

Advanced Sentiment Analysis Using Deep Learning: A Comprehensive Framework for High-Accuracy and Interpretable Models

Ata Amrullah

Department of Informatics, Darul Ulum Islamic University, East Java, Indonesia

Received: 2 December 2024

Accepted: 21 January 2025

Published: 24 January 2025

Keywords:

Sentiment Analysis;
Deep Learning;
Natural Language Processing;
Interpretability;
Hybrid Models.

Corresponding author:

Ata Amrullah
ata@unisda.ac.id

Abstract

Sentiment analysis has become a critical tool for understanding public opinion, customer feedback, and social media trends. Despite significant advancements in deep learning, existing models often struggle with accuracy, generalizability, and interpretability, particularly when applied to complex and noisy datasets. In this paper, we propose a novel deep learning framework for sentiment analysis that addresses these limitations by combining the strengths of convolutional neural networks (CNNs) and transformer-based architectures. Our framework leverages verified and high-quality datasets, including Twitter Sentiment140, IMDb movie reviews, and Amazon product reviews, to ensure robustness and reliability. We introduce a hybrid model that integrates multi-head attention mechanisms with hierarchical feature extraction, enabling the model to capture both local and global contextual information effectively. Additionally, we employ state-of-the-art interpretability techniques, such as SHAP and LIME, to provide transparent and human-understandable explanations for model predictions. Experimental results demonstrate that our framework achieves superior performance compared to existing state-of-the-art methods, with an accuracy of 94.3%, an F1-score of 93.8%, and an AUC-ROC score of 97.2%. Furthermore, our model's interpretability features offer valuable insights into decision-making processes, making it highly applicable for real-world applications such as brand monitoring, market analysis, and political sentiment tracking. This study not only advances the field of sentiment analysis but also provides a scalable and interpretable solution for future research in natural language processing.

1. INTRODUCTION

Sentiment analysis, also known as opinion mining, has become one of the most critical and widely studied areas in natural language processing (NLP) due to its extensive applications in understanding public opinion, customer feedback, and social media dynamics [1]. The exponential growth of user-generated content on platforms such as Twitter, Amazon, and IMDb has created an unprecedented demand for automated tools capable of extracting meaningful insights from textual data [2]. Sentiment analysis enables organizations to

gauge customer satisfaction, monitor brand reputation, and even predict market trends, making it a cornerstone of modern data-driven decision-making [3]. For instance, businesses use sentiment analysis to improve products and services based on customer reviews, while governments leverage it to understand public sentiment toward policies or events [4].

Despite the significant advancements in sentiment analysis, particularly with the advent of deep learning techniques, several challenges remain unresolved. Traditional methods, such as lexicon-based approaches and machine

learning models, often struggle to capture the nuanced and context-dependent nature of human language [5]. While deep learning models, including convolutional neural networks (CNNs) and recurrent neural networks (RNNs), have shown promise in improving accuracy, they still face limitations in handling complex datasets with noise, sarcasm, and ambiguity [6]. For example, sarcasm and irony in social media posts often lead to misclassification, as these linguistic features are difficult to detect using conventional methods [7]. Furthermore, the interpretability of these models remains a critical issue, as stakeholders often require transparent and explainable insights to trust and act upon the results [8].

The rise of transformer-based architectures, such as BERT (Bidirectional Encoder Representations from Transformers) and GPT (Generative Pre-trained Transformer), has revolutionized the field of NLP by achieving state-of-the-art performance on various tasks, including sentiment analysis [9]. However, these models are computationally expensive and often require large amounts of labeled data for training, which may not be available for specific domains or languages [10]. Additionally, while transformer-based models excel in capturing long-range dependencies and contextual information, their interpretability remains limited, making it difficult to understand how they arrive at specific predictions [11].

To address these challenges, this paper proposes a novel deep learning framework for sentiment analysis that combines the strengths of transformer-based architectures and hierarchical feature extraction mechanisms. Our framework is designed to achieve high accuracy while maintaining interpretability, making it suitable for real-world applications. The key contributions of this work are as follows:

1. **A Hybrid Deep Learning Model:** We introduce a hybrid architecture that integrates multi-head attention mechanisms with CNN-based feature extraction, enabling the model to capture both local and global contextual information effectively. This approach leverages the strengths of transformers in handling long-

range dependencies while utilizing CNNs for efficient local feature extraction.

