Rethink Rumor Detection in the Era of LLM: A Review

Anonymous ACL submission

Abstract

The rise of large language models (LLMs) has fundamentally reshaped the technological paradigm of rumor detection, offering transformative opportunities to construct adaptive detection systems while simultaneously ushering in new threats, such as "logically flawless" rumors. This paper focuses on modeling rumor detection in the era of LLMs by unifying existing methods in the field of rumor detection and uncovering their underlying logical mechanisms. From the perspective of complex systems, we innovatively propose a "Cognition-Interaction-Behavior" (CIB) trilevel framework for rumor detection based on collective intelligence and explore the synergistic relationship between LLMs and collective intelligence in rumor governance. Further-018 more, we analyze the core challenges in the 019 LLM era and outline future development pathways for social simulation agents. We hope this work lays a theoretical foundation for nextgeneration rumor detection paradigms and offers valuable insights for advancing the field.

1 Introduction

011

017

024

033

037

041

In the digital era, the widespread adoption of social media and the explosion of user-generated content has enabled rumors to threaten public safety and social trust at unprecedented speeds, scales, and levels of complexity (Shao et al., 2016; Kim and Dennis, 2019). Meanwhile, the rapid advancements in large language models (LLMs) have demonstrated remarkable performance across various fields (Tan et al., 2023; Poldrack et al., 2023), but it has also brought challenges that cannot be ignored. Models like GPT-4(Achiam et al., 2023) and DeepSeek(Guo et al., 2025), known for their deep semantic understanding and reasoning capabilities, can generate highly credible and logically coherent professional content. However, this ability can also be used to generate "logically perfect rumors" (such as false arguments based on chain reasoning),

which are far more concealed and misleading than traditional generation methods.(Bommasani et al., 2021; Kreps et al., 2022). For example, studies have shown that ChatGPT, when provided with malicious prompts, can not only optimize deceptive text but also proactively enhance their disguise by incorporating additional misleading details (Augenstein et al., 2024). Thus, leveraging the powerful capabilities of LLMs while addressing their inherent limitations has emerged as an urgent challenge in the field of rumor detection.

042

043

044

047

048

053

054

056

060

061

062

063

064

065

066

067

068

069

070

071

072

073

074

078

079

081

Existing rumor detection surveys primarily focus on the dissemination mechanisms of rumors on social media (Shu et al., 2017; Del Vicario et al., 2016; Johnson et al., 2020), the psychological mechanisms underlying belief in rumors (Roozenbeek et al., 2020), and effective intervention strategies (Zubiaga et al., 2015; Guess et al., 2020). However, most existing frameworks primarily rely on feature-based or technical classifications, which result in two primary issues: (1) the failure to thoroughly explore the theoretical and logical connections between detection methods and rumor propagation mechanisms, and (2) the inability to effectively reveal the intrinsic relationships between features, especially in the context of research on LLMs in this field (Chen and Shu, 2024).

To bridge this gap, we introduce a three-tiered "Cognition-Interaction-Behavior" (CIB) framework to systematically elucidate the underlying logic of rumor propagation and detection on social networks. The specific contributions of this work include: (1) A new theoretical paradigm for rumor detection. The construction of the CIB framework unifies existing rumor detection methods (as shown in Figure 4) and uncovers the multi-scale coupling mechanisms underlying rumor propagation, including collective knowledge emergence, interactive network evolution, and iterative behavioral patterns. (2) A systematic exploration of LLMs' multifaceted roles in rumor detection and their synergies



Figure 1: The three-layer architecture operates collaboratively. The cognition layer integrates multi-source evidence to provide informational support for the interaction layer. Through user interactions, the interaction layer facilitates the formation of the behavior layer. The behavior layer, in turn, continuously refines the cognition layer through accumulated experiences and collective cognitive feedback. (RD is rumor detection; CI is collective intelligence, driving information's dynamic reconstruction and optimization).

with collective intelligence, forming a more comprehensive adaptive governance system. (3) A summary of the core challenges of rumor detection in the LLM era and an outline of future development pathways for socially simulated agents.

2 Collective Intelligence-Based Rumor Detection Framework

In the social media ecosystem, user communities serve dual roles as both disseminators and evaluators, forming the self-organizing foundation of the networked information ecology. Studies have shown that through cross-validation among users and the interplay of opinions, social networks can facilitate collective cognitive correction (Ma et al., 2018). Compared to individual cognition, collective intelligence leverages the integration of diverse knowledge and dynamic interactions, demonstrating superior cognitive capabilities in addressing complex information (Castillo et al., 2011), thereby offering a novel approach to advancing rumor detection (Phan et al., 2023).

From the perspective of complex systems, the emergence of collective intelligence is essentially a self-organizing process driven by the reduction of information entropy. During this procesocial media users' the diverse cognition, social connections, and dynamic behaviers interact, facilitating information flow and collaborative evolution. Rumor diffusion, as a specific form of information dissemination, is often constrained by individuals' cognitive thresholds (e.g., cognitive abilities, emotional biases) and the topological structure of the social network. At its core, rumor diffusion can be viewed as a staged state of cognitive imbalance: it arises when users, driven by information uncertainty and emotional impetus, engage in social interactions to reduce uncertainty, which in turn drives the continuous evolution of network structures (Allcott and Gentzkow, 2017). This process generates macro-level dissemination behaviors (potentially unintentionally promoting rumor propagation). However, collective intelligence can dynamically correct such imbalanced states in social networks through multi-level knowledge sharing and interaction.

109

110

111

112

113

114

115

116

117

118

119

120

121

122

123

124

125

126

127

128

129

130

131

132

133

134

135

136

Based on the above theoretical construction, this study proposes a three-order model framework for rumor detection based on collective intelligence, as shown in Figure 1. The cognitive layer facilitates the construction of a collective knowledge system for rumor identification through knowledge sharing and evidence integration among users. It serves as the foundational support for detecting rumors and aggregating multidmulti-dimensionalnce. The in-

106

187

188

teraction layer analyzes users' social relationships 137 and interaction behaviors within social networks 138 to capture rumor signals. The behavior layer mod-139 els the evolution of information dissemination and 140 collective behavior. Finally, the feedback mech-141 anism based on collective intelligence optimizes 142 the dissemination path and reduces the spread of 143 rumors. 144

2.1 Cognitive Layer

145

146

147

149

150

151

153

154

155

157

158

159

160

161

163

164

165

166

167

169

170

171

172

173

174

175

176

177

178

179

181

182

183

186

The Cognitive Layer leverages **data-driven analysis** (§2.1.1) to explore the features of multi-source information on social networks deeply, while **knowledge-driven analysis** (§2.1.2) focuses on verifying information sources and effectively integrating multi-perspective evidence. Together, these processes construct collective knowledge as the cognitive foundation for rumor detection.

2.1.1 Data-Driven Analysis

Data-driven analysis focuses on extracting features from multi-source data within social networks, aiming to develop multidmulti-dimensionalaches for assessing information veracity. Research on linguistic semantics, visual content, and multimodal characteristics provides critical support for more in-depth rumor detection efforts.

In terms of semantic analysis, rumor texts often exhibit multi-multi-dimensionallies. Studies have revealed distinct patterns at various levels by integrating linguistic and psychological theories (Undeutsch, 1967; Zuckerman, 1981). Lexical Level: rumolevelts tend to avoid expressing indepth information (Aich et al., 2022; Horne and Adali, 2017), characterized by lower usage rates of first-person pronouns and higher proportions of adverbs and emotional words. Syntactic Level: rumolevelts exhibit a trend toward simplification (Pérez-Rosas et al., 2017), such as reduced lexical diversity and shorter average sentence lengths, abnormal frequency distributions of function words and punctuation, further highlight the "readability" characteristic of rumor texts. Stylistic Level: rumolevelts often feature exaggerated headlines, informal language, and frequent insertion of URLs or hashtags to enhance virality and attractiveness (Blass, 1984). These linguistic abnormalities also reveal rumors' strategic use of language to evoke negative emotions (e.g., anger, fear) in public, amplifying their dissemination, particularly in sensitive topics such as politics and health (Vosoughi et al., 2018). Studies have developed various detection methods through linguistic pattern analysis to leverage these characteristics to capture deeper semantic features.

In addition, deepfake technologies have significantly increased the deceptive capabilities and dissemination risks associated with image-based rumors (Vaccari and Chadwick, 2020). Early methods relied on spatial domain (Popescu and Farid, 2004) and frequency domain (Fridrich et al., 2003) analysis to extract pixel-level features for identifying common local manipulations in forged images but suffered from insufficient sensitivity to localized manipulations. Studies (Marra et al., 2019) have shown that statistical properties and spectral responses of GAN-generated content systematically differ from authentic images, providing a technological breakthrough for detection. The introduction of deep learning has further enhanced detection performance. Techniques such as local feature extraction (Bayar and Stamm, 2016; Wang et al., 2020a), temporal modeling (Güera and Delp, 2018), and the use of pre-trained models (Hao et al., 2021; Khan et al., 2022) have been effective in capturing complex visual forgery patterns. For video forgeries, beyond traditional frame-level classification methods (Montserrat et al., 2020), advancements have been made through inter-frame consistency analysis (Amerini et al., 2019), metadata validation (Huh et al., 2018), detection of visual artifacts (Matern et al., 2019), and leveraging biometric signal characteristics (Li et al., 2018). These approaches improve forgery detection performance by revealing inherent flaws in dynamic modeling. To more comprehensively capture the multimodal characteristics of rumor information, multimodal analysis focuses on feature fusion and consistency verification (Wang et al., 2018a; Jin et al., 2017). These methods balance modality-specific features and improve cross-modal detection performance by employing different fusion strategies (Singhal et al., 2019; Qian et al., 2021). Consistency verification is utilized to identify cross-modal information conflicts, including text-image inconsistencies (e.g., emotional conflicts) and audio-visual mismatches (e.g., forged videos) (Agarwal et al., 2020; Chugh et al., 2020). These approaches effectively enhance the performance of multimodal rumor detection.

2.1.2 Knowledge-Driven Analysis

The knowledge-driven analysis leverages external information resources to enrich and validate rumor content, offering critical support for social net-

241

242

243

244

247

248

249

251

252

260

261

262

264

270

272

273

275

276

281

works' noisy, concise content. Verification methods based on **knowledge graphs** and **evidence texts** are the core research content.

As structured data systems, knowledge graphs provide a networked organization of large-scale entities and relationships, supporting rumor detection through contextual verification and logical reasoning. By matching entities and relationships within the text, knowledge graphs can quickly validate content accuracy and identify potential contradictions (Cui et al., 2020). Several studies (Hu et al., 2021) further integrate semantic analysis and graph reasoning techniques to uncover implicit associations or supporting evidence, thereby improving verification reliability. For rumor content characterized by semantic ambiguity or missing information, knowledge graphs leverage their capabilities in semantic completion and reasoning to explore hidden, deeper-level entity associations. By incorporating contextual information (Dun et al., 2021) and multimodal data (Wang et al., 2020b), knowledge graphs can fill in critical gaps within implicit or ambiguous statements. Furthermore, they have demonstrated robust adaptability in cross-domain rumor detection tasks (Sun et al., 2022; Zhang et al., 2019).

On the other hand, evidence-based text verification focuses on fact-checking, aiming to validate the factuality of rumor content through authoritative information sources such as news articles, scientific literature, and fact-checking platforms. Traditionally, this task has relied on expertdriven manual verification. Internationally accredited organizations such as the International Fact-Checking Network (IFCN) (Porter and Wood, 2021), the European Fact-Checking Standards Network (EFCSN) (Wouters and Opgenhaffen, 2024), and government platforms (e.g., China's Internet Joint Rumor Debunking Platform, Tencent Fact-Check Platform) provide standardized evaluation procedures to assess the authenticity and timeliness of evidence, delivering high-quality validation services to the public (Vlachos and Riedel, 2014).

However, the expert-driven model faces efficiency bottlenecks and struggles to respond rapidly to large-scale, real-time information dissemination demands (Das et al., 2023; Guo et al., 2022). Automated fact-checking has increasingly become a focus of research to address these limitations. Key processes in these systems include multisource data retrieval, semantic alignment, and logical reasoning. By leveraging deep learning models combined with information retrieval methods (Hanselowski et al., 2019), semantic relevance between rumors and evidence is extracted from authoritative data sources such as Wikipedia and scientific literature (Schuster et al., 2021; Wadden et al., 2021). Additionally, semantic alignment techniques are employed to assess the extent to which the retrieved evidence supports or refutes the claims embedded in the rumor. For complex, multi-layered claims, deep learning methodologies (Zhong et al., 2019) can generate reliable verification results. Explainability-enhancing techniques are also applied to extract key logic and evidence chains, improving users' understanding of and trust in the verification outcomes (Lu and Li, 2020).

291

292

293

295

296

297

298

299

300

301

302

303

304

305

306

307

308

309

310

311

312

313

314

315

316

317

318

319

320

321

322

323

324

325

326

327

328

329

330

331

332

333

334

335

336

337

338

339

2.2 Interaction layer

The Interaction Layer emphasizes the analysis of **user feature (§2.2.1))** and **social context (§2.2.2)** to provide comprehensive contextual information, including individual behavior patterns, interactions between groups, and the dynamic formation of group consensus. This not only reveals the processes through which information spreads among users but also captures the collaboration and conflicts involved in users' efforts to discern and verify rumors.

2.2.1 User Feature Analysis

User groups within social networks serve as the core driving force behind rumor propagation. They are comprised of genuine users and social bots tasked with content generation and dissemination. Analyzing user characteristics and behavioral patterns can effectively uncover rumor dissemination's underlying mechanisms and risks.

As specialized accounts, social bots often manipulate public opinion through high-frequency content posting and synchronized interactions, significantly accelerating rumor dissemination. Bot detection methods can be broadly categorized into feature-based and network-based approaches. The feature-based analysis identifies non-human attributes, such as anomalous metadata (e.g., default profile pictures, short-lived accounts), highfrequency posting behaviors, and polarized content (e.g., repetitive sentence structures, simple semantics). Network-based detection detects bots by identifying abnormalities in dissemination structures. Social bots often amplify their influence by forming densely interconnected groups or fabricating community structures. In recent years, deep learning techniques, such as Graph Neural

Networks (GNNs), have been employed to model 340 the dependencies within network topologies (Guo 341 et al., 2021). When combined with content features, 342 pre-trained language models have further enhanced the generalization performance of social bot detection in heterogeneous environments (Haider et al., 2023).

Genuine users are the central driving force behind information dissemination (Aral and Walker, 2012). Among them, malicious users deliberately 349 spread rumors for personal or organizational gain (Kahneman, 1979), while ordinary users may inadvertently propagate rumors due to cognitive limitations or emotional triggers. Research has delved into the dissemination patterns of these users by 354 examining their static attributes (e.g., registration time, geographical location, number of friends) and dynamic behavioral characteristics (e.g., posting frequency, emotional traits such as anger and fear, 358 and account credibility scores) (Chu et al., 2012). For example, anomalies such as short-term, highfrequency interactions and highly repeated content 361 releases are often considered potential signals of rumor spreading (Zhao et al., 2014). Additionally, as 364 rumor dissemination increasingly transcends platform boundaries, cross-platform user identity association analysis has emerged as a research focus. For example, by matching user profiles (Iof-367 ciu et al., 2011), content geolocation (Riederer et al., 2016), and personalized traits such as posting style (Goga et al., 2013) and interest prefer-370 ences (Nie et al., 2016), researchers can identify attempts by malicious users to disguise their identities across platforms. Studies have increasingly applied unified analyses of static and dynamic be-374 haviors across broader network ecosystems by in-375 tegrating deep learning and network analysis techniques (Hamdi et al., 2020; Zhang et al., 2015; Zhou et al., 2015). 378

2.2.2 Social Context Analysis

Rumor propagation is influenced by individual be-

haviors and the deeper constraints imposed by so-

cial network structures and user interaction patterns.

By uncovering the synergies between user interac-

tions and network structures, social context analy-

sis provides critical support for understanding the

work structures exhibit distinct patterns in rumor

propagation. Due to a lack of supervision and infor-

mation verification mechanisms, Sparse networks

The interaction characteristics within social net-

mechanisms underlying the diffusion of rumors.

371

390

tend to form low-cohesion, flat diffusion structures, which accelerate the spread of false information (Vosoughi et al., 2018). In contrast, dense networks, with their strong connectivity, can partially filter or curb the propagation of rumors. Moreover, the heterogeneity of user roles within a network (e.g., "messengers" bridging multiple communities or "skeptics" constructing local verification networks) dynamically impacts interaction structures (Raponi et al., 2022). For example, user comment chains and resharing behaviors can be modeled as propagation tree structures (Kwon et al., 2013), which are utilized to capture both top-down and bottomup propagation patterns (Alrubaian et al., 2016; Ma et al., 2015). To more comprehensively represent such complex social contexts, multiple entities such as users, posts, and hashtags can be modeled as propagation graphs (Nguyen et al., 2020; Shu et al., 2020). Studies also introduce graph neural networks (Bian et al., 2020; Min et al., 2022), which aggregate node features to capture complex interaction relationships and diffusion dynamics among users. These significantly enhance the efficacy of modeling rumor propagation patterns and their underlying mechanisms.

391

392

393

394

395

396

397

398

399

400

401

402

403

404

405

406

407

408

409

410

411

412

413

414

415

416

417

418

419

420

421

422

423

424

425

426

427

428

429

430

431

432

433

434

435

436

437

438

439

440

441

442

User interaction patterns serve as critical clues for uncovering rumor social context. Studies have shown that genuine users, fake news producers, and hybrid users tend to form homogeneous clusters within collaborative networks. The polarization between different clusters reflects regional differences and political and ideological divisions. This phenomenon is particularly pronounced in rumor propagation, where emotion acts as a catalyst. Emotionally charged content, by evoking negative emotions such as anger and fear, triggers group polarization and amplifies the speed and scope of rumor dissemination (Zeng and Zhu, 2019; Pröllochs et al., 2021). According to the "two-step flow" theory in communication studies (Katz, 1957), information flows first from mass media to key opinion leaders (KOLs), who then influence broader audiences. KOLs often act as community bridges (Yang et al., 2018), playing a significant amplifying role in rumor dissemination while also having the potential to contribute effectively to rumor debunking. Modeling KOLs is crucial for understanding their influence mechanisms in rumor propagation (Wei and Meng, 2021). The most common approaches for KOL modeling involve social network analysis techniques, utilizing centrality metrics (Opsahl et al., 2010), statistical models (Amor et al., 2016),

network topology analysis (Zhao et al., 2016), and
deep learning methods (Shafiq et al., 2013). These
methods aid in designing effective intervention
strategies.

