# A Survey of Fake News Detection and Misinformation Detection Using NLP and Machine Learning

## Abstract

The proliferation of fake news and misinformation in the digital age poses significant challenges to societal trust and democratic processes, necessitating advanced detection methodologies. This survey explores the interdisciplinary approaches leveraging natural language processing (NLP) and machine learning (ML) to enhance the detection of fake news. It highlights the integration of content-based, social context-based, multimodal, and multilingual methods, which collectively improve the robustness and accuracy of detection systems. The survey emphasizes the role of advanced models like the three-level hierarchical attention network (3HAN) and the HERO method, which utilize sophisticated linguistic and semantic analyses to classify news articles effectively. The importance of comprehensive datasets, including non-English and multimodal data, is underscored to ensure the applicability and generalizability of detection models across diverse contexts. Challenges such as algorithmic vulnerabilities, explainability, and the evolving nature of misinformation are addressed, with a focus on developing adaptive and resilient frameworks. Future directions include enhancing dataset diversity, advancing model transparency, and fostering interdisciplinary collaboration to tackle ethical considerations. The survey concludes that continued research and innovation are crucial for developing effective detection systems that safeguard the integrity of information in an increasingly complex digital landscape.

## 1 Introduction

### 1.1 The Importance of Fake News Detection

The digital age has transformed information dissemination, with social media amplifying the spread of misleading content, including fake news [1]. This proliferation poses a significant societal threat, necessitating the development of artificial intelligence solutions to address the issue [2]. The prevalence of fake news, particularly during critical events like electoral campaigns, underscores its potential to manipulate public opinion and undermine democratic processes.

Detecting fake news is complicated by evolving writing styles and the dependence on large annotated datasets, which present challenges for current detection methods [3]. The complexity of distinguishing credible news from misinformation is intensified by intricate relationships among news articles, their creators, and subjects [4]. Stylistic modifications of fake news that mimic reliable sources further undermine existing detection strategies, highlighting the need for advanced techniques [5].

The societal implications of fake news are profound, threatening democratic integrity by influencing public perception and behavior [4]. The swift spread of misinformation across social networks endangers societal stability and public trust, emphasizing the urgent requirement for effective detection mechanisms [6]. Integrating machine learning and knowledge engineering has emerged as a promising approach to enhance detection capabilities and combat the growing problem of fake news on social media [7].
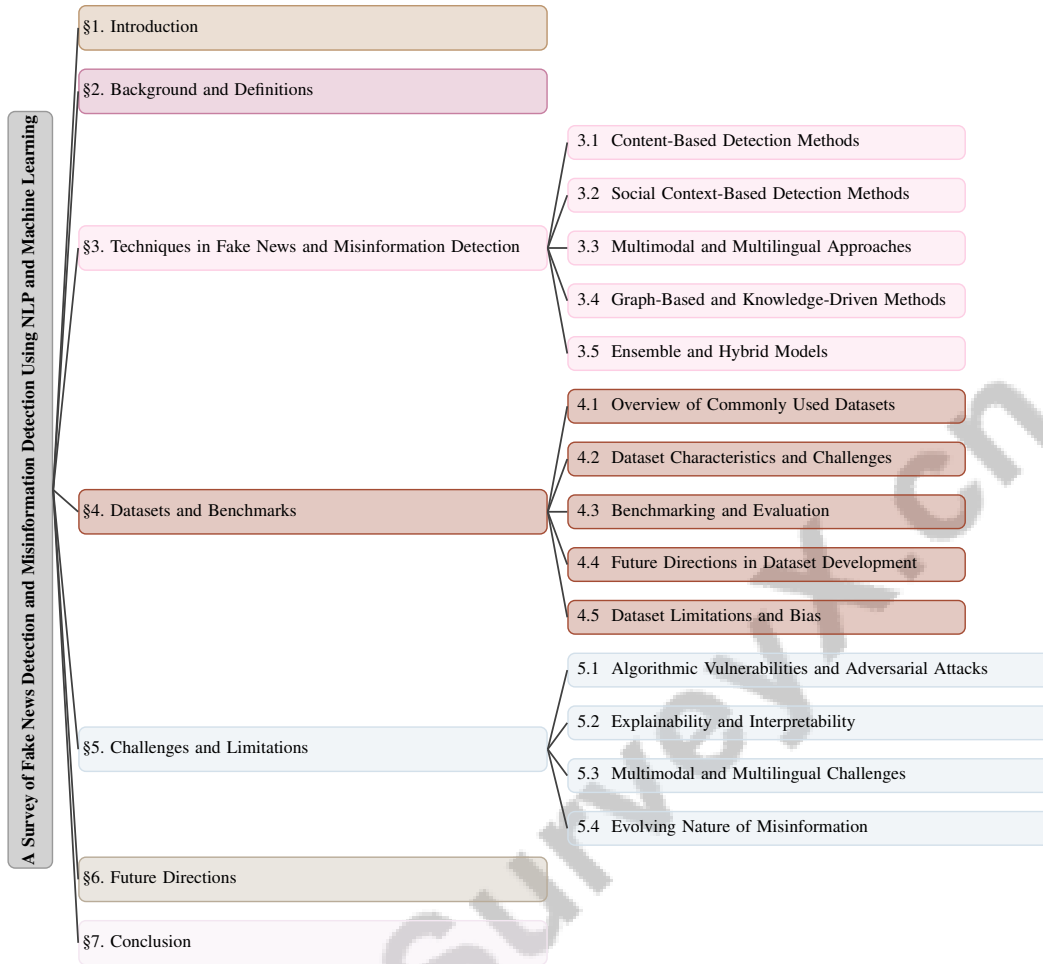
Figure 1: chapter structure

Detecting fake news is vital in today's digital landscape, necessitating a multidisciplinary approach that leverages advancements in computational linguistics, machine learning, and social sciences to safeguard information integrity and maintain public confidence in news content [8].

## 1.2 Role of NLP and Machine Learning

Natural Language Processing (NLP) and machine learning are crucial in developing sophisticated systems for detecting fake news and misinformation. These technologies provide a robust framework for analyzing and classifying textual data, employing advances in artificial intelligence, cognitive psychology, and computational linguistics to enhance detection capabilities. A notable approach is the Credibility-based Fake News Detection (CFND) method, which integrates source and content features to improve detection accuracy [9].

NLP and machine learning applications in misinformation detection include techniques that identify misinformation themes and keywords, demonstrating their effectiveness in addressing misinformation [10]. Existing methodologies are categorized into content-based and social context-based detection frameworks, reflecting the diverse strategies employed in this field [7].

Advanced models, such as the three-level hierarchical attention network (3HAN), effectively represent article semantics through news vectors, aiding in fake news detection [2]. The HERO method combines a hierarchical linguistic tree with a neural network to classify news articles based on their linguistic structures, capturing both local and global features [8].

Multimodal frameworks, which utilize NLP and deep learning techniques to identify fake news across various languages, demonstrate the adaptability of these technologies to different linguistic contexts

[11]. Hierarchical propagation networks enhance detection efficacy by capturing macro-level and micro-level features of news dissemination [12].

Deep neural networks play a vital role in detecting and categorizing fake news, enabling users to assess the credibility of information [13]. However, the emergence of Large Language Models (LLMs) that generate text resembling trustworthy sources presents a significant challenge to current detectors, necessitating continuous adaptation to maintain robustness and accuracy [5].

Benchmarking efforts evaluate machine learning models for their effectiveness in identifying and classifying fake news, providing standardized methods for performance comparison [1]. The integration of NLP and machine learning is essential for developing resilient systems capable of mitigating the societal impact of misinformation. By leveraging advanced models and techniques, researchers continue to enhance the efficacy of fake news detection frameworks.

## 1.3 Impact on Society

The societal ramifications of fake news are profound, as its rapid dissemination through social media can significantly undermine public trust and democratic processes [14]. The deceptive nature of fake news, often mimicking legitimate news content, complicates detection efforts and challenges the differentiation between true and false information [15]. This challenge is exacerbated by the scarcity of high-quality and extensive training data needed for effective model development [16].

The impact of fake news is particularly evident during critical events like elections, where misinformation can sway public opinion and influence political outcomes [4]. The 2016 US presidential election exemplifies how fake news can proliferate across social networks, misleading readers and affecting electoral integrity [6]. Additionally, the COVID-19 pandemic underscored the significant effects of misinformation on public health and safety, highlighting the urgent need for robust detection methods [17].

Experiments with models like Fakebox reveal vulnerabilities to adversarial attacks, emphasizing the necessity of integrating fact-checking into detection frameworks to enhance accuracy [18]. Furthermore, the societal consequences of fake news underscore the importance of developing models that are not only accurate but also interpretable, fostering better understanding and trust in detection processes [19].

The spread of misinformation through various media channels can lead to a misinformed public, emphasizing the critical need for effective detection methods to mitigate these negative societal impacts [20]. As fake news continues to challenge the integrity of information dissemination, developing advanced models capable of adapting to multimodal content is crucial for effectively countering its spread [5].

## 1.4 Structure of the Survey

This survey systematically explores the multifaceted domain of fake news and misinformation detection, leveraging insights from natural language processing and machine learning. We begin with an introduction that emphasizes the urgent need for effective detection of fake news and misinformation in today's digital landscape. This section elucidates the growing challenge posed by misleading information across various media platforms and underscores the vital role of NLP and machine learning technologies in developing computational tools that enhance the reliability assessment of online content. Recent research insights, including novel datasets for fake news detection and innovative methods leveraging authorship credibility to improve accuracy, are also discussed [9, 21, 22]. Additionally, the introduction addresses the societal implications of misinformation and outlines the paper's structure.

The second section provides essential background and definitions, clarifying key concepts such as fake news, misinformation, and disinformation, while discussing the interdisciplinary nature of the field. This section addresses the challenges inherent in defining and detecting fake news, setting the stage for the exploration of detection techniques.

