



# International Journal of Innovative Research in Science Engineering and Technology (IJIRSET)

*(A Monthly, Peer Reviewed, Refereed, Scholarly Indexed, Open Access Journal)*



**Impact Factor: 8.699**

**Volume 14, Issue 12, December 2025**

# Decoding Misinformation: How AI Detects Fake News across Social Media Platforms

Vaibhav Hariramani

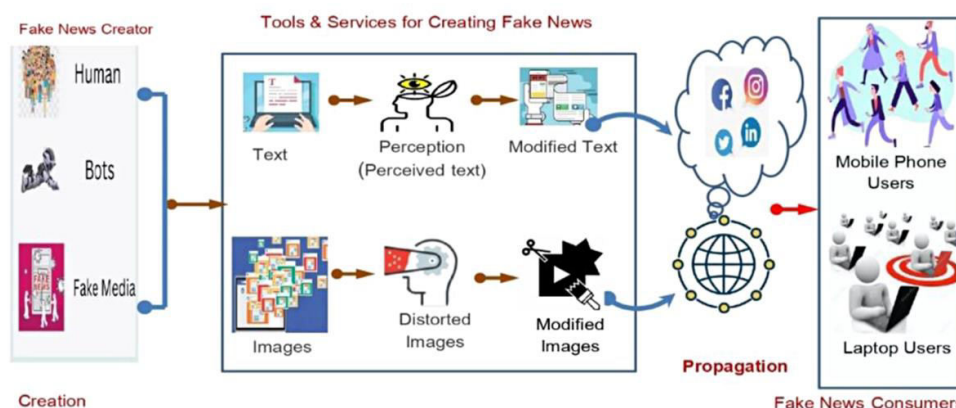
Student, Birla Institute of Technology & Science, Pilani, Dubai, UAE

**ABSTRACT:** The fast development of social media sites has enhanced the spread of fake information which poses major challenges to society, politics and the economy. This literature review critically analyzes how AI-based misinformation detection systems have evolved on five major themes namely traditional machine-learning models, deep learning developments, NLP-based fact-checking, multimodal fake news detection, and cross-platform ethical and technical issues. The review addresses how early feature-engineering systems have transformed to intricate transformer-based systems that can analyze linguistic pattern, contextual features and semantic meaning. It also examines the developments in multimodal detection incorporating text, image, and video analysis that are becoming increasingly popular on visually oriented social media like Tik Tok, Instagram, and YouTube. Although these developments have been made, some challenges remain, among them being bias in the datasets, multilinguality, adversarial manipulation, and the challenge of generalisation of models across different online ecosystems. Such ethical concerns as transparency, accountability, privacy and over-censorship risks are also examined critically. It is concluded that the lack of misinformation detection accuracy is caused by the necessity to introduce a hybridized framework of automated AI and human intelligence, data sharing across platforms, and more solid regulatory policies. Future studies are required on explainable AI, scalable multimodal architectures, and culturally inclusive datasets to enhance the global resilience to misinformation.

**KEYWORDS:** Artificial Intelligence; Misinformation Detection; Social Media Platforms.

## I. INTRODUCTION

The fast development of social media has made the process of accessing news and sharing opinions and interaction with information to be changed. Nevertheless, this free and dynamic online space has provided a rich soil in the propagation of misinformation and fake news (Nwaiwu et al., 2025). Inaccurate or untrue information can go viral in a few minutes, changing the minds of people, forming political beliefs and even leading to actual damage. Old methods of fact-checking are at least accurate, but too slow to compete with the volume and velocity of manufactured misinformation on the internet. Consequently, researchers and technology firms are looking to Artificial Intelligence (AI) in order to identify and prevent the dissemination of fake news on social networks like Facebook, X (formerly Twitter), Instagram, and Tik Tok (Rubin, 2022).