2.3 Behavior Layer

447

448

449

450

451

452

453

454

455

456

457

458

459

460

461

462

463

464

465

466

467

468

469

470

471

472

473

474 475

476

477

478

479

480

481

482

483

484

485

486

487

488

489

490

491

492

Rampant deepfakes and false rumors are often blamed as key culprits in influencing voter behavior. Studies suggest that while changing people's political opinions is challenging, influencing their actions is comparatively easier (Adam). The spread of rumors on social networks has been extensively studied from both macro and micro perspectives (Xuan et al., 2019). At the macro level, reselevelleverages **propagation pattern modeling** (§2.3.1) to analyze the dynamic processes of rumor propagation in complex networks (Zhu and Huang, 2019). At the micro level, studlevelocus on **behavioral pattern analysis** (§2.3.2) to analyze community characteristics and predict rumor dissemination (Alkhodair et al., 2020).

2.3.1 Propagation Pattern Modeling

Propagation modeling aims to uncover the diffusion patterns of rumors within the framework of a complex sociodynamic system, focusing on the coupled process of information dissemination and collective behaviors. Epidemic analogy models serve as the foundational approach in early studies, where state transition mechanisms (e.g., Susceptible-Infectious) are used to describe the spread of rumors across network topologies (e.g., SI (Kermack and McKendrick, 1927), SIS (Dong and Huang, 2018), SIR (Zhao et al., 2013), and their variants (Wan et al., 2017). Threshold-based diffusion models, such as the Linear Threshold Model (LT) (Chen et al., 2012) and the independent cascade model (IC), shift toward modeling the propagation process from the perspective of audience decision-making. These models are also used to design intervention strategies, such as node blocking or link disruption, to suppress diffusion (Yan et al., 2019). With the development of complex network theory, more studies focus on exploring multidmulti-dimensionalrs that influence rumor propagation in online social networks: temporal dimension (Tripathy et al., 2010), user dimension (Hosni et al., 2018, 2020), network dimension (Wang et al., 2018b), and information dimension (Xiao et al., 2019), providing a more systematic theoretical foundation for studying the diffusion of rumors in complex networks.

Rumor source detection based on propagation models plays a pivotal role in tracing the origins of information dissemination, providing crucial insights for intervention strategies. Early methods, grounded in centrality theory, estimate the importance of the nodes by traversing the global topology of social networks to identify the source nodes (Ali et al., 2020). However, these approaches often suffer from high computational complexity. Snapshot observation methods (Louni and Subbalakshmi, 2018) improve detection efficiency by extracting limited node infection states or propagation paths. These methods effectively perform on homogeneous and heterogeneous propagation models (Cai et al., 2018). Additionally, monitoring-based observation techniques (Qiu et al., 2022) leverage sensor nodes to capture real-time dissemination data, enhancing adaptability to dynamic propagation environments. In scenarios involving multi-source concurrent propagation (Zhu et al., 2022a), research focuses on decoupling infection networks into several independent regions. A divide-and-conquer strategy is then employed to locate the sources iteratively, reducing the computational complexity of detection. With the maturation of deep learning techniques (Wang et al., 2022; Ling et al., 2022) and Graph Neural Networks (GNNs) (Cheng et al., 2024), end-to-end frameworks have emerged as highly effective tools. By integrating propagation paths, temporal dynamics, and node characteristics, these approaches significantly enhance the robustness and accuracy of rumor source detection.

493

494

495

496

497

498

499

500

501

502

503

504

505

506

507

508

509

510

511

512

513

514

515

516

517

518

519

520

521

522

523

524

525

526

527

528

529

530

531

532

533

534

535

536

537

538

539

540

541

542

543

2.3.2 Behavioral Pattern Analysis

In social networks, nodes often cluster into tightlyknit communities, where rumors spread efficiently within communities but rely on bridge nodes or weak ties between communities for crosscommunity dissemination. Research shows that when bridge nodes are scarce, rumor dissemination remains localized, but when sufficiently abundant, it penetrates communities and reaches a broader audience (Zanette, 2001). To identify these critical nodes, community detection methods (Newman, 2004; Blondel et al., 2008) usually adopt heuristics to find closely related subgroups in the network (Zhang et al., 2018; Yang et al., 2016). Targeted interventions and immunization strategies can then be applied to these key nodes to minimize the spread of rumors.

Network immunization strategies are typically categorized into preventive and counteractive im-

munization: Preventive Immunization focuses 544 on proactively optimizing network structures and 545 key node distributions before harmful information 546 emerges. By analyzing topological properties (e.g., node centrality, spectral attributes) and community structure characteristics, high-risk nodes, and 549 cross-community bridge points can be identified to 550 disseminate accurate information (Petrescu et al., 2021). Counteractive immunization emphasizes real-time intervention during the propagation pro-553 cess. When initial sources or infected nodes are 554 known, dynamic detection and analysis of infected 555 nodes and their neighboring communities are per-556 formed. Key high-dissemination nodes are then 557 selected for removal or blocking to efficiently disrupt the propagation chain with minimal cost (Fan et al., 2013). It prioritizes local optimization under resource constraints, such as minimizing the dissemination of malicious information within the 562 network (Tariq et al., 2017), for instance, reducing the probability of some users sharing false content with their network connections.

3 Collective Intelligence-based Rumor Detection in the LLM Era

566

567

568

569

570

574

584

585

586

588

589

592

With the advancement of Natural Language Processing (NLP), rumor detection techniques have evolved from traditional statistical methods (Lim et al., 2017; Rayana and Akoglu, 2015) to deep learning models (Ma et al., 2016; Chen et al., 2019), and more recently to pre-trained language models (PLMs) (Kaliyar et al., 2021; Pelrine et al., 2021), progressively transitioning from static feature extraction to automated dynamic semantic analysis. However, traditional methods often suffer from limited adaptability in scenarios such as early-stage rumor detection and complex environments due to the scarcity of annotated data (He et al., 2021). LLMs further transform the rumor detection landscape with their extensive pre-training and contextual reasoning capabilities, enabling efficient reasoning in resource-poor environments and significantly improving the transparency and interpretability of detection results.

3.1 LLM-enhanced CIB Framework

LLMs play a multifaceted role in existing rumor detection approaches, deeply integrating into various roles ranging from knowledge analysis to adversarial defense, as shown in Figure 2. LLMs have brought revolutionary advancements to traditional methods, demonstrating significant technical potential and promising application prospects.



Figure 2: Multiple roles of LLM in rumor detection

At the cognitive layer, LLMs serve as highly efficient Key Evidence Extractors through largescale pretraining. Unlike traditional rumor detection systems that rely on static, structured approaches(such as knowledge graphs) to expand knowledge, LLMs leverage their implicit encoding capabilities to capture deep semantic associations in unstructured information, which is further used as evidence to enhance the generalization ability of traditional language models (Nan et al., 2024; Yang et al., 2023a). Additionally, LLMs operate as Scenario-Adaptive Decision Makers, leveraging their zero-shot reasoning abilities to address diversified rumor scenarios efficiently without requiring fine-tuning (Li et al., 2023c; Wu et al., 2023a). Combining Retrieval-Augmented Generation (RAG) with external knowledge bases (Peng et al., 2023; Niu et al., 2024), LLMs can dynamically integrate up-to-date knowledge, resolving limitations in knowledge coverage and timeliness inherent to traditional methods. This integration also effectively reduces the likelihood of hallucination phenomena, thereby enhancing the reliability of detection outputs (Ji et al., 2023; Rawte et al., 2023).

At the interaction layer, LLMs function as **Social Tool Coordinators**, coordinating external tools (e.g., search engines, deepfake detectors) through agents to expand rumor detection capabilities further (Chern et al., 2023; Wan et al., 2024; Li et al., 2024a). Unlike traditional, static social network analysis and modeling methods, LLM-based agents can perceive social environments by combining 595

596

598

599

600

601

602

603

604

605

606

607

608

610

611

612

613

614

615

616

617

618

619

620

621

622

623

624

625

626

715

716

717

718

719

720

721

722

723

724

725

726

679

680

681

short-term (contextual learning) and long-term (external knowledge retrieval) memory. These agents are capable of planning and calling external tools dynamically, improving analytical performance.
Furthermore, generative agents (Park et al., 2023) can simulate user interaction behaviors, driving a paradigm shift in rumor detection from static feature modeling to dynamic interaction simulation.

628

634

637

641

642

648

654

663

670

671

674

675

678

At the behavior layer, LLMs act as Rumor Analysis Experts with superior performance in tasks requiring advanced reasoning and cross-domain contextual knowledge. Traditional rumor detection approaches relied heavily on classification tasks, often requiring manually labeled large datasets. In contrast, the emergent abilities of LLMs, such as Chain of Thought (CoT) reasoning, can decompose complex problems into intermediate reasoning steps, thereby significantly improving logical transparency and explainability (Zhang and Gao, 2023a). LLMs also exhibit robust cross-domain transferability (Cao et al., 2023c,b), enabling unified reasoning across multimodal inputs, including text, images, and audio (Yao et al., 2023a). This addresses the limitations of traditional methods in multimodal fusion and shifts detection mechanisms from pattern classification to causal inference (Zhu et al., 2022b; Nan et al., 2021). Additionally, LLMs can act as Malicious Information Defenders, showcasing strong robustness in adversarial social network environments. By integrating adversarial training and red-teaming methods (Bhardwaj and Poria, 2023; OpenAI, 2023), LLMs can quickly adapt to evolving forgery techniques, overcoming the lag in model iteration and processing capabilities of traditional approaches (Wu et al., 2024b; Sun et al., 2024). For example, when tackling tasks such as detecting rumor dissemination, stylized language attacks, and deepfake content, this dynamic adaptability further enhances the robustness of rumor detection systems.

3.2 Collective Intelligence-driven LLM Detection

Network users are the primary agents of information dissemination and important participants in rumor detection. The propagation of rumors relies on user attention, trust, and further sharing behaviors. In contrast, user reports and feedback can effectively constrain this diffusion effect, which is critical in rumor detection. This user feedbackbased supervisory mechanism helps address the limitations of LLMs in adapting to dynamic scenarios by introducing ethical constraints and social consensus at the human cognitive level, thelevelinfusing greater flexibility and human-centricity into the rumor detection process (Kou et al., 2022).

Collective intelligence injects new technical implications into the development of LLM agents. Research demonstrates (Li et al., 2023b) that agents, by simulationing group behaviors modeled in sociology and economics theories, can exhibit emergent corrective mechanisms during collaborative tasks. These emergent behaviors provide theoretical support for rumor detection agents (Zhang et al., 2024b), demonstrating their significant potential in improving the simulation of social media ecosystems and other complex societal environments. For example, agents can produce socially simulated content that is indistinguishable from real-world community behavior (Park et al., 2022), simulate trust-building interactions in social dynamics (Xie et al., 2024), and facilitate harmless discussions that bridge biases and political divides, offering valuable insights into real-world phenomena (Törnberg et al., 2023).

In the field of rumor detection, existing research (Hu et al., 2025) further uses LLM-based multiagent simulation to explore the trend of rumor propagation and optimize intervention strategies and also uses tools to implement (Li et al., 2024a) realtime evaluation of information credibility based on shallow features such as language style and common sense rules. However, the complexity of the rumor detection task requires a more comprehensive approach that considers multi-level features. We discussed this in the future research route (Appendix B), paving the way for future exploration.

4 Conclusion

Based on the complex system characteristics of collective intelligence, we reconstructed the rumor detection paradigm adapted to the era of LLMs, the "Cognition-Interaction-Behavior" (CIB) framework. We emphasized the important role of LLM in rumor detection and its complementary relationship with collective intelligence. In addition, CIB can use cross-layer dynamic feedback to establish a "Macro-Micro Feedback Loop" to dynamically realize two-way rumor detection and intervention, providing a roadmap for applying LLM-based multi-agent social simulation in rumor detection.

741

742

743

745

746

747

748

749

750

751

752

753

754

755

757

762

763

765

767

772

775

776

5 Limitations

In the future research section of this paper, we propose a roadmap for collective intelligence-based 729 rumor detection agents under the CIB framework, 730 providing a comprehensive analysis of potential research challenges and corresponding directions. However, further exploration is needed to evaluate the practical application of this framework in 734 large-scale social media environments. Addition-735 ally, polarized contexts or anomalous interactions may introduce more significant complexities. To re-737 fine and optimize the framework, we will consider enhancing robustness and dynamic adaptability in complex scenarios.

References

- Alberto Acerbi and Joseph M Stubbersfield. 2023. Large language models show human-like content biases in transmission chain experiments. *Proceedings of the National Academy of Sciences*, 120(44):e2313790120.
- Josh Achiam, Steven Adler, Sandhini Agarwal, Lama Ahmad, Ilge Akkaya, Florencia Leoni Aleman, Diogo Almeida, Janko Altenschmidt, Sam Altman, Shyamal Anadkat, et al. 2023. Gpt-4 technical report. *arXiv preprint arXiv:2303.08774*.
- David Adam. Misinformation might sway elections-but not in the way that you think. *Nature*.
- Komi Afassinou. 2014. Analysis of the impact of education rate on the rumor spreading mechanism. *Physica A: Statistical Mechanics and Its Applications*, 414:43– 52.
- Shruti Agarwal, Hany Farid, Ohad Fried, and Maneesh Agrawala. 2020. Detecting deep-fake videos from phoneme-viseme mismatches. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition workshops*, pages 660–661.
- Muhammad Ahmad, Sarwan Ali, Juvaria Tariq, Imdadullah Khan, Mudassir Shabbir, and Arif Zaman.
 2020. Combinatorial trace method for network immunization. *Information Sciences*, 519:215–228.
- Ankit Aich, Souvik Bhattacharya, and Natalie Parde. 2022. Demystifying neural fake news via linguistic feature-based interpretation. In *Proceedings of the 29th International Conference on Computational Linguistics*, pages 6586–6599.
- Abdulrahman I Al-Ghadir, Aqil M Azmi, and Amir Hussain. 2021. A novel approach to stance detection in social media tweets by fusing ranked lists and sentiments. *Information Fusion*, 67:29–40.
- Syed Shafat Ali, Tarique Anwar, and Syed Afzal Murtaza Rizvi. 2020. A revisit to the infection source

identification problem under classical graph centrality measures. *Online Social Networks and Media*, 17:100061.

778

779

781

782

783

784

785

786

787

788

789

790

791

792

793

794

795

796

797

798

799

800

801

802

803

804

805

808

809

810

811

812

813

814

815

816

817

818

819

820

821

822

823

824

825

826

827

828

829

830

831

- Sarah A Alkhodair, Steven HH Ding, Benjamin CM Fung, and Junqiang Liu. 2020. Detecting breaking news rumors of emerging topics in social media. *Information Processing & Management*, 57(2):102018.
- Hunt Allcott and Matthew Gentzkow. 2017. Social media and fake news in the 2016 election. *Journal of economic perspectives*, 31(2):211–236.
- Gordon W Allport. 1947. The psychology of rumor. *Henry Holt.*
- Majed Alrubaian, Muhammad Al-Qurishi, Mohammad Mehedi Hassan, and Atif Alamri. 2016. A credibility analysis system for assessing information on twitter. *IEEE Transactions on Dependable and Secure Computing*, 15(4):661–674.
- Irene Amerini, Leonardo Galteri, Roberto Caldelli, and Alberto Del Bimbo. 2019. Deepfake video detection through optical flow based cnn. In *Proceedings of the IEEE/CVF international conference on computer vision workshops*, pages 0–0.
- Benjamin RC Amor, Sabine I Vuik, Ryan Callahan, Ara Darzi, Sophia N Yaliraki, and Mauricio Barahona. 2016. Community detection and role identification in directed networks: understanding the twitter network of the care. data debate. In *Dynamic networks and cyber-security*, pages 111–136. World Scientific.
- Markus Anderljung, Joslyn Barnhart, Anton Korinek, Jade Leung, Cullen O'Keefe, Jess Whittlestone, Shahar Avin, Miles Brundage, Justin Bullock, Duncan Cass-Beggs, et al. 2023. Frontier ai regulation: Managing emerging risks to public safety. *arXiv preprint arXiv:2307.03718*.
- Christopher MJ André, Helene FL Eriksen, Emil J Jakobsen, Luca CB Mingolla, and Nicolai B Thomsen. 2023. Detecting ai authorship: Analyzing descriptive features for ai detection.
- Fares Antaki, Samir Touma, Daniel Milad, Jonathan El-Khoury, and Renaud Duval. 2023. Evaluating the performance of chatgpt in ophthalmology: an analysis of its successes and shortcomings. *Ophthalmology science*, 3(4):100324.
- Dimosthenis Antypas, Jose Camacho-Collados, Alun Preece, and David Rogers. 2021. Covid-19 and misinformation: A large-scale lexical analysis on twitter. In Proceedings of the 59th Annual Meeting of the Association for Computational Linguistics and the 11th International Joint Conference on Natural Language Processing: Student Research Workshop, pages 119– 126.
- Elena-Simona Apostol, Özgur Coban, and Ciprian-Octavian Truică. 2023. Contain: A communitybased algorithm for network immunization. *arXiv preprint arXiv:2303.01934*.

938

939

940

941

889

890

Sinan Aral and Dylan Walker. 2012. Identifying influential and susceptible members of social networks. *Science*, 337(6092):337–341.

