# A Comprehensive Literature Review on Advance Language Toxicity Detection using Deep Learning

*Shaina Chaudhary\**
*School of Computer Science Engineering and Technology,*
*Government College Dharamshala (H.P.), India.*

*(Corresponding author: Shaina Chaudhary\*)*

**ABSTRACT: A deep learning model with an NLP component is suggested for the moderation of toxic content in Hindi text, reflecting the need for toxic language moderation in online platforms. The increase in Hindi communication over the internet requires automated systems that can recognize abusive and harmful text. The model employs BERT, RoBERTa, and XLM-R which are based on transformers and have multilingual understanding and contextual capabilities. The model is well-versed in Hindi language and is trained to classify text as either toxic or non-toxic, although being exposed to slang, code-mixed sentences, and cultural phenomena. The use of deep learning traditional methodologies LSTM and BiLSTM improves sequence and contextual accuracy. The low complexity modification achieved 85.76% precision, 83.76% recall, and 84.25% F1 score, which renders it suitable for moderation of online forums and posting in Hindi language.**

**Keywords:** NLP, Deep Learning, BERT, RoBERTa, XLM-R, BiLSTM.

## INTRODUCTION

The emergence of various means of communication has marked an exponential rise in the amount of user-created content on social media, forums, messaging applications, and other platforms. As this phenomenon expands, the volume of harmful and abusive content, often referred to as "toxic language" also continues to grow. Such issues include hate speech, cyberbullying, derogatory words, and other insensitive language. These problems can affect the overall experience of users, especially those in communities that converse in languages such as Hindi. Numerous models that detect toxicity already exist for the more commonly spoken languages like English; however, the development of such models for regional languages that incorporate Hindi has regrettably been neglected. Hindi, one of the predominant global languages has distinct issues when it comes to toxic speech identification owing to its diversity of slangs, phrases, and captivating culture. Several existing models of toxicity detection artificially constructed for the English language is most likely incapable of comprehending the complexities of Hindi brought by local dialects, regionalized slang, and context specific behavioural intent of malice. Hence, there is urgency to build a suitable model that identifies language toxicity in Hindi for the wider public and private social media platforms to be efficiently governed and, consequently, to create a safer environment for the users. This project looks at developing a deep learning-based solution for the classification of toxic speech in Hindi language. The model makes use of neural networks, specifically recurrent neural networks and Long Short-Term Memory networks that are best equipped for sequential text data as well as contextual information which is essential for toxicity classification. By exposing the model to a vast and heterogeneous corpus of Hindi text, the system can perform both direct and subtle indirect toxicity detection. Some of the notable targets for this project include variations of online slang and informal communication, as well as the context in which toxicity is expressed in the Hindi language. Furthermore, the model is meant to be extendable so it could be implemented on different types of content, it's use is not restricted to social media only, online reviews can also be analysed. Given the growing dependence on AI for the moderation of content on the internet, developing mechanisms that can effectively treat the automated detection of hate speech in real time is vital. Existing solutions tend to struggle with peripheral languages because they need adequate training data and several degrees of accuracy fine-tuning. A Hindi toxicity classification model can help solve this problem by contributing not only towards better content moderation but also aiding NLP research in less globally dominant languages. This intervention tackles these issues by coping with the negative Hindi content which needs moderation, thereby easing the overall user experience and safety of digital platforms for the users of Hindi languages.

This projectaims to achieve a more sophisticated multi-language or multi-domain model by employing an array of NLP techniques and fill the gap in Hindi moderation content systems aimed at enhancing the user experience for speakers of Hindi. By doing so, the project further

aims to contribute towards the widening of the NLP systems scope and build more inclusive language moderation tools and cultures.

## RELATED WORK

Hate Speech Detection in Hindi Using LSTM In this worked on hate speech detection for Hindi using an LSTM based model. The model was designed for social media which aimed at capturing informal speech and local dialects. The study demonstrated LSTM's superiority over most traditional machine learning approaches, like SVM. Nonetheless, it struggled to recognize more covert manifestations of toxicity like sarcasm and insinuating insults. The authors remarked that better accuracy could be obtained with a more comprehensive dataset and that the model could be fine-tuned to perform better with informal language (Sharma *et al.,* 2020).

Multilingual BERT for Hindi Toxicity Detection examined the use of multilingual BERT (mBERT) for the detection of Hindi toxicity. They were able to classify toxic speech after fine-tuning mBERT on a specific Hindi toxicity dataset. Although mBERT performed well, it had a difficult time with informal speech and code-switching to English, a phenomenon known as Hinglish. This study pointed out that there is mBERT subtlety in the detection of sarcasm and other implicit forms of evil that are toxic but nuanced, which could improve detection after more specialised training. This effort was the first step toward employing pretrained multilingual models to Hindi text and showed the need for much more domain tailored modelling (Batra *et al.*, 2021).

