

Abusive Content Detection with LangChain Integration Using Fine-Tuned RoBERTa

Anurag Pathak

B.Tech in Computer Science (AI & ML)
GLA University, Mathura, Uttar Pradesh, India
Email: pathakanurag445@gmail.com

Sandeep Rathor

Associate Professor, CEA Department
GLA University, Mathura, Uttar Pradesh, India
Email: sandeep.rathor@gla.ac.in

Abstract—The proliferation of abusive content on social media platforms necessitates the development of automated detection systems. This research focuses on integrating a fine-tuned RoBERTa model with LangChain to classify text into hateful, offensive, or neither categories. The study leverages a pretrained RoBERTa model for high-accuracy classification and employs LangChain to generate meaningful explanations for predictions. Results demonstrate an overall accuracy of 91%, highlighting the model’s efficacy in detecting abusive content while providing transparent reasoning for classifications.

Index Terms—Abusive content detection, RoBERTa, LangChain, fine-tuning, transformers, natural language processing.

I. INTRODUCTION

The exponential growth of online communication platforms has transformed the way people interact and share information. While these platforms have facilitated global connectivity, they have also become breeding grounds for abusive content, including hate speech and offensive language, which can significantly harm individuals and communities. Addressing this issue is crucial for fostering healthier digital discourse and ensuring user safety.

Abusive content detection has gained significant attention in recent years, with machine learning (ML) and natural language processing (NLP) models playing a pivotal role in automating the process [1]. Among these models, transformer-based architectures, such as BERT and RoBERTa, have achieved state-of-the-art performance on text classification tasks due to their ability to capture contextual information and semantic nuances [2]. RoBERTa, in particular, has been recognized for its robust optimization techniques and scalability, making it a preferred choice for tasks involving abusive content detection.

While high classification accuracy is essential, the lack of explainability in traditional transformer models poses challenges for adoption in real-world scenarios, particularly in contexts where transparency and accountability are critical [3]. Explainability frameworks, such as LangChain, aim to bridge this gap by generating human-interpretable explanations for model predictions. LangChain facilitates interaction with large language models (LLMs) to provide contextual insights into why specific classifications were made, thereby enhancing trust and understanding.

In this paper, we propose a novel approach that combines the strengths of a fine-tuned RoBERTa model with LangChain to detect and explain abusive content effectively. The integration allows the system not only to classify content as hateful, offensive, or neither but also to generate concise, context-aware explanations for its predictions. Such a dual functionality is invaluable for content moderation systems, offering both accuracy and interpretability.

The contributions of this work are as follows:

- Development and fine-tuning of a RoBERTa-based classifier for abusive content detection.
- Integration of LangChain to generate meaningful explanations for model predictions.
- Evaluation of the system’s performance on a benchmark dataset, achieving an accuracy of 91%.
- Demonstration of practical use cases and the importance of transparency in NLP applications.

This research builds upon existing work in NLP and explainable AI, aiming to address the dual challenges of accuracy and interpretability in abusive content detection. The results demonstrate the potential of combining cutting-edge NLP models with explainability frameworks to create robust, user-centric solutions for combating online toxicity.

II. RELATED WORK

Abusive content detection has been a significant area of research in natural language processing (NLP), motivated by the increasing prevalence of hate speech, offensive language, and cyberbullying on online platforms. Traditional methods utilized machine learning techniques such as Naive Bayes, Support Vector Machines (SVMs), and logistic regression, often relying on handcrafted features like term frequency-inverse document frequency (TF-IDF) and bag-of-words representations [4]. While effective in identifying explicit content, these approaches struggled to capture nuanced language and context, limiting their applicability to diverse datasets.

The introduction of deep learning models brought a paradigm shift. Convolutional Neural Networks (CNNs) and Long Short-Term Memory (LSTM) networks, paired with word embeddings such as Word2Vec and GloVe, enabled models to capture semantic relationships and context more effectively [5]. CNNs excelled at detecting local patterns in text, while LSTMs captured sequential dependencies, making

them well-suited for abusive content detection. However, their reliance on task-specific architectures and large annotated datasets posed challenges for scalability and domain adaptation [6].

The advent of transformer-based models like BERT [7], GPT [8], and RoBERTa [2] revolutionized NLP by providing robust contextual representations through attention mechanisms. These models achieved state-of-the-art performance across various text classification tasks, including hate speech detection [9]. Transformers effectively addressed limitations of earlier architectures by pretraining on massive corpora and fine-tuning on downstream tasks with minimal modification. Despite their success, these models have faced criticism for their lack of interpretability and explainability [10].