(Aimeur et al., 2023)



Detection based on AI is based on the methods of natural language processing, machine learning, and multimodal analysis to detect patterns that distinguish between false content and that which is true. These systems investigate language indicators, user behaviour, visualisation, and even how information dispersion occurs via networks (Aimeur et al., 2023). Despite a highly encouraging outlook, AI still faces certain obstacles, such as biased data, the changing tactics of misinformation, and the danger of over-moderation. The study will address the question of how AI models can be used to identify misinformation on the platforms and how well these methods perform, their drawbacks, and ethical implication (Waheed et al., 2025).

## II. RATIONALE AND OBJECTIVE

### Rationale

The common spread of fake news in social media has turned into a significant international issue, influencing the choice of the population in terms of their health, politics, and social confidence. Facebook, X (Twitter), Instagram, and Tik Tok have platforms that enable information to go viral, and users are not able to differentiate between quality content and manipulated and fake storytelling. The human fact-checking organisations strive to check the claims, but their work cannot match the rate and quantity of misinformation in the Internet (Ngueajio et al., 2025). Thus, AI-driven detection systems have become the key instruments that are capable of analysing patterns, classifying material, and detecting possibly harmful information on a large scale. It is important to study the functioning of such AI systems to determine their advantages and disadvantages, as well as the ways they can be improved. The reasoning justifies the necessity to explore AI-based misinformation detection to give information on more trustworthy, open, and moral systems capable of protecting the digital spaces (Kumar and Taylor, 2024).

### Objectives

1. To discuss the major AI methods of misinformation detection in the major social media, as well as natural language processing, machine learning, and multimodal analysis.
2. In order to assess the power and constraints of the AI-based fake news detectors models, emphasizing on the accuracy, dataset issues, and cross-platform contrasts.
3. To discuss the ethical, social, and technical challenges involved in automated misinformation detection, algorithmic bias, transparency, and over-moderation dangers.

## III. METHODOLOGY

The proposed research will employ the secondary research approach to explore how Artificial Intelligence can be used to identify misinformation on the social media. Secondary research can be used in this topic due to the fact that there are already a lot of scholarly researches, reports in the industry, and evaluations of technology, which offer credible and varied information sources. To conduct the research, it relies on peer-reviewed journal articles, conference papers, credible datasets, reports by fact-checking organisations, and technical commentaries on AI and social media sites. These sources provide information on AI detection methods like machine learning classifiers, multimodal systems, natural language processing and network-based models.

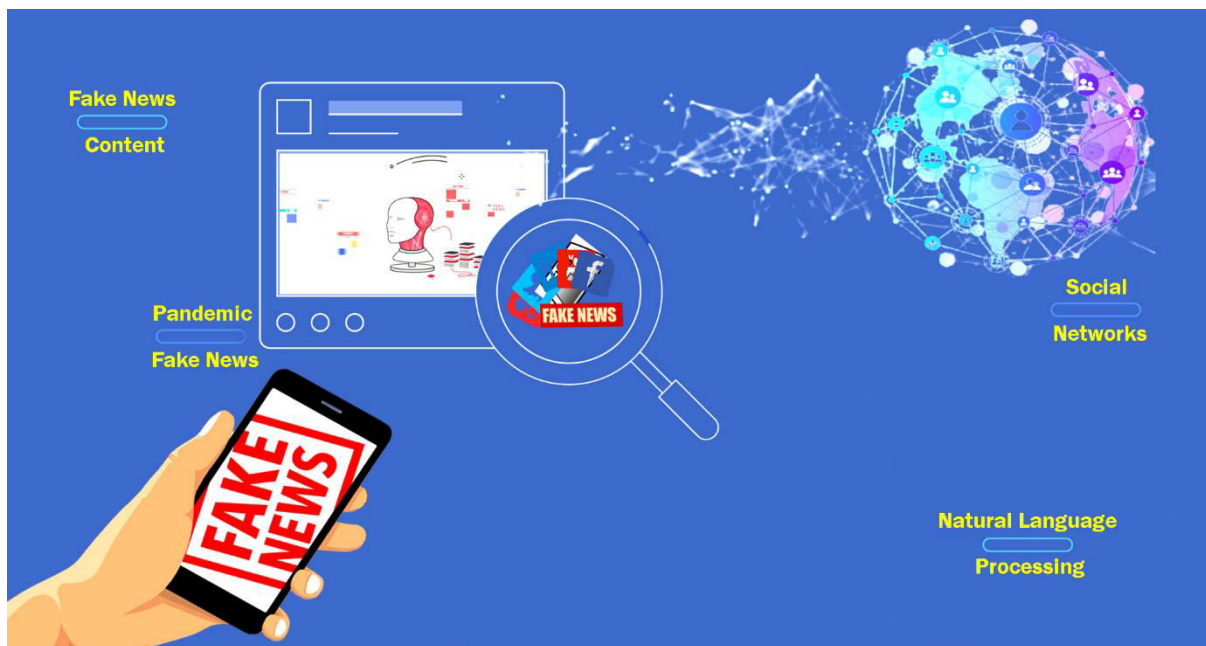
The process of review consists of a systematic identification of the relevant literature using database like Google Scholar. The keywords will be AI misinformation detection, fake news classification, social media algorithms, and multimodal fake news detection. The critical analysis of selected sources allows comparing methodologies, performance in detection, challenges, and ethical considerations.