833

834

836

839

843

848

849

853

854

855

858

859

870

871

874

875

876

879

885

- Isabelle Augenstein, Timothy Baldwin, Meeyoung Cha, Tanmoy Chakraborty, Giovanni Luca Ciampaglia, David Corney, Renee DiResta, Emilio Ferrara, Scott Hale, Alon Halevy, et al. 2024. Factuality challenges in the era of large language models and opportunities for fact-checking. *Nature Machine Intelligence*, 6(8):852–863.
- Anton Bakhtin, Sam Gross, Myle Ott, Yuntian Deng, Marc'Aurelio Ranzato, and Arthur Szlam. 2019. Real or fake? learning to discriminate machine from human generated text. *arXiv preprint arXiv:1906.03351*.
- Eytan Bakshy, Itamar Rosenn, Cameron Marlow, and Lada Adamic. 2012. The role of social networks in information diffusion. In *Proceedings of the 21st international conference on World Wide Web*, pages 519–528.
- Ritwik Banerjee, Song Feng, Jun Seok Kang, and Yejin Choi. 2014. Keystroke patterns as prosody in digital writings: A case study with deceptive reviews and essays. In *Proceedings of the 2014 conference on empirical methods in natural language processing* (*EMNLP*), pages 1469–1473.
- Yejin Bang, Samuel Cahyawijaya, Nayeon Lee, Wenliang Dai, Dan Su, Bryan Wilie, Holy Lovenia, Ziwei Ji, Tiezheng Yu, Willy Chung, Quyet V. Do, Yan Xu, and Pascale Fung. 2023. A multitask, multilingual, multimodal evaluation of ChatGPT on reasoning, hallucination, and interactivity. In *Proceedings of the 13th International Joint Conference on Natural Language Processing and the 3rd Conference of the Asia-Pacific Chapter of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 675–718, Nusa Dua, Bali. Association for Computational Linguistics.
 - Jawadul H Bappy, Amit K Roy-Chowdhury, Jason Bunk, Lakshmanan Nataraj, and BS Manjunath. 2017. Exploiting spatial structure for localizing manipulated image regions. In *Proceedings of the IEEE international conference on computer vision*, pages 4970– 4979.
- Anthony M Barrett, Dan Hendrycks, Jessica Newman, and Brandie Nonnecke. 2022. Actionable guidance for high-consequence ai risk management: Towards standards addressing ai catastrophic risks. *arXiv preprint arXiv:2206.08966*.
- Belhassen Bayar and Matthew C Stamm. 2016. A deep learning approach to universal image manipulation detection using a new convolutional layer. In *Proceedings of the 4th ACM workshop on information hiding and multimedia security*, pages 5–10.
- Yoshua Bengio, Geoffrey Hinton, Andrew Yao, Dawn Song, Pieter Abbeel, Yuval Noah Harari, Ya-Qin

Zhang, Lan Xue, Shai Shalev-Shwartz, Gillian Hadfield, et al. 2023. Managing ai risks in an era of rapid progress. *arXiv preprint arXiv:2310.17688*, page 18.

- Maciej Besta, Nils Blach, Ales Kubicek, Robert Gerstenberger, Michal Podstawski, Lukas Gianinazzi, Joanna Gajda, Tomasz Lehmann, Hubert Niewiadomski, Piotr Nyczyk, et al. 2024. Graph of thoughts: Solving elaborate problems with large language models. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 38, pages 17682–17690.
- Rishabh Bhardwaj and Soujanya Poria. 2023. Redteaming large language models using chain of utterances for safety-alignment. *arXiv preprint arXiv:2308.09662.*
- Xiuli Bi, Yang Wei, Bin Xiao, and Weisheng Li. 2019. Rru-net: The ringed residual u-net for image splicing forgery detection. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition Workshops*, pages 0–0.
- Tian Bian, Xi Xiao, Tingyang Xu, Peilin Zhao, Wenbing Huang, Yu Rong, and Junzhou Huang. 2020. Rumor detection on social media with bi-directional graph convolutional networks. In *Proceedings of the AAAI conference on artificial intelligence*, volume 34, pages 549–556.
- Marcel Binz and Eric Schulz. 2023. Using cognitive psychology to understand gpt-3. *Proceedings of the National Academy of Sciences*, 120(6):e2218523120.
- Thomas Blass. 1984. Social psychology and personality: Toward a convergence. *Journal of Personality and Social Psychology*, 47(5):1013.
- Vincent D Blondel, Jean-Loup Guillaume, Renaud Lambiotte, and Etienne Lefebvre. 2008. Fast unfolding of communities in large networks. *Journal of statistical mechanics: theory and experiment*, 2008(10):P10008.
- Rishi Bommasani, Drew A Hudson, Ehsan Adeli, Russ Altman, Simran Arora, Sydney von Arx, Michael S Bernstein, Jeannette Bohg, Antoine Bosselut, Emma Brunskill, et al. 2021. On the opportunities and risks of foundation models. *arXiv preprint arXiv:2108.07258*.
- Asgarali Bouyer, Hamid Ahmadi Beni, Bahman Arasteh, Zahra Aghaee, and Reza Ghanbarzadeh. 2023. Fip: A fast overlapping community-based influence maximization algorithm using probability coefficient of global diffusion in social networks. *Expert systems with applications*, 213:118869.
- Samuel R Bowman, Jeeyoon Hyun, Ethan Perez, Edwin Chen, Craig Pettit, Scott Heiner, Kamilė Lukošiūtė, Amanda Askell, Andy Jones, Anna Chen, et al. 2022. Measuring progress on scalable oversight for large language models. *arXiv preprint arXiv:2211.03540*.

Arastoo Bozorgi, Saeed Samet, Johan Kwisthout, and Todd Wareham. 2017. Community-based influence maximization in social networks under a competitive linear threshold model. *Knowledge-Based Systems*, 134:149–158.

942

943

945

946

947

951

952

953

954

955

956

957

958

960

961

962

963

964

965

966

967

968

969

970

971

972

974

979

982

983

985

986

987

991

993

994

997

- Jean-Flavien Bussotti, Luca Ragazzi, Giacomo Frisoni, Gianluca Moro, and Paolo Papotti. 2024. Unknown claims: Generation of fact-checking training examples from unstructured and structured data. In *Proceedings of the 2024 Conference on Empirical Methods in Natural Language Processing*, pages 12105– 12122.
- Kechao Cai, Hong Xie, and John CS Lui. 2018. Information spreading forensics via sequential dependent snapshots. *IEEE/ACM Transactions on Networking*, 26(1):478–491.
- Han Cao, Lingwei Wei, Mengyang Chen, Wei Zhou, and Songlin Hu. 2023a. Are large language models good fact checkers: A preliminary study. *arXiv preprint arXiv:2311.17355*.
- Rui Cao, Ming Shan Hee, Adriel Kuek, Wen-Haw Chong, Roy Ka-Wei Lee, and Jing Jiang. 2023b. Procap: Leveraging a frozen vision-language model for hateful meme detection. In *Proceedings of the 31st ACM International Conference on Multimedia*, pages 5244–5252.
- Rui Cao, Roy Ka-Wei Lee, Wen-Haw Chong, and Jing Jiang. 2023c. Prompting for multimodal hateful meme classification. *arXiv preprint arXiv:2302.04156*.
- Nicholas Carlini, Milad Nasr, Christopher A Choquette-Choo, Matthew Jagielski, Irena Gao, Pang Wei W Koh, Daphne Ippolito, Florian Tramer, and Ludwig Schmidt. 2023. Are aligned neural networks adversarially aligned? *Advances in Neural Information Processing Systems*, 36:61478–61500.
- Carlos Castillo, Marcelo Mendoza, and Barbara Poblete. 2011. Information credibility on twitter. In *Proceedings of the 20th international conference on World wide web*, pages 675–684.
- Chi-Min Chan, Weize Chen, Yusheng Su, Jianxuan Yu, Wei Xue, Shanghang Zhang, Jie Fu, and Zhiyuan Liu. 2023. Chateval: Towards better llm-based evaluators through multi-agent debate. *arXiv preprint arXiv:2308.07201*.
- Michael Chan, Francis LF Lee, and Hsuan-Ting Chen. 2021. Examining the roles of multi-platform social media news use, engagement, and connections with news organizations and journalists on news literacy: A comparison of seven democracies. *Digital Journalism*, 9(5):571–588.
- Samantha Chan, Pat Pataranutaporn, Aditya Suri, Wazeer Zulfikar, Pattie Maes, and Elizabeth F Loftus. 2024. Conversational ai powered by large language models amplifies false memories in witness interviews. arXiv preprint arXiv:2408.04681.
- Canyu Chen and Kai Shu. 2023. Can llm-generated 998 misinformation be detected? arXiv preprint 999 arXiv:2309.13788. Canyu Chen and Kai Shu. 2024. Combating misinfor-1001 mation in the age of llms: Opportunities and chal-1002 lenges. AI Magazine, 45(3):354-368. 1003 Hung-Hsuan Chen, Yan-Bin Ciou, and Shou-De Lin. 1004 2012. Information propagation game: A tool to ac-1005 quire humanplaying data for multiplayer influence 1006 maximization on social networks. In Proceedings of the 18th ACM SIGKDD international conference on 1008 Knowledge discovery and data mining, pages 1524-1009 1527. 1010 Junyi Chen, Leyuan Liu, and Fan Zhou. 2025. Do not 1011 wait: Preemptive rumor detection with cooperative 1012 llms and accessible social context. Information Pro-1013 cessing & Management, 62(3):103995. 1014 Lei Chen, Guanying Li, Zhongyu Wei, Yang Yang, Bao-1015 hua Zhou, Qi Zhang, and Xuanjing Huang. 2022. A 1016 progressive framework for role-aware rumor resolu-1017 tion. In Proceedings of the 29th International Con-1018 ference on Computational Linguistics, pages 2748-1019 2758, Gyeongju, Republic of Korea. International 1020 Committee on Computational Linguistics. 1021 Sanxing Chen, Yukun Huang, and Bhuwan Dhingra. 1022 2024. Real-time fake news from adversarial feed-1023 back. arXiv preprint arXiv:2410.14651. 1024 Weixuan Chen and Daniel McDuff. 2018. Deepphys: 1025 Video-based physiological measurement using con-1026 volutional attention networks. In Proceedings of the 1027 european conference on computer vision (ECCV), 1028 pages 349-365. 1029 Yimin Chen, Niall J Conroy, and Victoria L Rubin. 2015. 1030 Misleading online content: recognizing clickbait as" 1031 false news". In Proceedings of the 2015 ACM on 1032 workshop on multimodal deception detection, pages 1033 15 - 19. 1034 Yixuan Chen, Jie Sui, Liang Hu, and Wei Gong. 2019. 1035 Attention-residual network with cnn for rumor detec-1036 tion. In Proceedings of the 28th ACM international 1037 conference on information and knowledge management, pages 1121-1130. 1039 Le Cheng, Peican Zhu, Keke Tang, Chao Gao, and Zhen 1040 Wang. 2024. Gin-sd: source detection in graphs with 1041 incomplete nodes via positional encoding and atten-1042 tive fusion. In Proceedings of the AAAI Conference 1043 on Artificial Intelligence, volume 38, pages 55-63. 1044 I Chern, Steffi Chern, Shiqi Chen, Weizhe Yuan, Kehua 1045 1046
 - Feng, Chunting Zhou, Junxian He, Graham Neubig, Pengfei Liu, et al. 2023. Factool: Factuality detection in generative ai–a tool augmented framework for multi-task and multi-domain scenarios. *arXiv preprint arXiv:2307.13528*.

1047

1048

1049

Zi Chu, Steven Gianvecchio, Haining Wang, and Sushil Jajodia. 2012. Detecting automation of twitter accounts: Are you a human, bot, or cyborg? *IEEE Transactions on dependable and secure computing*, 9(6):811–824.

1051

1052

1053

1055

1056

1059

1060

1061

1062

1065

1066

1067

1068

1069

1070

1072

1073

1074

1078

1080

1081 1082

1083

1084

1086

1087

1088

1090

1092

1093

1094

1095

1096

1097

1099

1100

1101

1102

1103

1104

1105

1106

1107

- Lynn Chua, Badih Ghazi, Yangsibo Huang, Pritish Kamath, Ravi Kumar, Daogao Liu, Pasin Manurangsi, Amer Sinha, and Chiyuan Zhang. 2024. Mind the privacy unit! user-level differential privacy for language model fine-tuning. arXiv preprint arXiv:2406.14322.
- Komal Chugh, Parul Gupta, Abhinav Dhall, and Ramanathan Subramanian. 2020. Not made for each other-audio-visual dissonance-based deepfake detection and localization. In *Proceedings of the 28th ACM international conference on multimedia*, pages 439–447.
- Thomas H Costello, Gordon Pennycook, and David G Rand. 2024. Durably reducing conspiracy beliefs through dialogues with ai. *Science*, 385(6714):eadq1814.
- Chaoqun Cui and Caiyan Jia. 2024. Propagation tree is not deep: Adaptive graph contrastive learning approach for rumor detection. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 38, pages 73–81.
- Jian Cui, Kwanwoo Kim, Seung Ho Na, and Seungwon Shin. 2022. Meta-path-based fake news detection leveraging multi-level social context information. In Proceedings of the 31st ACM international conference on information & knowledge management, pages 325–334.
- Limeng Cui, Haeseung Seo, Maryam Tabar, Fenglong Ma, Suhang Wang, and Dongwon Lee. 2020. Deterrent: Knowledge guided graph attention network for detecting healthcare misinformation. In *Proceedings* of the 26th ACM SIGKDD international conference on knowledge discovery & data mining, pages 492– 502.
- Anubrata Das, Houjiang Liu, Venelin Kovatchev, and Matthew Lease. 2023. The state of human-centered nlp technology for fact-checking. *Information processing & management*, 60(2):103219.
- Soumita Das, Ravi Kishore Devarapalli, and Anupam Biswas. 2024. Leveraging cascading information for community detection in social networks. *Information Sciences*, 674:120696.
- Luigi De Angelis, Francesco Baglivo, Guglielmo Arzilli, Gaetano Pierpaolo Privitera, Paolo Ferragina, Alberto Eugenio Tozzi, and Caterina Rizzo. 2023. Chatgpt and the rise of large language models: the new ai-driven infodemic threat in public health. *Frontiers in public health*, 11:1166120.
- Michela Del Vicario, Alessandro Bessi, Fabiana Zollo, Fabio Petroni, Antonio Scala, Guido Caldarelli, H Eugene Stanley, and Walter Quattrociocchi. 2016. The spreading of misinformation online. *Proceedings of* the national academy of Sciences, 113(3):554–559.

Natalia Díaz-Rodríguez, Javier Del Ser, Mark Coeck-
elbergh, Marcos López de Prado, Enrique Herrera-
Viedma, and Francisco Herrera. 2023. Connecting
the dots in trustworthy artificial intelligence: From
ai principles, ethics, and key requirements to respon-
sible ai systems and regulation. Information Fusion,
99:101896.1108
1112

1115

1116

1117

1118

1119

1120

1121

1122

1123

1124

1125

1126

1127

1128

1129

1130

1131

1132

1133

1134

1135

1136

1137

1138

1139

1140

1141

1142

1143

1144

1145

1146

1147

1148

1149

1150

1151

1152

1153

1154

1155

1156

1157

1158

1159

1160

1161

1162

1163

- Peter Sheridan Dodds, Eric M Clark, Suma Desu, Morgan R Frank, Andrew J Reagan, Jake Ryland Williams, Lewis Mitchell, Kameron Decker Harris, Isabel M Kloumann, James P Bagrow, et al. 2015. Human language reveals a universal positivity bias. *Proceedings of the national academy of sciences*, 112(8):2389–2394.
- Ming Dong, Bolong Zheng, Nguyen Quoc Viet Hung, Han Su, and Guohui Li. 2019. Multiple rumor source detection with graph convolutional networks. In *Proceedings of the 28th ACM international conference on information and knowledge management*, pages 569–578.
- Suyalatu Dong and Yong-Chang Huang. 2018. Sis rumor spreading model with population dynamics in online social networks. In 2018 International Conference on Wireless Communications, Signal Processing and Networking (WiSPNET), pages 1–5. IEEE.
- John Dougrez-Lewis, Elena Kochkina, Maria Liakata, and Yulan He. 2024. Knowledge graphs for realworld rumour verification. In *Proceedings of the* 2024 Joint International Conference on Computational Linguistics, Language Resources and Evaluation (LREC-COLING 2024), pages 9843–9853.
- Yaqian Dun, Kefei Tu, Chen Chen, Chunyan Hou, and Xiaojie Yuan. 2021. Kan: Knowledge-aware attention network for fake news detection. In *Proceedings of the AAAI conference on artificial intelligence*, volume 35, pages 81–89.
- Yogesh K Dwivedi, Nir Kshetri, Laurie Hughes, Emma Louise Slade, Anand Jeyaraj, Arpan Kumar Kar, Abdullah M Baabdullah, Alex Koohang, Vishnupriya Raghavan, Manju Ahuja, et al. 2023. Opinion paper:"so what if chatgpt wrote it?" multidisciplinary perspectives on opportunities, challenges and implications of generative conversational ai for research, practice and policy. *International Journal of Information Management*, 71:102642.
- Owain Evans, Owen Cotton-Barratt, Lukas Finnveden, Adam Bales, Avital Balwit, Peter Wills, Luca Righetti, and William Saunders. 2021. Truthful ai: Developing and governing ai that does not lie. *arXiv preprint arXiv:2110.06674*.
- Kevin Eykholt, Ivan Evtimov, Earlence Fernandes, Bo Li, Amir Rahmati, Chaowei Xiao, Atul Prakash, Tadayoshi Kohno, and Dawn Song. 2018. Robust physical-world attacks on deep learning visual classification. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 1625–1634.

Yahya H Ezzeldin, Shen Yan, Chaoyang He, Emilio Ferrara, and A Salman Avestimehr. 2023. Fairfed: Enabling group fairness in federated learning. In Proceedings of the AAAI conference on artificial intelligence, volume 37, pages 7494–7502.

1165

1166

1167

1168 1169

1170

1171

1172

1173

1174

1175

1176

1177

1178

1179

1180

1181

1182

1183

1184

1185

1186

1187

1188

1189

1190

1191

1192

1193

1194

1195

1196

1197

1198

1199

1201

1202

1203

1204

1205

1206

1207

1208

1209

1210

1211

1212

1213

1214

1215

1216

1217

1218

1219

- Lidan Fan, Zaixin Lu, Weili Wu, Bhavani Thuraisingham, Huan Ma, and Yuanjun Bi. 2013. Least cost rumor blocking in social networks. In 2013 IEEE 33rd International Conference on Distributed Computing Systems, pages 540–549. IEEE.
- Mahmoud Fawzi and Walid Magdy. 2024. " pinocchio had a nose, you have a network!": On characterizing fake news spreaders on arabic social media. *Proceedings of the ACM on Human-Computer Interaction*, 8(CSCW1):1–20.
- Shangbin Feng, Zhaoxuan Tan, Herun Wan, Ningnan Wang, Zilong Chen, Binchi Zhang, Qinghua Zheng, Wenqian Zhang, Zhenyu Lei, Shujie Yang, et al. 2022. Twibot-22: Towards graph-based twitter bot detection. Advances in Neural Information Processing Systems, 35:35254–35269.
- Zhangyin Feng, Xiaocheng Feng, Dezhi Zhao, Maojin Yang, and Bing Qin. 2024. Retrieval-generation synergy augmented large language models. In *ICASSP* 2024-2024 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP), pages 11661–11665. IEEE.
- Santo Fortunato. 2010. Community detection in graphs. *Physics reports*, 486(3-5):75–174.
- Santo Fortunato and Marc Barthelemy. 2007. Resolution limit in community detection. *Proceedings of the national academy of sciences*, 104(1):36–41.
 - Santo Fortunato and Darko Hric. 2016. Community detection in networks: A user guide. *Physics reports*, 659:1–44.
 - Gabriel Freedman, Adam Dejl, Deniz Gorur, Xiang Yin, Antonio Rago, and Francesca Toni. 2024. Argumentative large language models for explainable and contestable decision-making. *arXiv preprint arXiv:2405.02079*.
- Jessica Fridrich, David Soukal, Jan Lukas, et al. 2003. Detection of copy-move forgery in digital images. In *Proceedings of digital forensic research workshop*, volume 3, pages 652–63. Cleveland, OH.
- Rinaldo Gagiano, Maria Myung-Hee Kim, Xiuzhen Jenny Zhang, and Jennifer Biggs. 2021. Robustness analysis of grover for machine-generated news detection. In *Proceedings of the 19th Annual Workshop of the Australasian Language Technology Association*, pages 119–127.
- ŁG Gajewski, Krzysztof Suchecki, and JA Hołyst. 2019. Multiple propagation paths enhance locating the source of diffusion in complex networks. *Physica A: Statistical Mechanics and its Applications*, 519:34– 41.