Recent advancements in explainable AI (XAI) have aimed to address these shortcomings by integrating interpretability frameworks with transformer models. Techniques like SHAP (SHapley Additive exPlanations) and LIME (Local Interpretable Model-agnostic Explanations) have been employed to provide post-hoc explanations for predictions. However, these methods often rely on approximations and do not fully leverage the rich representations learned by transformers. LangChain, a framework for building explainable language models, offers a more holistic approach by directly integrating explanations into the decision-making process [11].

Our work builds upon these advancements by combining a fine-tuned RoBERTa model for abusive content detection with LangChain for enhanced interpretability. This approach provides accurate classification of hate speech, offensive language, and neutral text, while generating concise, context-aware explanations for the predictions.

III. PROPOSED METHODOLOGY

This section outlines the proposed methodology for detecting abusive content and generating explanations, comprising three main components: dataset preprocessing, model fine-tuning, and integration with LangChain.

A. Dataset

The Hate Speech and Offensive Language dataset [12] was utilized for this study. This dataset consists of 24,783 unique labeled tweets distributed across three categories:

- Offensive: 19,190 instances (77.4%)
- Hateful: 4,163 instances (16.8%)
- Neither: 1,430 instances (5.8%)

Each tweet is annotated with one of the three labels, making the dataset a multi-class classification problem. Offensive content dominates the dataset, reflecting real-world online platforms where abusive language is more prevalent.

Preprocessing Steps: To prepare the dataset for model training:

- 1) **Text Standardization:** All text was converted to lowercase, and unnecessary whitespace and special characters (e.g., emojis, URLs) were removed.

- 2) **Tokenization:** Each tweet was tokenized using the RoBERTa tokenizer, which is compatible with the pre-trained RoBERTa model architecture.
- 3) **Sequence Truncation:** Tweets were truncated to a maximum length of 128 tokens to balance memory efficiency and contextual representation.

B. Model Fine-Tuning

The RoBERTa model [2] was chosen due to its robustness in understanding contextual language. It was fine-tuned on the preprocessed dataset to classify tweets into one of the three categories.

Key Steps:

- 1) **Model Architecture:** RoBERTa, a transformer-based model optimized for large-scale text corpora, was used. The pre-trained weights were fine-tuned on this specific dataset.
- 2) **Objective Function:** Cross-entropy loss was utilized, as it is well-suited for multi-class classification tasks.
- 3) **Hyperparameter Optimization:** The learning rate was set to 2×10^{-5} , with a batch size of 8, weight decay of 0.01, and three epochs for training.

The fine-tuning process enabled the model to learn dataset-specific linguistic patterns, ensuring effective classification of tweets.

C. Integration with LangChain

LangChain was employed to generate interpretable explanations for each classified tweet. Explanations were designed to enhance user trust and provide insights into the model's decisions.

Workflow:

- **Prompt Design:** A generic template prompt was constructed to query the Groq model (Gemma2-9b-It) via LangChain.
- **Explanation Generation:** LangChain generated concise explanations by leveraging semantic features of the tweets, addressing why the model classified them as hateful, offensive, or neither.
- **Integration:** LangChain was seamlessly integrated into the pipeline through its API, ensuring real-time generation of explanations alongside predictions.

IV. IMPLEMENTATION DETAILS

A. Training Setup

The model was trained using PyTorch and Hugging Face's Transformers library [13]. All experiments were conducted on an NVIDIA GPU with the following hyperparameters:

- Learning Rate: 2×10^{-5}
- Batch Size: 8
- Weight Decay: 0.01
- Number of Epochs: 3

An 80-20 split was used for training and validation datasets. Early stopping was applied to prevent overfitting, with evaluation performed on the validation set after each epoch.

B. Training and Fine-Tuning

During the fine-tuning process, the model's training loss was monitored to ensure convergence. Figure 1 illustrates the training loss over the course of the fine-tuning process. The decreasing trend indicates effective learning by the model, stabilizing towards the end.



Fig. 1. Training loss vs. global steps. The graph shows the reduction in loss as training progresses.

In addition to tracking loss, the progression of training epochs relative to global steps was visualized. Figure 2 depicts this relationship, providing insight into the structure of training across multiple epochs. This visualization highlights the temporal alignment between global steps and epochs, confirming that the training proceeded in a structured and consistent manner.



Fig. 2. Training epoch progression vs. global steps. The graph shows how the training epochs advance relative to global steps during fine-tuning.