In the third section, we delve into various techniques employed in fake news and misinformation detection, categorizing them into content-based, social context-based, multimodal, multilingual,

graph-based, and ensemble methods. This categorization aligns with existing research frameworks that organize detection methodologies based on task definitions and datasets [23].

The fourth section reviews the datasets and benchmarks pivotal to fake news detection research, providing an overview of commonly used datasets and discussing their characteristics, challenges, and limitations. This section emphasizes the importance of benchmarking in evaluating model performance and suggests future directions for dataset development.

The fifth section identifies the challenges and limitations faced in fake news detection, including algorithmic vulnerabilities, the need for explainability, and the evolving nature of misinformation. This section also discusses challenges specific to handling multimodal and multilingual data.

The penultimate section outlines potential future research directions, emphasizing the importance of developing robust and generalizable models, enhancing dataset quality, and fostering interdisciplinary collaboration. It also highlights the development of real-time detection systems and considers ethical implications.

The conclusion synthesizes the primary findings discussed throughout the document, emphasizing the critical need for ongoing interdisciplinary research to address the widespread challenges posed by fake news and misinformation. By reviewing detection methods from multiple perspectives—including source credibility, propagation patterns of false information, and content characteristics—the survey highlights the urgent demand for innovative solutions that not only enhance detection efficiency but also ensure the explainability of these methods. Collaborative efforts among experts in diverse fields such as computer science, social sciences, and journalism are essential for restoring public trust and safeguarding democratic processes [9, 24, 25].The following sections are organized as shown in Figure 1.

## 2 Background and Definitions

### 2.1 Defining Key Concepts

Fake news is characterized as fabricated content designed to deceive, often driven by financial or political motives, and presented as legitimate news, complicating its detection [6]. The linguistic nuances of fake news, distinct from authentic reporting, highlight the role of linguistic analysis in detection strategies [8]. The subtleties of language and deception challenge current methods, which often miss the nuanced features differentiating fake from real news [4].

Misinformation, the unintentional spread of false information, contrasts with disinformation, which involves deliberate falsehoods to manipulate audiences [3]. This distinction is crucial for understanding the broader landscape of information disorder, where both undermine digital informational integrity.

Natural Language Processing (NLP) is crucial for fake news detection, employing computational methods to analyze text and identify deceptive linguistic patterns, essential for classifying content as factual or fabricated [1]. However, machine learning models often face limitations due to the quality and availability of labeled datasets [4]. The rise of multimodal content, combining images and text, necessitates adaptive detection strategies to manage this complexity [5].

### 2.2 Interdisciplinary Nature

Fake news detection is inherently interdisciplinary, requiring insights from diverse fields to address misinformation and disinformation complexities. Detection methodologies are categorized into misinformation detection, source identification, and consumer interaction tools, underscoring the need for interdisciplinary integration [26]. This integration is vital for developing robust systems adaptable to the evolving nature of fake news.

Comprehensive surveys emphasize interdisciplinary approaches' necessity in enhancing fake news research efficacy [25]. By drawing from computer science, linguistics, psychology, and social sciences, researchers can create models that are technically proficient and sensitive to social and cognitive factors influencing fake news dissemination and reception. Incorporating structured linguistic features from social science into machine learning models has improved detection accuracy, illustrating cross-disciplinary collaboration's value [1].

4

Despite advancements, existing benchmarks focus predominantly on languages like English, creating a notable gap for languages such as Urdu, complicating automatic detection in non-English contexts [27]. Addressing these disparities is critical for developing globally applicable systems to combat misinformation across diverse linguistic landscapes.

The interdisciplinary nature of fake news detection is crucial for creating comprehensive solutions adaptable to misinformation challenges. Integrating insights from various disciplines enables researchers to develop robust frameworks that enhance fake news identification and mitigate societal impacts. For instance, analyzing news source credibility significantly improves detection accuracy, while advancements in retrieval-augmented large language models (LLMs) facilitate strategic evidence extraction for claim verification. Understanding human biases in news evaluation can inform tool development to better support truth discernment, fostering a more informed public [28, 9, 29].

## 2.3 Challenges in Definition and Detection

Defining and detecting fake news involves multifaceted challenges, including conceptual ambiguities and technical limitations. Distinguishing credible from fake news is complex, as both may contain factual elements, complicating detection efforts [9]. Varying definitions across studies lead to inconsistencies in conceptualizing and addressing fake news [7]. The intentional misleading nature of fake news, coupled with the vast scale and multimodal data from social media, further complicates detection [12].

Current methods struggle with feature complexity and dataset biases [20]. The rapid spread of misinformation and reliance on users who may inadvertently share false content complicate distinguishing real from fake news [30]. Challenges in obtaining high-quality annotated data and models' inability to generalize to new fake news styles hinder robust system development [3].

The abundance of unverified information online necessitates advanced methods capable of filtering extensive datasets [13]. Sophisticated techniques by fake news creators, mimicking legitimate sources, obscure the distinction between true and false news [6]. The lack of formal definitions and effective textual feature extraction presents ongoing challenges [31].

Current detection methods often fail to leverage articles' hierarchical structure, resulting in ineffective classification [2]. Limitations of existing datasets impact model generalizability, highlighting the need for comprehensive benchmarks reflecting fake news's multifaceted and evolving nature [7]. Advancements in theoretical understanding and practical implementation are required to enhance detection algorithms' interpretability and robustness.

## 3 Techniques in Fake News and Misinformation Detection

| Category | Feature | Method |
|---|---|---|
| **Content-Based Detection Methods** | Robustness to Style | SD[32] |
| **Social Context-Based Detection Methods** | Interaction Analysis | TriFN[33], NPD-FND[34], EF[35] |
| | Model Enhancement | CMTR-BERT[36] |
| | Feature Examination | FNDM[17] |
| | Multilingual Verification | MV[37] |
| **Multimodal and Multilingual Approaches** | Integration Strategies | LL-NLP[38], NT[39] |
| **Graph-Based and Knowledge-Driven Methods** | Explainable Graph Analysis | DISCO[40] |
| | Knowledge-Enhanced Graph Techniques | KEM-FND[41] |
| | Probabilistic Network Modeling | PIFM[42] |
| **Ensemble and Hybrid Models** | Adaptive Discrimination | DAAD[5] |
| | Transformer-Based Approaches | MR[43] |
| | Hybrid Learning Techniques | CMDM[10], FAD[31], TI-CNN[4], Bi-LSTM[6] |

Table 1: This table presents a comprehensive overview of various methodologies employed in fake news and misinformation detection. It categorizes the methods into content-based, social context-based, multimodal and multilingual, graph-based and knowledge-driven, and ensemble and hybrid approaches, highlighting specific features and techniques associated with each category. The table serves as a valuable resource for understanding the diverse strategies utilized in enhancing detection accuracy and effectiveness.

The surge in fake news and misinformation has prompted the creation of varied detection techniques to address its complex nature. Table 1 provides a detailed overview of the diverse methodologies employed in fake news and misinformation detection, categorizing them into distinct approaches

based on their unique features and techniques. Table 4 provides a comparative analysis of various methodologies employed in fake news and misinformation detection, highlighting the distinct features and techniques across different approaches. These techniques are categorized into content-based, social context-based, multimodal and multilingual, graph-based and knowledge-driven, and ensemble and hybrid approaches, each utilizing distinct facets of the information landscape to improve detection accuracy and effectiveness. Figure 2 illustrates the hierarchical categorization of these techniques, encompassing the aforementioned categories. Each category is further divided into specific methods and models, highlighting their unique contributions to enhancing detection accuracy and effectiveness. This visual representation not only aids in understanding the classification but also emphasizes the interconnectedness of the various approaches employed in combating misinformation.
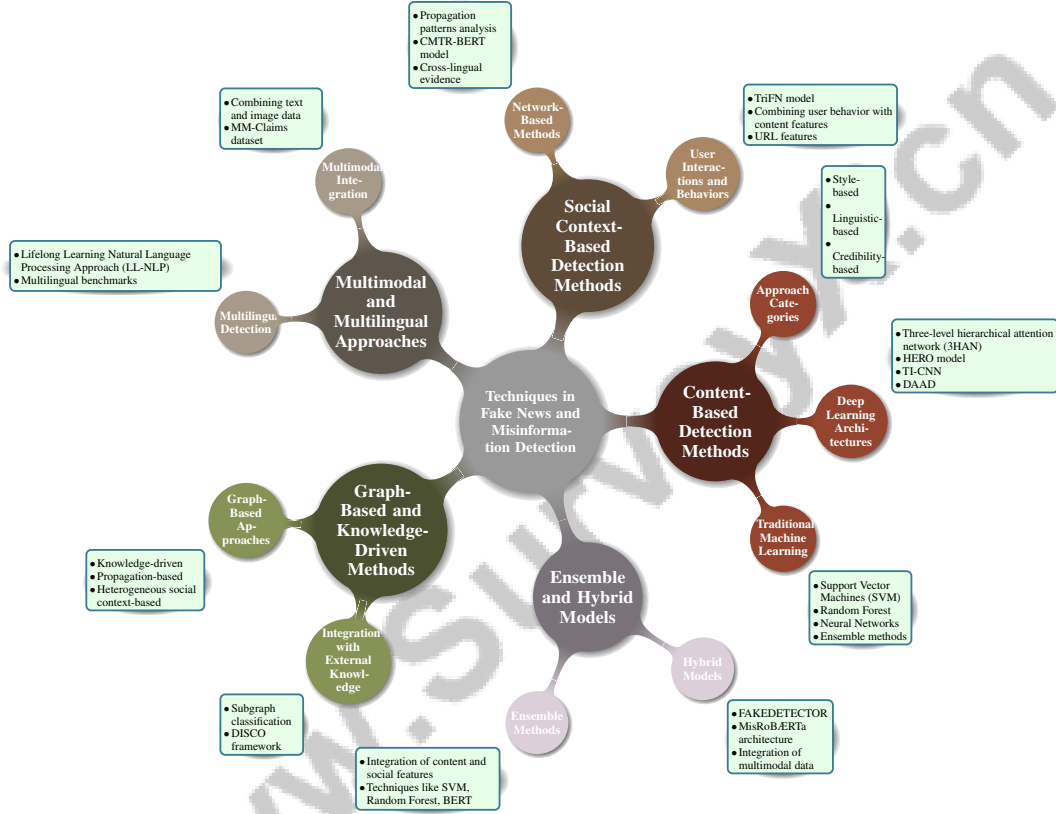


Figure 2: This figure illustrates the hierarchical categorization of techniques used in fake news and misinformation detection, encompassing content-based, social context-based, multimodal and multilingual, graph-based and knowledge-driven, and ensemble and hybrid approaches. Each category is further divided into specific methods and models, highlighting their unique contributions to improving detection accuracy and effectiveness.