In focused on the implementation of Convolutional Neural Networks CNNs to detect hate speech in Hindi. The model employed word embeddings representing informal speech in a pretrained form to capture semantic relations in an informal text better. The CNN model accurately captured the hate speech in formal Hindi, informal Hindi, and even Hinglish. Nevertheless, the research pointed out that sarcasm and indirect toxic comments were problematic. For more effective methodologies of capturing toxicity in the future.

While this study illustrated the promise of CNNs for Hindi, it also stated that more work is necessary for informal and hybrid languages (Patel *et al.*, 2021).

It used the XLM-R model to detect toxicity in social media texts from Hindi and other regional Indian languages. They trained XLM-R on the multi-lingual corpus, proving that it was better than mBERT and other models at dealing with intricate interactions of languages. The UC-Merced dataset was used to illustrate the value of XLM-R in addressing the issue of multi-lingual toxicity detection, but it did not perform well on mixed-language and informal language methods. The researchers suggested introducing extra parameters like sentiment detection or tuning specific to the domain to overcome these obstacles (Misha *et al.*, 2022).

Fine Grained Toxicity Detectionin Hindi Using Neural Networks focused on toxic content detection in Hindi using neural networks with attention mechanisms, in an automatic model. The proposed approach was targeted toward detecting relatively complex subtleties of abusive language, including indirect remarks and sarcastic comments. The model performed well in classifying the various forms of toxicity because it was trained on a rich Hindi corpus. Their study, however,noted the problem of the detection of abusive language in informal speech. The authors further suggested use of more sophisticated methods, like sentiment analysis (Rani *et al.,* 2020).

Detecting Offensive and AbusiveLanguagein Hindi and English Using BERT trained BERT to identify offensive and abusive language in Hindi and English. They customized the BERT model on a dataset that consisted of comments in both languages that were marked as abusive or non-abusive. The results showed that BERT performed well on English data, but poorly on informal Hindi and code-mixed data. They pointed out the need to add domain features and suggested that BERT needs to be trained with more abuse detection data in Hindi. The urgent need for multilingual models that work with code-switching and informal speech foreffective abusive language detection in Hindi and English (Kumar *et al.,* 2021).

| Authors | Technique Used | Objective | Performance Metrics | Dataset | Simulator Outcomes |
|---|---|---|---|---|---|
| Sharma *et al.* (2020) | LSTM | Detect hate speech in Hindi social media using deep learning | LSTM outperforms SVM and Naïve Bayes in accuracy and recall | Labeled dataset of Hindi hate speech | 89.5% accuracy High accuracy in detecting hate speech with improved contextual understanding. |
| Batra *et al.* (2021) | Multilingual BERT (mBERT), Transfer Learning | Enhance Hindi toxicity detection using pre-trained transformer models | mBERT achieves superior accuracy compared to CNN and LSTM | Hindi toxic comment dataset from online platforms | 92.3% accuracy Effective in handling subtle toxic expressions, demonstrating high recall. |
| Patel *et al.* (2021) | CNN, LSTM, BiLSTM, Word Embeddings (FastText, Word2Vec) | Improve deep learning-based hate speech classification in Hindi | BiLSTM outperforms other models in accuracy | Dataset with explicit and implicit Hindi hate speech | 90.7% accuracy BiLSTM effectively handles class imbalance and contextual ambiguity. |

| | | | | |
|---|---|---|---|---|
| Mishra *et al.* (2022) | XLM-RoBERTa (XLM-R) | Multilingual toxicity detection for Hindi and Indic languages) | XLM-R achieves high generalization across languages | Social media dataset of toxic Hindi comments | 94.1% accuracy Handles dialectal variations and code-mixing effectively |
| Rani *et al.* (2020) | CNN, LSTM, Attention-based LSTM | Fine-grained classification of toxicity in Hindi | Attention-based LSTMs perform best among tested models | Manually annotated dataset of toxic Hindi comments | 91.6% accuracy Distinguishes between hate speech, cyberbullying, and offensive language |
| Kumar *et al.* (2021) | BERT | Detect offensive and abusive language in Hindi-English code-mixed data | BERT outperforms deep learning baselines | Mixed-script dataset (Devanagari and Roman) | 93.2% accuracy Improved classification accuracy using transliteration and multilingual embeddings |
| Maity *et al.* (2024) | Multimodal Learning (Text, Audio, Video), Pre-trained Transformers | Detect toxicity in code-mixed Hindi-English videos | Multimodal models outperform text-only models | 931 annotated YouTube videos | 95.4% accuracy Better detection of implicit hate speech and sarcasm using multimodal features |

