifying discrepancies or potential fraud in resumes. This is particularly useful in high-stakes hiring scenarios or regulated industries.

8. Market Research and Labor Trend Analysis

Aggregated resume data can be used by analysts and policy makers to understand labor market trends, in-demand skills, and shifts in career patterns. This information is valuable for economic planning, training program development, and industry forecasting.

9. Custom Resume Databases

Organizations that regularly receive resumes (e.g., staffing agencies, freelance platforms) can use the parser to build searchable databases of candidate profiles. These databases support fast candidate discovery and long-term relationship management.



- arXiv Preprint, arXiv:1910.03089. Retrieved from: <https://arxiv.org/abs/1910.03089>
- [17]. Jakate, M., Lavangare, S., Bhoir, N., Das, A., and Kolhe, S.R. (2023). Resume Parser Using Hybrid Approach to Enhance the Efficiency of Automated Recruitment Processes, Authorea Preprints. DOI: 10.22541/au.168170278.82268853
- [18]. Zhang, S., Wang, F., and Zhao, T. (2023). Multilingual Resume Parsing Techniques, in Proc. International Conference on Natural Language Processing and Knowledge Engineering (NLP-KE), (pp. 1–6). Retrieved from: <https://www.researchgate.net/publication/355113696>
- [19]. Kumar, A., and Mehta, P.S. (2024). NLP-Based Resume Parser Using NER, in Proc. Computational Linguistics and Intelligent Systems (CLIS), (pp. 1–6). Retrieved from: <https://www.researchgate.net/publication/355113696>
- [20]. Gupta, M.R. (2023). Machine Learning Models (SVM, CRF) for Resume Extraction, in Proc. International Conference on Machine Learning and Data Engineering (ICMLDE), (pp. 1–6). Retrieved from: <https://www.researchgate.net/publication/355113696>
- [21]. Patel, S., Mehta, K., and Shah, A. (2021). Enhancing Resume Parsing Using NLP-Based Entity Extraction, in Proc. International Conference on Advances in Computing and Data Sciences, (pp. 254–264). DOI: 10.1007/978-3-030-79150-6_23

IJEAST

INTERNATIONAL JOURNAL
OF ENGINEERING APPLIED SCIENCE
AND TECHNOLOGY

ABOUT IJEAST

International Journal of Engineering Applied Science and Technology (IJEAST) is a peer-reviewed, open access journal that publishes high-quality research papers in the field of Engineering, Applied Science and Technology.

IJEAST aims to provide a platform for researchers, academicians, and professionals to share their innovative ideas, research findings, and practical experiences with the global scientific community.

FOCUS AREAS

- Engineering
- Applied Science
- Technology
- Innovation & Development
- Interdisciplinary Studies



PEER REVIEWED

All submissions are rigorously peer reviewed to ensure quality.



OPEN ACCESS

Free and unrestricted access to research for all.



GLOBAL REACH

Connecting researchers and professionals worldwide.



TIMELY PUBLICATION

We ensure a swift and efficient publication process.



For more information, visit our website

www.ijeast.com



INTERNATIONAL JOURNAL
OF ENGINEERING APPLIED SCIENCE
AND TECHNOLOGY

✉ editor@ijeast.com

🌐 www.ijeast.com

📍 India



2455-2143