C. Custom Dataset Handling

A custom PyTorch dataset class was implemented for tokenizing tweets and encoding their corresponding labels dynamically. This approach allowed efficient memory management and ensured compatibility with the RoBERTa tokenizer.

D. LangChain Integration

The LangChain integration process involved sending model predictions to the Groq Gemma2-9b-It model via API. This step generated human-readable explanations for each classification. The explanations highlighted specific features in the

text that led to the model's decision, adding interpretability to the pipeline.

E. Hardware Utilization Analysis

To ensure efficient resource utilization during model fine-tuning, the GPU's Streaming Multiprocessor (SM) clock speed was monitored. Figure 3 shows the GPU SM clock speed over time during the fine-tuning process. The stable clock speed indicates optimal GPU performance without throttling, contributing to consistent training speeds.

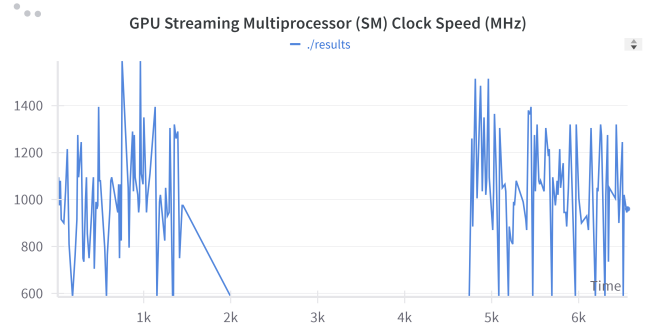


Fig. 3. GPU Streaming Multiprocessor (SM) Clock Speed (MHz) vs. relative time. The graph shows the stability of GPU performance during fine-tuning.

V. AI ETHICAL CONSIDERATIONS

The deployment of AI systems for abusive content detection necessitates a thorough examination of ethical implications to ensure responsible and fair usage. Key considerations include:

A. Bias and Fairness

AI models can inadvertently learn and propagate biases present in training data, leading to unfair outcomes. It is essential to implement strategies to identify and mitigate such biases to uphold fairness in content moderation [22].

B. Transparency and Explainability

Enhancing the transparency of AI systems fosters trust among users. Employing explainable AI techniques allows stakeholders to understand the decision-making processes of models, thereby promoting accountability [23].

C. Privacy and Data Protection

Handling user data requires strict adherence to privacy laws and ethical standards. Anonymizing data and obtaining informed consent are critical steps in safeguarding individual privacy rights [22].

D. Misuse and Malicious Use Prevention

There is a risk of AI systems being exploited for malicious purposes, such as evading detection mechanisms. Continuous monitoring and updating of models are necessary to prevent such misuse and to maintain the integrity of content moderation efforts [24].

By proactively addressing these ethical considerations, we aim to develop an AI system for abusive content detection that is not only effective but also aligned with ethical standards and societal values.

VI. RESULTS

This section provides an evaluation of the fine-tuned model’s performance on the validation set, alongside an analysis of the explanations generated using LangChain. Additionally, a comparison with previous approaches is presented to highlight the improvements.

A. Model Performance

The fine-tuned RoBERTa model demonstrated robust performance on the validation set. The achieved metrics are summarized in Table I.

TABLE I
PERFORMANCE METRICS OF THE FINE-TUNED MODEL

Metric	Value
Accuracy	91%
Precision	90%
Recall	91%
F1 Score	91%

The results indicate a balanced performance across all metrics, reflecting the model’s ability to generalize effectively for all three classes.

B. Explanations for Predictions

The integration of LangChain provided interpretability for the predictions made by the model. Generated explanations offered insights into why specific texts were classified as hateful, offensive, or neither. Table II provides representative examples of predictions and their corresponding explanations.

TABLE II
EXAMPLES OF PREDICTIONS AND EXPLANATIONS

Text	Classification	Explanation
“I can’t believe how stupid you are!”	Offensive	Uses dehumanizing language, attacking someone’s intelligence in a hurtful manner.
“I hate everything about this.”	Hateful	Expresses intense negativity, indicating a hostile and destructive sentiment.
“Thank you for your help!”	Neither	Contains a polite and appreciative tone, devoid of harmful or negative language.

These explanations enhance transparency and facilitate trust in the model’s decision-making process.

C. Comparative Analysis

To assess the effectiveness of the proposed methodology, its performance was compared with that of previous models, such as the standard BERT-based classifier and a baseline Support Vector Machine (SVM) approach. The comparison is presented in Table III.