## 3.1 Content-Based Detection Methods

Content-based detection methods focus on analyzing the textual content of news articles to identify fake news, employing linguistic, semantic, and syntactic features to detect inconsistencies [7]. Traditional machine learning techniques like Support Vector Machines (SVM), Random Forest, and Neural Networks are commonly used, often augmented by ensemble methods. Deep learning architectures have further advanced this field by capturing complex language relationships, enhancing classification accuracy. Models such as the three-level hierarchical attention network (3HAN) and the HERO model exemplify the use of sophisticated attention mechanisms and hierarchical trees to understand linguistic style and content structure [2, 8].

As illustrated in Figure 3, the categorization of content-based detection methods for fake news emphasizes the integration of various methodologies, including machine learning techniques, advanced models, and multimodal integration. This figure highlights the role of traditional techniques like

6

SVM and Random Forest alongside sophisticated models such as 3HAN, HERO, and TI-CNN, which concurrently analyzes text and image data, showcasing the potential of convolutional neural networks in fake news detection [4]. Additionally, DAAD employs prompt optimization and multiple discriminators for adaptive analysis [5]. The diversity of content-based methods is further illustrated by categorizing them into style-based, linguistic-based, and credibility-based approaches [7]. The SheepDog model, for instance, emphasizes substantive content over stylistic features [32].

Despite these advancements, the complexity of misinformation necessitates ongoing exploration of new strategies. Benchmarks highlighting human behavior in fake news consumption underscore the need for incorporating user interaction patterns in detection models [19]. These methods reveal the necessity for adaptive strategies to effectively counter the spread of fake news.
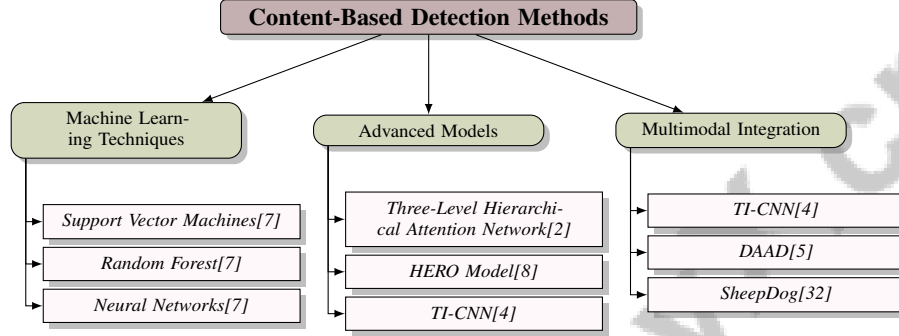


Figure 3: This figure illustrates the categorization of content-based detection methods for fake news, highlighting machine learning techniques, advanced models, and multimodal integration. It emphasizes the integration of various methodologies including Support Vector Machines, Random Forest, and Neural Networks, alongside sophisticated models like 3HAN, HERO, and TI-CNN, and multimodal approaches such as DAAD and SheepDog.

## 3.2   Social Context-Based Detection Methods

Social context-based detection methods leverage user interactions and behaviors within social networks to identify fake news. These approaches focus on the relationships among publishers, news content, and dissemination patterns to improve detection accuracy [44]. The TriFN model exemplifies this by analyzing the tri-relationship among publishers, news pieces, and users [33]. Combining user behavior analysis with content features demonstrates the potential of multifaceted approaches in combating misinformation [35]. URL features, such as length and domain age, also enhance detection precision [17].

Network-based methods analyze the propagation patterns of fake news within social networks, focusing on interactions among spreaders and their influence on information dissemination [34]. The CMTR-BERT model integrates social context with advanced NLP methods, employing an ensemble of BERT models and automatic text summarization techniques [36]. Additionally, cross-lingual evidence is utilized in methods like Multiverse, which analyzes news articles in multiple languages to enhance detection systems [37].

These methods underscore the significance of understanding user behavior and interactions in combating fake news. By synthesizing insights from social network analysis, linguistic features, and NLP techniques, they establish a comprehensive framework for detecting misinformation within the complex digital landscape.

## 3.3   Multimodal and Multilingual Approaches

Multimodal and multilingual approaches are essential for detecting misinformation across languages and formats, such as text and images. They enhance detection accuracy and address challenges posed by low-resource languages in social media [45, 46, 47, 48, 49]. These methods integrate data from various modalities to improve misinformation detection capabilities. For instance, the MM-Claims dataset emphasizes the importance of combining textual and visual data for claim detection across diverse topics [50].

7

In multilingual detection, approaches like the Lifelong Learning Natural Language Processing Approach (LL-NLP) combine classical feature extraction with deep learning models to facilitate fake news detection across languages [38]. This addresses linguistic diversity challenges in social media, where misinformation can transcend language barriers [39].

The introduction of multilingual and multimodal benchmarks enhances the robustness of deception detection frameworks beyond traditional single-language or single-modal content approaches [46]. By considering multiple languages and modalities, these methods provide a comprehensive framework for tackling misinformation in a globalized digital landscape.

The integration of multimodal and multilingual techniques underscores the importance of developing adaptable systems capable of countering misinformation across diverse platforms. These innovative detection approaches significantly improve accuracy while addressing cultural and linguistic variability. By leveraging advanced retrieval-augmented large language models (LLMs) and incorporating symbolic features, these systems enhance the reliability of verdicts and provide human-readable explanations, facilitating better interpretability across different domains [1, 29].

## 3.4 Graph-Based and Knowledge-Driven Methods

| Method Name | Structural Features | Knowledge Integration | Data Heterogeneity |
|---|---|---|---|
| KEM-FND[41] | Structural Information Subgraphs | Knowledge-based Approach | Multiple Modalities Information |
| PIFM[42] | Probabilistic Modeling | External Knowledge Bases | Diverse Data Sources |
| DISCO[40] | Word Graph | Graph Machine Learning | Heterogeneous Features |

Table 2: Comparison of graph-based and knowledge-driven methods for fake news detection, highlighting structural features, knowledge integration, and data heterogeneity. Each method utilizes unique structural characteristics and external knowledge to enhance misinformation detection accuracy across diverse data sources.

Graph-based and knowledge-driven methods offer a sophisticated framework for fake news detection by utilizing the structural characteristics of information dissemination networks and integrating external knowledge bases. This enables more accurate identification of misinformation through the analysis of relationships between entities and the context of news content [51, 52, 41]. These approaches capture intricate relationships among entities involved in misinformation spread. Table 2 provides a comparative analysis of various graph-based and knowledge-driven methods employed in fake news detection, emphasizing their structural features, knowledge integration, and handling of data heterogeneity.

Recent research categorizes graph-based fake news detection methods into knowledge-driven, propagation-based, and heterogeneous social context-based approaches [52]. The knowledge-driven approach transforms fake news detection into a subgraph classification task, enhancing identification accuracy through external knowledge integration [41]. Propagation-based methods focus on the dynamic spread of information within networks, classifying misinformation based on dissemination patterns [42].

Heterogeneous social context-based methods extend the graph-based framework by incorporating diverse data sources and user interactions. They analyze the co-sharing likelihood of mainstream and fake news articles, providing insights into narrative construction and sharing dynamics across platforms [53]. The integration of heterogeneous data enhances the model's ability to capture the complexities of social interactions influencing misinformation dissemination.

The effectiveness of graph-based methods is further improved by novel pre-training objectives tailored for heterogeneous graph structures, allowing for accurate modeling of diverse social networks [54]. The DISCO framework exemplifies the power of graph structures in enhancing detection accuracy, leveraging relationships between words to provide explainable predictions based on word contributions [40].

## 3.5 Ensemble and Hybrid Models

Ensemble and hybrid models are pivotal for enhancing the accuracy and robustness of fake news detection systems by synergizing diverse algorithms and methodologies. These models utilize the complementary capabilities of individual classifiers, forming a cohesive framework for misinfor-

| Method Name | Integration Techniques | Data Modalities | Model Frameworks |
|---|---|---|---|
| FAD[31] | Hybrid Feature Learning | Textual Information Relationships | Gated Diffusive Network |
| MR[43] | Ensemble OF Transformers | Content And Social | Transformer-based Architecture |
| TI-CNN[4] | Convolutional Neural Networks | Text And Images | Convolutional Neural Network |
| DAAD[5] | Prompt Optimization | Images And Text | Adaptive Discriminators |
| CMDM[10] | Machine Learning Algorithms | Content And Social | Ensemble And Hybrid |
| Bi-LSTM[6] | Hybrid Models | Content And Social | Bi-LSTM Architecture |

Table 3: Comparison of various ensemble and hybrid models used for fake news detection, detailing their integration techniques, data modalities, and model frameworks. The table highlights the diverse approaches employed to enhance detection accuracy and robustness by leveraging different algorithmic and data integration strategies.

mation identification. Ensemble methods integrating content and social features exhibit superior detection accuracy compared to standalone approaches [7]. They often incorporate advanced techniques, including Support Vector Machines, Random Forest, BERT, and RoBERTa, to optimize detection outcomes [55].