2. **Verified and High-Quality Datasets:** Our experiments are conducted on widely recognized datasets, including Twitter Sentiment140, IMDb movie reviews, and Amazon product reviews, ensuring the robustness and generalizability of our results. We also perform rigorous preprocessing to remove noise and bias, ensuring the reliability of our findings.
3. **State-of-the-Art Interpretability:** We employ advanced interpretability techniques, such as SHAP (SHapley Additive exPlanations) and LIME (Local Interpretable Model-agnostic Explanations), to provide transparent and human-understandable explanations for model predictions. This enhances the trustworthiness of our framework and makes it more accessible to non-technical stakeholders.
4. **Superior Performance:** Our framework achieves an accuracy of 94.3%, an F1-score of 93.8%, and an AUC-ROC score of 97.2%, outperforming existing state-of-the-art methods. These results demonstrate the effectiveness of our approach in handling complex and noisy datasets.

The remainder of this paper is organized as follows: Section 2 reviews related work in sentiment analysis and deep learning, highlighting recent advancements and limitations. Section 3 details the proposed methodology, including dataset preprocessing, model architecture, and evaluation metrics. Section 4 presents the experimental results and discussion, comparing our framework with existing methods. Finally, Section 5 concludes the paper and outlines future research directions.

2. RELATED WORK

Sentiment analysis has been a cornerstone of natural language processing (NLP) research, with applications spanning social media monitoring, customer feedback analysis, and market trend prediction. Over the past five years, advancements in deep learning and transformer-based architectures have significantly improved the performance of sentiment analysis models.

Table 1. Comparison of Literature Reviews

Study	Approach	Contribution	Limitations
Samonte et al. (2023) [12]	Hybrid CNN-LSTM with attention	Achieved state-of-the-art performance on tweets during the COVID-19 pandemic	Computationally expensive; struggles with long-range dependencies.
Wang et al. (2022) [13]	Graph-based model with GNNs	Improved performance on datasets with complex sentence structures (e.g., SemEval).	Requires dependency parsing, which may not generalize to all languages.
Shao et al. (2023) [14]	Capsule Networks (CapsNets)	Enhanced ability to handle polysemy and context-dependent sentiment.	Limited scalability for large datasets.
Nguyen et al. (2020) [15]	Fine-tuned BERT	Superior results on IMDb and Amazon reviews.	High computational cost and resource requirements.
Sanh et al. (2019) [16]	DistilBERT	Lightweight and efficient alternative to BERT.	Slight drop in performance compared to BERT.
Lan et al. (2020) [17]	ALBERT	Reduced parameter size while maintaining performance.	Requires extensive pre-training on large corpora.
Lee et al. (2020) [18]	BioBERT	Domain-specific model for biomedical sentiment analysis.	Limited to biomedical texts; not generalizable to other domains.
Liu et al. (2021) [19]	FinBERT	Optimized for financial sentiment analysis.	Requires domain-specific pre-training.
Córdova et al. (2022) [20]	Interpretable attention for BERT	Improved transparency of BERT predictions.	Adds computational overhead to the model.
Hemker et al. (2023) [21]	Hybrid deep learning + rule-based	Generated human-understandable explanations for predictions.	Rule-based component may not scale well to large datasets.
Liu et al. (2022) [22]	Domain-adaptive BERT for healthcare	Achieved significant improvements in healthcare sentiment analysis.	Requires domain-specific data for fine-tuning.
Yekrangi et al. (2023) [23]	BERT + domain-specific embeddings	Outperformed traditional methods on financial datasets.	Limited to financial texts; not generalizable to other domains.
Conneau et al. (2020) [24]	XLM-R (Cross-lingual Language Model)	State-of-the-art performance on multilingual sentiment analysis.	Struggles with code-switching and low-resource languages.
Agüero et al. (2021) [25]	Code-switching-aware model	Improved performance on multilingual social media data.	Requires annotated code-switching data, which is scarce.

However, challenges related to accuracy, interpretability, and domain adaptation persist. This section provides a comprehensive review of recent studies (2021–2025) in sentiment analysis, highlighting their contributions and limitations.

2.1. Deep Learning Approaches for Sentiment Analysis

Deep learning has revolutionized sentiment analysis by enabling models to automatically learn features from raw text data. Convolutional Neural Networks (CNNs) and Recurrent Neural Networks (RNNs) have been widely adopted for this task. For instance, Samonte et al. [12] proposed a hybrid CNN-LSTM model that leverages word embeddings and attention mechanisms to capture both local and global contextual information. Their model achieved state-of-the-art performance in tweets during the COVID-19 pandemic, demonstrating the effectiveness of combining CNNs and RNNs for sentiment analysis.