Deep Ganguli, Liane Lovitt, Jackson Kernion, Amanda Askell, Yuntao Bai, Saurav Kadavath, Ben Mann, Ethan Perez, Nicholas Schiefer, Kamal Ndousse, et al. 2022. Red teaming language models to reduce harms: Methods, scaling behaviors, and lessons learned. *arXiv preprint arXiv:2209.07858*. 1220

1221

1222

1224

1226

1227

1228

1229

1230

1231

1232

1233

1234

1235

1236

1237

1238

1239

1240

1241

1242

1243

1244

1245

1246

1247

1248

1249

1250

1251

1252

1253

1254

1255

1256

1258

1259

1261

1262

1263

1264

1265

1266

1267

1268

1269

1270

1271

1272

1273

1274

1275

- Maryanne Garry, Way Ming Chan, Jeffrey Foster, and Linda A Henkel. 2024. Large language models (llms) and the institutionalization of misinformation. *Trends in cognitive sciences*.
- Jiahui Geng, Fengyu Cai, Yuxia Wang, Heinz Koeppl, Preslav Nakov, and Iryna Gurevych. 2024. A survey of confidence estimation and calibration in large language models. In *Proceedings of the 2024 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies (Volume 1: Long Papers)*, pages 6577–6595.
- Moumita Ghosh, Samhita Das, and Pritha Das. 2022. Dynamics and control of delayed rumor propagation through social networks. *Journal of Applied Mathematics and Computing*, pages 1–30.
- David Glukhov, Ilia Shumailov, Yarin Gal, Nicolas Papernot, and Vardan Papyan. 2023. Llm censorship: A machine learning challenge or a computer security problem? *arXiv preprint arXiv:2307.10719*.
- Oana Goga, Howard Lei, Sree Hari Krishnan Parthasarathi, Gerald Friedland, Robin Sommer, and Renata Teixeira. 2013. Exploiting innocuous activity for correlating users across sites. In *Proceedings of the 22nd international conference on World Wide Web*, pages 447–458.
- Jian Guan, Jesse Dodge, David Wadden, Minlie Huang, and Hao Peng. 2023. Language models hallucinate, but may excel at fact verification. *arXiv preprint arXiv:2310.14564*.
- David Güera and Edward J Delp. 2018. Deepfake video detection using recurrent neural networks. In 2018 15th IEEE international conference on advanced video and signal based surveillance (AVSS), pages 1–6. IEEE.
- Andrew M Guess, Michael Lerner, Benjamin Lyons, Jacob M Montgomery, Brendan Nyhan, Jason Reifler, and Neelanjan Sircar. 2020. A digital media literacy intervention increases discernment between mainstream and false news in the united states and india. *Proceedings of the National Academy of Sciences*, 117(27):15536–15545.
- Bin Guo, Yasan Ding, Lina Yao, Yunji Liang, and Zhiwen Yu. 2020. The future of false information detection on social media: New perspectives and trends. *ACM Computing Surveys (CSUR)*, 53(4):1–36.
- Daya Guo, Dejian Yang, Haowei Zhang, Junxiao Song, Ruoyu Zhang, Runxin Xu, Qihao Zhu, Shirong Ma, Peiyi Wang, Xiao Bi, et al. 2025. Deepseek-r1: Incentivizing reasoning capability in llms via reinforcement learning. *arXiv preprint arXiv:2501.12948*.

Qinglang Guo, Haiyong Xie, Yangyang Li, Wen Ma, and Chao Zhang. 2021. Social bots detection via 1278 1279 fusing bert and graph convolutional networks. Symmetry, 14(1):30. 1281 Zhijiang Guo, Michael Schlichtkrull, and Andreas Vlachos. 2022. A survey on automated fact-checking. Transactions of the Association for Computational *Linguistics*, 10:178–206. Philipp Hacker. 2024. Sustainable ai regulation. Com-1285 1286 mon Market Law Review, 61(2). Philipp Hacker, Andreas Engel, and Marco Mauer. 2023. 1287 Regulating chatgpt and other large generative ai mod-1289 els. In Proceedings of the 2023 ACM Conference on 1290 Fairness, Accountability, and Transparency, pages 1112-1123. 1291 Thilo Hagendorff. 2024. Deception abilities emerged in 1293 large language models. Proceedings of the National Academy of Sciences, 121(24):e2317967121. 1294 Samar Haider, Luca Luceri, Ashok Deb, Adam Badawy, 1295 Nanyun Peng, and Emilio Ferrara. 2023. Detecting 1296 1297 social media manipulation in low-resource languages. 1298 In Companion Proceedings of the ACM Web Confer-1299 ence 2023, pages 1358-1364. 1300 Patrick Haller, Ansar Aynetdinov, and Alan Akbik. 2023. 1301 Opiniongpt: Modelling explicit biases in instruction-1302 tuned llms. arXiv preprint arXiv:2309.03876. 1303 Tarek Hamdi, Hamda Slimi, Ibrahim Bounhas, and 1304 Yahya Slimani. 2020. A hybrid approach for fake 1305 news detection in twitter based on user features and 1306 graph embedding. In Distributed Computing and Internet Technology: 16th International Conference, 1307 1308 ICDCIT 2020, Bhubaneswar, India, January 9–12, 1309 2020, Proceedings 16, pages 266-280. Springer. Ahmed Abdeen Hamed. 2023. Improving detection of 1310 chatgpt-generated fake science using real publication 1311 text: Introducing xfakebibs a supervised learning 1312 1313 network algorithm. Kateřina Haniková, David Chudán, Vojtěch Svátek, Pe-1314 ter Vajdečka, Raphaël Troncy, Filip Vencovskỳ, and 1315 Jana Syrovátková. 2024. Towards fact-check summa-1316 1317 rization leveraging on argumentation elements tied to entity graphs. In Companion Proceedings of the 1318 ACM on Web Conference 2024, pages 1473–1481. 1319 1320 Hans WA Hanley and Zakir Durumeric. 2024. Machinemade media: Monitoring the mobilization of 1321 machine-generated articles on misinformation and 1322 1323 mainstream news websites. In Proceedings of the 1324 International AAAI Conference on Web and Social 1325 Media, volume 18, pages 542-556. Andreas Hanselowski, Christian Stab, Claudia Schulz, 1326 Zile Li, and Iryna Gurevych. 2019. A richly anno-1327 tated corpus for different tasks in automated fact-1328 1329 checking. arXiv preprint arXiv:1911.01214.

1277

Jing Hao, Zhixin Zhang, Shicai Yang, Di Xie, and Shiliang Pu. 2021. Transforensics: image forgery localization with dense self-attention. In Proceedings of the IEEE/CVF International Conference on Computer Vision, pages 15055–15064.

1330

1331

1332

1334

1335

1336

1337

1338

1339

1340

1341

1342

1343

1344

1345

1346

1347

1348

1349

1350

1351

1352

1353

1354

1355

1356

1357

1358

1359

1360

1361

1362

1363

1364

1365

1366

1367

1368

1369

1370

1371

1372

1373

1374

1375

1376

1377

1378

1379

- Bing He, Mustaque Ahamad, and Srijan Kumar. 2023.Reinforcement learning-based countermisinformation response generation: a case study of covid-19 vaccine misinformation. In Proceedings of the ACM Web Conference 2023, pages 2698-2709.
- Zhenyu He, Ce Li, Fan Zhou, and Yi Yang. 2021. Rumor detection on social media with event augmentations. In Proceedings of the 44th international ACM SIGIR conference on research and development in information retrieval, pages 2020-2024.
- Da Silva Gameiro Henrique, Andrei Kucharavy, and Rachid Guerraoui. 2023. Stochastic parrots looking for stochastic parrots: Llms are easy to fine-tune and hard to detect with other llms. arXiv preprint arXiv:2304.08968.
- Benjamin Horne and Sibel Adali. 2017. This just in: Fake news packs a lot in title, uses simpler, repetitive content in text body, more similar to satire than real news. In Proceedings of the international AAAI conference on web and social media, volume 11, pages 759-766.
- Adil Imad Eddine Hosni, Kan Li, and Sadique Ahmad. 2020. Minimizing rumor influence in multiplex online social networks based on human individual and social behaviors. Information Sciences, 512:1458-1480.
- Adil Imad Eddine Hosni, Kan Li, Cangfeng Ding, and Sadique Ahmed. 2018. Least cost rumor influence minimization in multiplex social networks. In International Conference on Neural Information Processing, pages 93–105. Springer.
- Saghar Hosseini, Hamid Palangi, and Ahmed Hassan Awadallah. 2023. An empirical study of metrics to measure representational harms in pre-trained language models. arXiv preprint arXiv:2301.09211.
- Dongpeng Hou, Shu Yin, Chao Gao, Xianghua Li, and Zhen Wang. 2024. Propagation dynamics of rumor vs. non-rumor across multiple social media platforms driven by user characteristics. arXiv preprint arXiv:2401.17840.
- Beizhe Hu, Qiang Sheng, Juan Cao, Yuhui Shi, Yang Li, Danding Wang, and Peng Qi. 2024. Bad actor, good advisor: Exploring the role of large language models in fake news detection. In Proceedings of the AAAI Conference on Artificial Intelligence, volume 38, pages 22105–22113.
- Linmei Hu, Tianchi Yang, Luhao Zhang, Wanjun Zhong, 1381 Duyu Tang, Chuan Shi, Nan Duan, and Ming Zhou. 1382 2021. Compare to the knowledge: Graph neural fake 1383 news detection with external knowledge. In Proceed-1384 ings of the 59th Annual Meeting of the Association for 1385

Computational Linguistics and the 11th International Joint Conference on Natural Language Processing (Volume 1: Long Papers), pages 754–763.	Kathleen H Jamieson. 2008. Echo chamber: Rush Lim- baugh and the conservative media establishment. Ox- ford University Press.
Nan Hu, Zirui Wu, Yuxuan Lai, Xiao Liu, and Yansong Feng. 2022. Dual-channel evidence fusion for fact verification over texts and tables. In <i>Proceedings of</i> <i>the 2022 Conference of the North American Chap-</i> <i>ter of the Association for Computational Linguistics:</i> <i>Human Language Technologies</i> , pages 5232–5242.	Danish Javed, NZ Jhanjhi, Navid Ali Khan, Sayan Ku- mar Ray, Alanoud Al Mazroa, Farzeen Ashfaq, and Shampa Rani Das. 2024. Towards the future of bot detection: A comprehensive taxonomical review and challenges on twitter/x. <i>Computer Networks</i> , 254:110808.
Tianrui Hu, Dimitrios Liakopoulos, Xiwen Wei, Radu Marculescu, and Neeraja J Yadwadkar. 2025. Simu- lating rumor spreading in social networks using llm agents. <i>arXiv preprint arXiv:2502.01450</i> .	Ziwei Ji, Nayeon Lee, Rita Frieske, Tiezheng Yu, Dan Su, Yan Xu, Etsuko Ishii, Ye Jin Bang, Andrea Madotto, and Pascale Fung. 2023. Survey of halluci- nation in natural language generation. <i>ACM Comput-</i> <i>ing Surveys</i> , 55(12):1–38.
Di Huang, Jinbao Song, and Xingyu Zhang. 2025a. Semi-supervised social bot detection with relational graph attention transformers and characteristics of the social environment. <i>Information Fusion</i> , page 102956.	Zhengbao Jiang, Frank F Xu, Luyu Gao, Zhiqing Sun, Qian Liu, Jane Dwivedi-Yu, Yiming Yang, Jamie Callan, and Graham Neubig. 2023. Ac- tive retrieval augmented generation. <i>arXiv preprint</i> <i>arXiv:2305.06983</i> .
Linan Huang and Quanyan Zhu. 2023. An introduction of system-scientific approaches to cognitive security. <i>arXiv preprint arXiv:2301.05920</i> .	Zhiwei Jin, Juan Cao, Han Guo, Yongdong Zhang, and Jiebo Luo. 2017. Multimodal fusion with recurrent neural networks for rumor detection on microblogs. In Proceedings of the 25th ACM international con
Yue Huang and Lichao Sun. 2023. Harnessing the power of chatgpt in fake news: An in-depth explo- ration in generation, detection and explanation. <i>arXiv</i> <i>preprint arXiv:2310.05046</i> .	 In Proceedings of the 25th ACM International conference on Multimedia, pages 795–816. Zhiwei Jin, Juan Cao, Yongdong Zhang, and Jiebo Luo. 2016 News verification by exploiting conflicting
Yue Huang and Lichao Sun. 2024. Fakegpt: fake news generation, explanation and detection of large lan- guage models. <i>arxiv. org</i> .	social viewpoints in microblogs. In <i>Proceedings of</i> <i>the AAAI conference on artificial intelligence</i> , vol- ume 30.
Yuheng Huang, Jiayang Song, Zhijie Wang, Shengming Zhao, Huaming Chen, Felix Juefei-Xu, and Lei Ma. 2025b. Look before you leap: An exploratory study of uncertainty analysis for large language models. <i>IEEE Transactions on Software Engineering</i> .	 Matthew Jiwa, Patrick S Cooper, Trevor TJ Chong, and Stefan Bode. 2023. Hedonism as a motive for infor- mation search: biased information-seeking leads to biased beliefs. <i>Scientific Reports</i>, 13(1):2086.
Minyoung Huh, Andrew Liu, Andrew Owens, and Alexei A Efros. 2018. Fighting fake news: Image	The global landscape of ai ethics guidelines. <i>Nature</i> machine intelligence, 1(9):389–399.
splice detection via learned self-consistency. In <i>Proceedings of the European conference on computer vision (ECCV)</i> , pages 101–117.	Neil F Johnson, Nicolas Velásquez, Nicholas John- son Restrepo, Rhys Leahy, Nicholas Gabriel, Sara El Oud, Minzhang Zheng, Pedro Manrique, Stefan Wuchty, and Yonatan Lupu, 2020. The online com-
Hongwen Hui, Chengcheng Zhou, Xing Lü, and Jiarong Li. 2020. Spread mechanism and control strategy of social network rumors under the influence of covid-	petition between pro-and anti-vaccination views. <i>Na-</i> <i>ture</i> , 582(7811):230–233.
19. Nonlinear Dynamics, 101:1933–1949. Tereza Iofciu, Peter Fankhauser, Fabian Abel, and Ker-	Daniel Kahneman. 1979. Prospect theory: An analysis of decisions under risk. <i>Econometrica</i> , 47:278.
stin Bischoff. 2011. Identifying users across social tagging systems. In <i>Proceedings of the International AAAI Conference on Web and Social Media</i> , volume 5, pages 522–525.	Rohit Kumar Kaliyar, Anurag Goswami, and Pratik Narang. 2021. Fakebert: Fake news detection in so- cial media with a bert-based deep learning approach. <i>Multimedia tools and applications</i> , 80(8):11765– 11788.
Daphne Ippolito, Daniel Duckworth, Chris Callison- Burch, and Douglas Eck. 2020. Automatic detec- tion of generated text is easiest when humans are fooled. In <i>Proceedings of the 58th Annual Meeting of</i> <i>the Association for Computational Linguistics</i> , pages 1808–1822, Online. Association for Computational Linguistics.	Mohammadamin Kanaani. 2024. Triple-r: Automatic reasoning for fact verification using language mod- els. In Proceedings of the 2024 Joint International Conference on Computational Linguistics, Language Resources and Evaluation (LREC-COLING 2024), pages 16831–16840.

- 1495 1496 1497 1499 1500 1501 1505 1508 1510 1511 1513 1514 1516 1517 1518 1519 1520 1521 1522 1523 1524 1525 1526 1527 1528 1530 1531 1532 1533 1534 1536 1537 1538 1539 1540 1541 1542 1543 1544 1545

1547 1548

- Hamid Karimi and Jiliang Tang. 2019. Learning hierarchical discourse-level structure for fake news detection. arXiv preprint arXiv:1903.07389.
- Elihu Katz. 1957. The two-step flow of communication: An up-to-date report on an hypothesis. Public opinion quarterly, 21(1):61-78.
- William Ogilvy Kermack and Anderson G McKendrick. 1927. A contribution to the mathematical theory of epidemics. Proceedings of the royal society of london. Series A, Containing papers of a mathematical and physical character, 115(772):700-721.
- Nida Saddaf Khan, Maira Ata, and Quratulain Rajput. 2015. Identification of opinion leaders in social network. In 2015 International Conference on Information and Communication Technologies (ICICT), pages 1-6. IEEE.
 - Salman Khan, Muzammal Naseer, Munawar Hayat, Syed Waqas Zamir, Fahad Shahbaz Khan, and Mubarak Shah. 2022. Transformers in vision: A survey. ACM computing surveys (CSUR), 54(10s):1-41.
- Dhruv Khattar, Jaipal Singh Goud, Manish Gupta, and Vasudeva Varma. 2019. Mvae: Multimodal variational autoencoder for fake news detection. In The world wide web conference, pages 2915–2921.
- Antino Kim and Alan R Dennis. 2019. Says who? the effects of presentation format and source rating on fake news in social media. Mis quarterly, 43(3):1025-1039.
- Camille Koenders, Johannes Filla, Nicolai Schneider, and Vinicius Woloszyn. 2021. How vulnerable are automatic fake news detection methods to adversarial attacks? arXiv preprint arXiv:2107.07970.
- Ziyi Kou, Lanyu Shang, Yang Zhang, Zhenrui Yue, Huimin Zeng, and Dong Wang. 2022. Crowd, expert & ai: A human-ai interactive approach towards natural language explanation based covid-19 misinformation detection. In IJCAI, pages 5087-5093.
- Sarah Kreps, R Miles McCain, and Miles Brundage. 2022. All the news that's fit to fabricate: Aigenerated text as a tool of media misinformation. Journal of experimental political science, 9(1):104-117.
- Weirui Kuang, Bingchen Qian, Zitao Li, Daoyuan Chen, Dawei Gao, Xuchen Pan, Yuexiang Xie, Yaliang Li, Bolin Ding, and Jingren Zhou. 2024. Federatedscopellm: A comprehensive package for fine-tuning large language models in federated learning. In Proceedings of the 30th ACM SIGKDD Conference on Knowledge Discovery and Data Mining, pages 5260-5271.
- KP Krishna Kumar and G Geethakumari. 2014. Detecting misinformation in online social networks using cognitive psychology. Human-centric Computing and Information Sciences, 4(1):14.