Table 3 presents a comprehensive comparison of ensemble and hybrid models employed in fake news detection, illustrating the integration techniques, data modalities, and model frameworks utilized to improve detection efficacy. Hybrid models extend these capabilities by merging deep learning with heterogeneous data sources. The FAKEDETECTOR, for instance, employs a hybrid feature learning unit and a gated diffusive network model to effectively fuse disparate information, enhancing credibility inference accuracy [31]. The MisRoBÆRTa architecture showcases the efficacy of hybrid approaches in misinformation detection through transformer-based deep learning for multi-class classification [43].

The integration of multimodal data is a significant advantage of hybrid models, as demonstrated by TI-CNN, which excels in analyzing multiple data modalities for improved detection accuracy [4]. Additionally, the DAAD method enhances fake news detection by optimizing prompts and dynamically selecting discriminators based on news content characteristics [5].

Innovative hybrid models, combining insights from social psychology and machine learning, have been shown to improve detection accuracy, particularly concerning COVID-19 misinformation [10]. Future research could focus on expanding datasets and exploring hybrid models that integrate Bi-LSTM with other techniques to enhance accuracy and robustness [6].

Frameworks like BODEGA introduce a comprehensive evaluation framework that integrates multiple misinformation detection tasks and adversarial attack scenarios, providing a robust benchmark for assessing model performance across diverse challenges [56].

| Feature | Content-Based Detection Methods | Social Context-Based Detection Methods | Multimodal and Multilingual Approaches |
|---|---|---|---|
| Detection Focus | Textual Content | User Interactions | Text And Images |
| Techniques Used | Linguistic, Semantic, Syntactic | Network-based, Url Features | Multimodal, Multilingual |
| Model Examples | 3han, Hero, Ti-CNN | Trifn, Cmtr-BERT, Multiverse | Ll-NLP, Mm-Claims |

Table 4: Comparison of Detection Methods for Fake News and Misinformation: This table delineates the primary features and techniques employed by different detection methods, categorized into content-based, social context-based, and multimodal and multilingual approaches. It highlights the focus of each method, the techniques utilized, and provides examples of models employed within each category.

# 4 Datasets and Benchmarks

## 4.1 Overview of Commonly Used Datasets

The efficacy of fake news detection models hinges on the availability of diverse and comprehensive datasets that encapsulate the intricacies of misinformation. Notable datasets, such as the Kaggle dataset with a balanced mix of fake and real news articles, are crucial for training and evaluating detection methodologies [57]. Another dataset, comprising 20,015 articles divided into 70% training and 30% testing, supports various detection techniques [13]. Benchmark datasets like PolitiFact, GossipCop, and Labeled Unreliable News (LUN) provide essential resources for model evaluation

across different contexts [32]. A large-scale dataset with 20,372 fake and 20,932 genuine articles highlights the importance of extensive resources for enhancing model robustness [2].

Datasets such as Recovery and MM-COVID, containing labeled news documents of varying lengths, are vital for examining the linguistic styles in misinformation [8]. Moreover, datasets focusing on linguistic features from social media and web articles shed light on the symbolic traits distinguishing fake news [1]. Key datasets identified in systematic reviews, like FakeNewsNet, offer additional social context and image data, enriching the understanding of misinformation dynamics [7]. The Fake News Corpus, despite its substantial size, lacks manual labeling, underscoring the need for annotated datasets to improve detection accuracy [7].

Experiments with datasets like Weibo, Weibo-21, and GossipCop illustrate their widespread use in fake news detection research, providing valuable benchmarks for model evaluation [5]. A dataset with 20,015 articles, including 11,941 labeled as fake and 8,074 as real, exemplifies the necessity of well-labeled resources in refining detection techniques [4]. An experimental setup with 2,977 news articles from various sources, categorized as true, false, or partially false, further demonstrates the diversity of available datasets for fake news detection [6]. Collectively, these datasets form the foundation of fake news detection research, offering essential resources for developing effective detection models.

## 4.2 Dataset Characteristics and Challenges

The construction and use of datasets for fake news detection face inherent challenges that significantly affect model effectiveness. A primary limitation is the limited size and scope of existing benchmarks, often relying on crowdsourced data that may not accurately reflect real-world scenarios [58]. This issue is compounded by the rapid evolution of fake news, necessitating ongoing updates and adaptations in dataset construction [59].

Datasets like ISOT, which includes 21,417 real and 23,481 fake articles, offer a solid foundation for analysis but reveal difficulties in maintaining balanced and representative samples [16]. The LIAR dataset, featuring 12.8K labeled short statements with multiple truthfulness labels, attempts to capture the nuanced nature of misinformation [16]. However, the scalability of detection methods is hindered by the challenge of obtaining high-quality training data, as malicious actors continuously adapt their strategies to evade detection [58].

Class distribution imbalance presents another significant challenge, particularly in datasets comprising tweets where underrepresented classes can adversely affect model performance [60]. This imbalance is especially problematic for categories with smaller representation, such as 'disagree,' complicating accurate classification outcomes [59]. The MediaEval 2020 dataset, consisting of 5,842 tweets classified into three categories, exemplifies the importance of addressing class imbalance in dataset construction [60].

Comprehensive datasets have been developed, including articles employing various propaganda techniques alongside factual inaccuracies [61]. Additionally, datasets scraped from news articles across multiple countries provide a diverse representation of news content [15]. The characteristics and challenges of fake news detection datasets highlight the need for ongoing innovation and refinement in dataset construction. By addressing issues like class imbalance, noise, and the evolving nature of misinformation, researchers can enhance the robustness and applicability of detection models in diverse digital environments [59].

## 4.3 Benchmarking and Evaluation

Robust benchmarking practices are crucial for evaluating fake news detection models, providing a standardized framework for assessing performance across various datasets and methodologies. Table 5 presents a detailed overview of the benchmarks utilized in evaluating fake news detection models, highlighting their respective sizes, domains, task formats, and metrics employed for performance assessment. These benchmarks facilitate comparability and reproducibility of results, empowering researchers to evaluate the effectiveness of diverse approaches in detecting misinformation, as evidenced by advancements in retrieval-augmented large language models and credibility-based detection systems [64, 9, 29, 68].

10

| Benchmark | Size | Domain | Task Format | Metric |
|---|---|---|---|---|
| FND-BENCH[57] | 10,000 | Fake News Detection | Binary Classification | Accuracy, F1-score |
| FND[22] | 73,368 | Fake News | Classification | Accuracy, F1-score |
| FND-BA[62] | 25,886 | Fake News Detection | Binary Classification | F1-score |
| FND-USER[63] | 40 | Misinformation | Classification | User Accuracy, User Agreement |
| MISINFO[64] | 1,250,000 | Misinformation Detection | Content Classification | F1-score |
| LIAR[65] | 12,836 | Political Discourse | Text Classification | Accuracy |
| SOTA-FND[66] | 45,000 | Fake News Detection | Binary Classification | Accuracy, F1 |
| FND-BERT[67] | 1,055 | Fake News Detection | Classification | Accuracy, F1-score |

Table 5: Table summarizing various benchmarks used in the evaluation of fake news detection models, detailing the size, domain, task format, and performance metrics for each benchmark. These benchmarks provide a comprehensive framework for assessing model effectiveness across different datasets, contributing to the advancement of misinformation detection systems.

A crucial aspect of benchmarking involves cross-validation techniques that yield reliable performance estimates. Typically, models are evaluated using a fivefold cross-validation approach, repeated multiple times with shuffled dataset versions to ensure robustness and minimize performance variance. This methodology allows for a thorough assessment of a model's generalizability across different data splits [57].

Standard classification metrics, including accuracy and F1-score, are commonly used to gauge model performance. The F1-score is particularly critical for evaluating detection performance in imbalanced datasets, as it effectively balances precision and recall, providing a comprehensive assessment of model effectiveness, especially in contexts where misclassifications can have significant societal implications [56, 15]. These metrics offer insights into a model's ability to classify news articles accurately, with the F1-score being especially useful for balancing precision and recall in imbalanced datasets.

Furthermore, model evaluation typically involves splitting datasets into training and testing sets, training the models, and assessing their accuracy and F1-score on the test set. This process is vital for understanding a model's performance in practical scenarios and its capacity to generalize beyond the training data [57].

Benchmarks are fundamental to advancing detection systems, as they systematically evaluate various machine learning approaches, including advanced pre-trained models like BERT, which have demonstrated efficacy in fake news detection across diverse datasets. This comparative analysis enables researchers to identify the most effective methods, enhancing the reliability and accuracy of detection systems and supporting strategies to combat the pervasive issue of fake news [29, 69]. By employing a comprehensive set of metrics and evaluation strategies, researchers can ensure that models are robust, generalizable, and capable of effectively addressing misinformation in an increasingly complex digital landscape.

## 4.4 Future Directions in Dataset Development

Future dataset development in fake news detection aims to address current limitations and enhance detection system robustness. Expanding datasets to include a broader range of languages and cultural contexts is crucial for overcoming the existing bias toward English-centric resources [27]. This expansion is vital for creating detection models that are globally applicable across diverse linguistic landscapes.

Integrating multimodal data, combining textual, visual, and potentially audio elements, offers a more comprehensive representation of news content [70]. Developing datasets that incorporate these modalities will enable models to capture the full spectrum of misinformation characteristics, thereby improving detection accuracy and resilience against sophisticated fake news strategies.

Creating dynamic datasets that adapt to the evolving nature of misinformation is also essential. Such datasets would continuously update and incorporate new instances of fake news, reflecting the latest trends and tactics employed by misinformation creators [44]. This adaptability is crucial for maintaining the relevance and effectiveness of detection models in a rapidly changing digital environment.

11

Moreover, including detailed metadata, such as source credibility, publication date, and dissemination patterns, can enhance the contextual understanding of news articles, providing additional layers of information for model training [71]. This enriched context can improve models' ability to discern subtle cues indicative of fake news.

