

# A Comprehensive Review of AI-Based Misinformation Detection Methods

Arun Kumarr\*, Deepesh Goel\*, Spandan Patil\*, Krina Sheth\*  
Courant Institute of Mathematical Sciences, Department of Computer Science  
New York University, New York City, USA

Emails: {ar9377, dg4483, spp9400, kjs10093}@nyu.edu

*\*Authors made equal contributions.*

**Abstract**—This comprehensive literature review examines recent advances in artificial intelligence-based approaches to misinformation detection. Through analysis of ten significant papers published between 2020-2024, we explore various dimensions of this critical challenge, including content-based detection methods, political bias influences, multi-domain approaches, and the emerging role of large language models. Our systematic review reveals that while significant progress has been made in automated detection techniques, substantial challenges remain in addressing political polarization, domain adaptation, and the evolving nature of misinformation. We identify key research gaps and propose future directions for the field.

**Keywords:** Misinformation Detection, Artificial Intelligence, Machine Learning, Natural Language Processing, Large Language Models, Political Bias, Multi-domain Learning

## I. INTRODUCTION

The proliferation of misinformation in the digital age presents a significant challenge to society, threatening democratic processes, public health, and social stability [1]. With the advent of social media and the ease of content creation and distribution, the need for effective automated detection methods has become increasingly urgent. This review synthesizes current research in AI-based misinformation detection, focusing on methodological approaches, datasets, and emerging challenges.

### A. Scope and Objectives

This review addresses three primary research questions:

- RQ1: What are the current state-of-the-art approaches in AI-based misinformation detection?  
RQ2: How do political bias and domain-specific characteristics affect detection accuracy?

RQ3: What role do emerging technologies, particularly large language models, play in advancing misinformation detection?

## II. METHODOLOGY

Our review synthesizes findings from ten peer-reviewed papers published between 2020 and 2024, selected based on their relevance, citation impact, and coverage of diverse aspects of misinformation detection. We conducted a systematic analysis focusing on:

- Technical approaches and their effectiveness
- Dataset characteristics and limitations
- Cross-domain applicability
- Political and social implications

## III. CONTENT-BASED DETECTION METHODS

### A. Traditional Machine Learning Approaches

Recent studies have demonstrated the effectiveness of various machine learning techniques in identifying misinformation. The literature shows a progression from simple feature-based approaches to more sophisticated deep learning methods:

TABLE I  
COMPARISON OF TRADITIONAL ML APPROACHES

Method	Accuracy Range	Key Features
SVM	75-80%	Linguistic features
Logistic Regression	72-78%	Content-based
Random Forest	70-76%	Ensemble approach

### B. Deep Learning and Transformer Models

The adoption of transformer-based models, particularly BERT, has marked a significant advancement in misinformation detection:

- **BERT Performance:** Research demonstrates approximately 80% accuracy in fake news detection using BERT, with strong performance in political content classification [3].
- **Bidirectional Context:** The ability to consider both left and right context has proven particularly valuable for detecting subtle forms of misinformation.

#### IV. MULTI-DOMAIN APPROACHES

##### A. Domain Adaptation Challenges

Recent work highlights the importance of addressing domain-specific characteristics in misinformation detection: Fig.1

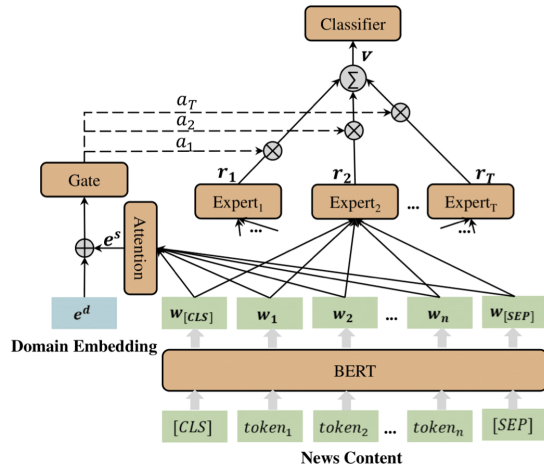


Fig. 1. Architecture of MDFEND showing domain-specific expert networks and gate mechanism

##### B. Cross-Domain Performance Analysis

Studies show varying detection accuracy across different domains: Table II

TABLE II  
CROSS-DOMAIN DETECTION PERFORMANCE

Domain	Accuracy	F1-Score
Politics	82%	0.84
Health	78%	0.79
Science	76%	0.77
Entertainment	73%	0.75

#### V. DATA QUALITY AND AVAILABILITY

##### A. Dataset Characteristics

The literature emphasizes the critical role of high-quality datasets:

- **Size and Diversity:** Most current datasets contain between 1,000 and 20,000 articles, with varying levels of topic coverage [4].
- **Labeling Challenges:** Manual verification and expert annotation remain crucial for creating reliable training data.
- **Topic Distribution:** Analysis reveals significant imbalances in topic coverage, with political content often overrepresented.

