

A Deep Learning-Driven Framework for Identifying Online Recruitment Scams

Syed Ishad Hussain¹ Mr. Rizwan Uz Zaman² Dr. Ruhiat Sultana³

¹ Research Scholar, Dept. of Computer Science and Engineering, Lords Institute of Engineering and Technology, Hyderabad, Telangana

² Assistant professor, Dept. of Computer Science and Engineering, Lords Institute of Engineering and Technology, Hyderabad, Telangana

³ Associate Professor, Dept. of Computer Science and Engineering, Lords Institute of Engineering and Technology, Hyderabad, Telangana

Article History:

Received: 12-01-2025

Revised: 15-02-2025

Accepted: 01-03-2025

Abstract

With the expansion of online recruitment platforms, the hiring process has become more efficient; however, this advancement has also led to a rise in fraudulent job postings, causing significant financial harm to job seekers. To address this growing concern, the study introduces a deep learning-based framework for detecting Online Recruitment Fraud (ORF). A novel, comprehensive dataset is constructed by integrating data from Fake Job Postings, Pakistan Job Postings, and US Job Postings. The proposed method leverages state-of-the-art Natural Language Processing (NLP) models—Bidirectional Encoder Representations from Transformers (BERT) and its optimized variant, RoBERTa—to encode job descriptions into meaningful vector representations. To resolve the challenge of class imbalance in the dataset, the Synthetic Minority Over-sampling Technique with Borderline Detection (SMOBD) is employed. These processed features are then input into a two-dimensional Convolutional Neural Network (CNN2D) for classification. Experimental results demonstrate that this integrated approach achieves an impressive classification accuracy of 98.68%. The proposed model not only surpasses existing techniques but also provides a robust and scalable solution for identifying fraudulent job listings, thereby contributing significantly to the prevention of online recruitment scams.

1. INTRODUCTION

In today's digital age, the internet has significantly transformed the way individuals search for jobs and how organizations manage their hiring processes. Traditional recruitment practices have been largely replaced by electronic recruitment (E-recruitment) platforms, which offer enhanced efficiency, accessibility, and convenience. These platforms enable employers to post detailed job openings—including descriptions, qualifications, salary expectations, and benefits—while allowing job seekers to browse and apply for roles that align with their skills and interests. The shift to E-recruitment became even more pronounced during the COVID-19 pandemic, as restrictions on physical interactions and the rise of remote work forced many organizations to adopt digital hiring methods. According to the World Economic Outlook

Report, the global unemployment rate rose sharply to 13% in 2020, up from 7.3% in 2019 and 3.9% in 2018, prompting companies to rely more heavily on online recruitment to maintain hiring continuity and address the increasing number of job seekers. However, the widespread use of online recruitment platforms has also created opportunities for cybercriminals to exploit the system through fraudulent job postings that promise attractive salaries and benefits. These scams have led to financial losses, identity theft, and emotional distress for many job seekers, highlighting the urgent need for reliable mechanisms to detect and prevent online recruitment fraud. To address this challenge, the present study proposes a deep learning-based approach that integrates natural language processing (NLP) techniques for the detection of fraudulent job postings. By incorporating intelligent models into recruitment systems, this approach aims to enhance the security and credibility of E-recruitment platforms, ultimately ensuring a safer environment for both employers and job seekers.

2. RELATED WORK

The increasing reliance on online recruitment platforms has been accompanied by a notable rise in fraudulent job postings, which pose significant threats to job seekers. In an effort to mitigate these risks, researchers have explored various techniques to effectively detect such fraudulent activity. Artificial Neural Networks (ANNs) have shown promise in this domain, as demonstrated by Nasser et al. [3], who utilized ANNs to identify intricate patterns in recruitment datasets. Their model effectively captured complex relationships and demonstrated strong performance in classifying fraudulent postings. In a similar vein, Lokku [4] employed machine learning algorithms to examine the textual and structural attributes of job advertisements, using supervised learning methods to assess the authenticity of listings with encouraging results. Habiba et al. [5] conducted a comparative study of different data mining approaches for detecting fake job postings, emphasizing the importance of selecting appropriate algorithms and applying effective preprocessing techniques—especially when dealing with imbalanced data—to enhance prediction accuracy.

