

Dimensionality-Reduced Frame Work for Predictive Credit Scoring

Dr. G.Ramasubba Reddy¹, Dudekula Sreya², Bogireddy venkata Sravani³, Guvvala Lasya⁴,
Lomada Sanjay Kumar Reddy⁵

¹Professor & Head, Department of CSE(AI&ML), Sai Rajeswari Institute of Technology, Proddatur-516362, Andhra Pradesh, India.

^{2,3,4}UG Students, Department of Computer Science & Engineering, Sai Rajeswari Institute of Technology, Proddatur-516362, Andhra Pradesh, India.

Abstract- The growing complexity and scale of financial and behavioral datasets have posed significant challenges to traditional credit scoring methods, as conventional models often fail to capture temporal dependencies and nonlinear feature interactions. To address this, a Hybrid Long Short- Term Memory (LSTM) network was designed, incorporating temporal features for accurate credit scoring prediction. The model integrates a hybrid loss function combining binary cross entropy for classification tasks and optimization techniques such as Min-Max normalization, Synthetic Minority Oversampling Technique (SMOTE) for imbalanced data handling, Recursive Feature Elimination (RFE) for feature selection, and Principal Component Analysis (PCA) for dimensionality reduction. Performance evaluation was conducted on benchmark datasets including the Credit Risk dataset, which provides both structured and unstructured financial data. Comparative analysis against baseline models such as Random Forest and XGBoost demonstrated the effectiveness of the approach. Furthermore, a self-attention mechanism was incorporated into the LSTM framework to enhance contextual learning by emphasizing critical input features, leading to improved predictive accuracy. Experimental results indicate that the Hybrid LSTM with self-attention achieved superior performance with 89.87% accuracy, outperforming existing machine learning and deep learning techniques in credit score prediction.

Keywords: LSTM, SMOTE, PCA.

I. INTRODUCTION

Credit scoring prediction plays a critical role in the financial sector, directly influencing risk management, loan approvals, and regulatory compliance. Traditional credit scoring methods, relying primarily on static rules and historical credit behaviors, struggle to adapt to the dynamic and nonlinear nature of modern consumer financial patterns. The rapid growth of digital financial services and alternative data sources, including transaction histories, social media behavior, and spending patterns, has created a demand for more sophisticated computational techniques capable of uncovering complex patterns. Modern financial datasets are increasingly high-dimensional and heterogeneous, containing both structured and unstructured information that conventional scoring systems cannot fully exploit. Additionally, temporal dependencies in user behavior, such as changes in spending habits or repayment trends over time,

present further challenges for traditional models. These limitations have highlighted the need for predictive frameworks that can automatically learn meaningful representations from diverse data sources while capturing temporal dynamics and nonlinear interactions. Accurate and interpretable credit scoring systems are essential for financial institutions to manage risk effectively, allocate credit responsibly, and maintain consumer trust. The increasing availability of rich financial data streams necessitates the exploration of advanced learning approaches that go beyond static modeling, offering robust, scalable, and adaptive solutions for contemporary credit assessment.

Objective

The main objective is to develop an advanced credit scoring prediction framework that enhances the accuracy and reliability of financial risk evaluation by integrating traditional and deep learning models. The system employs Random Forest and XGBoost to

establish baseline performance using structured financial data, capturing essential relationships between variables such as income, payment behavior, and loan history. To overcome the limitations of these models in handling sequential dependencies, a Hybrid Long Short-Term Memory (LSTM) model is introduced to process temporal and nonlinear patterns effectively. The approach incorporates preprocessing techniques such as Min-Max normalization, SMOTE, Recursive Feature Elimination (RFE), and Principal Component Analysis (PCA) to optimize data representation, reduce dimensionality, and ensure balanced, feature-rich learning for precise and efficient credit score prediction.