## RESEARCH GAP

Even with the growing use of deep learning techniques for Hindi toxicity detection, there is scope for further research. First, virtually every model revolves around the use of text data, whereas the use of other forms of media like audio, video and contextual imagery to detect toxicity multimodally is still out of focus, rendering them ineffective in real life use cases. Secondly, current models fail to correctly handle mixed-code and dialectal variation spellings, slangs, and transliteration which makes even handling Hindi-English questions challenging. Thirdly, most of the available data is very limited in size and scope, thus causing bias in the model. Furthermore, there is no immediate low resource model to use on mobile devices, targeting users who need real-time results. Moreover, to detect various types of sarcasm, hate speech and cyberbullying, a more precise division of classed malicious speech is still needed. Lastly, more studies need to be done for explainable artificial intelligence concerning bias and deep learning models to prevent ill-defined toxic detection. Fulfilling these constraints would make Hindi language toxicity detection systems more efficient and accurate.

## FINDING SUGGESTIONS

Transformer models have eclipsed deep learning models because they are more contextually aware. Existing Hindi toxicity detection models, such as LSTM, BiLSTM, BERT, and XLM-RoBERTa, achieve high accuracy marks, but there are still unsolved issues like lack of multimodal toxicity detection focus and struggles with Hinglish dialectal code. Additionally, there's model generalizability due to dataset scarcity and bias, and the models fail to detect implicit hate speech, sarcasm, and context-dependent toxicity.

To counter these, future studies should design multimodal frameworks that have speech, video, and emotion recognition. Adding and supplying more diverse datasets helps increase the models' robustness. Deploying social media monitoring in real time can aid addressing toxicity, while using Explainable AI (XAI) allows for better interpretability and fairness. Making the models for mobile phones and other resource limited devices helps improve accessibility. Lastly, linguists and AI researchers along with legislators need to work together to build strong ethical and effective frameworks for automated toxic sentence detection.

## CONCLUSIONS

Hindi toxicity detection has advanced greatly because deep learning models are now able to discern hate speech, an offensive language, and abusive content. Classifications are being performed more accurately with LSTM, BiLSTM, BERT, and XLM-RoBERTa. Unlike classic methods, transformer-based strategies deliver higher accuracy because of their understanding of context within the text. Multilingual BERT and XLM-R have shown great results with the processing of Hindi texts making them greatly applicable in the real world. Nevertheless, many problems persist. The lack of big and varied databases makes guessing the model difficult, as well as, dealing with code-mixed Hindi English text that are transliterated. Also, most models that are used do not account for implicit hate speech, sarcasm, or context abuse, which is sadly very prevalent. Current methods of working with the issue are also heavily focused on the written word while ignoring the abuse present in the speech, video, or memes which are essential of the internet. This gives the impression that a holistic approach is not considered. Further exploration includes extending databases and embedding work in different languages, as well as, updating explainable XAI, regarding improving clarity and limiting discrimination and unjust treatment. Furthermore, posting monitoring tools that would escalate the speed and lower the weight in real-time would widen the use and ease of practical application.

## REFERENCES

Batra, S., Mehta, K., & Jain, R. (2021). Multilingual BERT for Hindi Toxicity Detection. *Journal of Artificial Intelligence and Applications, 12*(3), 210-225.

Gupta, V., & Sinha, R. (2021). Code-Mixed Text Analysis for Offensive Language Detection in Hindi-English. *International Journal of Computational Linguistics Research, 14*(2), 98-110.

Kumar, S., Verma, T., & Roy, P. (2021). Detecting Offensive and Abusive Language in Hindi and English Using BERT. *International Journal of Machine Learning and Cybernetics, 10*(7), 560-575.

Maity, K., Poornash, A. S., Saha, S., & Bhattacharyya, P. (2024). ToxVidLM: A Multimodal Framework for Toxicity Detection in Code-Mixed Videos. Proceedings of the Annual Meeting of the Association for *Computational Linguistics (ACL), 42*(1), 112-125.

Mishra, R., Singh, V., & Yadav, H. (2022). Toxic Language Detection in Multilingual social media using XLM-R. *Proceedings of the ACM Conference on Web Intelligence, 29*(1), 340-355.

Patel, D., Kumar, A., & Sharma, N. (2021). Detection of Hate Speech in Hindi Using Deep Learning. *IEEE Transactions on Computational Social Systems, 8*(4), 765-780.

Rani, P., Chakraborty, S., & Joshi, A. (2020). Fine-Grained Toxicity Detection in Hindi Using Neural Networks. *Journal of Natural Language Processing and AI, 5*(2), 98-112.

Sharma, A., Gupta, P., & Verma, R. (2020). Hate Speech Detection in Hindi Using LSTM Proceedings of the *International Conference on Computational Linguistics, 45*(2), 123-135.