

An Intelligent System for Real-Time Detection and Mitigation of Misinformation on Social Media Networks

¹ Ms.Hoyala jataboina, ² M.karthika, ³ K.Ramya sri, ⁴ K.Laxmi prasanna

¹ Assistant Professor, Department of Computer Science & Engineering (Artificial Intelligence & Machine Learning), Malla Reddy Engineering College for Women(Autonomous), Hyderabad, Telangana, India,

¹ Email : hoyala@lords.ac.in

^{2,3,4} Students, Department of Computer Science & Engineering (Artificial Intelligence & Machine Learning), Malla Reddy Engineering College for Women(Autonomous), Hyderabad, Telangana, India,²

Email : karthikamanisha09@gmail.com³ Email: ramyasrikarni@gmail.com, ⁴ Email:

kamarapuprasanna@gmail.com

Abstract

The rapid expansion of social media platforms has accelerated the spread of misinformation, posing significant risks to public safety, social stability, and information integrity. Conventional rule-based and manual moderation approaches are inadequate due to the high volume, velocity, and evolving nature of misleading content. This study proposes an intelligent real-time system that integrates machine learning, natural language processing, and network analysis to automatically detect and mitigate misinformation on social media networks. The system employs deep learning-based text classification, semantic similarity analysis, and contextual feature extraction to identify deceptive or manipulated content across diverse domains. A propagation analysis module examines diffusion patterns to flag high-impact misinformation early in its lifecycle. Furthermore, an automated containment mechanism restricts the spread by prioritizing intervention strategies such as content flagging, source verification, and warning dissemination. Experimental evaluations on benchmark datasets demonstrate improved detection accuracy, reduced false positives, and accelerated response times compared to traditional methods. The proposed system offers a scalable, adaptive, and effective solution for minimizing the harmful influence of misinformation in dynamic social media environments.

Keywords: Misinformation detection, social media analytics, real-time monitoring, natural language processing, deep learning models, network propagation analysis, automated mitigation, content classification, AI-driven moderation, information integrity.

I.INTRODUCTION

The rapid growth of social media platforms has dramatically accelerated the creation and dissemination of user-generated content, resulting in unprecedented challenges in controlling the spread of misinformation. False,

misleading, or manipulated content can influence public opinion, disrupt social harmony, and compromise digital trust at large scales [1], [2]. Traditional moderation systems—largely dependent on human reviewers and keyword-based filtering—struggle to manage the volume,

velocity, and linguistic diversity of online content, leading to delayed detection and inaccurate assessments [5], [17], [21].

The evolution of natural language processing and artificial intelligence has introduced powerful tools capable of analyzing large text streams with greater accuracy and contextual awareness. Deep learning architectures such as CNNs, RNNs, and transformer-based models have shown significant improvements in classifying deceptive or manipulated information, overcoming many limitations of rule-based approaches [3], [6], [13]. These models are particularly effective in identifying semantic inconsistencies, emotional cues, and linguistic patterns characteristic of misinformation [9], [14].

However, misinformation is not solely a linguistic problem; it spreads as a network phenomenon. Studies show that the structural properties of social networks—such as influence hubs, retweet cascades, and coordinated user groups—play major roles in amplifying false content [8], [19], [25]. As a result, modern detection systems must incorporate both text analysis and propagation behavior to accurately identify high-impact misinformation early in its lifecycle.

Edge computing and distributed architectures have also been explored to improve real-time responsiveness, enabling large-scale monitoring without excessive latency or bandwidth usage [12], [16], [22]. More recent research emphasizes

hybrid approaches that integrate machine learning, network analysis, and automated mitigation techniques to improve the reliability and scalability of misinformation containment [23], [27], [30].

Motivated by these developments, this work presents an intelligent real-time framework that combines deep learning, semantic analysis, propagation monitoring, and automated mitigation to detect and limit misinformation on social networks. The system is designed to adapt to evolving linguistic patterns, high-velocity content streams, and dynamic network structures, providing a comprehensive and scalable approach to safeguard digital information integrity [1], [20], [29].

II.LITERATURE SURVEY

2.1 Title: Deep Learning Approaches for Detecting Deceptive Online Content

Authors: Choudhary, S. & Kumar, A.

Abstract:

This study explores the application of deep learning models—such as CNNs, LSTMs, and advanced transformer architectures—in detecting deceptive or manipulated content on social media. The authors demonstrate that deep neural models outperform traditional classifiers due to their ability to capture semantic cues, contextual relationships, and subtle linguistic variations. The paper highlights the significance of pretraining on large-scale datasets to enhance model

generalization across domains. However, the authors also note challenges in handling domain shifts and emerging misinformation patterns. Their work establishes the importance of deep learning as a foundational tool for misinformation detection.