Problem Statement

Traditional credit scoring systems rely on static rules and outdated modeling techniques that fail to capture temporal dependencies, nonlinear relationships, and dynamic financial behaviors, leading to inaccurate risk evaluations and inconsistent credit score predictions in modern financial environments. With the rapid growth of digital finance and diverse data sources, financial institutions now deal with high-dimensional, unstructured, and time-dependent data that traditional machine learning algorithms like Random Forest and XGBoost cannot efficiently process or interpret. Borrowers, lenders, and financial organizations are directly affected by inaccurate credit evaluations, resulting in unfair loan approvals, misjudged risk assessments, financial losses, and reduced trust in automated credit decision-making systems across banking and lending sectors. Inefficient credit scoring models lead to biased evaluations, higher default rates, and regulatory non-compliance, hindering financial inclusion and making it difficult for

- institutions to manage risk effectively while maintaining fairness and transparency in lending operations.
- To overcome these limitations, a Hybrid LSTM-based deep learning framework integrating temporal feature extraction, data balancing, and dimensionality reduction techniques is introduced to deliver more accurate, adaptive,

and interpretable credit scoring predictions for financial risk management.

II. FEASIBILITY STUDY

Feasibility study examines the viability or sustainability of an idea, project, or business. The study examines whether there are enough resources to implement it, and the concept has the potential to generate reasonable profits. In addition, it will demonstrate the benefits received in return for taking the risk of investing in the idea.

Types of Feasibility Study

There are several different kinds of feasibility studies. Understanding the types of feasibility studies and the technicalities of the concept is important for any business. They are elaborated below:

Technical Feasibility

Technical feasibility study checks for accessibility of technical resources in the organization. In case technological resources exist, the study team will conduct assessments to check whether the technical team can customize or update the existing technology to suit the new method of workings for the project by properly checking the health of the hardware and software. Many factors need to be taken into consideration here, like staffing requirements, transportation, and technological competency.

Financial Feasibility

Financial feasibility allows an organization to determine cost-benefit analysis. It gives details about the investment that has to go in to get the desired level of benefit (profit). Factors such as total cost and expenses are considered to arrive simultaneously. With this data, the companies know their present state of financial affairs and anticipate future monetary requirements and the sources from which the company can acquire them. Investors can largely benefit from the economic analysis done. Assessing the return on investment of a particular asset or acquisition can be a financial feasibility study example.

Market Feasibility

It assesses the industry type, the existing marketing characteristics and improvements to make it better, the growth evident and needed, competitive environment of the company's products and services. Preparations of sales projections can thus be a good market feasibility study example.

Organization Feasibility

Organization feasibility focuses on the organization's structure, including the legal system, management team's competency, etc. It checks whether the existing conditions will suffice to implement the business idea.

III. LITERATURE SURVEY

Yunus Emre Gür, et. Al., 2025[1] has been developed the Integration of CNN Models and Machine Learning Methods in Credit Score Classification: 2D Image Transformation and Feature Extraction.

The problem of accurately classifying credit scores is critical for financial institutions to assess individual creditworthiness and effectively manage credit risk. Traditional methods often face limitations when processing large datasets, resulting in lower accuracy and longer processing time. To address this issue, this paper proposes a novel approach to credit score classification by integrating convolutional neural networks (CNN) with machine learning methods. First, a 1D dataset of sequential text data is transformed into 2D greyscale images to use 2D CNN models for feature extraction and classification. Six CNN architectures—DenseNet201, GoogLeNet, MobileNetV2, ResNet18, ShuffleNet, and SqueezeNet—are implemented, and the features in the last layer (1000 features) of each CNN are classified using the softmax method. To further improve the performance, the two best CNN models were selected, and a new fully connected layer (NewFC) was added.