Synthesising results across studies is also part of the methodology, where the objective is to bring out similar trends, gaps, and contradictions in the current studies on AI detection. This research relies on the secondary data only to obtain a fairly wide and comparative perspective of the functioning of AI systems without any primary experiments or gathering of original data.

## IV. LITERATURE REVIEW

**Theme 1: Development and Nature of Misinformation on Social media**

The development of misinformation on social media has been conditioned by the active development of the digital communication, an emergence of the user-generated content and the power of the platform algorithms. The initial misinformation mainly spread on blogs, forums, and email chains but the migration to social networking sites altered its magnitude, velocity and magnitude (Gondwe, 2025). These platforms (Facebook, X (ex Twitter), Instagram, YouTube, and Tik Tok) are currently the main sources of information used by millions of people and this situation provides a more favorable environment than ever before where fake information may be shared more extensively and faster than previously.



(Hussain et al., 2025)

All the platforms have unique features that determine the way misinformation spreads. Long-form posts, image-based accounts, and closed groups on Facebook allow forming a close-knit community that strengthens the shared beliefs (Hussain et al., 2025). Its short text format and ability to retweet makes X an instantaneous sharing platform, which makes it a hotbed of breaking-news fake news and politically charged stories. Instagram and TikTok are very visual based, and photographic manipulations, memes, and brief videos may lead to the misunderstanding of the viewers without using complicated written descriptions. YouTube, in turn, allows very long misleading content, conspiracy theories, and algorithmically suggested contents that can lead users to the harmful content loops (Giansiracusa, 2021). The virality of misinformation is based on user engagement. Studies constantly indicate that fake news invariably creates more emotional reactions, e.g., fear, anger, or surprise, and thus users are more likely to repost them. This appeal to the emotions, in combination with algorithmic amplification, enhances the appearance of misleading content (Joshi et al., 2023). The social media algorithms give priority to posts which have high engagement, and this boosts content that attracts reactions, whether they are accurate or not. Echo chambers and filter bubbles only exacerbate the issue by showing people only the information that confirms their preconceived notions of information, with little room to clearly reconsider the information.