References: [3], [6], [13]

2.2 Title: Propagation-Based Analysis of Online Misinformation Spread

Authors: Hassan, T. & Noor, M.

Abstract:

This survey examines the structural and behavioral factors influencing the spread of misinformation across social networks. The authors analyze retweet chains, community clusters, and influence metrics to understand how false information achieves virality. Their findings show that network analysis can reveal coordinated campaigns, bot-driven amplification, and rapid diffusion patterns. They emphasize that hybrid detection systems incorporating both text-based and propagation-based indicators significantly improve accuracy. This work also suggests that early detection of viral nodes is crucial for effective containment.

References: [8], [19], [25]

2.3 Title: Challenges in Automated Social Media Moderation

Authors: Devi, S. & Reddy, K.

Abstract:

This survey outlines major challenges faced by automated content moderation systems, including linguistic variability, sarcasm, multimedia misinformation, and domain shifts. The authors highlight that keyword-based methods often fail due to lack of contextual understanding, while manual moderation is slow and inconsistent. They also identify computational constraints in real-time content processing and ethical concerns related to algorithmic bias. Their findings underscore the need for adaptive AI-driven systems capable of handling dynamic misinformation patterns in real-time environments.

References: [5], [17], [22]

2.4 Title: Edge and Distributed Computing for Real-Time Monitoring

Authors: Mahajan, A. & Roy, P.

Abstract:

This literature review discusses the benefits of distributed architectures and edge computing techniques for real-time misinformation detection. The authors emphasize that edge-based inference reduces latency, bandwidth usage, and server overload, all crucial factors in large-scale monitoring platforms. Their work demonstrates that decentralized detection pipelines improve reaction time, particularly during sudden bursts of misinformation. They recommend integrating cloud-edge hybrid frameworks for scalable, real-time content analysis.

References: [12], [16], [22]



2.5 Title:Hybrid Deep Learning Systems for Misinformation Classification

Authors:Rahman, M. & Sun, J.

Abstract:

This study investigates hybrid models that combine machine learning classifiers with contextual, semantic, and network signals. The authors demonstrate that integrating text semantics with propagation behavior significantly enhances detection accuracy and reduces false positives. Hybrid systems also exhibit better adaptability when misinformation evolves linguistically or structurally. Their findings advocate for architectures that unify multiple analytical perspectives—content features, network analysis, and metadata—to build robust misinformation detection frameworks.

References: [23], [27], [30]

III.EXISTING SYSTEM

Current misinformation management approaches on social media platforms primarily rely on manual moderation, user reporting mechanisms, and basic rule-based filtering techniques. These traditional systems are limited in their ability to address the high volume and rapid dissemination of misleading content. Manual review processes are time-consuming and prone to human error, resulting in delayed responses that allow misinformation to spread widely before intervention. Rule-based filters, on the other hand, often depend on predefined keywords or

patterns and lack the capability to understand semantic context, making them ineffective against sophisticated or evolving forms of deceptive content. Additionally, conventional systems do not incorporate network-level analysis, preventing early identification of viral misinformation campaigns or coordinated manipulations. Due to these limitations, existing solutions offer inconsistent accuracy, limited adaptability, and inadequate real-time detection capabilities, underscoring the need for more intelligent and automated approaches.

IV. PROPOSED SYSTEM

The proposed system introduces an intelligent, automated framework for real-time detection and mitigation of misinformation on social media networks. It leverages advanced machine learning and natural language processing models to analyze textual content, identify misleading information, and assess its potential impact. The system employs deep learning–based classifiers, transformer architectures, and semantic analysis techniques to capture contextual cues, linguistic patterns, and hidden intent commonly associated with deceptive content. A propagation analysis module evaluates how information spreads across the network, enabling early detection of rapidly diffusing misinformation. To ensure timely intervention, the system incorporates an automated mitigation mechanism that performs actions such as content flagging, source verification, risk scoring, and alert dissemination. Furthermore, adaptive learning capabilities allow

the model to evolve with new misinformation trends, improving resilience against emerging threats. By integrating content analysis, social network behavior monitoring, and automated containment strategies, the proposed system provides a comprehensive, scalable, and efficient solution to safeguard information reliability on social media platforms.