Andry Alamsyah, et al., 2026([2] has been developed the Innovative Credit Risk Assessment: Leveraging Social Media Data for Inclusive Credit Scoring in Indonesia's Fintech Sector. This study introduces a novel approach leveraging social media analytics

and advanced machine learning techniques to assess the creditworthiness of individuals without traditional credit histories and collateral assets. Conventional credit scoring methods tend to rely heavily on central bank credit information, especially traditional collateral assets such as property or savings accounts.

young individuals lacking traditional credit histories or collateral assets—as either good or bad credit risks based on expert judgment thresholds. This innovative approach questions conventional financial evaluation methods and enhances access to credit for marginalized communities. The research question addressed in this study is how to develop a credit scoring mechanism using social media data.

Isaac Adinoyi Salami, et al., 2025[3] has been developed the Addressing-Bias-and-Data-Privacy-Concerns-in-AI-Driven-Credit-Scoring-Systems-Through-Cybersecurity-Risk-Assessment. This study investigates the role of cybersecurity risk assessment in mitigating these risks, utilizing multiple datasets, including the Home Mortgage Disclosure Act (HMDA) dataset, the Equifax Data Breach Report, the Financial Cybersecurity Incidents Database, and the MITRE ATT&CK Financial Sector Threat Intelligence Dataset.

The increasing reliance on artificial intelligence (AI) in credit scoring has raised concerns about algorithmic bias and data privacy, necessitating robust cybersecurity risk assessment frameworks. This study investigates the role of cybersecurity risk assessment in mitigating these risks, utilizing multiple datasets, including the Home Mortgage Disclosure Act (HMDA) dataset, the Equifax Data Breach Report, the Financial Cybersecurity Incidents Database, and the MITRE ATT&CK Financial Sector Threat Intelligence Dataset. We employ statistical fairness metrics, Bayesian Probability Modeling, Markov Chain Analysis, and Monte Carlo Simulations to evaluate the extent of bias, privacy risks, and cybersecurity vulnerabilities. Findings reveal significant disparities in loan approvals, with Black applicants receiving approval rates 28% lower than White applicants ($\chi^2 = 59.83$, $p < 0.001$), highlighting systemic bias in AI-driven credit scoring.

Neepta Biswas. Et al.,2022[4] has been developed the Automated credit assessment framework using ETL process and machine learning.

The automated ETL (extraction, transformation, and load) process ensures data ingestion into the data warehouse in near real-time, and insights are generated through the BI process based on real-time data. In this paper, we have concentrated on automated credit risk assessment in the financial domain based on the machine learning approach. The machine learning-based classification techniques can furnish a self-regulating process to categorize data. Establishing an automated credit decision-making system helps the lending institution to manage the risks, increase operational efficiency and comply with regulators. In this paper, an empirical approach is taken for credit risk assessment using logistic regression and neural network classification method in compliance with Basel II standards. Here, Basel II standards are adopted to calculate the expected loss.

Derick Kazimoto, et al.,2025[5] has been developed the Empowering Africa's Disfranchised SMEs: Machine Learning-Based Credit Scoring for Informal African Merchants.

This study addresses this challenge by developing tailored credit scoring models for informal merchants using supervised learning techniques. Leveraging data from a financial technology company in Lesotho (South Africa), we applied logistic regression and support vector machines to predict the likelihood of loan defaults among merchants. Our methodology involved the evaluation of six logistic regression models and twelve support vector machine models, assessing their effectiveness in default prediction. The results provide a robust tool for more accurate assessment of creditworthiness, reducing the risk of lending to potential defaulters. The study underscores the potential of supervised learning methods to create impactful financial solutions and suggests a pathway towards narrowing the financial inclusion gap in the informal African economy. This approach not only aids in risk reduction for lenders but also empowers

a critical segment of the economy by enabling better financial support and growth opportunities for informal sector merchants.