Improving annotation quality and consistency in datasets is another key focus area. High-quality annotated datasets are essential for training models capable of accurately classifying news articles, and standardizing annotation practices will contribute to the reliability and comparability of detection systems [7].

To effectively combat misinformation, collaborations among academic institutions, industry stakeholders, and fact-checking organizations are crucial. Such partnerships can lead to the development of comprehensive and authoritative datasets that enhance fake news detection and fact verification efforts. By leveraging diverse expertise, these collaborations can tackle data quality and relevance challenges, ultimately improving misinformation mitigation strategies [72, 73, 74, 1, 29].

As illustrated in Figure 4, the future of dataset development in fake news detection centers around creating diverse, multimodal, and dynamic datasets that incorporate high-quality annotations and rich contextual metadata. This evolution is essential for enhancing detection models' effectiveness and robustness, as underscored by recent surveys emphasizing the critical role of dataset quality and diversity. Comprehensive datasets should integrate multiple information sources, including news content, social context, and dynamic propagation patterns, to address the complexities of fake news detection. Furthermore, ongoing research identifies challenges in dataset construction and offers opportunities to improve the depth and applicability of fake news studies, ultimately facilitating more effective detection and intervention strategies [72, 75, 59, 22]. These advancements will play a pivotal role in enhancing detection models' capabilities, enabling them to effectively counter the spread of misinformation in an increasingly complex digital landscape.
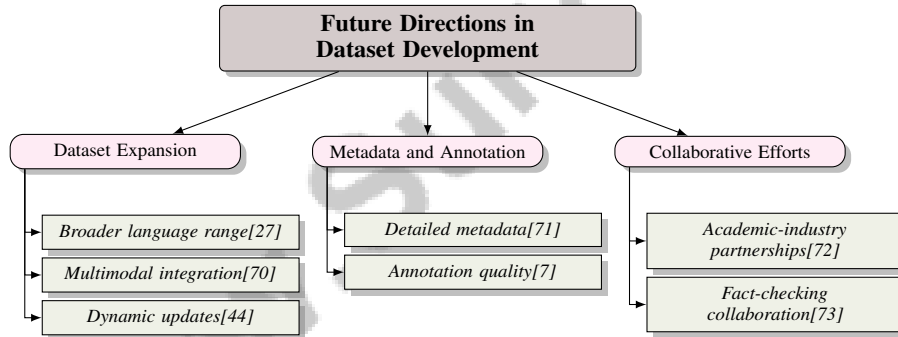


Figure 4: This figure illustrates future directions in dataset development for fake news detection, focusing on expanding datasets to include diverse languages and multimodal data, enhancing metadata and annotation quality, and fostering collaborations among academic, industry, and fact-checking organizations.

## 4.5 Dataset Limitations and Bias

The effectiveness of fake news detection models is often constrained by inherent limitations and biases present in existing datasets. A significant concern is selection bias, which occurs when datasets fail to represent the wide array of misinformation encountered in real-world scenarios. This bias can lead to models that perform well on specific datasets but struggle to generalize across different contexts and platforms [59].

Labeling bias is another critical issue, as the annotation process is susceptible to subjective interpretations, resulting in inconsistencies in classifying news articles as fake or real. Such biases can undermine the reliability of training data, ultimately affecting the accuracy and robustness of detection models [59]. The reliance on crowdsourced labeling further exacerbates this problem, as the quality and consistency of annotations can vary significantly.

Moreover, many datasets predominantly focus on English-language content, leading to a linguistic bias that limits the applicability of detection models in non-English contexts. This bias restricts the

global applicability of these models and overlooks the nuances and cultural differences inherent in misinformation across diverse linguistic landscapes [59].

The evolving nature of misinformation presents additional challenges, as datasets may quickly become outdated, failing to capture the latest trends and tactics employed by fake news creators. This temporal bias necessitates continuous updates and expansions of datasets to maintain their relevance and effectiveness in training detection models [59].

Addressing these limitations requires a concerted effort to develop more diverse, balanced, and dynamic datasets that can accurately reflect the complex and multifaceted nature of misinformation. By tackling biases and enhancing dataset quality, researchers can significantly bolster the generalizability and reliability of fake news detection systems. This improvement is crucial for ensuring their effectiveness in a rapidly evolving digital environment, where information accuracy is paramount. Employing advanced frameworks like retrieval-augmented large language models (LLMs) not only facilitates the strategic extraction of relevant evidence but also enhances interpretability through human-readable explanations. Furthermore, incorporating credibility assessments based on authorship and source history can strengthen detection methods, particularly in low-resource scenarios where traditional approaches may falter. Collectively, these strategies can lead to more robust systems capable of navigating the complexities of misinformation in today's information landscape [28, 9, 29, 76].

# 5 Challenges and Limitations

## 5.1 Algorithmic Vulnerabilities and Adversarial Attacks

Fake news detection models are frequently challenged by adversarial attacks that exploit algorithmic weaknesses. A notable issue is the over-reliance on stylistic features, making systems vulnerable to manipulations that alter these traits [8]. The static nature of traditional models limits their adaptability to the rapidly evolving misinformation landscape [5]. This issue is exacerbated by a dependence on manually crafted features, which undermines effectiveness against sophisticated fake news [5].

The quality and diversity of training datasets critically influence detection frameworks' generalizability and robustness [4]. The demand for extensive labeled data, as seen in the FAKE DETECTOR's limitations, presents a significant challenge for practical application [6]. Additionally, reliance on social propagation information, often unavailable during early detection, highlights the need for models that effectively leverage content-based features [7].

Training data quality is a crucial limitation, as it may not encompass all fake news types, necessitating the inclusion of more recent and diverse data sources [4]. The narrow focus on textual content without adequately addressing multimodal aspects limits current studies' effectiveness [1]. Challenges in assessing the credibility of single-author articles further complicate detection efforts [7].

Algorithmic vulnerabilities also stem from small sample sizes and potential biases in manual annotations, underscoring the need for innovative approaches to enhance detection frameworks' robustness and adaptability [3]. By improving dataset quality, incorporating multimedia data, and ensuring model transparency, researchers can significantly advance fake news detection systems' effectiveness.

## 5.2 Explainability and Interpretability

Developing explainable and interpretable models is essential for fostering trust and understanding in fake news detection. Current AI models often lack transparency, complicating predictions' comprehension and potentially leading to mistrust [77]. This issue is compounded by reliance on datasets that may not fully represent fake news, affecting model generalizability [78].

While the DISCO framework offers a comprehensive approach, it encounters challenges with complex disinformation forms, indicating a need for further refinement to improve interpretability [40]. Additionally, the quality of automatic summarization in some models may not match human-generated summaries, impacting explanations' clarity and effectiveness [36].

Dataset biases, particularly those from specific domains or languages, complicate models' interpretability and generalizability across diverse contexts [79]. Existing studies often lack a comprehensive approach to model interpretability and algorithm accountability in news distribution, highlighting the need for more transparent frameworks [80].

Some methods, despite offering explainability, may exhibit slightly lower accuracy than state-of-the-art deep learning models, limiting their applicability in precision-critical scenarios [15]. Limited sample sizes in certain datasets may hinder findings' generalizability, necessitating larger and more diverse datasets for improved reliability [81].

The complexity of some models can obscure explanation interpretability, leading to potential misunderstandings of underlying behaviors [19]. Enhancing transparency and interpretability in fake news detection systems is crucial to ensure users can understand and trust these models' decisions. Prioritizing explainability and interpretability can improve detection efficacy and foster greater confidence in system outputs.

Figure 5 illustrates the key challenges, methods, and future directions in enhancing explainability and interpretability for fake news detection models. This figure highlights specific frameworks and emphasizes the need for improved dataset and feature integration, reinforcing the discussion on the importance of transparency in this domain.
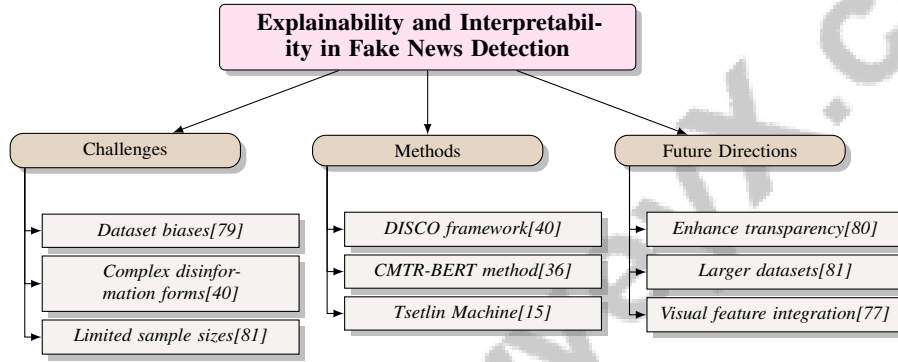


Figure 5: This figure illustrates the key challenges, methods, and future directions in enhancing explainability and interpretability for fake news detection models, highlighting specific frameworks and the need for improved dataset and feature integration.

## 5.3 Multimodal and Multilingual Challenges

Detecting fake news in multimodal and multilingual contexts presents unique challenges that require sophisticated approaches to address diverse data types and languages effectively. Integrating modalities such as text, images, and videos necessitates advanced techniques for accurate interpretation and analysis. The complexity and time consumption of preprocessing steps, including EDU segmentation and dependency graph construction, further complicate efficient model development [82].

In the multilingual domain, language diversity and linguistic nuances introduce additional difficulties for detection systems. Models must process multiple languages, each with distinct syntactic and semantic characteristics. This necessitates developing extensive datasets covering a broad spectrum of languages and cultural contexts, ensuring detection methods are robust across various linguistic landscapes. This approach is crucial for addressing misinformation challenges in low-resource languages and improving accuracy in distinguishing human-written from machine-generated text [48, 83].