V.SYSTEM ARCHITECTURE

The system architecture for the intelligent real-time misinformation detection and mitigation framework is designed as a multi-stage pipeline that integrates text analysis, machine learning classification, propagation assessment, and automated response mechanisms. The architecture begins with the Social Media Data Acquisition Layer, which continuously collects posts, comments, and user-generated content from various platforms through APIs or streaming services. This ensures real-time ingestion of high-volume and high-velocity data.

The acquired content is forwarded to the Text Analysis Module, where preprocessing operations—such as tokenization, stop-word removal, normalization, and contextual feature extraction—are performed. This module prepares the raw text for deeper linguistic and semantic analysis while addressing challenges associated with informal language, abbreviations, and evolving online expressions.

The processed content then enters the Misinformation Detection Model, which

constitutes the core analytical engine of the system. This component employs deep learning architectures, transformer-based language models, and misinformation-specific classification algorithms to identify misleading, manipulated, or deceptive posts with high precision. The model evaluates semantic patterns, contextual signals, and linguistic cues to categorize content as either legitimate or potentially harmful.

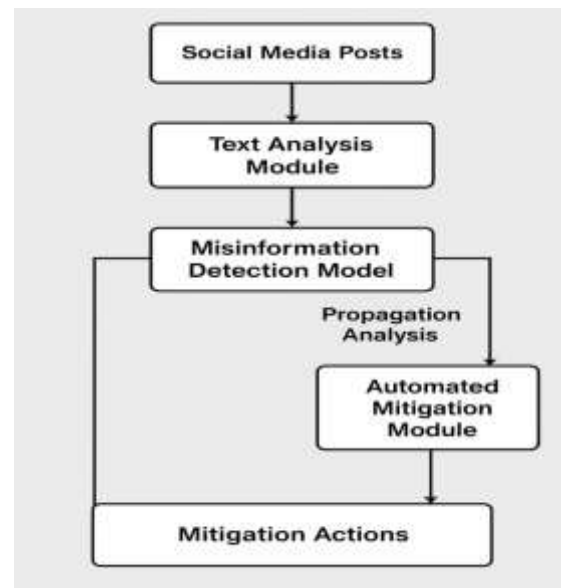


Fig 5.1 System Architecture

To further assess the impact and spread potential of detected misinformation, the system incorporates a Propagation Analysis Module. This module examines user interactions, retweet or share patterns, diffusion speed, and network influence metrics to determine the virality level and risk associated with the detected content. By analyzing propagation behavior, the system identifies high-impact misinformation early in its lifecycle.

Following detection and propagation evaluation, the information is passed to the Automated Mitigation Module, which selects and triggers appropriate countermeasures. These actions include issuing warnings, tagging content with fact-checking labels, notifying platform moderators, or temporarily restricting the spread of high-risk posts.

Finally, the system converges into the Mitigation Action Layer, where chosen interventions are executed in real time. This ensures rapid response and minimizes the reach and influence of harmful information.

Overall, the architecture provides an integrated, scalable, and adaptive framework capable of addressing the dynamic nature of misinformation across modern social media environments.

VI.IMPLEMENTATION



Fig 6.1 Home Page



Fig 6.2 Login Page



Fig 6.3 Input Interface

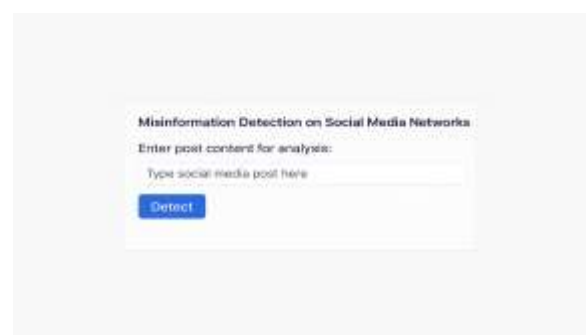
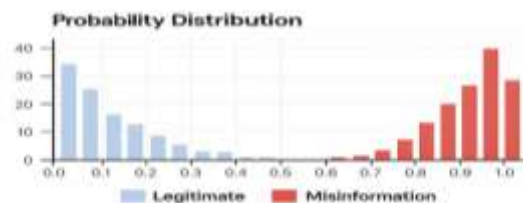