IV. EXISTING SYSTEM

Existing credit scoring systems primarily rely on traditional machine learning algorithms such as Random Forest, XGBoost, and feedforward neural networks to evaluate borrower risk. These models utilize structured financial data, including income, payment history, and employment stability, to predict creditworthiness. While effective for basic pattern recognition, they often struggle to capture temporal dependencies and complex non-linear relationships inherent in financial behaviors. Conventional methods typically require extensive feature engineering, data preprocessing, and manual tuning to handle imbalanced datasets and missing values. Random Forest and XGBoost provide reasonable predictive performance through ensemble learning and gradient boosting techniques but are limited in processing sequential data. Feedforward neural networks can model non-linear interactions but do not inherently capture time-dependent trends in credit activity. Overall, these approaches are constrained by their inability to fully leverage both structured and unstructured data or to adapt to evolving consumer financial patterns, creating gaps in predictive accuracy and robustness for modern credit scoring applications.

Disadvantages of Existing System:

- Struggles to capture temporal dependencies, limiting prediction accuracy for sequential financial data.
- Requires extensive manual feature engineering, preprocessing, and tuning, increasing implementation complexity.
- Cannot effectively handle unstructured data, reducing model adaptability to modern heterogeneous datasets.
- Limited scalability for evolving consumer behaviors and complex nonlinear relationships, constraining robustness in dynamic financial environments.

SYSTEM ARCHITECTURE



Fig.4.1 System Architecture

IMPLEMENTATION MODULES

- **Importing the Packages:** This module loads all required Python libraries and classes for data handling, preprocessing, visualization, and model building. Essential packages like NumPy, Pandas, Scikit-learn, TensorFlow, and Matplotlib are imported to enable smooth execution of subsequent modules and ensure compatibility across all processing and modeling steps.
- **Exploring the Dataset:** This module loads the Credit Risk dataset and displays its structure, values, and statistical summary. It provides an initial understanding of the data, highlighting features, data types, and potential inconsistencies, helping to plan preprocessing and analysis strategies.
- **Data Processing:** Handles missing values, duplicates, and inconsistencies in the dataset. It ensures all entries are valid, complete, and suitable for model training. Proper preprocessing guarantees data quality, enhancing predictive accuracy and minimizing errors during model execution.
- **Visualization:** Generates graphical representations of data distributions, class labels, and categorical feature relationships. This module helps in understanding patterns such as class imbalance, loan purposes, and applicant behavior, enabling informed decisions for feature engineering and model preparation.
- **Label Encoding:** Converts non-numeric categorical variables into numeric form using encoding techniques. This step ensures compatibility with machine learning and deep learning models, allowing algorithms to process and interpret categorical information effectively.
- **Feature Selection:** Applies Recursive Feature Elimination (RFE) and Principal Component Analysis (PCA) to select relevant features and reduce dimensionality. This module improves model performance, computational efficiency, and focuses on features contributing most to credit score prediction.
- **Split the Data into Train & Test:** Divides the processed dataset into training and testing subsets, typically 80% for training and 20% for testing. This module ensures unbiased evaluation of models and effective learning on representative data.
- **Model Training:** Trains multiple models, including Random Forest, XGBoost, Propose Hybrid LSTM, and Extension Hybrid LSTM with Attention. Each model learns patterns from financial and behavioral data to predict creditworthiness accurately, handling both structured and sequential features.
- **Evaluation:** Assesses trained models using metrics like accuracy, precision, recall, and F1-score. This module determines model performance, compares algorithms, and identifies the most effective approach for credit scoring prediction.
- **SHAP:** Applies SHAP analysis to interpret model predictions by identifying feature importance. This module explains which factors influence credit score predictions, enhancing transparency and interpretability of deep learning and machine learning results.
- **Flask Server:** Initializes a web server for running the credit scoring system online. This module allows users to interact with the trained models via a web interface, providing real-time credit score predictions.
- **User Login:** Authenticates users by verifying credentials before accessing the system. This module ensures secure access, restricting operations to authorized personnel for data privacy and system integrity.
- **Credit Score Prediction:** Processes uploaded test data to generate credit scores as Good or Bad. It displays predictions clearly for users,