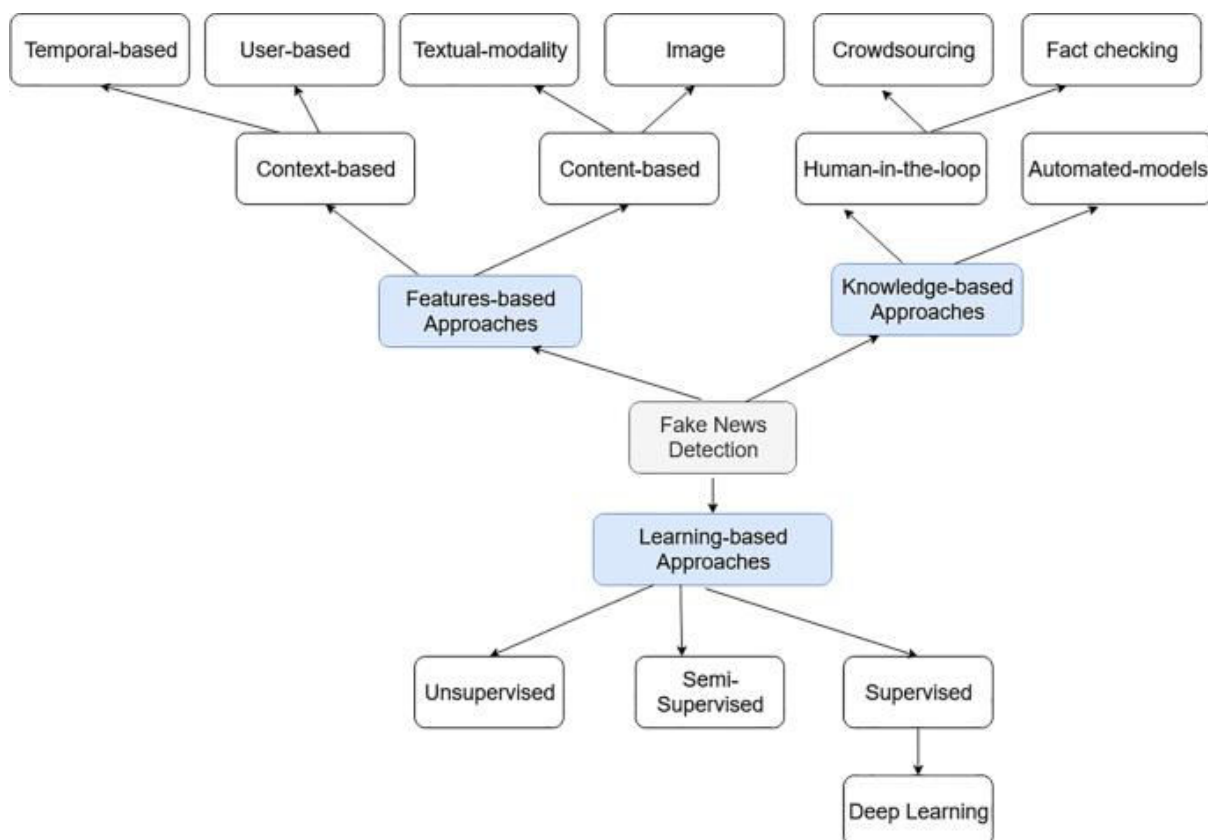
In general, misinformation on the internet is influenced by the interaction of the design, behaviour, and content dynamics. These patterns are critical to the analysis of the possibilities of AI tools to mediate effectively in various social media ecosystems.

## Theme 2: Text-Based Detection Techniques (NLP and Machine Learning) (AI-Driven)

Text-based detection is one of the most commonly studied and used methods in detecting misinformation in the social media platforms. Initial techniques were based on ancient Natural Language Processing (NLP) and machine learning algorithms that aimed at identifying linguistic, lexical, and stylistic attributes of textual data (Mouratidis et al., 2025). Text was modeled numerically using the bag-of-words (BoW) and term frequency-inverse document frequency (TF-IDF) models, and the use of classifiers, including Support Vector Machines (SVM), Naive Bayes, and logistic regression, was used to distinguish between genuine and fake content. These methods reasonably worked in controlled data, especially when misinformation had a characteristic linguistic signature, such as exaggerated affect or sensational wording or reduced readability (Hashmi et al., 2024).