Further, Vidros et al. [7] introduced an automated detection framework that leveraged the behavioral characteristics of fraudulent job ads. Their use of a public dataset and machine learning models demonstrated the effectiveness of extracting and analyzing key textual features to identify deceptive content. Likewise, Dutta and Bandyopadhyay [8] explored various machine learning models for fraud detection, stressing the crucial role of feature engineering in boosting model performance. They also addressed practical challenges such as data noise and imbalance, proposing solutions to improve detection accuracy. Alghamdi and Alharby [9] advanced this line of research by developing an intelligent fraud detection system based on linguistic analysis. Their model focused on identifying subtle inconsistencies and patterns in job descriptions, and demonstrated how combining multiple feature types and algorithms could significantly enhance overall system reliability. Collectively, these studies underscore the growing sophistication in online recruitment fraud detection and highlight the potential of machine learning-based approaches. However, challenges such as dataset imbalance, real-world noise, and optimal feature selection remain areas for further improvement, providing a strong foundation for the development of more robust detection models.

3. MATERIALS AND METHODS

The proposed system aims to detect Online Recruitment Fraud (ORF) by utilizing a novel and comprehensive dataset compiled from three key sources: the Fake Job Posting dataset [16], the Pakistan Job Posting dataset [18], and the US Job Posting dataset [17]. These datasets contain job advertisements labeled as genuine or fraudulent, along with associated metadata such as job descriptions, company names, and other relevant attributes. This diverse data serves as the foundation for training and evaluating the fraud detection models. To convert textual job descriptions into a machine-readable format, the system employs state-of-the-art deep learning models—Bidirectional Encoder Representations from Transformers (BERT) [15] and its optimized variant, RoBERTa [12]. These models are utilized to extract contextualized embeddings, transforming raw textual data into dense numerical vectors that capture the semantic relationships within the job postings. These vectorized features are then used as input for downstream machine learning models. A major challenge in fraudulent job detection is the issue of class imbalance, where genuine job listings far outnumber fraudulent ones. To mitigate this, the system incorporates the Synthetic Minority Over-sampling Technique with Borderline Detection (SMOBD), a variant of the traditional SMOTE algorithm [14]. SMOBD generates synthetic samples near the decision boundary of the minority class, helping the model learn more effectively from fraudulent examples and reducing classification bias. The processed data—including the SMOBD-balanced features derived from BERT and RoBERTa embeddings—is then passed into a Two-Dimensional Convolutional Neural Network (CNN2D). CNN2D is chosen for its ability to capture local and hierarchical feature patterns, enhancing the system’s capacity to distinguish between legitimate and fraudulent postings. This combination of contextual embeddings, oversampling techniques, and deep learning classification results in a robust and scalable fraud detection system.

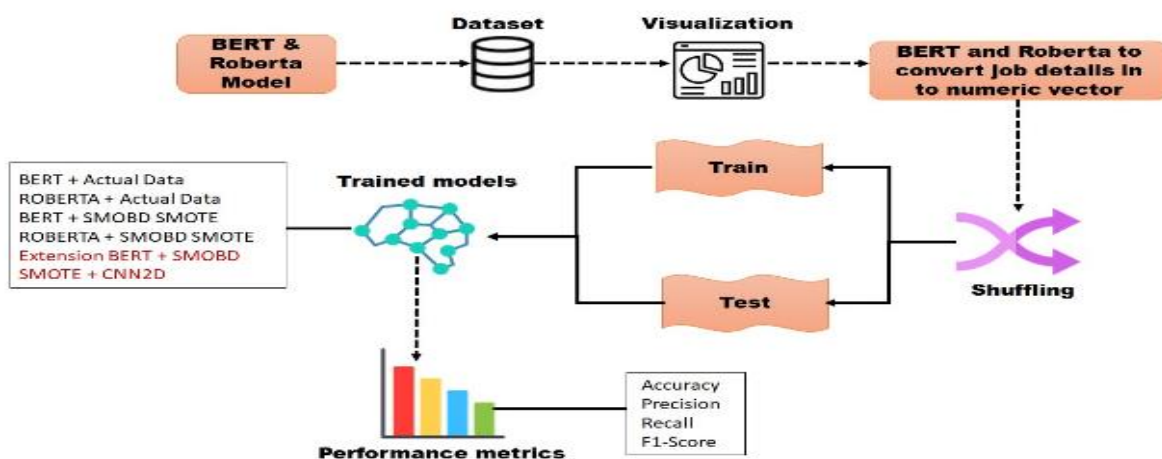


Figure 3.1: Proposed Architecture

The proposed architecture involves multiple stages: BERT and RoBERTa are first used to encode job descriptions into vector representations. These vectors are augmented through SMOBD to balance class distribution. The enhanced dataset is then used to train models such

as CNN2D, as well as other variants including fig 3.1 SMOTE-CNN2D and SMORD-based classifiers. Shuffling is applied to the data before training to improve model generalization and prevent overfitting. Model performance is evaluated using standard metrics, including accuracy, precision, recall, and F1-score.

i) Dataset Collection

The datasets used in this research comprise labeled job postings collected from reliable online recruitment sources. Each posting is classified as either legitimate or fraudulent and includes rich information such as job descriptions, company details, and other metadata necessary for feature extraction [13]. The data is stored in CSV format and processed using Python-based tools to extract relevant features for classification. This preprocessing step is essential for preparing the data to train the fraud detection models and assess their effectiveness in identifying fraudulent job postings.