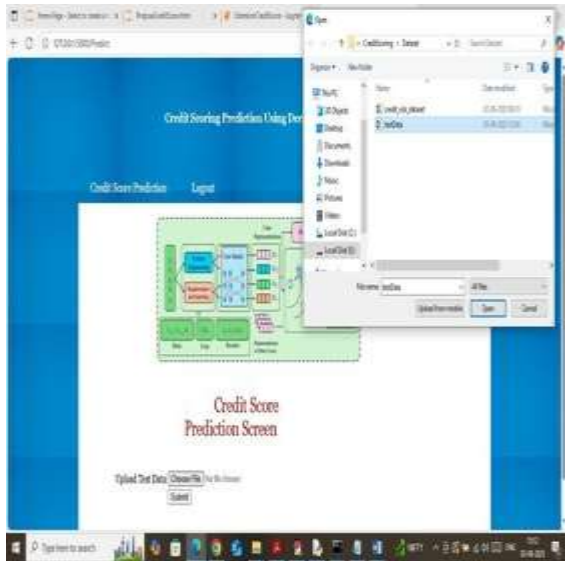


Fig:5.4

In above screen selecting and uploading test data file and then click on buttons to get below page

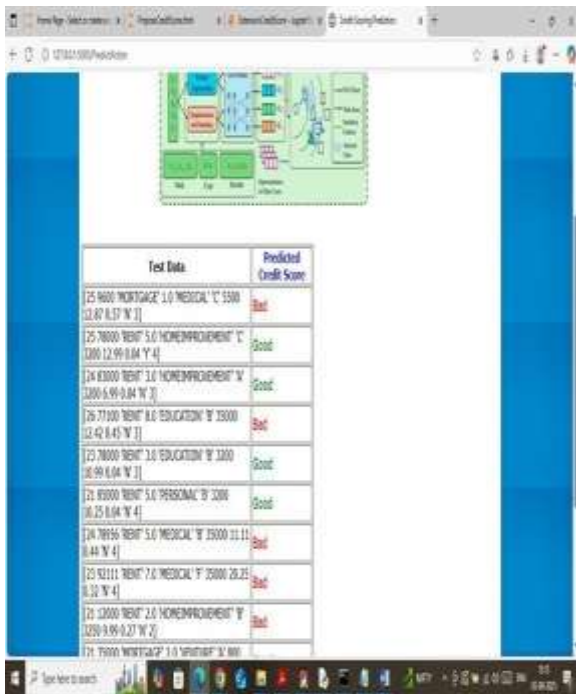


Fig:5.5

In above screen in first column can see test data values and in second column can see predicted Credit score as 'Good or Bad'. This score can be utilized by banks to approved or reject loans

VI. CONCLUSION

The study demonstrates that integrating temporal features with advanced deep learning architectures significantly enhances credit scoring prediction in the financial sector. By employing a Hybrid LSTM model optimized with Min-Max normalization, SMOTE for imbalanced data, RFE feature selection, and PCA for dimensionality reduction, the model effectively captures complex relationships in structured and unstructured financial datasets. Comparative evaluations show that traditional machine learning models, while competent, are limited in handling sequential patterns and nonlinear dependencies inherent in financial behavior data.

The introduction of a self-attention mechanism further improves the LSTM's ability to focus on the most relevant features, capturing contextual dependencies and enhancing prediction accuracy. Experimental results on the Credit Risk dataset reveal that the Hybrid LSTM with self-attention achieved 89.87% accuracy, outperforming other models and demonstrating robustness in distinguishing between good and bad credit scores. The outcome validates that attention-enhanced temporal modeling can provide more reliable insights for financial institutions, supporting informed decision-making in loan approvals and risk assessment. This approach highlights the importance of combining temporal feature extraction, data preprocessing, and attention mechanisms in deep learning models to address challenges in credit scoring applications effectively.