Raghvendra Kumar, Bhargav Goddu, Sriparna Saha, and Adam Jatowt. 2025. Silver lining in the fake news cloud: Can large language models help detect misinformation? IEEE Transactions on Artificial Intelligence, 6(1):14–24.

1549

1550

1551

1553

1554

1555

1556

1557

1558

1559

1560

1561

1562

1563

1564

1565

1566

1567

1568

1569

1570

1571

1572

1573

1574

1575

1576

1577

1578

1579

1580

1581

1582

1583

1584

1585

1586

1587

1588

1591

1593

1594

1595

- Tharindu Kumarage and Huan Liu. 2023. Neural authorship attribution: Stylometric analysis on large language models. In 2023 International Conference on Cyber-Enabled Distributed Computing and Knowledge Discovery (CyberC), pages 51-54. IEEE.
- Sejeong Kwon, Meeyoung Cha, Kyomin Jung, Wei Chen, and Yajun Wang. 2013. Prominent features of rumor propagation in online social media. In 2013 IEEE 13th international conference on data mining, pages 1103-1108. IEEE.
- An Lao, Qi Zhang, Chongyang Shi, Longbing Cao, Kun Yi, Liang Hu, and Duoqian Miao. 2024. Frequency spectrum is more effective for multimodal representation and fusion: A multimodal spectrum rumor detector. In Proceedings of the AAAI Conference on Artificial Intelligence, volume 38, pages 18426-18434.
- Stephan Lewandowsky. 2022. Fake news and participatory propaganda. In Cognitive Illusions, pages 324-340. Routledge.
- Stephan Lewandowsky, Ullrich KH Ecker, Colleen M Seifert, Norbert Schwarz, and John Cook. 2012. Misinformation and its correction: Continued influence and successful debiasing. Psychological science in *the public interest*, 13(3):106–131.
- Bingxin Li and Linhe Zhu. 2024. Turing instability analvsis of a reaction-diffusion system for rumor propagation in continuous space and complex networks. Information Processing & Management, 61(3):103621.
- Bo Li, Peng Qi, Bo Liu, Shuai Di, Jingen Liu, Jiquan Pei, Jinfeng Yi, and Bowen Zhou. 2023a. Trustworthy ai: From principles to practices. ACM Computing Surveys, 55(9):1-46.
- Chen Li, Hao Peng, Jianxin Li, Lichao Sun, Lingjuan Lyu, Lihong Wang, S Yu Philip, and Lifang He. 2021. Joint stance and rumor detection in hierarchical heterogeneous graph. IEEE transactions on neural networks and learning systems, 33(6):2530–2542.
- Guohao Li, Hasan Hammoud, Hani Itani, Dmitrii Khizbullin, and Bernard Ghanem. 2023b. Camel: Communicative agents for" mind" exploration of large language model society. Advances in Neural Information Processing Systems, 36:51991–52008.
- Jiarong Li, Haijun Jiang, Xuehui Mei, Cheng Hu, and 1597 Guoliang Zhang. 2020. Dynamical analysis of ru-1598 mor spreading model in multi-lingual environment 1599 and heterogeneous complex networks. Information 1600 Sciences, 536:391-408. 1601

1602

- 1617 1618
- 1619 1620
- 1621 1622
- 1623 1624 1625
- 1626 1627 1628 1629
- 1630 1631 1632
- 163 163
- 1637 1638 1639

1640

- 1641 1642 1643
- 1644 1645
- 1646 1647
- 1648 1649 1650 1651

1652 1653

1656 1655 1656

- Miaoran Li, Baolin Peng, Michel Galley, Jianfeng Gao, and Zhu Zhang. 2023c. Self-checker: Plug-and-play modules for fact-checking with large language models. *arXiv preprint arXiv:2305.14623*.
- Sha Li, Ruining Zhao, Manling Li, Heng Ji, Chris Callison-Burch, and Jiawei Han. 2023d. Opendomain hierarchical event schema induction by incremental prompting and verification. *arXiv preprint arXiv:2307.01972*.
 - Weimin Li, Xiaokang Zhou, Chao Yang, Yuting Fan, Zhao Wang, and Yanxia Liu. 2022. Multi-objective optimization algorithm based on characteristics fusion of dynamic social networks for community discovery. *Information Fusion*, 79:110–123.
- Xinyi Li, Yongfeng Zhang, and Edward C Malthouse. 2024a. Large language model agent for fake news detection. *arXiv preprint arXiv:2405.01593*.
- Yuezun Li, Ming-Ching Chang, and Siwei Lyu. 2018. In ictu oculi: Exposing ai created fake videos by detecting eye blinking. In 2018 IEEE International workshop on information forensics and security (WIFS), pages 1–7. Ieee.
- Yupeng Li, Haorui He, Jin Bai, and Dacheng Wen. 2024b. Mcfend: a multi-source benchmark dataset for chinese fake news detection. In *Proceedings of the ACM on Web Conference 2024*, pages 4018–4027.
- Zhuohang Li, Jiaxin Zhang, Chao Yan, Kamalika Das, Sricharan Kumar, Murat Kantarcioglu, and Bradley A Malin. 2024c. Do you know what you are talking about? characterizing query-knowledge relevance for reliable retrieval augmented generation. *arXiv preprint arXiv:2410.08320*.
- Gang Liang, Wenbo He, Chun Xu, Liangyin Chen, and Jinquan Zeng. 2015. Rumor identification in microblogging systems based on users' behavior. *IEEE Transactions on Computational Social Systems*, 2(3):99–108.
- Wee Yong Lim, Mong Li Lee, and Wynne Hsu. 2017. Ifact: An interactive framework to assess claims from tweets. In Proceedings of the 2017 ACM on Conference on Information and Knowledge Management, pages 787–796.
- Stephanie Lin, Jacob Hilton, and Owain Evans. 2022. Teaching models to express their uncertainty in words. *arXiv preprint arXiv:2205.14334*.
- Chen Ling, Junji Jiang, Junxiang Wang, and Zhao Liang. 2022. Source localization of graph diffusion via variational autoencoders for graph inverse problems. In *Proceedings of the 28th ACM SIGKDD conference on knowledge discovery and data mining*, pages 1010– 1020.
- Aiwei Liu, Qiang Sheng, and Xuming Hu. 2024a. Preventing and detecting misinformation generated by large language models. In *Proceedings of the 47th International ACM SIGIR Conference on Research*

and Development in Information Retrieval, pages 3001–3004.

1657

1658

1659

1660

1661

1662

1663

1665

1666

1667

1668

1669

1672

1673

1674

1678

1679

1680

1682

1683

1684

1685

1686

1688

1689

1690

1691

1692

1693

1694

1695

1696

1697

1698

1699

1700

1701

1702

1703

1704

1705

1706

- Hui Liu, Wenya Wang, Haoru Li, and Haoliang Li. 2024b. Teller: A trustworthy framework for explainable, generalizable and controllable fake news detection. *arXiv preprint arXiv:2402.07776*.
- Xuannan Liu, Peipei Li, Huaibo Huang, Zekun Li, Xing Cui, Jiahao Liang, Lixiong Qin, Weihong Deng, and Zhaofeng He. 2024c. Fka-owl: Advancing multimodal fake news detection through knowledge-augmented lvlms. In *Proceedings of the 32nd ACM International Conference on Multimedia*, pages 10154–10163.
- Yi Liu, Gelei Deng, Yuekang Li, Kailong Wang, Zihao Wang, Xiaofeng Wang, Tianwei Zhang, Yepang Liu, Haoyu Wang, Yan Zheng, et al. 2023. Prompt injection attack against llm-integrated applications. *arXiv preprint arXiv:2306.05499*.
- Alireza Louni and KP Subbalakshmi. 2018. Who spread that rumor: Finding the source of information in large online social networks with probabilistically varying internode relationship strengths. *IEEE transactions on computational social systems*, 5(2):335–343.
- Yi-Ju Lu and Cheng-Te Li. 2020. Gcan: Graphaware co-attention networks for explainable fake news detection on social media. *arXiv preprint arXiv:2004.11648*.
- Qing Lyu, Shreya Havaldar, Adam Stein, Li Zhang, Delip Rao, Eric Wong, Marianna Apidianaki, and Chris Callison-Burch. 2023. Faithful chain-ofthought reasoning. In Proceedings of the 13th International Joint Conference on Natural Language Processing and the 3rd Conference of the Asia-Pacific Chapter of the Association for Computational Linguistics (Volume 1: Long Papers), pages 305–329, Nusa Dua, Bali. Association for Computational Linguistics.
- Jing Ma, Wei Gao, Shafiq Joty, and Kam-Fai Wong. 2019. Sentence-level evidence embedding for claim verification with hierarchical attention networks. Association for Computational Linguistics.
- Jing Ma, Wei Gao, Prasenjit Mitra, Sejeong Kwon, Bernard J Jansen, Kam-Fai Wong, and Meeyoung Cha. 2016. Detecting rumors from microblogs with recurrent neural networks.
- Jing Ma, Wei Gao, Zhongyu Wei, Yueming Lu, and Kam-Fai Wong. 2015. Detect rumors using time series of social context information on microblogging websites. In *Proceedings of the 24th ACM international on conference on information and knowledge management*, pages 1751–1754.
- Jing Ma, Wei Gao, and Kam-Fai Wong. 2018. Rumor1708detection on twitter with tree-structured recursive1709neural networks. Association for Computational Linguistics.1710

1823

Zihan Ma, Minnan Luo, Hao Guo, Zhi Zeng, Yiran Hao, and Xiang Zhao. 2024. Event-radar: Event-driven multi-view learning for multimodal fake news detection. In Proceedings of the 62nd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers), pages 5809–5821.

1712

1713

1714

1717

1718

1719

1720

1799

1723

1724

1725

1726

1727

1728

1729

1730

1731

1732

1733

1734

1735

1736

1737

1738

1739

1740

1741

1742

1743

1744

1745

1746

1747

1748

1749

1750

1752

1753

1754

1755 1756

1757

1758

1759

1760

1761

1762

1763

1764

1765

1766

1767

- Peihua Mai, Ran Yan, Zhe Huang, Youjia Yang, and Yan Pang. 2023. Split-and-denoise: Protect large language model inference with local differential privacy. *arXiv preprint arXiv:2310.09130*.
- Shrikant Malviya and Stamos Katsigiannis. 2024. Evidence retrieval for fact verification using multi-stage reranking. In *Findings of the Association for Computational Linguistics: EMNLP 2024*, pages 7295– 7308, Miami, Florida, USA. Association for Computational Linguistics.
 - Francesco Marra, Diego Gragnaniello, Luisa Verdoliva, and Giovanni Poggi. 2019. Do gans leave artificial fingerprints? In 2019 IEEE conference on multimedia information processing and retrieval (MIPR), pages 506–511. IEEE.
- David Martín-Gutiérrez, Gustavo Hernández-Peñaloza, Alberto Belmonte Hernández, Alicia Lozano-Diez, and Federico Álvarez. 2021. A deep learning approach for robust detection of bots in twitter using transformers. *IEEE Access*, 9:54591–54601.
- Falko Matern, Christian Riess, and Marc Stamminger. 2019. Exploiting visual artifacts to expose deepfakes and face manipulations. In 2019 IEEE Winter Applications of Computer Vision Workshops (WACVW), pages 83–92. IEEE.
- SC Matz, JD Teeny, Sumer S Vaid, H Peters, GM Harari, and M Cerf. 2024. The potential of generative ai for personalized persuasion at scale. *Scientific Reports*, 14(1):4692.
- Mohit Mayank, Shakshi Sharma, and Rajesh Sharma. 2022. Deap-faked: Knowledge graph based approach for fake news detection. In 2022 IEEE/ACM International Conference on Advances in Social Networks Analysis and Mining (ASONAM), pages 47–51. IEEE.
- Nikhil Mehta, María Leonor Pacheco, and Dan Goldwasser. 2022. Tackling fake news detection by continually improving social context representations using graph neural networks. In *Proceedings of the* 60th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers), pages 1363–1380.
- Erxue Min, Yu Rong, Yatao Bian, Tingyang Xu, Peilin Zhao, Junzhou Huang, and Sophia Ananiadou. 2022. Divide-and-conquer: Post-user interaction network for fake news detection on social media. In *Proceedings of the ACM web conference 2022*, pages 1148–1158.
- Daniel Mas Montserrat, Hanxiang Hao, Sri K Yarlagadda, Sriram Baireddy, Ruiting Shao, János Horváth, Emily Bartusiak, Justin Yang, David Guera,

Fengqing Zhu, et al. 2020. Deepfakes detection with automatic face weighting. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition workshops*, pages 668–669.

- Xuemei Mou, Wei Xu, Yangfu Zhu, Qian Li, and Yunpeng Xiao. 2022. A social topic diffusion model based on rumor, anti-rumor, and motivation-rumor. *IEEE Transactions on Computational Social Systems*, 10(5):2644–2659.
- Qiong Nan, Juan Cao, Yongchun Zhu, Yanyan Wang, and Jintao Li. 2021. Mdfend: Multi-domain fake news detection. In *Proceedings of the 30th ACM International Conference on Information & Knowledge Management*, pages 3343–3347.
- Qiong Nan, Qiang Sheng, Juan Cao, Beizhe Hu, Danding Wang, and Jintao Li. 2024. Let silence speak: Enhancing fake news detection with generated comments from large language models. In *Proceedings of the 33rd ACM International Conference on Information and Knowledge Management*, pages 1732–1742.
- Mark EJ Newman. 2004. Fast algorithm for detecting community structure in networks. *Physical Review E—Statistical, Nonlinear, and Soft Matter Physics*, 69(6):066133.
- Lynnette Hui Xian Ng and Kathleen M. Carley. 2022. Online coordination: Methods and comparative case studies of coordinated groups across four events in the united states. In *Proceedings of the 14th ACM Web Science Conference 2022*, WebSci '22, page 12–21, New York, NY, USA. Association for Computing Machinery.
- Van-Hoang Nguyen, Kazunari Sugiyama, Preslav Nakov, and Min-Yen Kan. 2020. Fang: Leveraging social context for fake news detection using graph representation. In *Proceedings of the 29th ACM international conference on information & knowledge management*, pages 1165–1174.
- Raymond S Nickerson. 1998. Confirmation bias: A ubiquitous phenomenon in many guises. *Review of general psychology*, 2(2):175–220.
- Yuanping Nie, Yan Jia, Shudong Li, Xiang Zhu, Aiping Li, and Bin Zhou. 2016. Identifying users across social networks based on dynamic core interests. *Neurocomputing*, 210:107–115.
- Cheng Niu, Yang Guan, Yuanhao Wu, Juno Zhu, Juntong Song, Randy Zhong, Kaihua Zhu, Siliang Xu, Shizhe Diao, and Tong Zhang. 2024. Veract scan: Retrieval-augmented fake news detection with justifiable reasoning. *arXiv preprint arXiv:2406.10289*.
- Cheng Niu, Yuanhao Wu, Juno Zhu, Siliang Xu, Kashun Shum, Randy Zhong, Juntong Song, and Tong Zhang. 2023. Ragtruth: A hallucination corpus for developing trustworthy retrieval-augmented language models. *arXiv preprint arXiv:2401.00396*.
- R OpenAI. 2023. Gpt-4 technical report. arxiv 2303.08774. *View in Article*, 2(5).

Tore Opsahl, Filip Agneessens, and John Skvoretz. 2010. Node centrality in weighted networks: Generalizing degree and shortest paths. Social networks, 32(3):245-251.

1824

1825

1826

1828

1830

1835 1836

1837 1838

1839

1840

1843

1844

1845

1846

1847

1848

1849

1850

1851 1852

1853

1854

1855

1856

1857 1858

1859

1862

1865 1866

1867

1868 1869

1870

1871

1872

1873 1874

1875

1876

1877

1878

1879

Liangming Pan, Xiaobao Wu, Xinyuan Lu, Anh Tuan Luu, William Yang Wang, Min-Yen Kan, and Preslav Nakov. 2023a. Fact-checking complex claims with program-guided reasoning. In Proceedings of the 61st Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers), pages 6981-7004, Toronto, Canada. Association for Computational Linguistics.

- Yikang Pan, Liangming Pan, Wenhu Chen, Preslav Nakov, Min-Yen Kan, and William Wang. 2023b. On the risk of misinformation pollution with large language models. In Findings of the Association for Computational Linguistics: EMNLP 2023, pages 1389–1403, Singapore. Association for Computational Linguistics.
- Yikang Pan, Liangming Pan, Wenhu Chen, Preslav Nakov, Min-Yen Kan, and William Yang Wang. 2023c. On the risk of misinformation pollution with large language models. arXiv preprint arXiv:2305.13661.
- Joon Sung Park, Joseph O'Brien, Carrie Jun Cai, Meredith Ringel Morris, Percy Liang, and Michael S Bernstein. 2023. Generative agents: Interactive simulacra of human behavior. In Proceedings of the 36th annual acm symposium on user interface software and technology, pages 1-22.
- Joon Sung Park, Lindsay Popowski, Carrie Cai, Meredith Ringel Morris, Percy Liang, and Michael S Bernstein. 2022. Social simulacra: Creating populated prototypes for social computing systems. In Proceedings of the 35th Annual ACM Symposium on User Interface Software and Technology, pages 1–18.
- Peter S Park, Simon Goldstein, Aidan O'Gara, Michael Chen, and Dan Hendrycks. 2024. Ai deception: A survey of examples, risks, and potential solutions. Patterns, 5(5).
- E Parliament. 2023. Artificial intelligence act: deal on comprehensive rules for trustworthy ai. Pressemitteilung vom, 9.
- Kellin Pelrine, Jacob Danovitch, and Reihaneh Rabbany. 2021. The surprising performance of simple baselines for misinformation detection. In *Proceedings* of the Web Conference 2021, pages 3432-3441.
- Baolin Peng, Michel Galley, Pengcheng He, Hao Cheng, Yujia Xie, Yu Hu, Qiuyuan Huang, Lars Liden, Zhou Yu, Weizhu Chen, et al. 2023. Check your facts and try again: Improving large language models with external knowledge and automated feedback. arXiv preprint arXiv:2302.12813.
- Matjaž Perc, Mahmut Ozer, and Janja Hojnik. 2019. Social and juristic challenges of artificial intelligence. Palgrave Communications, 5(1).