Nevertheless, traditional models had weaknesses in context, sarcasm, and subtle manipulations capture as the misinformation became more advanced. This prompted the use of deep learning methods, such as recurrent neural networks (RNNs) and convolutional neural networks (CNNs) that were more effective at acquiring semantic relationships in text. Nevertheless, these models needed big labelled datasets and were not able to handle long-range dependencies (Hangloo and Arora, 2022).

Transformer based architecture became another breakthrough in the field of fake news as a text. BERT, RoBERTa, XLNet, and DistilBERT were used as models that enabled the incorporation of contextual knowledge in the tasks of classification with the help of the bidirectional attention mechanism. These models are very effective in finding subtle patterns, sentence construction and detecting misinformation in situations where the linguistic clues are subtle (Gifu, 2023). The misinformation datasets (LIAR or FakeNewsNet) fine-tuned with transformers have achieved significantly higher accuracy and generalisability than the previous approaches.



(Xu et al., 2023)

Transformer models have challenges in spite of their strong points. They can be computationally costly and might need large amounts of training data and can be biased by whatever is included in datasets, which can result in unjust or misleading classifications (Xu et al., 2023). Also, text-only models can have difficulties with the content that does not have obvious linguistic clues, e.g. short posts, coded language, or culturally specific phenomena. However, AI-based text classification is still a fundamental part of misinformation detection on which more complex and multimodal systems are created.

### Theme 3: Network and Propagation-Based Fake News Detection

The network and propagation-based methods are aimed at studying how misinformation disseminates on the social media as opposed to assessing the content itself. These approaches do not rely on linguistic indicators as it is done in text-based models but analyzes structural patterns in social networks, user interactions, and diffusion behaviours (Roy et al., 2025). The assumption here is that fake news tends to have recognizable propagation patterns that may not be similar to the information that can be considered credible.

The key component of this approach is based on graph, i.e. the interactions in social media are expressed in the form of a network with users as nodes and relationships between users like follows, shares or retweets as edges. Misinformation tends to propagate in small, highly dynamic communities and has a characteristic pattern of diffusion, such as early resharing and a concentration in ideologically consistent groups (Jahanbakhsh et al., 2023). Mapping such network dynamics, AI models are able to identify anomalies that indicate coordinated behaviour or even deceptive behaviour.

Diffusion modelling also has the advantage of improving the detection through the temporal stream of information. Models do not only monitor the speed of a post spreading, and multiple layers of resharing, but also whether an amplifying tendency corresponds to organic or artificially enhanced activity (Abdali et al., 2024). It was found out that fake news is more likely to spread more quickly and worldwide at the early stages of development, so propagation-based indicators can help identify it at an early stage.

Another technique that is used commonly is user credibility score. Through historical behaviour, posting patterns, follower networks, and verification status, algorithms are able to predict the probability of a user publishing reliable or unreliable posts (Raza et al., 2025). Tell-tale behavioural signatures of low-credibility accounts (bots, troll farms, new accounts, etc.) can be identified by machine learning models.

The network-based detection is also a key in the detection of coordinated inauthentic behaviour. This encompasses groups of accounts sharing the same content, coordinated actions, or fake boosting with an aim of affecting platform recommendations. Even the content may seem innocuous but such coordinated campaigns can be identified.

Comprehensively, propagation and network-based methods offer a strong supplementary view of the content analysis (Bashaddadh et al., 2025). They are very good at cross-platform monitoring and are capable of detecting patterns regardless of language, which makes them very useful in tracking the early-phase detection of emerging misinformation campaigns.

### Theme 4: Multimodal Fake News Detection (Text + images + Videos)

With the visualization of social media, multimodal fake news detection has become one of the most important fields of study. Social networks like Instagram, Tik Tok, and YouTube are more dependent on photos and videos, and it seems easier to ascertain misinformation through text only (Liu et al., 2025). To overcome this issue, multimodal AI systems consider the analysis of a mix of text, pictures, audio, and video to can produce a more detailed interpretation of content authenticity.