Despite their success, CNNs and RNNs face limitations in handling long-range dependencies and complex linguistic structures. To address this, recent studies have explored the use of graph neural networks (GNNs) for sentiment analysis. Wang et al. [13] introduced a graph-based model that constructs a dependency tree from text and uses GNNs to capture syntactic and semantic relationships. Their approach outperformed traditional CNN and RNN models on datasets with complex sentence structures, such as SemEval-2014.

Another emerging trend is the use of capsule networks (CapsNets) for sentiment analysis. CapsNets, which were originally designed for image processing, have shown promise in capturing hierarchical relationships in text. Shao et al. [14] proposed a CapsNet-based model that uses dynamic routing to group words into higher-level semantic units, improving the model's ability to handle polysemy and context-dependent sentiment. Their results demonstrated significant improvements over baseline methods on multi-domain sentiment analysis tasks.

2.2. Transformer-Based Models

The introduction of transformer-based architectures, such as BERT (Bidirectional Encoder Representations from Transformers)

[9], has marked a significant breakthrough in NLP. BERT's ability to capture bidirectional context has led to substantial improvements in sentiment analysis. For example, Nguyen et al. [15] fine-tuned BERT for sentiment classification and achieved superior results on multiple datasets, including IMDb and Amazon product reviews. Their work highlighted the importance of pre-training on large corpora for achieving high accuracy in sentiment analysis tasks.

However, BERT's computational complexity and high resource requirements remain a challenge, particularly for real-time applications. To address this, recent studies have explored lightweight variants of transformers. For instance, DistilBERT [16] and ALBERT [17] have been proposed as more efficient alternatives, achieving comparable performance with fewer parameters. These models have shown promise in sentiment analysis tasks, particularly in scenarios where computational resources are limited.

Another notable development is the use of domain-specific transformers for sentiment analysis. For example, BioBERT [18] and FinBERT [19] have been pre-trained on biomedical and financial texts, respectively, and fine-tuned for sentiment analysis in these domains. These models have demonstrated significant improvements over general-purpose transformers, highlighting the importance of domain adaptation in sentiment analysis.

2.3. Interpretability in Sentiment Analysis

Interpretability has emerged as a critical concern in sentiment analysis, as stakeholders often require transparent and explainable models. Techniques such as SHAP (SHapley Additive exPlanations) [20] and LIME (Local Interpretable Model-agnostic Explanations) [8] have been widely adopted to provide insights into model predictions. For example, Ribeiro et al. [8] demonstrated how LIME can be used to explain the predictions of black-box models, making them more accessible to non-technical users.

Recent work has focused on improving the interpretability of transformer-based models. Córdova et al. [20] proposed an interpretable attention mechanism for BERT, enabling users to understand how the model assigns importance to different words in a sentence.

Their approach has shown potential in improving the transparency of sentiment analysis models, particularly in high-stakes applications such as healthcare and finance.

Another promising direction is the use of rule-based explanations for deep learning models. For instance, Hemker et al. [21] developed a hybrid model that combines deep learning with rule-based reasoning to generate human-understandable explanations for sentiment predictions. Their approach has been successfully applied to customer feedback analysis, demonstrating its potential for real-world applications.

2.4. Domain-Specific Sentiment Analysis

Domain-specific sentiment analysis has gained attention in recent years, as models trained on general datasets often struggle to perform well in specialized domains. For example, medical sentiment analysis requires understanding complex terminology and context, which is not adequately captured by general-purpose models. To address this, Liu et al. [22] developed a domain-adaptive BERT model for healthcare sentiment analysis, achieving significant improvements over baseline methods.

Similarly, financial sentiment analysis has been explored to analyze market trends and investor sentiment. Yekrangi et al. [23] proposed a hybrid model combining BERT with domain-specific embeddings, which outperformed traditional methods on financial datasets. These studies highlight the importance of domain adaptation in sentiment analysis and the need for specialized models.

2.5. Multilingual Sentiment Analysis

With the increasing globalization of digital content, multilingual sentiment analysis has become a critical area of research. Traditional models often fail to generalize across languages due to differences in syntax, semantics, and cultural context. Recent work by Conneau et al. [24] introduced XLM-R (Cross-lingual Language Model-RoBERTa), a transformer-based model pre-trained on multiple languages, which achieved state-of-the-art performance on multilingual sentiment analysis tasks.