The rapid evolution of misinformation, especially with Large Language Models (LLMs), adds complexity. LLMs' adaptability in generating text that mimics legitimate sources necessitates continuous evolution of detection methods to keep pace [84]. This dynamic environment requires benchmarks that effectively assess detection systems' performance in adapting to new modalities and languages, ensuring continued efficacy in identifying fake news.

To address misinformation complexities in an increasingly multimodal digital landscape, innovative methodologies integrating diverse data forms—such as text, images, and videos—and addressing linguistic diversity are essential. Developing advanced detection systems capable of identifying cross-modal discrepancies and leveraging synthetic data generation can enhance accuracy and robustness in misinformation detection [49, 85]. Employing advanced techniques and comprehensive benchmarks

14

can enhance fake news detection systems' adaptability and accuracy, ensuring resilience in an ever-changing digital landscape.

## 5.4 Evolving Nature of Misinformation

The evolving nature of misinformation poses significant challenges for detection systems, necessitating adaptive models capable of responding to new trends and contexts. Current studies often struggle with the rapidly changing landscape of fake news, complicating detection methods' effectiveness across various settings [86]. A primary obstacle is the reliance on web-sourced articles' quality and evolving fake news tactics, which can undermine detection frameworks' robustness [87].

The use of real text and images in misleading contexts complicates misinformation detection, as systems must discern subtle cues indicating deceptive intent [71]. Significant distribution shifts between events create challenges, leading to ineffective knowledge transfer and feature application between historical and upcoming events [88]. This underscores the need for models capable of adapting to new events and contexts, ensuring continued relevance and effectiveness.

Misinformation often transcends language barriers, complicating detection efforts in multilingual environments [37]. Detecting truth and falsity alone is insufficient; it is crucial to focus on misleadingness and content intent [89]. Dataset imbalances can affect model performance generalizability [55].

Potential biases in user demographics and the inability to track users' fact-checking behavior outside the test environment may impact findings' validity, highlighting the need for comprehensive evaluation frameworks [28]. Benchmarks may face challenges related to fake news' evolving nature, necessitating continuous updates to datasets and model evaluations to maintain utility and accuracy [57].

Future research directions include refining models for improved real-time detection capabilities, exploring additional features from user behaviors, and integrating deep learning techniques to enhance performance [12]. Developing innovative approaches responsive to new trends and capable of integrating diverse data sources can enhance fake news detection systems' resilience and efficacy in a rapidly changing digital landscape.

# 6 Future Directions

The future of combating misinformation relies heavily on the advancement of fake news detection models. This section delves into essential areas, highlighting the importance of enhancing generalization and robustness within detection systems. As misinformation continues to evolve, models must adapt to maintain consistent performance across varied contexts, thereby addressing the challenges posed by fake news in a dynamic digital environment.

## 6.1 Generalization and Robustness

Improving the generalization and robustness of fake news detection models is crucial due to the volatile nature of misinformation. Future research should focus on enhancing model adaptability to new misinformation trends, improving dataset reliability, and considering the ethical implications of automated detection technologies [90]. This involves integrating additional features and applying benchmarks across diverse social media platforms to strengthen detection capabilities [91]. Understanding the psychological aspects of fake news dissemination is vital for developing effective strategies [92]. Optimizing Large Language Model (LLM) inference efficiency will also enhance model generalization, enabling better handling of misinformation complexities [93]. Incorporating knowledge graphs can improve reasoning capabilities, while exploring linguistic features and alternative NLP techniques can increase detection accuracy, as evidenced by recent studies achieving over 0.80 F1-scores [51, 1]. These advancements will enhance model transferability across domains, adapting to thematic and contextual shifts.

As illustrated in Figure 6, refining evidence retrieval and integrating diverse evidence sources can significantly bolster a model's robustness in complex informational environments. The figure highlights key areas of focus for enhancing generalization and robustness in fake news detection, emphasizing model improvement, detection techniques, and evidence retrieval as critical components.

15

The quality and relevance of evidence are critical for accurate verdict predictions, particularly in fake news detection. Traditional methods often rely on outdated data sources, limiting effectiveness for emerging claims. In contrast, retrieval-augmented LLMs utilize multi-round retrieval strategies to extract current and relevant evidence from various web sources, enhancing performance and providing interpretable explanations [94, 29]. Focusing on generalization and robustness will enable the development of resilient models capable of effectively addressing the evolving challenges of misinformation.
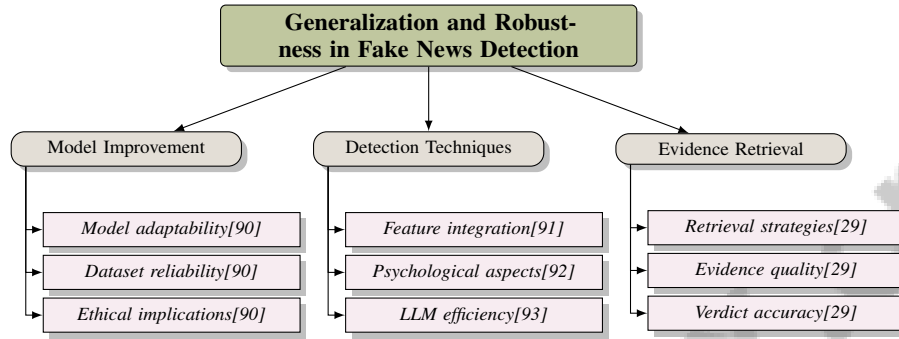
Figure 6: This figure illustrates the key areas of focus for enhancing generalization and robustness in fake news detection, highlighting model improvement, detection techniques, and evidence retrieval as critical components.

## 6.2 Enhancing Dataset Quality and Diversity

Improving dataset quality and diversity is vital for advancing fake news detection models, directly impacting their robustness and accuracy. Future research should prioritize integrating multimodal data—text, images, and other media types—to enrich datasets and capture the nuanced features of misinformation [30]. Expanding datasets to include a broader range of languages and topics will enhance the generalizability of detection systems across cultural and linguistic contexts [4]. Additionally, incorporating social media interactions and indirect relationships among heterogeneous nodes can refine detection capabilities [39]. Enhancing dataset quality also involves implementing explainable AI methods to provide insights into detection models' decision-making processes, addressing ethical considerations and fostering trust [7]. Developing adaptive datasets that evolve with misinformation patterns is essential for maintaining model relevance and effectiveness over time. Furthermore, leveraging synthetic data generation techniques can enhance dataset diversity, contributing to more effective detection capabilities [30]. By focusing on these areas, researchers can significantly advance dataset quality and diversity, leading to more robust and reliable fake news detection systems.

## 6.3 Advanced Model Development

Advancing fake news detection models requires a multifaceted approach that incorporates cutting-edge techniques and interdisciplinary insights. Future research should focus on enhancing model transparency and interpretability through the integration of external knowledge sources, improving counterfactual explanations and clarity of model decisions [19]. This approach builds trust in AI systems and enhances adaptability across diverse contexts. Integrating multimodal and social-context information is essential for refining detection models, as exemplified by efforts to enhance HERO's performance [8]. By leveraging multimodal data, including text, images, and user interactions, models can achieve a comprehensive understanding of misinformation, improving detection accuracy and robustness.

Future research could explore automating deceit pattern discovery, extending methods like DAAD to domains such as sarcasm detection [5]. This automation streamlines the detection process, enabling models to adapt to new forms of misinformation. Additionally, integrating symbolic knowledge into language models presents a promising avenue for enhancing adaptability and generalizability, allowing for capturing nuanced linguistic patterns indicative of fake news [1]. Enhancements to models like SheepDog, particularly in multimodal contexts, can improve detection accuracy by

16

capturing the complexity of misinformation [32]. Prioritizing the development of models leveraging Bi-LSTM architectures, which have shown significant effectiveness in fake news detection, is also crucial [6]. By concentrating on integrating advanced retrieval-augmented LLMs and sophisticated machine learning techniques, researchers can create highly effective and adaptable misinformation detection systems. These systems will source relevant, high-quality evidence from diverse web resources, enhancing resilience against evolving misinformation tactics. Incorporating symbolic features and linguistic attributes specific to misleading content will enable superior performance, providing accurate verdicts and interpretable results essential for combating the pervasive issue of fake news and disinformation [95, 1, 29]. This comprehensive approach will ensure detection systems remain robust and effective, safeguarding information integrity in a complex media environment.

## 6.4 Real-time Detection Systems

The increasing prevalence of misinformation across digital platforms necessitates the development of real-time detection systems capable of swiftly and accurately identifying fake news. These systems must leverage advanced computational techniques and large-scale datasets for efficient processing. Future research should focus on expanding datasets to include non-English content, facilitating cross-language comparisons and enhancing the generalizability of detection systems across diverse linguistic contexts [81]. Incorporating multimodal data, such as text and images, can further improve claim detection accuracy, as demonstrated by datasets like MM-Claims [50].

Developing real-time detection systems should also prioritize integrating automatic fact-checking mechanisms and exploring additional topics beyond politics and gossip to comprehensively understand misinformation dynamics [10]. The potential for real-time application in automatic hoax detection has been demonstrated, achieving remarkable accuracy in classifying posts as hoaxes or non-hoaxes [96]. This highlights the feasibility of deploying systems that dynamically respond to the evolving misinformation landscape. Moreover, deploying web applications offering fake news detection as a service, with real-time learning from new manually fact-checked articles, represents a promising direction for enhancing accessibility and usability [2]. Implementing hierarchical propagation networks (HPN) has shown significant improvements over existing state-of-the-art methods, achieving an average F1 score greater than 0.80 across various datasets, underscoring the efficacy of advanced methodologies in real-time detection [12]. Focusing on developing advanced computational tools for fake news detection will enable researchers to create sophisticated systems that leverage novel datasets and credibility assessment methods. This proactive approach enhances the accuracy of fake news detection through integrating linguistic, visual, and metadata features and facilitates timely interventions to combat misinformation, thereby safeguarding information integrity in digital environments [90, 9, 22].