Image forensics, which is the analysis of visual evidence, inconsistencies, or evidence of manipulation is a central element of multimodal detection. Error-level analysis, noise pattern detection, and pixel-level anomaly detection are some of the techniques used to detect manipulated or fake images (Kwon and Jang, 2025). As the popularity of AI-generated visuals and deepfakes has increased, more sophisticated procedures have been created, such as analysis of face-movement, lip-sync detection, and temporal consistency checks by frame.



The current AI models can improve multimodal analysis with embeddings that combine the text and visual data. The CLIP model of the OpenAI, in particular, represents images and text in a common set of vectors, allowing the algorithm to detect discrepancies between captions and images (Sadiq et al., 2023). Likewise, vision transformers (ViT) and convolutional neural networks (CNNs) are high-level visual feature extractors that can be used together with textual embeddings produced by models like BERT or RoBERTa.

Another important role of metadata analysis is multimodal detection. Contextual information about authentic or fake content can be given by time stamps, geolocation tags, camera signatures and user posting behaviours (Bashardoust et al., 2024). On video-based networks, speech recognition algorithms can be used to retrieve audio transcripts to analyse the text, and computer vision methods can be used to analyse video frames to identify discrepancies or abnormal motion patterns that could be evidence of artificial alteration.

The multimodal systems are especially useful, since they are able to identify misinformation in cases where one modality does not have any clear cues. As an example, a non-veridical image with a neutral text can be flagged as a result of anomalies at the image level (Yildirim-Yayilgan et al., 2024). Moreover, multimodal models are more robust multi-culturally and multi-lingually because visual imagery tends to go beyond the limits of the text.

Altogether, multimodal fake news detection is the new horizon in the fight against misinformation that offers more accurate and engaging information that can be used in relation to the dynamism of online information.

#### **Theme 5: Cross-platform, Ethical and technical issues in AI Misinformation Detection**

Misinformation detection using AI has much potential, but it is also offered along with a set of ethical, technical, and cross-platform issues which make its use in the real world a difficult task. A significant ethical issue is the bias in datasets (i.e. their overrepresentation of particular languages, cultures, political opinions, or geographic locations). The majority of the benchmark datasets are based on Western cultures and using English as the language, and it makes it hard to expect the AI models to perform well in generalizing to the world (Su et al., 2024). This is the absence of multilingual and multicultural diversity, which makes the misclassification more risky, especially in case of the analysis of such culturally specific manifestations, satire, or informal dialects.

The limitation of privacy also affects the creation of effective detecting systems. In social media, user data is not easily accessible and researchers can hardly use social media to collect large and representative datasets. Moreover, the process of gathering and processing user behaviour data provokes the question of surveillance and the possibility of abuse of personal data (Gong et al., 2024). Another ethical concern is transparency because most AI-powered models, particularly deep learning models, can be seen as black boxes that are difficult to explain their actions. This interpretability is a shortcoming that casts doubt on the reliability of the users and question accountability particularly when AI is employed in content-moderation.

On the technical side, generalisation across platforms is one of the largest issues. The model that was trained under misinformation on Twitter might not work on YouTube or Tik Tok because of the content format, behaviour of the community, and the platform algorithms (Sarkar and Ghosh, 2024). Adversarial manipulation also makes it more difficult to be detected: malicious participants are more likely to dodge automated systems by making use of coded language, AI-generated pictures, deepfakes and coordinated inauthentic behaviour.

Finally, over-censorship is a very real danger to the extent that something that truly is not unethical (e.g. political discussion etc.) is wrongly considered misinformation. This highlights the necessity of cautious human-AI cooperation, where automated systems are used to offer first-line screening and human professionals with the final check (Nwaiwu et al., 2025). The lack of a balance between accuracy, fairness, and freedom of expression is a major concern since AI tools are increasingly becoming part of the misinformation control mechanisms.