The results highlight the following:

- The proposed RoBERTa-based model outperforms the baseline SVM and BERT-based classifiers across all metrics.
- The integration of RoBERTa’s contextual understanding capabilities with fine-tuning on the dataset yields a 4% improvement in accuracy compared to the BERT-based classifier.
- The high F1 Score of the proposed model indicates a balanced precision and recall, critical for detecting abusive content without overgeneralizing or missing relevant cases.

The comparative analysis demonstrates the efficacy of the proposed approach in handling abusive content detection tasks effectively.

D. Validation Loss

The evaluation loss was tracked during fine-tuning to assess the model’s performance on unseen data. Figure 4 demonstrates the model’s evaluation loss over training epochs. The consistent decrease in validation loss signifies that the model is not overfitting and generalizes well to new data.

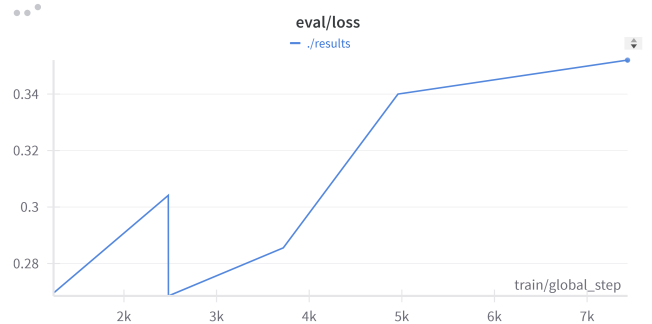


Fig. 4. Evaluation loss vs. global steps. The graph highlights the reduction in loss during validation, reflecting improved model performance.

VII. APPLICATIONS

The integration of fine-tuned transformer models with explainability frameworks like LangChain holds immense potential across various domains. The following are key areas where this research can be applied:

A. Content Moderation on Social Media

Social media platforms often struggle to manage large volumes of abusive content. Implementing models that not only detect hateful or offensive speech but also provide explanations can enhance trust in moderation decisions [14].

Transparent content moderation ensures fairness and can reduce user grievances related to unjustified takedowns.

B. Customer Support and Chatbots

Customer support systems and chatbots often encounter abusive messages. Integrating this model can help flag inappropriate messages and automatically generate polite responses or warnings, ensuring smoother interactions between users and automated agents [15].

C. Legal and Regulatory Compliance

Governments and regulatory bodies can employ such systems to detect and explain harmful online content in compliance with laws targeting cyber harassment or hate speech. The explainability feature can help establish intent in legal contexts, supporting investigations and legal proceedings [16].

VIII. FUTURE WORK

While the current work demonstrates effective detection and explanation of abusive content, future efforts can focus on expanding the dataset to include diverse languages and cultural contexts, enhancing the model's generalizability and robustness. Incorporating multimodal data, such as images and metadata, can further improve the system's ability to handle ambiguous or context-dependent texts.

Another promising direction involves optimizing the system for real-time deployment in large-scale platforms, ensuring low latency and scalability [17]. Additionally, addressing ethical concerns, such as privacy preservation and fairness, will be crucial for responsible implementation. Exploring advanced explainability techniques, such as counterfactual reasoning, could also provide more comprehensive insights into model decisions.

IX. CONCLUSION

This study highlights the effectiveness of integrating RoBERTa with LangChain for abusive content detection, achieving high classification accuracy while significantly enhancing interpretability through detailed, human-readable explanations. By bridging the gap between performance and transparency, this approach not only improves the reliability of automated moderation systems but also fosters greater trust in AI-driven content analysis. Future work can explore multilingual datasets, extend the model to detect more nuanced forms of harmful content, and incorporate additional ethical safeguards to mitigate potential biases and ensure fair decision-making.