4. RESULTS AND DISCUSSION

Accuracy: The accuracy of a test is its ability to differentiate the patient and healthy cases correctly. To estimate the accuracy of a test, we should calculate the proportion of true positive and true negative in all evaluated cases. Mathematically, this can be stated as:

$$\text{Accuracy} = \frac{TP + TN}{TP + FP + TN + FN} \quad (1)$$

Precision: Precision evaluates the fraction of correctly classified instances or samples among the ones classified as positives. Thus, the formula to calculate the precision is given by:

$$\text{Precision} = \frac{\text{True Positive}}{\text{True Positive} + \text{False Positive}} \quad (2)$$

Recall: Recall is a metric in machine learning that measures the ability of a model to identify all relevant instances of a particular class. It is the ratio of correctly predicted positive observations to the total actual positives, providing insights into a model's completeness in capturing instances of a given class.

$$\text{Recall} = \frac{TP}{TP + FN} \quad (3)$$

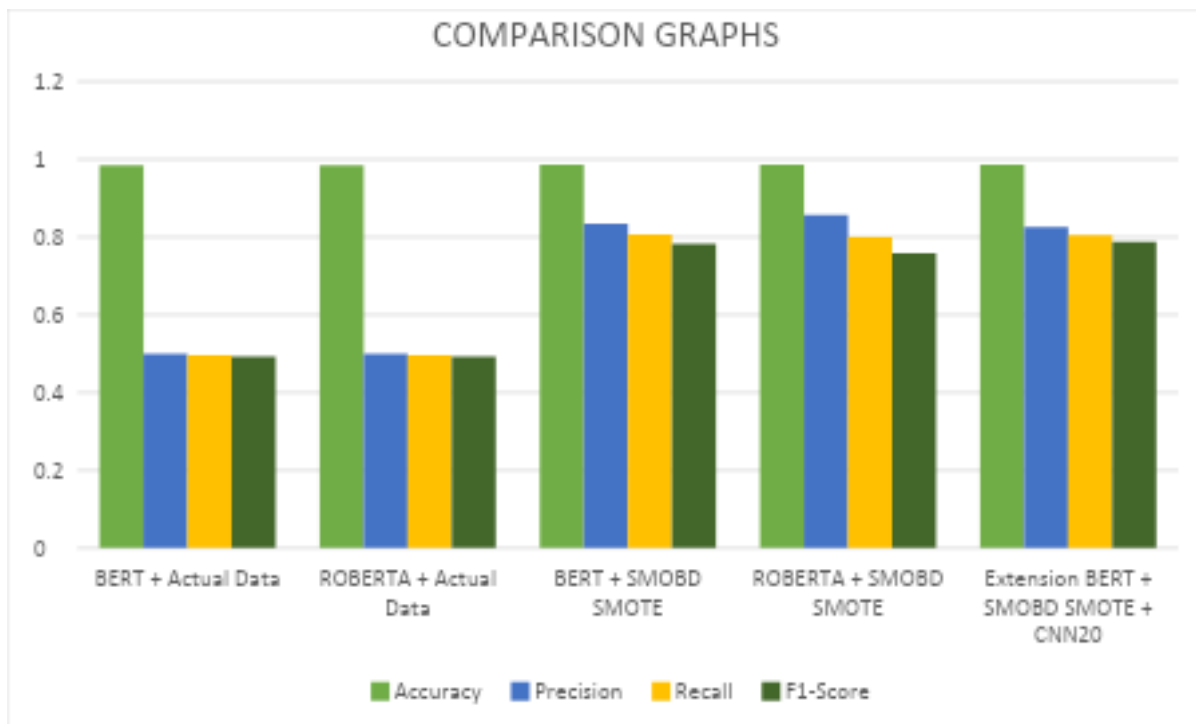
F1-Score: F1 score is a machine learning evaluation metric that measures a model's accuracy. It combines the precision and recall scores of a model. The accuracy metric computes how many times a model made a correct prediction across the entire dataset.

$$\text{F1 Score} = \frac{2 * \text{Recall} * \text{Precision}}{\text{Recall} + \text{Precision}} * 100 \quad (4)$$

We evaluate the performance metrics—accuracy, precision, recall, and F1-score—for each algorithm in Table 1. The BERT + SMOBD SMOTE + CNN2D achieves the highest scores. The table below also presents the metrics of other algorithms for comparison.