## 6.5 Interdisciplinary Collaboration and Ethical Considerations

Interdisciplinary collaboration is crucial in addressing the multifaceted challenges posed by fake news and misinformation. Integrating insights from computer science, social sciences, and journalism is essential for developing comprehensive strategies that effectively counter misinformation [24]. Such collaboration fosters the creation of advanced methodologies that consider psychological and social dynamics, leading to more effective consumer-facing tools [26]. This approach is further supported by collaboration between linguists and technologists, enhancing misinformation detection systems' efficacy [48]. The importance of dataset accessibility is paramount, as it prevents redundancy in dataset construction and facilitates broader research efforts [72]. Ensuring datasets are readily available can significantly enhance robust detection model development. Moreover, user education regarding algorithmic processes is crucial, as interpretable algorithms build trust and accountability in news dissemination [97]. Ethical considerations are equally critical in developing and deploying fake news detection systems. The ethical implications of controlled access to models and data must be carefully considered to prevent misuse and ensure responsible technology use [93]. Maintaining user privacy and avoiding censorship are paramount in interdisciplinary collaborations aimed at combating misinformation [10]. Transparency and accountability in dataset usage are essential to address privacy and consent concerns, thereby minimizing negative societal impacts [59].

# 7 Conclusion

This survey underscores the significant advancements in fake news detection, driven by the integration of natural language processing and machine learning. Techniques such as the three-level hierarchical attention network (3HAN) have enhanced classification accuracy, demonstrating the potential of these technologies in combating misinformation. The development of early detection models highlights the importance of timely intervention, achieving notable accuracy and underscoring the need for further refinement of these approaches.

Document encoding methods, including the DOCEMB approach, illustrate the effectiveness of encoding techniques that rival complex deep learning models, emphasizing the ongoing need for research into adaptable and robust detection methods. Fact-based detection strategies have emerged as crucial for delivering reliable and interpretable results, fostering trust in these systems. The introduction of diverse datasets, particularly those encompassing non-English languages, broadens the applicability of detection models, promoting inclusivity and efficacy across different linguistic landscapes.

However, the challenges of developing sophisticated datasets and models that transcend simple binary classifications remain. Addressing these complexities requires sustained interdisciplinary collaboration to devise comprehensive strategies that effectively meet the evolving challenges of fake news and misinformation in the digital age.

# References

[1] Flavio Merenda and José Manuel Gómez-Pérez. Capturing pertinent symbolic features for enhanced content-based misinformation detection, 2024.

[2] Sneha Singhania, Nigel Fernandez, and Shrisha Rao. 3han: A deep neural network for fake news detection, 2023.

[3] Yuzhou Yang, Yangming Zhou, Qichao Ying, Zhenxing Qian, Dan Zeng, and Liang Liu. Fact-checking based fake news detection: a review, 2024.

[4] Yang Yang, Lei Zheng, Jiawei Zhang, Qingcai Cui, Zhoujun Li, and Philip S. Yu. Ti-cnn: Convolutional neural networks for fake news detection, 2023.

[5] Xinqi Su, Yawen Cui, Ajian Liu, Xun Lin, Yuhao Wang, Haochen Liang, Wenhui Li, and Zitong Yu. Daad: Dynamic analysis and adaptive discriminator for fake news detection, 2024.

[6] A proposed bi-lstm method to fake news detection.

[7] Nathaniel Hoy and Theodora Koulouri. A systematic review on the detection of fake news articles, 2021.

[8] Xinyi Zhou, Jiayu Li, Qinzhou Li, and Reza Zafarani. Linguistic-style-aware neural networks for fake news detection, 2023.

[9] Niraj Sitaula, Chilukuri K. Mohan, Jennifer Grygiel, Xinyi Zhou, and Reza Zafarani. Credibility-based fake news detection, 2019.

[10] Mir Mehedi A. Pritom, Rosana Montanez Rodriguez, Asad Ali Khan, Sebastian A. Nugroho, Esra'a Alrashydah, Beatrice N. Ruiz, and Anthony Rios. Case study on detecting covid-19 health-related misinformation in social media, 2021.

[11] Adhish S. Sujan, Ajitha. V, Aleena Benny, Amiya M. P., and V. S. Anoop. Malfake: A multimodal fake news identification for malayalam using recurrent neural networks and vgg-16, 2023.

[12] Kai Shu, Deepak Mahudeswaran, Suhang Wang, and Huan Liu. Hierarchical propagation networks for fake news detection: Investigation and exploitation, 2019.

[13] Álvaro Ibrain Rodríguez and Lara Lloret Iglesias. Fake news detection using deep learning, 2019.

[14] Sajjad Ahmed, Knut Hinkelmann, and Flavio Corradini. Combining machine learning with knowledge engineering to detect fake news in social networks-a survey, 2022.

[15] Bimal Bhattarai, Ole-Christoffer Granmo, and Lei Jiao. Explainable tsetlin machine framework for fake news detection with credibility score assessment, 2021.

[16] Mohammad Hadi Goldani, Saeedeh Momtazi, and Reza Safabakhsh. Detecting fake news with capsule neural networks, 2020.

[17] V. Mazzeo, A. Rapisarda, and G. Giuffrida. Detection of fake news on covid-19 on web search engines, 2021.

[18] Zhixuan Zhou, Huankang Guan, Meghana Moorthy Bhat, and Justin Hsu. Fake news detection via nlp is vulnerable to adversarial attacks, 2019.

[19] Julien Delaunay, Luis Galárraga, and Christine Largouët. Does it make sense to explain a black box with another black box?, 2024.

[20] Lovedeep Singh. Hybrid ensemble for fake news detection: An attempt, 2022.

[21] Verónica Pérez-Rosas, Bennett Kleinberg, Alexandra Lefevre, and Rada Mihalcea. Automatic detection of fake news. *arXiv preprint arXiv:1708.07104*, 2017.

[22] Verónica Pérez-Rosas, Bennett Kleinberg, Alexandra Lefevre, and Rada Mihalcea. Automatic detection of fake news, 2017.

[23] Ray Oshikawa, Jing Qian, and William Yang Wang. A survey on natural language processing for fake news detection. *arXiv preprint arXiv:1811.00770*, 2018.

[24] Xinyi Zhou and Reza Zafarani. A survey of fake news: Fundamental theories, detection methods, and opportunities, 2020.

[25] Xinyi Zhou and Reza Zafarani. A survey of fake news: Fundamental theories, detection methods, and opportunities. *ACM Computing Surveys (CSUR)*, 53(5):1–40, 2020.

[26] Alireza Karduni. Human-misinformation interaction: Understanding the interdisciplinary approach needed to computationally combat false information, 2019.

[27] Maaz Amjad, Grigori Sidorov, Alisa Zhila, Alexander Gelbukh, and Paolo Rosso. Overview of the shared task on fake news detection in urdu at fire 2020, 2022.

[28] Giancarlo Ruffo and Alfonso Semeraro. Fakenewslab: Experimental study on biases and pitfalls preventing us from distinguishing true from false news, 2022.

[29] Guanghua Li, Wensheng Lu, Wei Zhang, Defu Lian, Kezhong Lu, Rui Mao, Kai Shu, and Hao Liao. Re-search for the truth: Multi-round retrieval-augmented large language models are strong fake news detectors, 2024.

[30] Sakshini Hangloo and Bhavna Arora. Fake news detection tools and methods – a review, 2021.

[31] Jiawei Zhang, Bowen Dong, and S Yu Philip. Fakedetector: Effective fake news detection with deep diffusive neural network. In *2020 IEEE 36th international conference on data engineering (ICDE)*, pages 1826–1829. IEEE, 2020.

[32] Jiaying Wu, Jiafeng Guo, and Bryan Hooi. Fake news in sheep's clothing: Robust fake news detection against llm-empowered style attacks, 2024.

[33] Kai Shu, Suhang Wang, and Huan Liu. Beyond news contents: The role of social context for fake news detection, 2018.

[34] Xinyi Zhou and Reza Zafarani. Network-based fake news detection: A pattern-driven approach, 2019.

[35] Mesay Gemeda Yigezu, Melkamu Abay Mersha, Girma Yohannis Bade, Jugal Kalita, Olga Kolesnikova, and Alexander Gelbukh. Ethio-fake: Cutting-edge approaches to combat fake news in under-resourced languages using explainable ai, 2024.

[36] Philipp Hartl and Udo Kruschwitz. Applying automatic text summarization for fake news detection, 2022.

[37] Daryna Dementieva, Mikhail Kuimov, and Alexander Panchenko. Multiverse: Multilingual evidence for fake news detection, 2022.

[38] Jędrzej Kozal, Michał Leś, Paweł Zyblewski, Paweł Ksieniewicz, and Michał Woźniak. Lifelong learning natural language processing approach for multilingual data classification, 2022.

[39] Xinyi Zhou, Reza Zafarani, and Emilio Ferrara. From fake news to fakenews: Mining direct and indirect relationships among hashtags for fake news detection, 2022.

[40] Dongqi Fu, Yikun Ban, Hanghang Tong, Ross Maciejewski, and Jingrui He. Disco: Comprehensive and explainable disinformation detection, 2022.

[41] Yi Han, Amila Silva, Ling Luo, Shanika Karunasekera, and Christopher Leckie. Knowledge enhanced multi-modal fake news detection, 2021.

[42] Daniel Toma and Wasim Huleihel. Sequential classification of misinformation, 2024.

[43] Ciprian-Octavian Truică and Elena-Simona Apostol. Misrobærta: Transformers versus misinformation, 2023.

[44] Bin Guo, Yasan Ding, Lina Yao, Yunji Liang, and Zhiwen Yu. The future of misinformation detection: New perspectives and trends, 2019.

[45] Hao Chen, Hui Guo, Baochen Hu, Shu Hu, Jinrong Hu, Siwei Lyu, Xi Wu, and Xin Wang. A self-learning multimodal approach for fake news detection, 2024.