Despite these advancements, challenges such as code-switching and low-resource languages remain. For instance, Agüero et al.

[25] proposed a code-switching-aware model for sentiment analysis in multilingual social media data, demonstrating the potential of leveraging linguistic diversity to improve model performance.

The systematic summary of related work, highlighting key contributions and limitations, is presented in Table 1.

3. RESEARCH METHODOLOGY

This section outlines the methodology employed in this study, including the dataset selection, preprocessing steps, model architecture, training process, and evaluation metrics. The proposed framework aims to achieve high accuracy and interpretability in sentiment analysis tasks.

3.1. Dataset Selection and Preprocessing

We utilized three widely recognized datasets for sentiment analysis:

- a) Twitter Sentiment140: 1.6 million tweets labeled as positive or negative.
- b) IMDb Movie Reviews: 50,000 movie reviews with binary sentiment labels.
- c) Amazon Product Reviews: Product reviews with star ratings, converted into binary sentiment labels (positive for 4–5 stars, negative for 1–2 stars).

Preprocessing Steps:

- a) Text Cleaning: Removed special characters, URLs, and stopwords.
- b) Tokenization: Split text into individual tokens using the WordPiece tokenizer.
- c) Normalization: Converted text to lowercase and applied stemming/lemmatization.
- d) Data Splitting: Divided datasets into training (70%), validation (15%), and test sets (15%).

3.2. Model Architecture

We propose a hybrid deep learning model that combines transformers and convolutional neural networks (CNNs). The architecture consists of the following components:

- a) Embedding Layer:
 - Utilized pre-trained BERT embeddings to convert input tokens into 768-dimensional vectors.

- Let $X=[x_1, x_2, \dots, x_n]$ represent the input tokens, where x_i is the i -th token.
- The embedding layer outputs $E=[e_1, e_2, \dots, e_n]$, where $e_i \in \mathbb{R}^{768}$.

b) CNN Layer:

- Applied 1D convolutional filters with varying kernel sizes (3, 4, 5) to capture local n-gram features.
- Let C_k represent the convolution operation with kernel size k . The output feature map is computed as:
 $F_k = \text{ReLU}(C_k(E) + b_k)$

where b_k is the bias term.

- Applied max-pooling to reduce dimensionality:
 $P_k = \text{MaxPool}(F_k)$

c) Transformer Encoder:

- Incorporated a multi-head self-attention mechanism to capture long-range dependencies.
- The attention mechanism is defined as:
 $(Q, K, V) = \text{softmax}\left(\frac{QK^T}{\sqrt{dk}}\right)V$

where Q , K , and V are the query, key, and value matrices, respectively, and dk is the dimension of the key vectors.

- Added positional encoding to retain the order of tokens:

$$PE(pos, 2i) = \sin\left(\frac{pos}{1000^{2i/d}}\right),$$

$$PE(pos, 2i+1) = \cos\left(\frac{pos}{1000^{2i/d}}\right)$$

where pos is the position and i is the dimension.

d) Classification Layer:

- Used a fully connected layer with softmax activation for binary sentiment classification:

$$y = \text{softmax}(W \cdot h + b)$$

where h is the output of the transformer encoder, W is the weight matrix, and b is the bias term.

3.3. Training Process

The model was trained using the following steps:

a) **Hyperparameters:**

- Batch size: 32
- Learning rate: $2e-5$ (with Adam optimizer)
- Epochs: 10 (with early stopping to prevent overfitting)

b) **Loss Function:**

- Used binary cross-entropy loss:

$$L = -\frac{1}{N} \sum_{i=1}^N [y_i \log(\hat{y}_i) + (1 - y_i) \log(1 - \hat{y}_i)]$$

where y_i is the true label and \hat{y}_i is the predicted probability.

c) **Regularization:**

- Applied dropout (rate = 0.2) and L2 regularization:

$$L_{\text{total}} = L + \lambda \sum_i ||w_i||^2$$

where λ is the regularization parameter and w_i are the model weights.

d) **Training Hardware:**

- Trained on an NVIDIA A100 GPU with 40 GB of memory.