Ethan Perez, Saffron Huang, Francis Song, Trevor Cai,	1880
Roman Ring, John Aslanides, Amelia Glaese, Nat	1881
McAleese, and Geoffrey Irving. 2022. Red team-	1882
ing language models with language models. arXiv	1883
preprint arXiv:2202.03286.	1884
Verónica Pérez-Rosas, Bennett Kleinberg, Alexandra	1885
Lefevre, and Rada Mihalcea. 2017. Automatic detec-	1886
tion of fake news. arXiv preprint arXiv:1708.07104.	1887
Hainrich Paters and Sandra Matz 2024. Large language	1000
models can infer psychological dispositions of social	1880
media users. <i>PNAS Nexus</i> , page pgae231.	1890
inean asersi i inis rienss, page parezo i	
Alexandru Petrescu, Ciprian-Octavian Truică, Elena-	1891
Simona Apostol, and Panagiotis Karras. 2021.	1892
Sparse shield: Social network immunization vs.	1893
harmful speech. In Proceedings of the 30th ACM	1894
International Conference on Information & Knowl-	1895
edge Management, pages 1426–1436.	1896
Huyen Trang Phan, Ngoc Thanh Nguyen, and Dosam	1897
Hwang. 2023. Fake news detection: A survey of	1898
graph neural network methods. Applied Soft Comput-	1899
ing, 139:110235.	1900
Moritz Pilarski, Kirill Olegovich Solovev, and Nico-	1901
las Pröllochs. 2024. Community notes vs. snoping:	1902
How the crowd selects fact-checking targets on so-	1903
cial media. In Proceedings of the International AAAI	1904
Conference on Web and Social Media, volume 18,	1905
pages 1262–1275.	1906
Sanjaikanth E Vadakkethil Somanathan Pillai. 2024. En-	1907
hancing misinformation detection through semantic	1908
analysis and knowledge graphs. In 2024 4th Interna-	1909
tional Conference on Data Engineering and Commu-	1910
nication Systems (ICDECS), pages 1–5. IEEE.	1911
Clara Pizzuti. 2017. Evolutionary computation for com-	1912
munity detection in networks: A review. IEEE Trans-	1913
actions on Evolutionary Computation, 22(3):464-	1914
483.	1915
Russell A Poldrack, Thomas Lu, and Gašper Beguš.	1916
2023. Ai-assisted coding: Experiments with gpt-4.	1917
arXiv preprint arXiv:2304.13187.	1918
Pascal Pons and Matthieu Latapy. 2005. Computing	1919
communities in large networks using random walks.	1920
In Computer and Information Sciences-ISCIS 2005:	1921
20th International Symposium, Istanbul, Turkey, Oc-	1922
tober 26-28, 2005. Proceedings 20, pages 284-293.	1923
Springer.	1924
Alin C Popescu and Hany Farid. 2004. Exposing digital	1925
forgeries by detecting duplicated image regions.	1926
Ethan Porter and Thomas J Wood. 2021. The global ef-	1927
fectiveness of fact-checking: Evidence from simulta-	1928
neous experiments in argentina, nigeria, south africa,	1929
and the united kingdom. Proceedings of the National	1930
Academy of Sciences, 118(37):e2104235118	1931

1932 Martin Potthast, Johannes Kiesel, Kevin Reinartz, Janek Shebuti Rayana and Leman Akoglu. 2015. Collective 1985 1933 Bevendorff, and Benno Stein. 2017. A stylometopinion spam detection: Bridging review networks 1986 and metadata. In Proceedings of the 21th acm sigkdd 1934 ric inquiry into hyperpartisan and fake news. arXiv preprint arXiv:1702.05638. international conference on knowledge discovery and 1988 data mining, pages 985–994. 1989 1936 Ofir Press, Muru Zhang, Sewon Min, Ludwig Schmidt, Noah A Smith, and Mike Lewis. 2022. Measuring Veronica Red, Eric D Kelsic, Peter J Mucha, and Maand narrowing the compositionality gap in language son A Porter. 2011. Comparing community structure models. arXiv preprint arXiv:2210.03350. to characteristics in online collegiate social networks. SIAM review, 53(3):526–543. Nicolas Pröllochs, Dominik Bär, and Stefan Feuerriegel. 1940 1941 2021. Emotions in online rumor diffusion. EPJ Data Zhicheng Ren, Zhiping Xiao, and Yizhou Sun. 2024. 1994 Science, 10(1):51. 1942 Do we trust what they say or what they do? a multi-1995 modal user embedding provides personalized expla-1996 Peng Qi, Zehong Yan, Wynne Hsu, and Mong Li Lee. nations. arXiv preprint arXiv:2409.02965. 1944 2024. Sniffer: Multimodal large language model for 1945 explainable out-of-context misinformation detection. Christopher Riederer, Yunsung Kim, Augustin Chain-1998 *Preprint*, arXiv:2403.03170. treau, Nitish Korula, and Silvio Lattanzi. 2016. Linking users across domains with location data: Theory Xiangyu Qi, Yi Zeng, Tinghao Xie, Pin-Yu Chen, and validation. In Proceedings of the 25th international conference on world wide web, pages 707-719. 2002 1948 Ruoxi Jia, Prateek Mittal, and Peter Henderson. 2023. Fine-tuning aligned language models compromises 1949 1950 safety, even when users do not intend to! Preprint, Hamid Roghani, Asgarali Bouyer, and Esmaeil Nourani. arXiv:2310.03693. 1951 2021. Pldls: A novel parallel label diffusion and label selection-based community detection algorithm 1952 Shengsheng Qian, Jinguang Wang, Jun Hu, Quan Fang, based on spark in social networks. Expert Systems 2006 1953 and Changsheng Xu. 2021. Hierarchical multi-modal with Applications, 183:115377. 2007 1954 contextual attention network for fake news detection. 1955 In Proceedings of the 44th international ACM SIGIR Jon Roozenbeek, Claudia R Schneider, Sarah Dryhurst, 2008 John Kerr, Alexandra LJ Freeman, Gabriel Recchia, 1956 conference on research and development in informa-1957 tion retrieval, pages 153-162. Anne Marthe Van Der Bles, and Sander Van Der Lin-2010 den. 2020. Susceptibility to misinformation about 2011 1958 Liqing Qiu, Shiqi Sai, and Moji Wei. 2022. Bpsl: a new covid-19 around the world. Royal Society open sci-2012 1959 rumor source location algorithm based on the timeence, 7(10):201199. 2013 1960 stamp back propagation in social networks. Applied 1961 Intelligence, pages 1–13. Martin Rosvall and Carl T Bergstrom. 2008. Maps of 2014 random walks on complex networks reveal commu-1962 Ori Ram, Yoav Levine, Itay Dalmedigos, Dor Muhlgay, nity structure. Proceedings of the national academy 2016 1963 Amnon Shashua, Kevin Leyton-Brown, and Yoav of sciences, 105(4):1118–1123. 2017 Shoham. 2023. In-context retrieval-augmented lan-1964 guage models. Transactions of the Association for Shailendra Sahu and T Sobha Rani. 2022. A neighbour-1965 1966 Computational Linguistics, 11:1316–1331. similarity based community discovery algorithm. Ex-2019 pert Systems with Applications, 206:117822. Yuan Rao and Jianggun Ni. 2016. A deep learning ap-1967 1968 proach to detection of splicing and copy-move forg-Chiman Salavati, Alireza Abdollahpouri, and Zhaleh 2021 1969 eries in images. In 2016 IEEE international work-Manbari. 2019. Ranking nodes in complex networks shop on information forensics and security (WIFS), based on local structure and improving closeness 1970 2023 pages 1-6. IEEE. centrality. Neurocomputing, 336:36-45. 1971 2024 Simone Raponi, Zeinab Khalifa, Gabriele Oligeri, and Amine Sallah, Said Agoujil, Mudasir Ahmad Wani, Mo-1972 2025 Roberto Di Pietro. 2022. Fake news propagation: A hamed Hammad, Ahmed A Abd El-Latif, Yassine 1973 review of epidemic models, datasets, and insights. Maleh, et al. 2024. Fine-tuned understanding: En-1974 ACM Transactions on the Web (TWEB), 16(3):1-34. hancing social bot detection with transformer-based 1975 classification. IEEE Access. Hannah Rashkin, Eunsol Choi, Jin Yea Jang, Svitlana 1976 Volkova, and Yejin Choi. 2017. Truth of varying Dietram A Scheufele and Nicole M Krause. 2019. 2030 1977 1978 shades: Analyzing language in fake news and polit-Science audiences, misinformation, and fake news. 2031 1979 ical fact-checking. In Proceedings of the 2017 con-Proceedings of the National Academy of Sciences, 2032 1980 ference on empirical methods in natural language 116(16):7662-7669. 2033 1981 processing, pages 2931–2937. Kaylyn Jackson Schiff, Daniel S Schiff, and Natália S 2034 Vipula Rawte, Amit Sheth, and Amitava Das. 2023. A Bueno. 2023. The liar's dividend: Can politicians 2035 1982 survey of hallucination in large foundation models. claim misinformation to evade accountability? Amer-1983 arXiv preprint arXiv:2309.05922. ican Political Science Review, pages 1-20. 1984

2038

2039

- 2087
- 2092

- Tal Schuster, Adam Fisch, and Regina Barzilay. 2021. Get your vitamin c! robust fact verification with contrastive evidence. arXiv preprint arXiv:2103.08541.
- M Zubair Shafiq, Muhammad U Ilyas, Alex X Liu, and Hayder Radha. 2013. Identifying leaders and followers in online social networks. IEEE Journal on Selected Areas in Communications, 31(9):618–628.
- Chintan Shah, Nima Dehmamy, Nicola Perra, Matteo Chinazzi, Albert-László Barabási, Alessandro Vespignani, and Rose Yu. 2020. Finding patient zero: Learning contagion source with graph neural networks. arXiv preprint arXiv:2006.11913.
- Chengcheng Shao, Giovanni Luca Ciampaglia, Alessandro Flammini, and Filippo Menczer. 2016. Hoaxy: A platform for tracking online misinformation. In Proceedings of the 25th international conference companion on world wide web, pages 745-750.
 - Chengcheng Shao, Giovanni Luca Ciampaglia, Onur Varol, Kai-Cheng Yang, Alessandro Flammini, and Filippo Menczer. 2018. The spread of low-credibility content by social bots. Nature communications, 9(1):1-9.
 - Junming Shao, Zhichao Han, Qinli Yang, and Tao Zhou. 2015. Community detection based on distance dynamics. In Proceedings of the 21th ACM SIGKDD International Conference on Knowledge Discovery and Data Mining, KDD '15, page 1075-1084, New York, NY, USA. Association for Computing Machinery.
 - Kai Shu, Limeng Cui, Suhang Wang, Dongwon Lee, and Huan Liu. 2019. defend: Explainable fake news detection. In Proceedings of the 25th ACM SIGKDD international conference on knowledge discovery & data mining, pages 395-405.
 - Kai Shu, Yichuan Li, Kaize Ding, and Huan Liu. 2021. Fact-enhanced synthetic news generation. In Proceedings of the AAAI Conference on Artificial Intelligence, volume 35, pages 13825–13833.
 - Kai Shu, Deepak Mahudeswaran, Suhang Wang, and Huan Liu. 2020. Hierarchical propagation networks for fake news detection: Investigation and exploitation. In Proceedings of the international AAAI conference on web and social media, volume 14, pages 626-637.
 - Kai Shu, Amy Sliva, Suhang Wang, Jiliang Tang, and Huan Liu. 2017. Fake news detection on social media: A data mining perspective. ACM SIGKDD explorations newsletter, 19(1):22-36.
- Ronit Singhal, Pransh Patwa, Parth Patwa, Aman Chadha, and Amitava Das. 2024. Evidence-backed fact checking using rag and few-shot in-context learning with llms. arXiv preprint arXiv:2408.12060.
- Shivangi Singhal, Rajiv Ratn Shah, Tanmoy Chakraborty, Ponnurangam Kumaraguru, and Shin'ichi Satoh. 2019. Spotfake: A multi-modal

framework for fake news detection. In 2019 IEEE fifth international conference on multimedia big data (BigMM), pages 39-47. IEEE.

2096

2097

2098

2099

2100

2101

2102

2103

2104

2105

2106

2107

2108

2110

2111

2112

2113

2114

2115

2116

2117

2118

2119

2120

2121

2122

2123

2124

2125

2126

2127

2128

2129

2130

2131

2132

2133

2134

2135

2136

2137

2138

2139

2140

2141

- Xing Su, Shan Xue, Fanzhen Liu, Jia Wu, Jian Yang, Chuan Zhou, Wenbin Hu, Cecile Paris, Surya Nepal, Di Jin, et al. 2022. A comprehensive survey on community detection with deep learning. IEEE Transactions on Neural Networks and Learning Systems.
- Mengzhu Sun, Xi Zhang, Jiaqi Zheng, and Guixiang Ma. 2022. Ddgcn: Dual dynamic graph convolutional networks for rumor detection on social media. In Proceedings of the AAAI conference on artificial intelligence, volume 36, pages 4611-4619.
- Yanshen Sun, Jianfeng He, Limeng Cui, Shuo Lei, and Chang-Tien Lu. 2024. Exploring the deceptive power of llm-generated fake news: A study of real-world detection challenges. arXiv preprint arXiv:2403.18249.
- Tetsuro Takahashi and Nobuyuki Igata. 2012. Rumor detection on twitter. In The 6th International Conference on Soft Computing and Intelligent Systems, and The 13th International Symposium on Advanced Intelligence Systems, pages 452–457. IEEE.
- Shulong Tan, Ziyu Guan, Deng Cai, Xuzhen Qin, Jiajun Bu, and Chun Chen. 2014. Mapping users across networks by manifold alignment on hypergraph. In Proceedings of the AAAI Conference on Artificial Intelligence, volume 28.
- Yiming Tan, Dehai Min, Yu Li, Wenbo Li, Nan Hu, Yongrui Chen, and Guilin Qi. 2023. Can chatgpt replace traditional kbga models? an in-depth analysis of the question answering performance of the gpt llm family. In International Semantic Web Conference, pages 348–367. Springer.
- Zhen Tan, Chengshuai Zhao, Raha Moraffah, Yifan Li, Song Wang, Jundong Li, Tianlong Chen, and Huan Liu. 2024. Glue pizza and eat rocks - exploiting vulnerabilities in retrieval-augmented generative models. In Proceedings of the 2024 Conference on Empirical Methods in Natural Language Processing, pages 1610-1626, Miami, Florida, USA. Association for Computational Linguistics.
- Juvaria Tariq, Muhammad Ahmad, Imdadullah Khan, and Mudassir Shabbir. 2017. Scalable approximation algorithm for network immunization. arXiv preprint arXiv:1711.00784.
- Kassym-Jomart Tokayev. 2023. Ethical implications of large language models a multidimensional exploration of societal, economic, and technical concerns. International Journal of Social Analytics, 8(9):17-33.
- Guangmo Tong, Weili Wu, and Ding-Zhu Du. 2018. 2143 Distributed rumor blocking with multiple positive 2144 cascades. IEEE Transactions on Computational So-2145 cial Systems, 5(2):468-480. 2146

2249

2250

2251

2252

Petter Törnberg, Diliara Valeeva, Justus Uitermark, and Christopher Bail. 2023. Simulating social media using large language models to evaluate alternative news feed algorithms. *arXiv preprint arXiv:2310.05984*.

2147

2148

2149

2150 2151

2153

2154

2155

2156

2157

2158

2159

2160

2161

2162

2163

2164

2165

2166 2167

2168

2169

2171

2172 2173

2174

2175

2176

2178

2179

2180

2181

2182 2183

2184

2185

2187 2188

2189

2190

2191

2192

2193

2194

2195

2196

2197 2198

2199

- Rudra M Tripathy, Amitabha Bagchi, and Sameep Mehta. 2010. A study of rumor control strategies on social networks. In *Proceedings of the 19th ACM international conference on Information and knowledge management*, pages 1817–1820.
- U Undeutsch. 1967. Beurteilung der glaubhaftigkeit von aussagen [veracity assessment of statements]. undeutsch. *Handbuch der Psychologie*, 11:26–181.
- Cristian Vaccari and Andrew Chadwick. 2020. Deepfakes and disinformation: Exploring the impact of synthetic political video on deception, uncertainty, and trust in news. *Social media+ society*, 6(1):2056305120903408.
- Karthik Valmeekam, Matthew Marquez, Alberto Olmo, Sarath Sreedharan, and Subbarao Kambhampati.
 2023. Planbench: An extensible benchmark for evaluating large language models on planning and reasoning about change. *Advances in Neural Information Processing Systems*, 36:38975–38987.
- Neeraj Varshney, Wenlin Yao, Hongming Zhang, Jianshu Chen, and Dong Yu. 2023. A stitch in time saves nine: Detecting and mitigating hallucinations of llms by validating low-confidence generation. *arXiv* preprint arXiv:2307.03987.
- Nikhita Vedula and Srinivasan Parthasarathy. 2021. Face-keg: Fact checking explained using knowledge graphs. In *Proceedings of the 14th ACM International Conference on Web Search and Data Mining*, pages 526–534.
- Andreas Vlachos and Sebastian Riedel. 2014. Fact checking: Task definition and dataset construction. In Proceedings of the ACL 2014 workshop on language technologies and computational social science, pages 18–22.
- Soroush Vosoughi, Deb Roy, and Sinan Aral. 2018. The spread of true and false news online. *science*, 359(6380):1146–1151.
- Ivan Vykopal, Matúš Pikuliak, Ivan Srba, Robert Moro, Dominik Macko, and Maria Bielikova. 2023. Disinformation capabilities of large language models. *arXiv preprint arXiv:2311.08838*.
- David Wadden, Kyle Lo, Lucy Lu Wang, Arman Cohan, Iz Beltagy, and Hannaneh Hajishirzi. 2021. Multivers: Improving scientific claim verification with weak supervision and full-document context. *arXiv preprint arXiv:2112.01640*.
- Nathan Walter and Riva Tukachinsky. 2020. A metaanalytic examination of the continued influence of misinformation in the face of correction: How powerful is it, why does it happen, and how to stop it? *Communication research*, 47(2):155–177.