[46] Maria Glenski, Ellyn Ayton, Josh Mendoza, and Svitlana Volkova. Multilingual multimodal digital deception detection and disinformation spread across social platforms, 2019.

[47] Santiago Alonso-Bartolome and Isabel Segura-Bedmar. Multimodal fake news detection, 2021.

[48] Xinyu Wang, Wenbo Zhang, and Sarah Rajtmajer. Monolingual and multilingual misinformation detection for low-resource languages: A comprehensive survey, 2024.

[49] Sara Abdali, Sina shaham, and Bhaskar Krishnamachari. Multi-modal misinformation detection: Approaches, challenges and opportunities, 2024.

[50] Gullal S. Cheema, Sherzod Hakimov, Abdul Sittar, Eric Müller-Budack, Christian Otto, and Ralph Ewerth. Mm-claims: A dataset for multimodal claim detection in social media, 2022.

[51] Jeff Z Pan, Siyana Pavlova, Chenxi Li, Ningxi Li, Yangmei Li, and Jinshuo Liu. Content based fake news detection using knowledge graphs. In *The Semantic Web–ISWC 2018: 17th International Semantic Web Conference, Monterey, CA, USA, October 8–12, 2018, Proceedings, Part I 17*, pages 669–683. Springer, 2018.

[52] Shuzhi Gong, Richard O. Sinnott, Jianzhong Qi, and Cecile Paris. Fake news detection through graph-based neural networks: A survey, 2023.

[53] Pranav Goel, Jon Green, David Lazer, and Philip Resnik. Mainstream news articles co-shared with fake news buttress misinformation narratives, 2023.

[54] Gregor Donabauer and Udo Kruschwitz. Challenges in pre-training graph neural networks for context-based fake news detection: An evaluation of current strategies and resource limitations, 2024.

[55] Haoming Guo, Tianyi Huang, Huixuan Huang, Mingyue Fan, and Gerald Friedland. Detecting covid-19 conspiracy theories with transformers and tf-idf, 2022.

[56] Piotr Przybyła, Alexander Shvets, and Horacio Saggion. Verifying the robustness of automatic credibility assessment, 2024.

[57] Kahlil bin Abdul Hakim and Sathishkumar Veerappampalayam Easwaramoorthy. Impact of fake news on social media towards public users of different age groups, 2024.

[58] William Yang Wang. "liar, liar pants on fire": A new benchmark dataset for fake news detection, 2017.

[59] Soveatin Kuntur, Anna Wróblewska, Marcin Paprzycki, and Maria Ganzha. Fake news detection: It's all in the data!, 2024.

[60] Chahat Raj and Priyanka Meel. People lie, actions don't! modeling infodemic proliferation predictors among social media users, 2021.

[61] Kung-Hsiang Huang, Kathleen McKeown, Preslav Nakov, Yejin Choi, and Heng Ji. Faking fake news for real fake news detection: Propaganda-loaded training data generation, 2023.

[62] Camille Koenders, Johannes Filla, Nicolai Schneider, and Vinicius Woloszyn. How vulnerable are automatic fake news detection methods to adversarial attacks?, 2021.

[63] Bruno Tafur and Advait Sarkar. User perceptions of automatic fake news detection: Can algorithms fight online misinformation?, 2023.

[64] Kellin Pelrine, Jacob Danovitch, and Reihaneh Rabbany. The surprising performance of simple baselines for misinformation detection, 2021.

[65] William Yang Wang. " liar, liar pants on fire": A new benchmark dataset for fake news detection. *arXiv preprint arXiv:1705.00648*, 2017.

[66] Matthew Iceland. How good are sota fake news detectors, 2023.

[67] Shafna Fitria Nur Azizah, Hasan Dwi Cahyono, Sari Widya Sihwi, and Wisnu Widiarto. Performance analysis of transformer based models (bert, albert and roberta) in fake news detection, 2023.

[68] Ronald Denaux and Jose Manuel Gomez-Perez. Linked credibility reviews for explainable misinformation detection, 2020.

[69] Junaed Younus Khan, Md. Tawkat Islam Khondaker, Sadia Afroz, Gias Uddin, and Anindya Iqbal. A benchmark study of machine learning models for online fake news detection, 2021.

[70] Gullal S. Cheema, Sherzod Hakimov, Eric Müller-Budack, and Ralph Ewerth. On the role of images for analyzing claims in social media, 2021.

[71] Yizhou Zhang, Loc Trinh, Defu Cao, Zijun Cui, and Yan Liu. Interpretable detection of out-of-context misinformation with neural-symbolic-enhanced large multimodal model, 2024.

[72] Taichi Murayama. Dataset of fake news detection and fact verification: A survey, 2021.

[73] Bing He, Yibo Hu, Yeon-Chang Lee, Soyoung Oh, Gaurav Verma, and Srijan Kumar. A survey on the role of crowds in combating online misinformation: Annotators, evaluators, and creators, 2024.

[74] Nguyen Vo and Kyumin Lee. Learning from fact-checkers: Analysis and generation of fact-checking language, 2019.

[75] Kai Shu, Deepak Mahudeswaran, Suhang Wang, Dongwon Lee, and Huan Liu. Fakenewsnet: A data repository with news content, social context and spatialtemporal information for studying fake news on social media, 2019.

[76] Ye Liu, Jiajun Zhu, Xukai Liu, Haoyu Tang, Yanghai Zhang, Kai Zhang, Xiaofang Zhou, and Enhong Chen. Detect, investigate, judge and determine: A knowledge-guided framework for few-shot fake news detection, 2025.

[77] Athira A B, S D Madhu Kumar, and Anu Mary Chacko. Towards smart fake news detection through explainable ai, 2022.

[78] Ahmed Akib Jawad Karim, Kazi Hafiz Md Asad, and Aznur Azam. Strengthening fake news detection: Leveraging svm and sophisticated text vectorization techniques. defying bert?, 2024.

[79] Anusua Trivedi, Alyssa Suhm, Prathamesh Mahankal, Subhiksha Mukuntharaj, Meghana D. Parab, Malvika Mohan, Meredith Berger, Arathi Sethumadhavan, Ashish Jaiman, and Rahul Dodhia. Defending democracy: Using deep learning to identify and prevent misinformation, 2021.

[80] Sina Mohseni, Eric Ragan, and Xia Hu. Open issues in combating fake news: Interpretability as an opportunity, 2019.

[81] Taichi Murayama, Shohei Hisada, Makoto Uehara, Shoko Wakamiya, and Eiji Aramaki. Annotation-scheme reconstruction for "fake news" and japanese fake news dataset, 2022.

[82] Yuhang Wang, Li Wang, Yanjie Yang, and Yilin Zhang. Detecting fake news by enhanced text representation with multi-edu-structure awareness, 2022.

[83] Ganesh Jawahar, Muhammad Abdul-Mageed, and Laks V. S. Lakshmanan. Automatic detection of machine generated text: A critical survey, 2020.

[84] Jinyan Su, Claire Cardie, and Preslav Nakov. Adapting fake news detection to the era of large language models, 2024.

[85] Fatma Shalabi, Huy H. Nguyen, Hichem Felouat, Ching-Chun Chang, and Isao Echizen. Image-text out-of-context detection using synthetic multimodal misinformation, 2024.

[86] Fake news: Fundamental theories.

[87] Luiz Giordani, Gilsiley Darú, Rhenan Queiroz, Vitor Buzinaro, Davi Keglevich Neiva, Daniel Camilo Fuentes Guzmán, Marcos Jardel Henriques, Oilson Alberto Gonzatto Junior, and Francisco Louzada. fakenewsbr: A fake news detection platform for brazilian portuguese, 2023.

[88] Yasan Ding, Bin Guo, Yan Liu, Yunji Liang, Haocheng Shen, and Zhiwen Yu. Metadetector: Meta event knowledge transfer for fake news detection, 2021.

[89] Algorithmic detection of misinfo.

[90] Shaily Bhatt, Sakshi Kalra, Naman Goenka, and Yashvardhan Sharma. Fake news detection: Experiments and approaches beyond linguistic features, 2021.

[91] Inna Vogel and Meghana Meghana. Fake news spreader detection on twitter using character n-grams. notebook for pan at clef 2020, 2020.

[92] Tanveer Khan, Antonis Michalas, and Adnan Akhunzada. Sok: Fake news outbreak 2021: Can we stop the viral spread?, 2021.

[93] Herun Wan, Shangbin Feng, Zhaoxuan Tan, Heng Wang, Yulia Tsvetkov, and Minnan Luo. Dell: Generating reactions and explanations for llm-based misinformation detection, 2024.

[94] Hao Liao, Jiaohao Peng, Zhanyi Huang, Wei Zhang, Guanghua Li, Kai Shu, and Xing Xie. Muser: A multi-step evidence retrieval enhancement framework for fake news detection, 2023.

[95] Michal Choras, Konstantinos Demestichas, Agata Gielczyk, Alvaro Herrero, Pawel Ksieniewicz, Konstantina Remoundou, Daniel Urda, and Michal Wozniak. Advanced machine learning techniques for fake news (online disinformation) detection: A systematic mapping study, 2020.

[96] Eugenio Tacchini, Gabriele Ballarin, Marco L. Della Vedova, Stefano Moret, and Luca de Alfaro. Some like it hoax: Automated fake news detection in social networks, 2017.

[97] Sina Mohseni and Eric Ragan. Combating fake news with interpretable news feed algorithms, 2018.

**Disclaimer:**

SurveyX is an AI-powered system designed to automate the generation of surveys. While it aims to produce high-quality, coherent, and comprehensive surveys with accurate citations, the final output is derived from the AI's synthesis of pre-processed materials, which may contain limitations or inaccuracies. As such, the generated content should not be used for academic publication or formal submissions and must be independently reviewed and verified. The developers of SurveyX do not assume responsibility for any errors or consequences arising from the use of the generated surveys.