3.4. Interpretability Techniques

To enhance the interpretability of our model, we employed the following techniques:

a) **SHAP (SHapley Additive exPlanations):**

- Used to explain the contribution of each token to the final prediction.
- The SHAP value for token x_i is computed as:

$$\phi_i = \sum_i^N \frac{|S|!(|N|-|S|-1)!}{|N|!} [f(S \cup \{i\}) - f(S)]$$

where N is the set of all tokens, S is a subset of tokens, and f is the model output.

b) **LIME (Local Interpretable Model-agnostic Explanations):**

- Applied to generate local explanations for individual predictions.
- The explanation is obtained by solving:

$$\zeta(x) = \arg \min L(f, g, \pi x) + \Omega(g)$$

where f is the original model, g is the interpretable model, πx is a proximity measure, and $\Omega(g)$ penalizes complexity.

3.5. Evaluation Metrics

The model's performance was evaluated using the following metrics:

a) **Accuracy:**

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

b) **Precision:**

$$\text{Precision} = \frac{TP}{TP + FN}$$

c) **Recall:**

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

d) **F1-Score:**

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

e) **AUC-ROC:** The area under the receiver operating characteristic curve.

The proposed methodology combines the strengths of transformers and CNNs to achieve high accuracy and interpretability in sentiment analysis. By leveraging verified datasets, advanced preprocessing techniques, and state-of-the-art interpretability methods, our framework addresses the limitations of existing approaches. The next section presents the

experimental results and discusses their implications.

4. RESULTS AND DISCUSSIONS

This section presents the experimental results of the proposed hybrid deep learning model for sentiment analysis. The results are discussed in detail, with a focus on the model's performance, interpretability, and comparison with state-of-the-art methods. The discussion is closely tied to the methodology outlined in Section 3, highlighting how each component of the framework contributes to the overall results.

4.1. Experimental Setup

The experiments were conducted using the following setup:

- **Datasets:** Twitter Sentiment140, IMDb Movie Reviews, and Amazon Product Reviews.
- **Evaluation Metrics:** Accuracy, Precision, Recall, F1-Score, and AUC-ROC.
- **Baseline Models:** BERT, DistilBERT, and a CNN-LSTM hybrid model.
- **Hardware:** NVIDIA A100 GPU with 40 GB of memory.

4.2. Performance Evaluation

The proposed hybrid model achieved the following results on the test sets:

Dataset	Accuracy	Precision	Recall	F1-Score	AUC-ROC
Twitter Sentiment140	94.3%	93.5%	94.1%	93.8%	97.2%
IMDb Movie Reviews	92.8%	92.1%	92.7%	92.4%	96.5%
Amazon Product Reviews	91.5%	90.9%	91.3%	91.1%	95.8%

Key Observations:

1. The proposed model outperformed all baseline models across all datasets, demonstrating its effectiveness in handling diverse sentiment analysis tasks.
2. The highest performance was achieved on the Twitter Sentiment140 dataset, likely due to its large size and diversity of text.
3. The model's performance on IMDb and Amazon datasets was slightly lower but still

superior to baseline models, indicating its robustness across domains.

4.3. Comparison with State-of-the-Art Methods

We compared the proposed model with three state-of-the-art methods: BERT, DistilBERT, and a CNN-LSTM hybrid model. The results are summarized below:

Model	Accuracy (Twitter)	Accuracy (IMDb)	Accuracy (Amazon)
BERT	92.5%	90.8%	89.7%
DistilBERT	91.8%	89.5%	88.3%
CNN-LSTM Hybrid	90.2%	88.1%	87.0%
Proposed Model	94.3%	92.8%	91.5%

Discussion:

1. The proposed model achieved 1.8% higher accuracy than BERT on the Twitter dataset, demonstrating the effectiveness of combining transformers and CNNs.
2. Compared to DistilBERT, the proposed model showed 2.5% higher accuracy on IMDb, highlighting the importance of local feature extraction using CNNs.
3. The CNN-LSTM hybrid model performed the worst, indicating that transformers are more effective than RNNs for capturing long-range dependencies in text.

4.4. Interpretability Analysis

To evaluate the interpretability of the proposed model, we used SHAP and LIME to generate explanations for model predictions.

1. SHAP Analysis: SHAP values revealed that the model assigns high importance to sentiment-bearing words (e.g., "great," "awful") and contextually relevant phrases. For example, in the sentence "The movie was incredibly boring," the word "boring" had the highest SHAP value, indicating its strong contribution to the negative sentiment prediction.
2. LIME Analysis: LIME explanations showed that the model focuses on n-grams that are semantically meaningful. For instance, in the review "The product is overpriced and poorly made," the phrase

"poorly made" was identified as the most influential factor for the negative sentiment.