### **V. CONCLUSION**

The increased size and complexity of misinformation on social media has rendered AI-based detection a critical part of the contemporary content moderation. In this review, it is demonstrated that misinformation develops differently on the platforms, and it is influenced by user behaviour, visual trends, and algorithm amplification. In order to deal with

this complexity, AI systems apply various techniques. Text-based models which are driven by NLP, machine learning and transformer architectures have high capabilities of analysing linguistic patterns, but have issues with context, cultural variation, and bias on data set. Network and propagation-based methods give the complementary view of studying the propagation of misinformation, and therefore, they allow detecting and identifying coordinated or inauthentic behavior early. Meanwhile, multimodal solutions have gained more significance as platforms are shifting more towards visual output, applying sophisticated image forensics, CLIP embeddings and video analysis to identify manipulated or AI generated content.

The efficiency of these systems is however curtailed by strong ethical and technical problems such as the absence of transparency, privacy issues, adversarial manipulation, multilingual issues, and the danger of over-censorship. All these problems underscore the necessity of more balanced datasets, interpretable AI models, and robust human-AI partnership. By and large, although AI has evolved significantly in identifying false information, there is a need to refine the system to guarantee accuracy, justice, and credibility in various online settings.

### REFERENCES

- 1 Abdali, S., Shaham, S. and Krishnamachari, B., 2024. Multi-modal misinformation detection: Approaches, challenges and opportunities. *ACM Computing Surveys*, 57(3), pp.1-29.
- 2 Aïmeur, E., Amri, S. and Brassard, G., 2023. Fake news, disinformation and misinformation in social media: a review. *Social Network Analysis and Mining*, 13(1), p.30.
- 3 Bashaddadh, O., Omar, N., Mohd, M. and Khalid, M.N.A., 2025. Machine Learning and Deep Learning Approaches for Fake News Detection: A Systematic Review of Techniques, Challenges, and Advancements. *IEEE Access*.
- 4 Bashardoust, A., Feuerriegel, S. and Shrestha, Y.R., 2024. Comparing the willingness to share for human-generated vs. AI-generated fake news. *Proceedings of the ACM on Human-Computer Interaction*, 8(CSCW2), pp.1-21.
- 5 Giansiracusa, N., 2021. How algorithms create and prevent fake news (pp. 17-39). Berkeley, CA: Apress.
- 6 Gifu, D., 2023. An intelligent system for detecting fake news. *Procedia Computer Science*, 221, pp.1058-1065.
- 7 Gondwe, G., 2025. Can AI outsmart fake news? Detecting misinformation with AI models in real-time. *Emerging Media*, p.27523543251325902.
- 8 Gong, Y., Shang, L. and Wang, D., 2024. Integrating social explanations into explainable artificial intelligence (XAI) for combating misinformation: Vision and challenges. *IEEE Transactions on Computational Social Systems*, 11(5), pp.6705-6726.
- 9 Hangloo, S. and Arora, B., 2022. Combating multimodal fake news on social media: methods, datasets, and future perspective. *Multimedia systems*, 28(6), pp.2391-2422.
- 10 Hashmi, E., Yayilgan, S.Y., Yamin, M.M., Ali, S. and Abomhara, M., 2024. Advancing fake news detection: Hybrid deep learning with fasttext and explainable ai. *IEEE Access*, 12, pp.44462-44480.
- 11 Hussain, F.G., Wasim, M., Hameed, S., Rehman, A., Asim, M.N. and Dengel, A., 2025. Fake News Detection Landscape: Datasets, Data Modalities, AI Approaches, their Challenges, and Future Perspectives. *IEEE Access*.
- 12 Jahanbakhsh, F., Katsis, Y., Wang, D., Popa, L. and Muller, M., 2023, April. Exploring the use of personalized AI for identifying misinformation on social media. In *Proceedings of the 2023 CHI Conference on Human Factors in Computing Systems* (pp. 1-27).
- 13 Joshi, G., Srivastava, A., Yagnik, B., Hasan, M., Saiyed, Z., Gabralla, L.A., Abraham, A., Walambe, R. and Kotecha, K., 2023. Explainable misinformation detection across multiple social media platforms. *IEEE Access*, 11, pp.23634-23646.
- 14 Kumar, A. and Taylor, J.W., 2024. Feature importance in the age of explainable AI: Case study of detecting fake news & misinformation via a multi-modal framework. *European Journal of Operational Research*, 317(2), pp.401-413.
- 15 Kwon, S. and Jang, B., 2025. A comprehensive survey of fake text detection on misinformation and lm-generated texts. *IEEE Access*.
- 16 Liu, Y., Shen, X., Zhang, Y., Wang, Z., Tian, Y., Dai, J. and Cao, Y., 2025. A systematic review of machine learning approaches for detecting deceptive activities on social media: Methods, challenges, and biases. *International Journal of Data Science and Analytics*, pp.1-26.
- 17 Mouratidis, D., Kanavos, A. and Kermanidis, K., 2025. From Misinformation to Insight: Machine Learning Strategies for Fake News Detection. *Information* (2078-2489), 16(3).