Algorithm Name	Accuracy	Precision	Recall	F1-Score
BERT + Actual Data	0.9839	0.5000	0.4959	0.4919
ROBERTA + Actual Data	0.9850	0.5000	0.4962	0.4925
BERT + SMOBD SMOTE	0.9866	0.8339	0.8065	0.7834
ROBERTA + SMOBD SMOTE	0.9858	0.8577	0.7992	0.7577
Extension BERT + SMOBD SMOTE + CNN2D	0.9868	0.8256	0.8052	0.7872

Table. 4.1 Performance Evaluation Metrics



Graph.4.1 Comparison Graphs

Graph 1 displays accuracy in light green, precision in blue, recall in light yellow, and the F1 score in green. The BERT + SMOBD SMOTE + CNN2D outperforms the other algorithms in all metrics, with the highest values compared to the remaining models. The above graph visually represents these details.

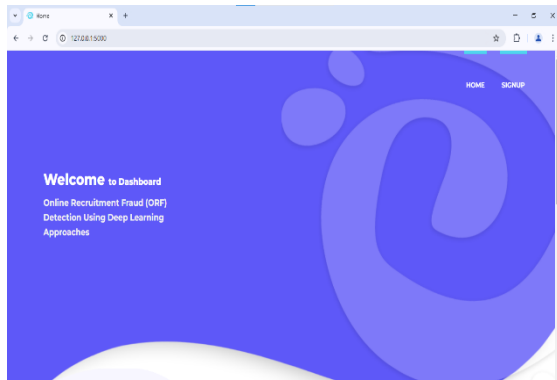


Fig.5 Home Page

In above fig.5 user interface dashboard with navigation and a welcome message.

New Account

username	Username
name	Name
email	Email
number	Mobile Number
password	Password

Remember me [Forgot Password](#)

Register

Already have an account? [Sign in](#)

Fig.6 Registration Page

In above fig.6 sign-up form with fields for username, name, email, mobile number, and password buttons.

The login form consists of a white box with a blue border. Inside, there are two input fields: 'username' with the text 'admin' and 'password' with masked characters '.....'. Below the fields are two links: 'Remember me' with a checked checkbox and 'Forgot Password'. At the bottom is a large blue button labeled 'Log In'. Below the form is a link 'Register here! Sign Up'.

Fig.7 Login Page

In above fig.7 Sign-in form with username and password fields, "Remember Me," "Forgot Password,".

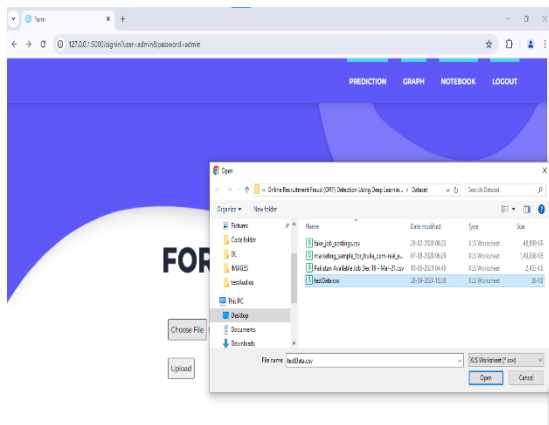


Fig.8 Upload Input Page

In above Fig.8 form with coordinate input field and upload button.

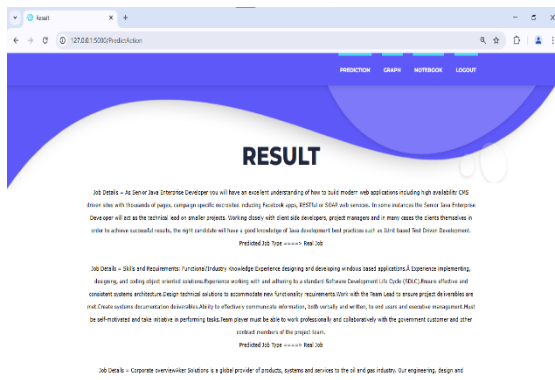


Fig.9 Predict Result for given input

In above Fig.9 Predicted result based on the input test data.