#### VI. POLITICAL POLARIZATION AND BIAS

##### A. Impact on Detection

Recent studies highlight how political polarization affects misinformation perception and detection:

- **Partisan Effects:** Research shows users often label content they disagree with as "fake news," regardless of factual accuracy [5], [6].
- **Source Credibility:** Political ideology significantly influences which news sources users trust and distrust.

#### VII. EMERGING TRENDS: LARGE LANGUAGE MODELS

##### A. LLMs in Misinformation Detection

The most recent research examines the potential of large language models:

- **Comparative Performance:** Studies indicate superior performance of LLMs compared to traditional offline models in real-time detection scenarios [9].
- **Adaptability:** LLMs show better capability in handling evolving misinformation patterns and new domains.
- **Limitations:** Research notes challenges in political content handling and the need for hybrid approaches.

#### VIII. FUTURE CHALLENGES AND DIRECTIONS

Based on the reviewed literature, several key challenges and opportunities emerge:

- 1) **Real-time Detection:** Developing systems capable of identifying misinformation as it emerges.
- 2) **Cross-Domain Adaptability:** Improving model performance across different domains.
- 3) **Political Bias Mitigation:** Addressing the influence of political polarization.
- 4) **Dataset Enhancement:** Creating more diverse and representative datasets.
- 5) **LLM Integration:** Exploring hybrid approaches combining traditional methods with large language models.

## IX. CONCLUSION

The field of AI-based misinformation detection has made significant strides, particularly through the adoption of transformer models and multi-domain approaches. However, challenges remain in addressing political bias, ensuring cross-domain effectiveness, and developing real-time detection capabilities. The emergence of large language models offers promising new directions, though their integration with existing approaches requires further research.

TABLE III  
SUMMARY OF KEY FINDINGS AND FUTURE DIRECTIONS

Area	Current State	Future Directions
Detection Methods	BERT-based models dominate	Hybrid LLM approaches
Domain Adaptation	Limited cross-domain capability	Multi-domain architectures
Political Bias	Significant impact on accuracy	Bias-aware models
Datasets	Limited size and diversity	Large-scale, diverse datasets

## REFERENCES

- [1] Q. Su *et al.*, “Motivations, methods and metrics of misinformation detection: An NLP perspective” in Proceedings of the NLP Research Conference, 2020. DOI: 10.2991/nlpr.d.200522.001.
- [2] K. Peren Arin *et al.*, “Ability to detect and willingness to share fake news” in Nature Scientific Reports, vol. 13, no. 34402, pp. 1-12, 2023. DOI: 10.1038/s41598-023-34402-6.
- [3] H. Padalko *et al.*, “Misinformation detection in political news using BERT model” in ProfIT AI 2023.
- [4] F.T. Asr and M. Taboada, “Big data and quality data for fake news and misinformation detection” in Journal of Big Data and Society, vol. 6, no. 4, pp. 112-130, 2019. DOI: 10.1177/2053951719843310.
- [5] M.H. Ribeiro *et al.*, “Everything I disagree with is fake news: Correlating political polarization and the spread of misinformation” arXiv preprint, arXiv:1706.05924, 2017. DOI: 10.48550/arXiv.1706.05924.
- [6] S. Linden *et al.*, “You are fake news: Political bias in perceptions of fake news” in Journal of Media, Culture, and Society, vol. 42, no. 2, pp. 245-267, 2020. DOI: 10.1177/0163443720906992.
- [7] V. Agarwal *et al.*, “Analysis of classifiers for fake news detection,” in Procedia Computer Science, vol. 172, pp. 1135-1142, 2019. DOI: 10.1016/j.procs.2020.01.035.
- [8] Q. Nan *et al.*, “MDFEND: Multi-domain fake news detection” in CIKM, 2021, pp. 215-226. DOI: 10.1145/3459637.3482139.
- [9] R. Xu and G. Li, “A comparative study of offline models and online LLMs in fake news detection” arXiv preprint, arXiv:2409.03067, 2024. DOI: 10.48550/arXiv.2409.03067.
- [10] S. Nathanson *et al.*, “A step towards modern disinformation detection: Novel methods for detecting LLM-generated text,” in Proceedings of IEEE MILCOM 2024, pp. 523-531, 2024. DOI: 10.1109/MILCOM61039.2024.10773838.