- 18 Ngueajio, M.K., Aryal, S., Atemkeng, M., Washington, G. and Rawat, D., 2025. Decoding fake news and hate speech: A survey of explainable ai techniques. *ACM Computing Surveys*, 57(7), pp.1-37.
- 19 Nwaiwu, S., Jongsawat, N. and Tungkasthan, A., 2025. Decoding Disinformation: A Feature-Driven Explainable AI Approach to Multi-Domain Fake News Detection. *Applied Sciences*, 15(17), p.9498.
- 20 Raza, S., Paulen-Patterson, D. and Ding, C., 2025. Fake news detection: comparative evaluation of BERT-like models and large language models with generative AI-annotated data. *Knowledge and Information Systems*, 67(4), pp.3267-3292.
- 21 Roy, S., Bhanu, M., Priya, S., Chandra, J. and Dandapat, S.K., 2025. Beyond just saying it's false: explainable AI for multimodal misinformation detection. *Applied Intelligence*, 55(11), p.781.
- 22 Rubin, V.L., 2022. *Misinformation and Disinformation: Detecting Fakes with the Eye and AI*. Springer Nature.
- 23 Sadiq, S., Aljrees, T. and Ullah, S., 2023. Deepfake detection on social media: leveraging deep learning and fasttext embeddings for identifying machine-generated tweets. *IEEE Access*, 11, pp.95008-95021.
- 24 Sarkar, S. and Ghosh, A., 2024. Leveraging artificial intelligence to enhance media literacy and combat misinformation. *JNRID*, ISSN, pp.2984-8687.
- 25 Su, J., Cardie, C. and Nakov, P., 2024, June. Adapting fake news detection to the era of large language models. In *Findings of the association for computational linguistics: NAACL 2024* (pp. 1473-1490).
- 26 Waheed, A., Azfar, S., Ali, A. and Soomro, M., 2025. NEURAL NETWORKS FOR DETECTING FAKE NEWS AND MISINFORMATION: AN AI-POWERED FRAMEWORK FOR SECURING DIGITAL MEDIA AND SOCIAL PLATFORMS. *Kashf Journal of Multidisciplinary Research*, 2(02), pp.90-111.
- 27 Xu, D., Fan, S. and Kankanhalli, M., 2023, October. Combating misinformation in the era of generative AI models. In *Proceedings of the 31st ACM International Conference on Multimedia* (pp. 9291-9298).
- 28 Yildirim-Yayilgan, S., Yamin, M.M., Ali, S. and Abomhara, M.A.S., 2024. Advancing Fake News Detection: Hybrid Deep Learning With FastText and Explainable AI.



INTERNATIONAL  
STANDARD  
SERIAL  
NUMBER  
INDIA



# INTERNATIONAL JOURNAL OF INNOVATIVE RESEARCH

IN SCIENCE | ENGINEERING | TECHNOLOGY

 9940 572 462  6381 907 438  [ijirset@gmail.com](mailto:ijirset@gmail.com)



[www.ijirset.com](http://www.ijirset.com)

Scan to save the contact details