5. CONCLUSIONS AND FUTURE WORK

Conclusion:

The proposed system for detecting Online Recruitment Fraud (ORF) offers an effective solution to the increasing problem of fraudulent job advertisements on digital hiring platforms. By integrating cutting-edge deep learning models such as Bidirectional Encoder Representations from Transformers (BERT) and its improved version, RoBERTa, the system enhances the accuracy and efficiency of fraud detection. These models successfully capture the contextual and semantic meaning within job descriptions, allowing for precise identification of subtle indicators of fraud. Additionally, the application of the SMOTE SMOBD technique addresses the issue of class imbalance by generating synthetic samples for the underrepresented class, ensuring that the model is trained fairly and is not biased toward legitimate postings. One of the key highlights of this research is the performance of the combined model that integrates BERT embeddings, SMOBD oversampling, and a two-dimensional Convolutional Neural Network (CNN2D). This hybrid architecture achieved a high classification accuracy of 98.68%, showcasing the effectiveness of combining contextual understanding, synthetic data balancing, and spatial feature extraction for fraud detection. The system thus provides a reliable and scalable framework that can help online recruitment platforms protect job seekers from scams and enhance the credibility of job postings

Future Work:

Looking forward, the system can be further enhanced by incorporating additional machine learning techniques such as ensemble learning, which could improve generalization by combining the strengths of multiple models. Future enhancements may also include advanced feature extraction methods, the use of recurrent neural networks (RNNs), and attention mechanisms to capture the sequential and contextual flow of information in job postings. Moreover, applying transfer learning from large pre-trained models can improve performance, especially when dealing with smaller or domain-specific datasets. These future directions aim to increase the system's adaptability, accuracy, and reliability, contributing to a safer and more trustworthy online recruitment environment.

REFERENCES

- 1) R. L. Siegel, K. D. Miller, H. E. Fuchs, and A. Jemal, "Cancer Statistics, 2021.," CA. Cancer J. Clin., vol. 71, no. 1, pp. 7–33, Jan. 2021, doi: 10.3322/caac.21654.
- 2) S. Lei et al., "Global patterns of breast cancer incidence and mortality: A population-based cancer registry data analysis from 2000 to 2020," Cancer Commun., vol. 41, no. 11, pp. 1183–1194, 2021.
- 3) M. M. Cackowski et al., "The absence of lymph nodes removed (pNx status) impacts survival in patients with lung cancer treated surgically," Surg. Oncol., vol. 48, p. 101941, 2023.

- 4) S. S. Raoof, M. A. Jabbar, and S. A. Fathima, "Lung Cancer Prediction using Machine Learning: A Comprehensive Approach," in 2020 2nd International Conference on Innovative Mechanisms for Industry Applications (ICIMIA), 2020, pp. 108–115. doi: 10.1109/ICIMIA48430.2020.9074947.
- 5) Y. Chen, E. Zitello, R. Guo, and Y. Deng, "The function of LncRNAs and their role in the prediction, diagnosis, and prognosis of lung cancer," *Clin. Transl. Med.*, vol. 11, no. 4, p. e367, 2021.
- 6) J. A. Barta, C. A. Powell, and J. P. Wisnivesky, "Global epidemiology of lung cancer," *Ann. Glob. Heal.*, vol. 85, no. 1, 2019.
- 7) M. A. Heuvelmans et al., "Lung cancer prediction by Deep Learning to identify benign lung nodules," *Lung Cancer*, vol. 154, no. November 2020, pp. 1–4, 2021, doi: 10.1016/j.lungcan.2021.01.027.
- 8) J. Jumanto, M. F. Mardiansyah, R. Pratama, M. F. Al Hakim, and B. Rawat, "Optimization of breast cancer classification using feature selection on neural network," *J. Soft Comput. Explor.*, vol. 3, no. 2, pp. 105–110, 2022, doi: 10.52465/josce.v3i2.78.
- 9) B. He et al., "Image segmentation algorithm of lung cancer based on neural network model," *Expert Syst.*, vol. 39, no. 3, p. e12822, 2022.
- 10) Q. Wang, Y. Zhou, W. Ding, Z. Zhang, K. Muhammad, and Z. Cao, "Random forest with self-paced bootstrap learning in lung cancer prognosis," *ACM Trans. Multimed. Comput. Commun. Appl.*, vol. 16, no. 1s, pp. 1–12, 2020.
- 11) S. Wankhade and S. Vigneshwari, "A novel hybrid deep learning method for early detection of lung cancer using neural networks," *Healthc. Anal.*, p. 100195, 2023.