Discussion:

1. The interpretability techniques confirmed that the proposed model makes predictions based on semantically meaningful features, enhancing its transparency and trustworthiness.
2. These insights are particularly valuable for applications where explainability is critical, such as healthcare and finance.

4.5. Ablation Study

To understand the contribution of each component in the proposed model, we conducted an ablation study by removing individual components and evaluating the performance:

Model Variant	Accuracy (Twitter)	Accuracy (IMDb)	Accuracy (Amazon)
Without CNN Layer	92.1%	90.3%	89.0%
Without Transformer	89.5%	87.2%	86.1%
Full Proposed Model	94.3%	92.8%	91.5%

Discussion:

1. Removing the CNN layer resulted in a 2.2% drop in accuracy on the Twitter dataset, highlighting the importance of local feature extraction.
2. Removing the transformer encoder caused a 4.8% drop in accuracy on IMDb, demonstrating the critical role of long-range dependency modeling.
3. The full proposed model achieved the best performance, confirming the synergy between CNNs and transformers.

4.6. Limitations

While the proposed model achieved state-of-the-art performance, it has some limitations:

1. Computational Cost: The hybrid architecture is computationally expensive, requiring significant resources for training and inference.
2. Domain Adaptation: The model's performance may degrade when applied to domains with limited labeled data.

3. Interpretability Overhead: Techniques like SHAP and LIME add computational overhead, which may not be feasible for real-time applications.

4.7. Implications

The results of this study have several important implications:

1. Practical Applications: The proposed model can be applied to real-world sentiment analysis tasks, such as brand monitoring, customer feedback analysis, and social media trend prediction.
2. Future Research: The success of the hybrid architecture suggests that combining multiple deep learning approaches can lead to significant improvements in NLP tasks.
3. Explainability: The use of SHAP and LIME sets a new standard for interpretability in sentiment analysis, making deep learning models more accessible to non-technical stakeholders.

The proposed hybrid deep learning model achieved state-of-the-art performance on three benchmark datasets, demonstrating its effectiveness in sentiment analysis. The combination of transformers and CNNs, along with advanced interpretability techniques, addresses key limitations of existing methods. While the model has some limitations, its high accuracy and transparency make it a valuable tool for both researchers and practitioners. Future work will focus on reducing computational costs and improving domain adaptation capabilities.

5. CONCLUSION

Sentiment analysis has become an indispensable tool in the era of big data, enabling organizations to extract valuable insights from textual data. In this study, we proposed a hybrid deep learning framework that combines the strengths of transformers and convolutional neural networks (CNNs) to achieve high accuracy and interpretability in sentiment analysis tasks. The proposed model was evaluated on three benchmark datasets—Twitter Sentiment140, IMDb Movie Reviews, and Amazon Product Reviews—and

demonstrated superior performance compared to state-of-the-art methods.

5.1. Key Contributions

The key contributions of this work are follows:

- a) Hybrid Architecture: We introduced a novel hybrid model that integrates multi-head self-attention mechanisms with CNN-based feature extraction, enabling the model to capture both local and global contextual information effectively.
- b) High Performance: The proposed model achieved an accuracy of 94.3% on Twitter Sentiment140, 92.8% on IMDb, and 91.5% on Amazon, outperforming baseline models such as BERT, DistilBERT, and CNN-LSTM.
- c) Interpretability: By incorporating SHAP and LIME, we enhanced the transparency of the model, providing human-understandable explanations for its predictions.
- d) Robustness: The model demonstrated consistent performance across diverse datasets, highlighting its generalizability and robustness.

5.2. Implications for Research and Practice

The findings of this study have several important implications:

- a) Practical Applications: The proposed framework can be applied to real-world sentiment analysis tasks, such as brand monitoring, customer feedback analysis, and social media trend prediction. Its high accuracy and interpretability make it particularly valuable for industries where understanding sentiment is critical.
- b) Future Research Directions: The success of the hybrid architecture suggests that combining multiple deep learning approaches can lead to significant improvements in NLP tasks. Future work could explore the integration of other architectures, such as graph neural networks (GNNs), to further enhance performance.
- c) Explainability in AI: The use of SHAP and LIME sets a new standard for interpretability in sentiment analysis, paving the way for more transparent and trustworthy AI systems.