REFERENCES

- Gür, Y. E., Toğaçar, M., & Solak, B. (2025). Integration of CNN models and machine learning methods in credit score classification: 2D image transformation and feature extraction. *Computational Economics*, 65(5), 2991-3035.
- Alamsyah, A., Hafidh, A. A., & Mulya, A. D. (2025). Innovative Credit Risk Assessment: Leveraging Social Media Data for Inclusive Credit Scoring in Indonesia's Fintech Sector. *Journal of Risk and Financial Management*, 18(2), 74.
- Salami, I. A., Adesokan-Imran, T. O., Tiwo, O. J., Metibemu, O. C., Olutimehin, A. T., & Olaniyi, O.

- O. (2025). Addressing bias and data privacy concerns in AI-driven credit scoring systems through cybersecurity risk assessment. *Asian Journal of Research in Computer Science*, 18(4), 59-82.
4. Biswas, N., Mondal, A. S., Kusumastuti, A., Saha, S., & Mondal, K. C. (2025). Automated credit assessment framework using ETL process and machine learning. *Innovations in Systems and Software Engineering*, 21(1), 257-270.
 5. Kazimoto, D., & Baadel, S. (2025). Empowering Africa's Disfranchised SMEs: Machine Learning-Based Credit Scoring for Informal African Merchants. *Human-Centric Intelligent Systems*, 1-9.
 6. T. Schick and H. Schütze, "Exploiting cloze-questions for few-shot text classification and natural language inference," in Proc. 16th Conf. Eur. Chapter Assoc. Comput. Linguistics, Main Volume, 2021. [Online]. Available: <https://arxiv.org/abs/2001.07676>
 7. S. Hu, N. Ding, H. Wang, Z. Liu, J. Wang, J. Li, W. Wu, and M. Sun, "Knowledgeable prompt-tuning: Incorporating knowledge into prompt verbalizer for text classification," in Proc. Annu. Meeting Assoc. Comput. Linguistics, 2022. [Online]. Available: <https://arxiv.org/abs/2108.02035>
 8. X. Sun, X. Li, J. Li, F. Wu, S. Guo, T. Zhang, and G. Wang, "Text classification via large language models," in Proc. Findings Assoc. Comput. Linguistics (EMNLP), 2023. [Online]. Available: <https://arxiv.org/abs/2305.08377>
 9. X. Han, W. Zhao, N. Ding, Z. Liu, and M. Sun, "PTR: Prompt tuning with rules for text classification," *AI Open*, vol. 3, pp. 182-192, Jan. 2022.
 10. Z. Wen and Y. Fang, "Augmenting low-resource text classification with graph-grounded pre-training and prompting," in Proc. 46th Int. ACM SIGIR Conf. Res. Develop. Inf. Retr., Jul. 2023, pp. 506-516.
 11. S. Minaee, N. Kalchbrenner, E. Cambria, N. Nikzad-Khasmakh, M. Chenaghlu, and J. Gao, "Deep learning-based text classification," *ACM Comput. Surveys*, vol. 54, no. 3, pp. 1-40, 2021.
 12. Z. Jiang, M. Y. R. Yang, M. Tsirlin, R. Tang, Y. Dai, and J. J. Lin, "'Low resource' text classification: A parameter-free classification method with compressors," in Proc. Annu. Meeting Assoc. Comput. Linguistics, 2023. [Online]. Available: https://aclanthology.org/2023.findings_acl.426/?trk=public_post_comment-text
 13. S. Min, M. Lewis, H. Hajishirzi, and L. Zettlemoyer, "Noisy channel language model prompting for few-shot text classification," in Proc. 60th Annu. Meeting Assoc. Comput. Linguistics (Long Papers), vol. 1, 2022. [Online]. Available: <https://arxiv.org/abs/2108.04106>
 14. A. Karimi, L. Rossi, and A. Prati, "AEDA: An easier data augmentation technique for text classification," in Proc. Findings Assoc. Comput. Linguistics: EMNLP, 2021. [Online]. Available: <https://arxiv.org/abs/2108.13230>
 15. S. Molitorys, B. Pilot, and K. Smyrak, "Text classification," in Proc. Int. Conf. Inf. Technol., 2024. [Online]. Available: https://www.researchgate.net/profile/V-Tampakas/publication/234812606_Text_classification_a_recent_overview/links/0c96051ee1dfd2a9f8000000/Text-classification-a-recent-overview.pdf
 16. A. Gasparetto, M. Marcuzzo, A. Zangari, and A. Albarelli, "A survey on text classification algorithms: From text to predictions," *Information*, vol. 13, no. 2, p. 83, Feb. 2022.
 17. R. Qasim, W. H. Bangyal, M. A. Alqarni, and A. Ali Almazroi, "A fine tuned BERT-based transfer learning approach for text classification," *J. Healthcare Eng.*, vol. 2022, pp. 1-17, Jan. 2022.
 18. Z. Wang, P. Wang, L. Huang, X. Sun, and H. Wang, "Incorporating hierarchy into text encoder: A contrastive learning approach for hierarchical text classification," in Proc. 60th Annu. Meeting Assoc. Comput. Linguistics (Long Papers), vol. 1, 2022. [Online]. Available: <https://arxiv.org/abs/2203.03825>
 19. V. Dogra, S. Verma, Kavita, P. Chatterjee, J. Shafi, J. Choi, and M. F. Ijaz, "A complete process of text classification system using state-of-the-art NLP models," *Comput. Intell. Neurosci.*, vol. 2022, pp. 1-26, Jun. 2022.
 20. M. Bayer, M.-A. Kaufhold, and C. Reuter, "A survey on data augmentation for text classification," *ACM Comput. Surv.*, vol. 55, no. 7, pp. 1-39, Jul. 2023.