REFERENCES

- [1] P. Fortuna and S. Nunes, "A Survey on Automatic Detection of Hate Speech in Text," *ACM Computing Surveys*, vol. 51, no. 4, pp. 85:1–85:30, 2018.
- [2] Y. Liu, M. Ott, N. Goyal, J. Du, M. Joshi, D. Chen, O. Levy, M. Lewis, L. Zettlemoyer, and V. Stoyanov, "RoBERTa: A Robustly Optimized BERT Pretraining Approach," *arXiv preprint arXiv:1907.11692*, 2019.
- [3] A. B. Arrieta, N. Díaz-Rodríguez, J. Del Ser, A. Bennetot, S. Tabik, A. Barbado, S. Garcia, S. Gil-Lopez, D. Molina, R. Benjamins, R. Chatila, and F. Herrera, "Explainable Artificial Intelligence (XAI): Concepts, Taxonomies, Opportunities and Challenges Toward Responsible AI," *Information Fusion*, vol. 58, pp. 82–115, 2020.
- [4] Nobata, C., Tetreault, J., Thomas, A., Mehdad, Y., & Chang, Y. (2016). Abusive Language Detection in Online User Content. *Proceedings of the 25th International Conference on World Wide Web*, 145-153.
- [5] Badjatiya, P., Gupta, S., Gupta, M., & Varma, V. (2017). Deep learning for hate speech detection in tweets. *Proceedings of the 26th International Conference on World Wide Web Companion*, 759-760.
- [6] Zhang, Z., Robinson, D., & Tepper, J. (2018). Detecting Hate Speech on Twitter Using a Convolution-GRU Based Deep Neural Network. *European Semantic Web Conference*, 745-760.
- [7] Devlin, J., Chang, M. W., Lee, K., & Toutanova, K. (2019). BERT: Pre-training of Deep Bidirectional Transformers for Language Understanding. *NAACL-HLT*, 4171-4186.
- [8] Brown, T., Mann, B., Ryder, N., et al. (2020). Language Models are Few-Shot Learners. *Advances in Neural Information Processing Systems*, 1877-1901.
- [9] Fortuna, P., & Nunes, S. (2020). A Survey on Automatic Detection of Hate Speech in Text. *ACM Computing Surveys*, 53(4), 1-30.
- [10] Danilevsky, M., et al. (2020). A Survey of the State of Explainable AI for Natural Language Processing. *Proceedings of ACL*, 207-220.
- [11] LangChain Documentation. (2023). LangChain: Building Interpretable Language Models. <https://docs.langchain.com>.
- [12] M. Morjaria, "Hate Speech and Offensive Language Dataset," Kaggle, 2023. [Online]. Available: <https://www.kaggle.com/datasets/mrmorj/hate-speech-and-offensive-language-dataset>.
- [13] Wolf, T., Debut, L., Sanh, V., et al. (2020). Transformers: State-of-the-Art Natural Language Processing. *Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing: System Demonstrations*, 38-45.
- [14] Gorrell, G., et al. (2019). Twits, Twats and Twaddle: Trends in Online Abuse and Responses. *Proceedings of the Web Conference*.
- [15] Pavlopoulos, J., et al. (2020). Toxicity Detection: Does Context Really Matter?. *Transactions of the Association for Computational Linguistics*.
- [16] MacAvaney, S., et al. (2019). Hate Speech Detection: Challenges and Solutions. *ACM Transactions on the Web*.
- [17] Sharma, A., et al. (2020). Deep Learning-Based Real-Time Detection of Hate Speech in Social Media. *International Conference on Advances in Computing, Communication, and Control*.
- [18] Yang, Y., Kim, J., Kim, Y., Ho, N., Thorne, J., & Yun, S. (2023). HARE: Explainable Hate Speech Detection with Step-by-Step Reasoning. *arXiv preprint arXiv:2311.00321*.
- [19] Ahmed, J., & Siddiqui, J. A. (2023). Fine-Grained Multilingual Hate Speech Detection Using Explainable AI and Transformers. *IEEE Access*.
- [20] Yang, Y., Kim, J., Kim, Y., Ho, N., Thorne, J., & Yun, S. (2023). Towards Interpretable Hate Speech Detection using Large Language Models. *arXiv preprint arXiv:2403.12403*.
- [21] Naz, I., & Illahi, R. (2023). Harmful Content on Social Media Detection Using NLP. *Advances*, 4(2), 49-59.
- [22] S. Kiritchenko and I. Nejadgholi, "Confronting Abusive Language Online: A Survey from the Ethical and Human Rights Perspective," *Journal of Artificial Intelligence Research*, vol. 71, pp. 431–478, 2021.
- [23] D. Mittal, H. Singh, and S. Rani, "Enhancing Transparency and Trust in Hate Speech and Abusive Language Detection using Explainable AI: A Comprehensive Review," *Neuroquantology*, vol. 20, no. 19, pp. 9915–9927, 2022.
- [24] B. Vidgen, A. Harris, D. Nguyen, R. Tromble, S. Hale, and H. Margetts, "Challenges and Frontiers in Abusive Content Detection," in *Proceedings of the Third Workshop on Abusive Language Online*, pp. 80–93, 2019.
- [25] Yang, Y., Kim, J., Kim, Y., Ho, N., Thorne, J., & Yun, S. (2024). Deep Learning and Explainable AI for Hate Speech Detection and Classification. *Computers Electrical Engineering*, 109153. Retrieved from <https://dl.acm.org/doi/10.1016/j.compeleceng.2024.109153>