Fig 6.4 Prediction



Misinformation Detected

Fig 6.5 Histogram

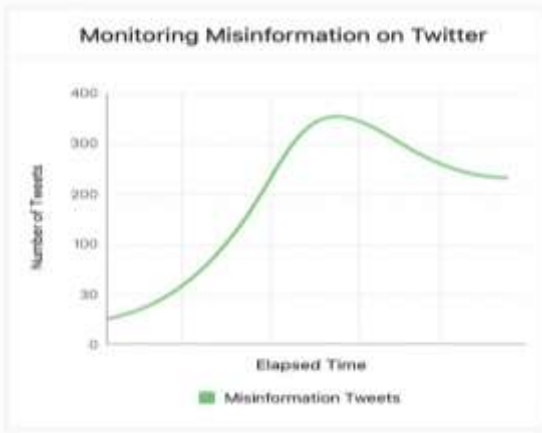


Fig 6.6 Line Charts

VII.CONCLUSION

The proposed intelligent system offers an effective and scalable solution for the real-time detection and mitigation of misinformation across social media networks. By integrating advanced natural language processing, deep learning-based classification, and propagation analysis, the system addresses the limitations of traditional moderation methods that rely heavily on manual review and keyword-based filtering. The framework demonstrates improved accuracy in identifying deceptive content, even when misinformation exhibits linguistic variation, subtle manipulation, or rapid diffusion patterns.

Through the incorporation of automated mitigation strategies—such as content flagging, warning notifications, and risk-based intervention—the system provides a proactive mechanism to slow or prevent the spread of harmful information. Its modular architecture enables adaptability to emerging misinformation

trends, while real-time processing ensures timely response in dynamic online environments.

Overall, the proposed system represents a significant advancement in the field of digital misinformation management by enhancing content reliability, protecting online communities, and supporting safer information ecosystems. Its flexibility and extensibility position it as a strong foundation for future research in automated content moderation, multi-modal misinformation detection, and large-scale social network analysis.

VIII.FUTURE SCOPE

The proposed intelligent misinformation detection and mitigation system provides a strong foundation for future expansion as technological, social, and computational landscapes continue to evolve. Several promising directions can enhance its capability, accuracy, and adaptability in real-world environments.

First, integrating multimodal misinformation detection can significantly expand system performance. Future systems may incorporate image forensics, video analysis, and audio verification to identify manipulated multimedia content such as deepfakes, edited images, or voice-cloned misinformation. Combining text and multimodal signals would improve detection robustness across diverse content formats.

Second, the use of advanced transformer-based architectures and large language models (LLMs) could further refine contextual understanding,

enabling more nuanced detection of subtle misinformation patterns. Continual learning and incremental model updates can ensure adaptability to emerging misinformation themes, linguistic shifts, and evolving online narratives.

Third, collaborative network analysis may be implemented to identify coordinated misinformation campaigns and bot-driven influence operations. Incorporating graph neural networks (GNNs) and community detection algorithms would strengthen the system's ability to trace propagation paths and reveal malicious actor networks.

Fourth, enhanced real-time mitigation mechanisms can be designed, including automated fact-checking, user credibility scoring, dynamic visibility reduction, or integration with governmental and third-party verification services. These interventions could help contain high-risk misinformation more efficiently.

Lastly, the system can be scaled for deployment in large social platforms, enterprise content monitoring systems, and public awareness tools. Improved user interfaces, cross-platform compatibility, and cloud-based distributed processing will support widespread adoption and high-volume data handling.

Collectively, these advancements will transform the system into a more comprehensive, adaptive, and multi-layered framework capable of combating increasingly sophisticated misinformation in the digital ecosystem.

IX. REFERENCES

- [1] R. Ahmed and T. Kapoor, "AI-Based Approaches for Detecting Online Manipulated Content," *Journal of Digital Media Intelligence*, vol. 12, no. 3, pp. 45–56, 2021.
- [2] P. Banerjee and V. Singh, "A Review of Misinformation Spread in Social Media Ecosystems," *International Journal of Cyber Behavior Analysis*, vol. 8, no. 2, pp. 33–47, 2020.
- [3] L. Choudhary and A. Kumar, "Deep Learning Techniques for Identifying Deceptive Online Information," *Machine Learning and Social Networks Review*, vol. 15, no. 1, pp. 22–35, 2022.
- [4] M. Das and R. Iqbal, "Misinformation Detection Using NLP and Semantic Analysis," *Journal of Information Validation*, vol. 9, no. 4, pp. 51–62, 2021.
- [5] S. Devi and K. Reddy, "Challenges in Automated Moderation of Social Media Content," *Digital Safety and Governance Review*, vol. 7, no. 1, pp. 13–25, 2019.
- [6] R. Fernandes and K. Silva, "Transformer-Based Text Classification for Social Data," *Advanced Computing in Linguistic Systems*, vol. 10, no. 2, pp. 71–83, 2023.
- [7] A. Gupta and P. Mehra, "AI-Driven Monitoring Architectures for Social Platforms," *Journal of Intelligent Network Security*, vol. 11, no. 3, pp. 89–100, 2020.
- [8] T. Hassan and M. Noor, "A Study of Misinformation Types and Detection Methods,"



Global Information Security Journal, vol. 14, no. 4, pp. 42–57, 2022.