21. S. Soni, S. S. Chouhan, and S. S. Rathore, "TextConvoNet: A convolutional neural network based architecture for text classification," *Appl. Intell.*, vol. 53, no. 11, pp. 14249–14268, Jun. 2023.
22. M. Umer, Z. Imtiaz, M. Ahmad, M. Nappi, C. Medaglia, G. S. Choi, and A. Mehmood, "Impact of convolutional neural network and FastText embedding on text classification," *Multimedia Tools Appl.*, vol. 82, no. 4, pp. 5569–5585, Feb. 2023.
23. J. Chen, R. Zhang, Y. Mao, and J. Xu, "ContrastNet: A contrastive learning framework for few-shot text classification," in *Proc. AAAI Conf. Artif. Intell.*, vol. 36, 2022, pp. 10492–10500.
24. A. S. Alammery, "BERT models for Arabic text classification: A systematic review," *Appl. Sci.*, vol. 12, no. 11, p. 5720, Jun. 2022.
25. X. Chen, P. Cong, and S. Lv, "A long-text classification method of Chinese news based on BERT and CNN," *IEEE Access*, vol. 10, pp. 34046–34057, 2022.
26. Y. Lin, Y. Meng, X. Sun, Q. Han, K. Kuang, J. Li, and F. Wu, "BertGCN: Transductive text classification by combining GNN and BERT," in *Proc. Findings Assoc. Comput. Linguistics (ACL-IJCNLP)*, 2021. [Online]. Available: <https://arxiv.org/abs/2105.05727>
27. D. Shao, C. Li, C. Huang, Y. Xiang, and Z. Yu, "A news classification applied with new text representation based on the improved LDA," *Multimedia Tools Appl.*, vol. 81, no. 15, pp. 21521–21545, Jun. 2022.
28. M. Sadat and C. Caragea, "Hierarchical multi-label classification of scientific documents," 2022, arXiv:2211.02810.
29. Q. Li, "Research on text sentiment classification and recognition based on improved BiLSTM and TextCNN," in *Proc. 2nd Int. Conf. Mechatronics, IoT Ind. Informat. (ICMIII)*, Jun. 2024, pp. 210–219.
30. H.H.Luong, L.T.T.Le, and H.T.Nguyen, "An approach for web content classification with FastText," in *Proc. Int. Conf. Comput. Data Social Netw.*, 2024, pp. 13.