- Chen Wan, Tao Li, and Zhicheng Sun. 2017. Global stability of a seir rumor spreading model with demographics on scale-free networks. *Advances in Difference Equations*, 2017(1):253.
- Herun Wan, Shangbin Feng, Zhaoxuan Tan, Heng Wang, Yulia Tsvetkov, and Minnan Luo. 2024. Dell: Generating reactions and explanations for Ilm-based misinformation detection. *arXiv preprint arXiv:2402.10426*.
- Biao Wang, Ge Chen, Luoyi Fu, Li Song, and Xinbing Wang. 2017. Drimux: Dynamic rumor influence minimization with user experience in social networks. *IEEE Transactions on Knowledge and Data Engineering*, 29(10):2168–2181.
- Bo Wang, Jing Ma, Hongzhan Lin, Zhiwei Yang, Ruichao Yang, Yuan Tian, and Yi Chang. 2024a. Explainable fake news detection with large language model via defense among competing wisdom. In *Proceedings of the ACM on Web Conference 2024*, pages 2452–2463.
- Bo Wang, Jing Ma, Hongzhan Lin, Zhiwei Yang, Ruichao Yang, Yuan Tian, and Yi Chang. 2024b. Explainable fake news detection with large language model via defense among competing wisdom. In *Proceedings of the ACM Web Conference 2024*, WWW '24, page 2452–2463, New York, NY, USA. Association for Computing Machinery.
- Chenxi Wang, Zongfang Liu, Dequan Yang, and Xiuying Chen. 2024c. Decoding echo chambers: Llmpowered simulations revealing polarization in social networks. *arXiv preprint arXiv:2409.19338*.
- Junxiang Wang, Junji Jiang, and Liang Zhao. 2022. An invertible graph diffusion neural network for source localization. In *Proceedings of the ACM Web Conference 2022*, pages 1058–1069.
- Sheng-Yu Wang, Oliver Wang, Richard Zhang, Andrew Owens, and Alexei A Efros. 2020a. Cnn-generated images are surprisingly easy to spot... for now. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 8695–8704.
- Yaqing Wang, Fenglong Ma, Zhiwei Jin, Ye Yuan, Guangxu Xun, Kishlay Jha, Lu Su, and Jing Gao. 2018a. Eann: Event adversarial neural networks for multi-modal fake news detection. In *Proceedings* of the 24th acm sigkdd international conference on knowledge discovery & data mining, pages 849–857.
- Yi Wang, Jinde Cao, Xiaodi Li, and Ahmed Alsaedi. 2018b. Edge-based epidemic dynamics with multiple routes of transmission on random networks. *Nonlinear Dynamics*, 91:403–420.
- Youze Wang, Shengsheng Qian, Jun Hu, Quan Fang,
and Changsheng Xu. 2020b. Fake news detection via
knowledge-driven multimodal graph convolutional
networks. In Proceedings of the 2020 international
conference on multimedia retrieval, pages 540–547.2253
2254
2256

2258 2259 2260 2261	Zecong Wang, Jiaxi Cheng, Chen Cui, and Chenhao Yu. 2023. Implementing bert and fine-tuned roberta to detect ai generated news by chatgpt. <i>arXiv preprint arXiv:2306.07401</i> .	Yang Xia, Haijun Jiang, Shuzhen Yu, and Zhiyong Yu. 2024. The dynamic analysis of the rumor spreading and behavior diffusion model with higher-order inter- actions. <i>Communications in Nonlinear Science and</i> <i>Numerical Simulation</i> , 138:108186.
2262 2263 2264	Alicia Wanless and Michael Berk. 2020. The audience is the amplifier: Participatory propaganda. <i>The SAGE</i> <i>handbook of propaganda</i> , pages 85–104.	Bin Xiao, Yang Wei, Xiuli Bi, Weisheng Li, and Jian- feng Ma. 2020. Image splicing forgery detection combining coarse to refined convolutional neural net-
2265	Jason Wei, Xuezhi Wang, Dale Schuurmans, Maarten	work and adaptive clustering. <i>Information Sciences</i> , 511:172–191
2266	Bosma, Fei Xia, Ed Chi, Quoc V Le, Denny Zhou,	
2207	soning in large language models. Advances in neural	Yunpeng Xiao, Diqiang Chen, Shihong Wei, Qian Li,
2269	information processing systems, 35:24824–24837.	tion dynamic model based on evolutionary game and anti-rumor. <i>Nonlinear Dynamics</i> , 95:523–539.
2270	Jianliang Wei and Fei Meng. 2021. How opinion dis-	Vennen Vice linear Ven Werling 7hos Oier Li
2271	tortion appears in super-influencer dominated so-	runpeng Xiao, Jinsong Yang, Wanjing Zhao, Qian Li, and Yucai Pang. 2024. Cross domain social rumor
2272	115:542–552.	propagation model based on transfer learning. <i>IEEE</i> <i>Transactions on Neural Networks and Learning Sys-</i>
2274	Teresa Weikmann and Sophie Lecheler. 2023. Visual	tems.
2275	disinformation in a digital age: A literature synthe-	Chengxing Xie, Canyu Chen, Feiran Jia, Ziyu Ye,
2276 2277	25(12):3696–3713.	Kai Shu, Adel Bibi, Ziniu Hu, Philip Torr, Bernard Ghanem, and Guohao Li. 2024. Can large language model agents simulate human trust behaviors? <i>arXiv</i>
2278	Ferre Wouters and Michaël Opgenhaffen. 2024. Re-	preprint arXiv:2402.04559.
2279	sub-state fact-checking initiatives in europe <i>Media</i>	Jierui Xie, Stephen Kelley, and Boleslaw K Szymanski
2281	and Communication, 12.	2013. Overlapping community detection in networks: The state-of-the-art and comparative study. <i>Acm com</i> -
2282	Feijie Wu, Zitao Li, Yaliang Li, Bolin Ding, and Jing	puting surveys (csur), $45(4)$:1–35.
2283	ated learning without full model. In <i>Proceedings of</i>	Yueqi Xie, Jingwei Yi, Jiawei Shao, Justin Curl,
2285 2286	the 30th ACM SIGKDD Conference on Knowledge Discovery and Data Mining, pages 3345–3355.	Lingjuan Lyu, Qifeng Chen, Xing Xie, and Fangzhao Wu. 2023. Defending chatgpt against jailbreak at-
		5(12):1486–1496
2287	Guangyang Wu, Weijie Wu, Xiaohong Liu, Kele Xu,	5(12).1400 1490.
2288	detection with llm using prompt engineering. In 2023	Miao Xiong, Zhiyuan Hu, Xinyang Lu, Yifei Li, Jie
2290 2291	<i>IEEE International Conference on Multimedia and Expo Workshops (ICMEW)</i> , pages 105–109. IEEE.	express their uncertainty? an empirical evaluation of confidence elicitation in llms. <i>arXiv preprint</i> <i>arXiv:2306.13063</i> .
2292	Jiaying Wu, Jiafeng Guo, and Bryan Hooi. 2024b. Fake	Junhaa Vu Langdi Vian Zaning Liu Mingliang Chan
2293	news in sheep's clothing: Robust fake news detection	Oiuvang Vin and Fenghua Song 2024a The future
2294	of the 30th ACM SIGKDD conference on knowledge	of combating rumors? retrieval, discrimination, and
2296	discovery and data mining, pages 3367–3378.	generation. arXiv preprint arXiv:2403.20204.
		Rongwu Xu, Brian Lin, Shujian Yang, Tiangi Zhang
2297	Jun Wu, Xuesong Ye, and Chengjie Mou. 2023b. Bot-	Weiyan Shi, Tianwei Zhang, Zhixuan Fang, Wei Xu,
2298 2299	havioral patterns. <i>arXiv preprint arXiv:2303.10214</i> .	and Han Qiu. 2024b. The earth is flat because: Investigating LLMs' belief towards misinformation
2300	Ke Wu, Song Yang, and Kenny Q Zhu. 2015. False	Via persuasive conversation. In <i>Proceedings of the</i> 62nd Annual Meeting of the Association for Compu-
2301	rumors detection on sina weibo by propagation struc-	tational Linguistics (Volume 1: Long Papers), pages
2302	tures. In 2015 IEEE 31st international conference on	16259–16303, Bangkok, Thailand. Association for
2303	aata engineering, pages 651–662. IEEE.	Computational Linguistics.
2304	Siyuan Wu, Yue Huang, Chujie Gao, Dongping Chen,	Rongwu Xu, Brian S Lin, Shujian Yang, Tiangi Zhang,
2305	Qihui Zhang, Yao Wan, Tianyi Zhou, Xiangliang	Weiyan Shi, Tianwei Zhang, Zhixuan Fang, Wei
2306	Zhang, Jianfeng Gao, Chaowei Xiao, et al. 2024c.	Xu, and Han Qiu. 2023. The earth is flat be-
2307	Unigen: A unified framework for textual dataset gen- eration using large language models arXiv preprint	cause: investigating lims belief towards misinfor- mation via persuasive conversation arXiv preprint
2309	arXiv:2406.18966.	arXiv:2312.09085.
	2	3

2311

2312

2313 2314

2315

2316

2317

2318 2319

2320 2321

2322

2323

2324

2325

2326

2327 2328

2329

2330

2331

2332 2333 2334

2335

2336 2337

2338

2339

2340

2341 2342

2343

2344

2345

2346 2347

2348

2349

2350

2351

2352

2353 2354

2355

2356

2357

2358

2359 2360

2361

2362

2363

2364

2365 2366

Ziwei Xu, Sanjay Jain, and Mohan Kankanhalli. 2024c. Hallucination is inevitable: An innate limitation of large language models. *arXiv preprint arXiv:2401.11817*.

2367

2370

2371

2376

2379

2380

2381

2388

2389

2390

2391

2395

2396

2397

2398

2399

2402

2404

2406

2407

2408

2409

2410

2411

2412

2413

2414

2415

2416

2417

2418

- Qi Xuan, Xincheng Shu, Zhongyuan Ruan, Jinbao Wang, Chenbo Fu, and Guanrong Chen. 2019. A self-learning information diffusion model for smart social networks. *IEEE Transactions on Network Science and Engineering*, 7(3):1466–1480.
- Ruidong Yan, Yi Li, Weili Wu, Deying Li, and Yongcai Wang. 2019. Rumor blocking through online link deletion on social networks. *ACM Transactions on Knowledge Discovery from Data (TKDD)*, 13(2):1– 26.
- Chang Yang, Xia Yu, JiaYi Wu, BoZhen Zhang, and HaiBo Yang. 2024a. Graph-aware multi-feature interacting network for explainable rumor detection on social network. *Expert Systems with Applications*, 249:123687.
- Chang Yang, Peng Zhang, Hui Gao, and Jing Zhang. 2024b. Deciphering rumors: A multi-task learning approach with intent-aware hierarchical contrastive learning. In *Proceedings of the 2024 Conference on Empirical Methods in Natural Language Processing*, pages 4471–4483.
- Chang Yang, Peng Zhang, Wenbo Qiao, Hui Gao, and Jiaming Zhao. 2023a. Rumor detection on social media with crowd intelligence and chatgpt-assisted networks. In *Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing*, pages 5705–5717.
- Fan Yang, Yang Liu, Xiaohui Yu, and Min Yang. 2012. Automatic detection of rumor on sina weibo. In Proceedings of the ACM SIGKDD Workshop on Mining Data Semantics, MDS '12, New York, NY, USA. Association for Computing Machinery.
- Fan Yang, Shiva K Pentyala, Sina Mohseni, Mengnan Du, Hao Yuan, Rhema Linder, Eric D Ragan, Shuiwang Ji, and Xia Hu. 2019. Xfake: Explainable fake news detector with visualizations. In *The world wide web conference*, pages 3600–3604.
- Jaewon Yang and Jure Leskovec. 2013. Overlapping community detection at scale: a nonnegative matrix factorization approach. In *Proceedings of the sixth ACM international conference on Web search and data mining*, pages 587–596.
- Kaijun Yang, Qing Bao, and Hongjun Qiu. 2023b. Identifying multiple propagation sources with motifbased graph convolutional networks for social networks. *IEEE Access*, 11:61630–61645.
- Li Yang, Yafeng Qiao, Zhihong Liu, Jianfeng Ma, and Xinghua Li. 2018. Identifying opinion leader nodes in online social networks with a new closeness evaluation algorithm. *Soft Computing*, 22:453–464.

- Liang Yang, Xiaochun Cao, Dongxiao He, Chuan Wang, Xiao Wang, and Weixiong Zhang. 2016. Modularity based community detection with deep learning. In *IJCAI*, volume 16, pages 2252–2258.
- Ruichao Yang, Jing Ma, Wei Gao, and Hongzhan Lin. 2025. Llm-enhanced multiple instance learning for joint rumor and stance detection with social context information. *ACM Transactions on Intelligent Systems and Technology*.
- Tianbao Yang, Rong Jin, Yun Chi, and Shenghuo Zhu. 2009. Combining link and content for community detection: a discriminative approach. In *Proceedings of the 15th ACM SIGKDD international conference on Knowledge discovery and data mining*, pages 927– 936.
- Xianjun Yang, Xiao Wang, Qi Zhang, Linda Petzold, William Yang Wang, Xun Zhao, and Dahua Lin. 2023c. Shadow alignment: The ease of subverting safely-aligned language models. *arXiv preprint arXiv:2310.02949*.
- Zhengyuan Yang, Linjie Li, Kevin Lin, Jianfeng Wang, Chung-Ching Lin, Zicheng Liu, and Lijuan Wang. 2023d. The dawn of lmms: Preliminary explorations with gpt-4v (ision). *arXiv preprint arXiv:2309.17421*, 9(1):1.
- Barry Menglong Yao, Aditya Shah, Lichao Sun, Jin-Hee Cho, and Lifu Huang. 2023a. End-to-end multimodal fact-checking and explanation generation: A challenging dataset and models. In *Proceedings of the 46th International ACM SIGIR Conference on Research and Development in Information Retrieval*, pages 2733–2743.
- Shunyu Yao, Dian Yu, Jeffrey Zhao, Izhak Shafran, Tom Griffiths, Yuan Cao, and Karthik Narasimhan. 2023b. Tree of thoughts: Deliberate problem solving with large language models. *Advances in neural information processing systems*, 36:11809–11822.
- Xiaopeng Yao, Yue Gu, Chonglin Gu, and Hejiao Huang. 2022. Fast controlling of rumors with limited cost in social networks. *Computer Communications*, 182:41–51.
- Yining Ye, Xin Cong, Shizuo Tian, Jiannan Cao, Hao Wang, Yujia Qin, Yaxi Lu, Heyang Yu, Huadong Wang, Yankai Lin, et al. 2023. Proagent: From robotic process automation to agentic process automation. *arXiv preprint arXiv:2311.10751*.
- Shuzhen Yu, Zhiyong Yu, Haijun Jiang, and Shuai Yang. 2021. The dynamics and control of 2i2sr rumor spreading models in multilingual online social networks. *Information Sciences*, 581:18–41.
- Damián H Zanette. 2001. Critical behavior of propagation on small-world networks. *Physical Review E*, 64(5):050901.
- Runxi Zeng and Di Zhu. 2019. A model and simulation2474of the emotional contagion of netizens in the process2475of rumor refutation. Scientific reports, 9(1):14164.2476

- 2477 2478 2479 2480 2481 2482 2483 2484 2485 2486 2490 2492 2493 2494 2495 2496 2497 2499 2503 2504 2505 2506 2507 2508 2509 2512 2514 2515 2516
- 2515 2516 2517 2518 2519 2520
- 2521 2522 2523 2524
- 2524 2525 2526 2527
- 2528 2529 2530
- 25 25

- Huaiwen Zhang, Quan Fang, Shengsheng Qian, and Changsheng Xu. 2019. Multi-modal knowledgeaware event memory network for social media rumor detection. In *Proceedings of the 27th ACM international conference on multimedia*, pages 1942–1951.
- Litian Zhang, Xiaoming Zhang, Ziyi Zhou, Feiran Huang, and Chaozhuo Li. 2024a. Reinforced adaptive knowledge learning for multimodal fake news detection. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 38, pages 16777– 16785.
- Litian Zhang, Xiaoming Zhang, Ziyi Zhou, Xi Zhang, S Yu Philip, and Chaozhuo Li. 2025. Knowledgeaware multimodal pre-training for fake news detection. *Information Fusion*, 114:102715.
- Mingqing Zhang, Haisong Gong, Qiang Liu, Shu Wu, and Liang Wang. 2024b. Breaking event rumor detection via stance-separated multi-agent debate. *arXiv preprint arXiv:2412.04859*.
- Qi Zhang, Yuan Li, Jialing Zou, Jianming Zhu, Dingning Liu, and Jianbin Jiao. 2024c. Unifying multimodal interactions for rumor diffusion prediction with global hypergraph modeling. *Knowledge-Based Systems*, 301:112246.
- Si Zhang and Hanghang Tong. 2016. Final: Fast attributed network alignment. In *Proceedings of the* 22nd ACM SIGKDD international conference on knowledge discovery and data mining, pages 1345– 1354.
- Xuan Zhang and Wei Gao. 2023a. Towards llm-based fact verification on news claims with a hierarchical step-by-step prompting method. *arXiv preprint arXiv:2310.00305*.
- Xuan Zhang and Wei Gao. 2023b. Towards LLM-based fact verification on news claims with a hierarchical step-by-step prompting method. In Proceedings of the 13th International Joint Conference on Natural Language Processing and the 3rd Conference of the Asia-Pacific Chapter of the Association for Computational Linguistics (Volume 1: Long Papers), pages 996–1011, Nusa Dua, Bali. Association for Computational Linguistics.
- Yao Zhang and B Aditya Prakash. 2015. Data-aware vaccine allocation over large networks. *ACM Transactions on Knowledge Discovery from Data (TKDD)*, 10(2):1–32.
- Yuan Zhang, Tianshu Lyu, and Yan Zhang. 2018. Cosine: Community-preserving social network embedding from information diffusion cascades. In *Proceedings of the AAAI conference on artificial intelligence*, volume 32.
- Yutao Zhang, Jie Tang, Zhilin Yang, Jian Pei, and Philip S Yu. 2015. Cosnet: Connecting heterogeneous social networks with local and global consistency. In *Proceedings of the 21th ACM SIGKDD international conference on knowledge discovery and data mining*, pages 1485–1494.