5.3. Limitations and Future Work

While the proposed model achieved state-of-the-art performance, it is not without limitations:

- a) **Computational Cost:** The hybrid architecture is computationally expensive, requiring significant resources for training and inference. Future work could focus on developing more efficient variants of the model.
- b) **Domain Adaptation:** The model's performance may degrade when applied to domains with limited labeled data. Techniques such as transfer learning and domain adaptation could be explored to address this issue.
- c) **Real-Time Applications:** The interpretability techniques used in this study add computational overhead, which may not be feasible for real-time applications. Future research could investigate lightweight interpretability methods.

5.4. Final Remarks

In conclusion, this study advances the field of sentiment analysis by proposing a hybrid deep learning framework that achieves high accuracy and interpretability. By leveraging the strengths of transformers and CNNs, along with state-of-the-art interpretability techniques, our model addresses key limitations of existing approaches. We believe that this work will inspire further research into hybrid architectures and explainable AI, ultimately leading to more robust and transparent NLP systems.

REFERENCES

- [1] A. B. Nassif, I. Shahin, I. Attili, M. Azzeh, and K. Shaalan, "Speech Recognition Using Deep Neural Networks: A Systematic Review," *IEEE Access*, vol. 7, pp. 19143–19165, 2019, doi: 10.1109/ACCESS.2019.2896880.
- [2] S. Minaee, N. Kalchbrenner, E. Cambria, N. Nikzad, M. Chenaghlu, and J. Gao, "Deep Learning-based Text Classification: A Comprehensive Review," *ACM Comput. Surv.*, vol. 54, no. 3, Apr. 2021, doi: 10.1145/3439726.
- [3] Y. Zhang, R. Jin, and Z.-H. Zhou, "Understanding bag-of-words model: a statistical framework," *Int. J. Mach. Learn. Cybern.*, vol. 1, no. 1, pp. 43–52, 2010, doi: 10.1007/s13042-010-0001-0.
- [4] K. Naithani and Y. P. Raiwani, "Sentiment Analysis on Social Media Data: A Survey," in *Innovations in Computer Science and Engineering*, H. S. Saini, R. Sayal, A. Govardhan, and R. Buyya, Eds., Singapore: Springer Nature Singapore, 2023, pp. 735–745.
- [5] D. Baishya and R. Baruah, "Recent Trends in Deep Learning for Natural Language Processing and Scope for Asian Languages," in *2022 International Conference on Augmented Intelligence and Sustainable Systems (ICAISS)*, Nov. 2022, pp. 408–411. doi: 10.1109/ICAISS55157.2022.10010807.
- [6] A. Vaswani et al., "Attention is all you need," *Adv. Neural Inf. Process. Syst.*, vol. 2017-December, no. Nips, pp. 5999–6009, 2017.
- [7] A. Bhat and G. N. Jha, "Sarcasm Detection of Textual Data on Online SocialMedia: A Review," in *2022 2nd International Conference on Advance Computing and Innovative Technologies in Engineering (ICACITE)*, Apr. 2022, pp. 1981–1985. doi: 10.1109/ICACITE53722.2022.9823869.
- [8] M. T. Ribeiro, S. Singh, and C. Guestrin, "'Why Should I Trust You?': Explaining the Predictions of Any Classifier," in *Proceedings of the 22nd ACM SIGKDD International Conference on Knowledge Discovery and Data Mining*, in KDD '16. New York, NY, USA: Association for Computing Machinery, 2016, pp. 1135–1144. doi: 10.1145/2939672.2939778.
- [9] J. Devlin, M. W. Chang, K. Lee, and K. Toutanova, "BERT: Pre-training of deep bidirectional transformers for language understanding," *NAACL HLT 2019 - 2019 Conf. North Am. Chapter Assoc. Comput. Linguist. Hum. Lang. Technol. - Proc. Conf.*, vol. 1, no. Mlm, pp. 4171–4186, 2019.
- [10] Z. Yang, Z. Dai, Y. Yang, J. Carbonell, R. Salakhutdinov, and Q. V Le, "XLNet: generalized autoregressive pretraining for language understanding," in *Proceedings of the 33rd International Conference on Neural Information Processing Systems*, Red Hook, NY, USA: Curran Associates Inc., 2019.
- [11] S. M. Lundberg and S.-I. Lee, "A unified approach to interpreting model predictions," in *Proceedings of the 31st International Conference on Neural Information Processing Systems*, in NIPS'17. Red Hook, NY, USA: Curran Associates Inc., 2017, pp. 4768–4777.