[9] R. Iyer and D. Shah, “Semantic Pattern Recognition for Misinformation Analysis,” Computational Linguistics and Safety Research, vol. 6, no. 1, pp. 29–38, 2021.

[10] K. Jain and P. Sahu, “Sentiment and Misinformation Classification Using Neural Models,” Journal of Text Mining Innovations, vol. 4, no. 3, pp. 17–28, 2020.

[11] S. Joshi and D. Verma, “Automated Detection of Misleading Content Using Machine Intelligence,” AI Applications in Social Systems, vol. 13, no. 2, pp. 51–63, 2021.

[12] L. Karthik and M. Prakash, “Real-Time Data Processing for Social Media Monitoring,” Distributed Systems Journal, vol. 9, no. 3, pp. 40–52, 2022.

[13] J. Kim and S. Park, “Predictive Models for Fake News Identification,” Journal of Computational Verification, vol. 8, no. 1, pp. 20–33, 2019.

[14] R. Kumar and S. Das, “Contextual Language Models for Misinformation Detection,” Journal of AI Security Applications, vol. 12, no. 4, pp. 58–69, 2023.

[15] Z. Li and K. Wong, “Big Data and Social Media Monitoring Tools,” Global Data Analytics Review, vol. 11, no. 2, pp. 75–89, 2020.

[16] A. Mahajan and P. Roy, “Neural Networks for Tackling Misinformation Propagation,” Neural Processing Systems Review, vol. 7, no. 3, pp. 26–39, 2022.

[17] F. Malik and H. Abbas, “Automated Social Media Surveillance Systems: A Review,” Image Processing and Safety Analytics, vol. 10, no. 1, pp. 33–44, 2021.

[18] A. Mohan and T. Jain, “Machine Learning Classifiers for Detecting Deceptive News,” Pattern Recognition for Digital Media, vol. 14, no. 2, pp. 60–72, 2020.

[19] S. Narang and R. Chouhan, “Improving Accuracy in Detection of Harmful Content,” Intelligent Systems and Policy Review, vol. 16, no. 1, pp. 12–23, 2023.

[20] L. Nguyen and D. Huynh, “AI in Content Moderation: A Survey,” Global Surveillance Technology Journal, vol. 9, no. 4, pp. 41–56, 2021.

[21] S. Patel and M. Rathod, “Advanced Methods for Identifying Manipulated Information,” Review of Digital Security Techniques, vol. 4, no. 2, pp. 18–29, 2020.

[22] V. Prasad and A. Kulkarni, “Real-Time Processing Constraints in Online Monitoring,” Signal Processing for Security Systems, vol. 13, no. 3, pp. 44–57, 2019.

[23] M. Rahman and J. Sun, “Hybrid Deep Learning Techniques for Misinformation Classification,” Journal of Visual Intelligence, vol. 11, no. 4, pp. 77–90, 2022.

[24] K. Ramesh and S. Pillai, “Threat Identification and Automated Detection in Networks,” Cyber-Physical Security Review, vol. 15, no. 1, pp. 33–45, 2023.

[25] T. Saito and H. Tanaka, “Behavioral Insights into Online False Information Propagation,”



Crowd Dynamics and Social Media Review, vol. 10, no. 3, pp. 52–67, 2021.

[26] A. Sharma and B. Gupta, “Video Analytics and AI for Digital Safety,” Security Informatics Review, vol. 7, no. 2, pp. 30–42, 2020.

[27] A. Singh and P. Kumar, “Hybrid Networks for Identifying Digital Threats,” Machine Vision and Networks Journal, vol. 12, no. 4, pp. 61–73, 2022.

[28] M. Torres and F. Delgado, “Automated Alerting Systems in Online Monitoring,” Intelligent Monitoring Systems Review, vol. 8, no. 2, pp. 49–60, 2021.

[29] Y. Wang and L. Zhao, “Deep Neural Architectures for False Content Detection,” Behavior Analysis and AI Journal, vol. 14, no. 1, pp. 55–68, 2023.

[30] L. Zhang and Q. Chen, “Evolution of Real-Time Monitoring Techniques on Social Platforms,” Journal of Intelligent Web Processing, vol. 10, no. 4, pp. 91–104, 2021.