Jian Zhao, Nan Cao, Zhen Wen, Yale Song, Yu-Ru Lin, and Christopher Collins. 2014. # fluxflow: Visual analysis of anomalous information spreading on social media. *IEEE transactions on visualization and computer graphics*, 20(12):1773–1782. 2534

2537

2539

2540

2541

2542

2543

2546

2548

2552

2555

2557

2558

2560

2561

2563

2564

2565

2566

2570

2571

2572

2573

2574

2575

2577

2578

2579

2581

2583

2584

2586

- Laijun Zhao, Hongxin Cui, Xiaoyan Qiu, Xiaoli Wang, and Jiajia Wang. 2013. Sir rumor spreading model in the new media age. *Physica A: Statistical Mechanics and its Applications*, 392(4):995–1003.
- Ruochen Zhao, Hailin Chen, Weishi Wang, Fangkai Jiao, Xuan Long Do, Chengwei Qin, Bosheng Ding, Xiaobao Guo, Minzhi Li, Xingxuan Li, et al. 2023a. Retrieving multimodal information for augmented generation: A survey. *arXiv preprint arXiv:2303.10868*.
- Wayne Xin Zhao, Kun Zhou, Junyi Li, Tianyi Tang, Xiaolei Wang, Yupeng Hou, Yingqian Min, Beichen Zhang, Junjie Zhang, Zican Dong, et al. 2023b. A survey of large language models. *arXiv preprint arXiv:2303.18223*.
- Yiran Zhao, Jinghan Zhang, I Chern, Siyang Gao, Pengfei Liu, Junxian He, et al. 2023c. Felm: Benchmarking factuality evaluation of large language models. *Advances in Neural Information Processing Systems*, 36:44502–44523.
- Yuxin Zhao, Shenghong Li, and Feng Jin. 2016. Identification of influential nodes in social networks with community structure based on label propagation. *Neurocomputing*, 210:34–44.
- Wanjun Zhong, Jingjing Xu, Duyu Tang, Zenan Xu, Nan Duan, Ming Zhou, Jiahai Wang, and Jian Yin. 2019. Reasoning over semantic-level graph for fact checking. arXiv preprint arXiv:1909.03745.
- Jiawei Zhou, Yixuan Zhang, Qianni Luo, Andrea G Parker, and Munmun De Choudhury. 2023. Synthetic lies: Understanding ai-generated misinformation and evaluating algorithmic and human solutions. In *Proceedings of the 2023 CHI Conference on Human Factors in Computing Systems*, pages 1–20.
- Jie Zhou, Xu Han, Cheng Yang, Zhiyuan Liu, Lifeng Wang, Changcheng Li, and Maosong Sun. 2019a. Gear: Graph-based evidence aggregating and reasoning for fact verification. *arXiv preprint arXiv:1908.01843*.
- Peng Zhou, Xintong Han, Vlad I Morariu, and Larry S Davis. 2018. Learning rich features for image manipulation detection. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 1053–1061.
- Xiaoping Zhou, Xun Liang, Haiyan Zhang, and Yuefeng Ma. 2015. Cross-platform identification of anonymous identical users in multiple social media networks. *IEEE transactions on knowledge and data engineering*, 28(2):411–424.

Xinyi Zhou and Reza Zafarani. 2018. Fake news: A survey of research, detection methods, and opportunities. arXiv preprint arXiv:1812.00315, 2.

2591

2592

2594

2595

2597

2599

2600

2601

2603

2604

2605

2606

2607

2609

2610

2611

2613

2614

2615

2616

2617

2618

2619

2621

2628

2629

2630

2631

2632

2634

2636

2640

- Zhixuan Zhou, Huankang Guan, Meghana Moorthy Bhat, and Justin Hsu. 2019b. Fake news detection via nlp is vulnerable to adversarial attacks. *arXiv preprint arXiv:1901.09657*.
- Biru Zhu, Xingyao Zhang, Ming Gu, and Yangdong Deng. 2021. Knowledge enhanced fact checking and verification. *IEEE/ACM Transactions on Audio*, *Speech, and Language Processing*, 29:3132–3143.
- Lingxuan Zhu, Weiming Mou, and Peng Luo. 2024. Potential of large language models as tools against medical disinformation. *JAMA Internal Medicine*, 184(4):450–450.
- Linhe Zhu and Xiaoyuan Huang. 2019. Sis model of rumor spreading in social network with time delay and nonlinear functions. *Communications in Theoretical Physics*, 72(1):015002.
- Peican Zhu, Le Cheng, Chao Gao, Zhen Wang, and Xuelong Li. 2022a. Locating multi-sources in social networks with a low infection rate. *IEEE Transactions* on Network Science and Engineering, 9(3):1853– 1865.
- Yongchun Zhu, Qiang Sheng, Juan Cao, Shuokai Li, Danding Wang, and Fuzhen Zhuang. 2022b. Generalizing to the future: Mitigating entity bias in fake news detection. In *Proceedings of the 45th International ACM SIGIR Conference on Research and Development in Information Retrieval*, pages 2120– 2125.
- Linlin Zong, Jiahui Zhou, Wenmin Lin, Xinyue Liu, Xianchao Zhang, and Bo Xu. 2024. Unveiling opinion evolution via prompting and diffusion for short video fake news detection. In *Findings of the Association for Computational Linguistics ACL 2024*, pages 10817–10826.
- Arkaitz Zubiaga, Maria Liakata, Rob Procter, Kalina Bontcheva, and Peter Tolmie. 2015. Towards detecting rumours in social media. In Workshops at the Twenty-Ninth AAAI conference on artificial intelligence.
- M Zuckerman. 1981. Verbal and nonverbal communication of deception. Advances in experimental social psychology/Academic Press.

A Rumor Definition and Related Tasks

The core characteristic of rumors lies in their "unverified ambiguity and uncertainty," which makes them highly prone to misinterpretation or misuse during the dissemination process. Unlike debunked false information (misinformation) (Scheufele and Krause, 2019; Kumar and Geethakumari, 2014), deliberately fabricated falsehoods (disinformation) (Guo et al., 2020), or fake news that adopts the 2641 form of journalistic reporting to deliberately mis-2642 lead the public (Shu et al., 2017, 2019) (Detecting 2643 fake news with NLP)], the uniqueness of rumors lies in the dynamic evolution of their verification status. Currently, rumor detection in a broad sense 2646 largely focuses on the verification of rumor verac-2647 ity, emphasizing the description of the potential 2648 risks posed by false rumors to societal trust (Takahashi and Igata, 2012; Wu et al., 2015; Liang et al., 2650 2015). From a narrower perspective, studies on 2651 rumors also consider their dissemination character-2652 istics and societal impacts (Allport, 1947; Zubiaga 2653 et al., 2015). This provides theoretical support for 2654 uncovering the deeper logic underpinning rumor 2655 propagation while laying the foundational frame-2656 work for research in rumor detection.

2658

2659

2660

2661

2662

2663

2664

2667

2670

2671

2673

2674

2675

2676

2677

2678

2679

2681

2682

2686

2687

2690

B Future Research Directions

B.1 LLM-based Multi-agent Social Simulation in Rumor Detection

While existing studies have shed light on the potential of rumor detection agents, there remains a lack of in-depth research into the complexities of the task. Critical gaps persist in addressing the key aspects of handling the complexity of rumor propagation, which can be summarized into the following three areas:

Lack of a deep understanding of the mechanisms of rumor propagation. Current research often relies on shallow features to classify information but fails to account for rumor dissemination's intricate cognitive behavioral patterns and sociodynamic characteristics. For instance, individuals tend to process information in ways that align with their pre-existing beliefs(cognitive consistency theory) while exhibiting significant non-objective and selective cognitive biases (Nickerson, 1998). Emotional drivers (such as fear and anger)(Chan et al., 2021), and social pressures (Blass, 1984) amplify the effects of rumor propagation. Moreover, some users spread false information not purely based on judgments of its truthfulness but rather due to hedonistic motives (Jiwa et al., 2023) or a sense of group identity (Lewandowsky, 2022; Wanless and Berk, 2020). Social networks' distributed and decentralized nature further reduces individuals' capacity to discern false information. It increases the likelihood of social cascades, forming extreme attitudes (Jamieson, 2008). This aspect has also received attention, e.g., simulating social phenomena

2731

2730

2735

2732 2733

2736

2737 2738

2740

2741

like echo chambers (Wang et al., 2024c).

Over-simplified modeling of rumor propaga-

tion and intervention. Current models are often restricted to static propagation roles (e.g., spreaders, bystanders, fact-checkers) and fail to comprehensively capture the dynamic changes in individual behaviors during the dissemination process and the dynamic role shifts in group interactions. In contrast, the information propagation dynamics extensively studied in information epidemiology offer valuable inspiration, such as integrating complex social variables (educational levels, forgetting mechanisms) to model dynamic processes more accurately. At the same time, intervention measures also lack flexibility and dynamic optimization regarding accurate verification and influencer blocking. This limits their ability to adapt to the gradual and complex evolution of rumor propagation, which often requires intervention strategies that are adaptable and sensitive to contextual changes.

Insufficient global modeling of the dynamic attributes of rumors. Existing studies are often confined to a specific dimension of rumors, such as content, propagation, or interaction, without systematically integrating these interdependent elements within a logical framework. Furthermore, the combined effect of multimodal information on rumor dissemination and the critical role of social bots as key propagation drivers have not been sufficiently considered. Additionally, active user behaviors, such as retrieval and learning of controversial information, remain underexplored. Although individual differences exist, such behaviors (such as actively verifying the authenticity of information and subsequently choosing to share, report, or ignore it) significantly impact information perception and dissemination.

Multi-Agent Systems Under the CIB **B.1.1** Framework

To tackle existing challenges, we propose a research roadmap for multi-agent systems within the CIB framework, as illustrated in Figure 3. This approach utilizes cross-layer dynamic feedback to establish a Macro-Micro Feedback Loop, enabling the modeling of macro-level information dissemination based on micro-level individual cognition. It facilitates dynamic, bidirectional interventions between macro and micro levels, fostering evolution driven by deeper cognitive insights.

At the cognitive layer, agents can utilize LLMs and multimodal analysis tools to achieve a deep semantic understanding of text, images, videos, and other content associated with rumors, also performing real-time monitoring and dynamic analysis of content flow on social media platforms. By incorporating psychological models, agents can dynamically assess the authenticity of information and precisely quantify user cognitive biases (e.g., selective processing or emotion-driven behaviors). For instance, agents can infer malicious intents behind false information through context-aware reasoning and identify the potential motivations of target users for spreading such information. This enables agents to provide information support for subsequent collaborative actions.

2742

2743

2744

2745

2746

2747

2748

2749

2750

2751

2752

2753

2754

2755

2756

2757

2758

2759

2760

2761

2762

2763

2764

2765

2766

2767

2768

2769

2770

2771

2772

2773

2774

2775

2776

2777

2778

2779

2780

2781

2782

2783

2784

2785

2786

2787

2788

2789

2790

2791

2792

2793

At the interaction layer, multi-agent systems leverage sophisticated communication protocols to collaborate efficiently, simulating the diversity of user behaviors in social networks, such as information sharing, commenting, and reporting. Moreover, they can also simulate external factors such as social bots, constructing dynamic environments that better reflect real-world propagation patterns. By comprehensively modeling collective knowledge and social interactions, these dynamic simulations help address critical questions, such as: How can we generate more targeted intervention strategies with greater accuracy? How can we improve collaborative efficiency when users participate in rumor reporting or debunking efforts?

At the behavior layer, multi-agent systems can capture the nonlinear propagation paths of collective behavior and dynamically generate personalized debunking content and intervention strategies through cognitive modeling. For example, By analyzing a target audience's cognitive and behavioral characteristics, agents can produce debunking content that is more persuasive and tailored to the audience. In response to the dynamic evolution of rumor propagation, agents can adaptively adjust intervention models, enabling more efficient and precise countermeasures.

Furthermore, in addition to modeling information dissemination of cognition to behavior, the framework enables a dynamic rumor detection and intervention mechanism in reverse. At the behavior layer, agents detect anomalous communities as initial targets. Subsequently, agents at the interaction layer analyze the suspicious users and interaction structures within these communities. Building on this, agents in the cognition layer perform collective knowledge analysis and social context reasoning on suspicious conversational threads to identify



Figure 3: The CIB framework establishes a Macro-Micro Feedback Loop that integrates cross-layer dynamic feedback to bridge macro-level information dissemination with micro-level individual cognition. This enables bidirectional interventions, wherein macro-level propagation dynamics inform micro-level behavior modeling, while individual insights refine system-wide strategies. Beyond modeling the flow from cognition to behavior, the framework supports a reverse feedback mechanism for dynamic rumor detection and intervention, driving continuous adaptation and iterative improvement in complex, evolving misinformation environments.

critical evidence. Finally, the behavior layer intervenes promptly, generating personalized debunking content tailored to the cognitive characteristics of the target audience as part of belief-based interventions. This feedback mechanism allows rumor detection and intervention models to continuously improve and optimize themselves, enhancing their adaptability to the dynamically evolving nature of rumor propagation.

2794

2795

2796

2797 2798

2799

2804

2809

2810

2811

2812

2813 2814

2815

2818

B.2 Cross-disciplinary Collaborative Optimization

LLMs also face significant challenges across three dimensions: content credibility, cognitive alignment, and technology adaptability (Liu et al., 2024a; De Angelis et al., 2023). First, the intertwining of hallucinated content generated by LLMs with malicious misinformation substantially complicates assessing content reliability (Pan et al., 2023c; Chen and Shu, 2024; Shu et al., 2021). Technologies such as deepfakes mislead the public and stigmatize genuine content, further eroding transparency in public discourse (e.g., real but negative content dismissed as synthetic and subsequently discredited (Schiff et al., 2023)). Second, limitations in LLMs' semantic alignment and ability to perform complex reasoning may reinforce users' preexisting cognitive biases (Xu et al., 2023; Garry et al., 2024; Hosseini et al., 2023). In complex scenes, the robustness and dynamic adaptability of the model (Carlini et al., 2023; Haller et al., 2023), as well as early benchmarks (Zhou et al., 2023; Chen and Shu, 2023), show weak performance. (He et al., 2023; Li et al., 2024b). To address these challenges, future research must promote crossdisciplinary collaborative optimization: 2819

2821

2823

2824

2827

2831

2832

2833

2836

2839

2841

2843

2844

Cognition layer: Enhancing content credibility and knowledge attribution. Future research should focus on building dual mechanisms that combine technological forensics and social validation to improve content credibility, particularly for attributing and explaining the trustworthiness of LLM-generated content (André et al., 2023; Kumarage and Liu, 2023). Another area of importance is addressing "technology-amplified cognitive biases" and their impact on information credibility perception and cognitive schema activation. For instance, The high fluency of LLM-generated content amplifies the halo effect (Augenstein et al., 2024); Multimodal synthetic content, such as deepfakes, enhances emotional infiltration, making misinformation harder to detect.

2901

2902

2903

2904

2905

2906

2907

2908

2909

2910

2911

2912

2913

2914

2915

2917

2918

2919

2921

2922

2925

Interaction layer: strengthening cognitive alignment and belief intervention in social interactions. Integrating psychology and behavioral science insights can enable belief interventions using LLM-generated personalized and persuasive debunking content (Costello et al., 2024; Matz et al., 2024). For example, Employing dynamic intervention strategies, such as the "Friction Strategy," can suppress blind adherence and impulsive information sharing by increasing the cognitive processing cost of user decisions; LLMs could also provide real-time knowledge enhancement services for low-education user groups, helping mitigate the continued influence effect (Lewandowsky et al., 2012; Walter and Tukachinsky, 2020) (where misinformation persists even after being debunked) and the recognition gap driven by educational disparities (Afassinou, 2014; Hui et al., 2020). This approach fosters a synergistic effect between educational compensation and behavioral interventions.

2846

2847

2851

2858

2864

2868

2870

2871

2873

2875

2879

2881

2884

2887

2891

Behavior layer: Enhancing robustness and adaptability in dynamic environments. To improve the adaptability of LLMs in dynamic environments, future research can leverage LLMagent architectures that perform multi-phase reasoning chains (e.g., planning-execution-reflection), enabling zero-shot automatic feature annotation and latent rumor identification in real-time. Additionally, a more comprehensive rumor baseline assessment should be established to address the evolving characteristics of misinformation in complex environments.

By integrating beneficial insights from across disciplines into a rumor detection framework based on collective intelligence, this research line can promote the development of efficient and highly adaptable systems, pushing the boundaries of rumor detection methods.

B.3 Information Ecosystem Governance Under Multi-Multi-dimensionalraints

In the context of rumor governance, the synergistic governance of legal, ethical, and technological constraints emerges as a necessary approach.

Legal measures should focus on regulating data usage while ensuring privacy protection. Privacypreserving technologies(such as differential privacy (Chua et al., 2024; Mai et al., 2023) and federated learning](Kuang et al., 2024; Wu et al., 2024a; Ezzeldin et al., 2023)), combined with compliance frameworks(General Data Protection Regulation (GDPR) and the Artificial Intelligence Act (AI Act) (Parliament, 2023)), enhance model performance while safeguarding data security. These measures serve as a foundation for responsible rumor detection and governance in the digital age.

LLMs have been shown to possess the capability of inferring psychological tendencies from user-generated texts (Peters and Matz, 2024; Perc et al., 2019), potentially influencing users' false memories (Chan et al., 2024; Acerbi and Stubbersfield, 2023). Furthermore, LLM-generated content could be weaponized for privacy infringements, cognitive attacks, and social media manipulation (Huang and Zhu, 2023; Park et al., 2024). To address these challenges, platforms, and developers should proactively disclose algorithm designs, ensure data sources and security measures, and establish transparent accountability chains to enhance transparency and responsibility allocation (Dwivedi et al., 2023).

At the social governance level, advancing multistakeholder collaborative mechanisms is essential. This involves building a governance ecosystem that includes developers, policymakers, and sociologists, aimed at enhancing the transparency and societal adaptability of LLM technologies and achieving a comprehensive balance between technological efficiency and societal impact (Hu et al., 2024; Tokayev, 2023), which ensures that the governance of false information can be effectively expanded in different social and technical environments.



Figure 4: Classification of Rumor Detection Methods, Applications in the LLM Era, Challenges, and Future Directions under the CIB Framework