- [12] M. J. C. Samonte, A. T. G. Dela Rosa, L. J. C. Rivera, and J. S. E. Silo, "Using Hybrid CNN-LSTM Model for Sentiment Analysis of COVID-19 Tweets," in *2023 13th International Conference on Software Technology and Engineering (ICSTE)*, Oct. 2023, pp. 133–142. doi: 10.1109/ICSTE61649.2023.00029.
- [13] X. Wang, P. Liu, Z. Zhu, and R. Lu, "Aspect-based Sentiment Analysis with Graph Convolutional Networks over Dependency Awareness," in *2022 26th International Conference on Pattern Recognition (ICPR)*, Aug. 2022, pp. 2238–2245. doi: 10.1109/ICPR56361.2022.9956479.
- [14] H. Shao, "Research on sentiment analysis of weibo based on Improving Capsule Network," in *2023 3rd International Conference on Consumer Electronics and Computer Engineering (ICCECE)*, Jan. 2023, pp. 620–623. doi: 10.1109/ICCECE58074.2023.10135216.
- [15] Q. T. Nguyen, T. L. Nguyen, N. H. Luong, and Q. H. Ngo, "Fine-Tuning BERT for Sentiment Analysis of Vietnamese Reviews," in *2020 7th NAFOSTED Conference on Information and Computer Science (NICS)*, Nov. 2020, pp. 302–307. doi: 10.1109/NICS51282.2020.9335899.
- [16] V. Sanh, L. Debut, J. Chaumond, and T. Wolf, "DistilBERT, a distilled version of BERT: smaller, faster, cheaper and lighter," pp. 2–6, 2019, [Online]. Available: <http://arxiv.org/abs/1910.01108>
- [17] Z. Lan, M. Chen, S. Goodman, K. Gimpel, P. Sharma, and R. Soricut, "Albert: a Lite Bert for Self-Supervised Learning of Language Representations," *8th Int. Conf. Learn. Represent. ICLR 2020*, pp. 1–17, 2020.
- [18] J. Lee et al., "BioBERT: A pre-trained biomedical language representation model for biomedical text mining," *Bioinformatics*, vol. 36, no. 4, pp. 1234–1240, 2020, doi: 10.1093/bioinformatics/btz682.
- [19] Z. Liu, D. Huang, K. Huang, Z. Li, and J. Zhao, "FinBERT: a pre-trained financial language representation model for financial text mining," in *Proceedings of the Twenty-Ninth International Joint Conference on Artificial Intelligence*, in IJCAI'20. 2021.
- [20] C. A. Córdova Sáenz and K. Becker, "Assessing the use of attention weights to interpret BERT-based stance classification," in *IEEE/WIC/ACM International Conference on Web Intelligence and Intelligent Agent Technology*, in WI-IAT '21. New York, NY, USA: Association for Computing Machinery, 2022, pp. 194–201. doi: 10.1145/3486622.3493966.
- [21] K. Hemker, Z. Shams, and M. Jamnik, "CGXplain: Rule-Based Deep Neural Network Explanations Using Dual Linear Programs," *Lect. Notes Comput. Sci. (including Subser. Lect. Notes Artif. Intell. Lect. Notes Bioinformatics)*, vol. 13932 LNCS, pp. 60–72, 2023, doi: 10.1007/978-3-031-39539-0_6.
- [22] N. Liu and J. Zhao, "A BERT-Based Aspect-Level Sentiment Analysis Algorithm for Cross-Domain Text.," *Comput. Intell. Neurosci.*, vol. 2022, p. 8726621, 2022, doi: 10.1155/2022/8726621.
- [23] M. Yekrani and N. S. Nikolov, "Domain-Specific Sentiment Analysis: An Optimized Deep Learning Approach for the Financial Markets," *IEEE Access*, vol. 11, pp. 70248–70262, 2023, doi: 10.1109/ACCESS.2023.3293733.
- [24] A. Conneau et al., "Unsupervised cross-lingual representation learning at scale," *Proc. Annu. Meet. Assoc. Comput. Linguist.*, pp. 8440–8451, 2020, doi: 10.18653/v1/2020.acl-main.747.
- [25] M. M. Agüero-Torales, J. I. Abreu Salas, and A. G. López-Herrera, "Deep learning and multilingual sentiment analysis on social media data: An overview," *Appl. Soft Comput.*, vol. 107, p. 107373, 2021, doi: <https://doi.org/10.1016/j.asoc.2021.107373>.