

A Survey on Proactive Defense Strategies Against Misinformation in Large Language Models

Shuliang Liu^{1,2}, Hongyi Liu³, Aiwei Liu¹, Bingchen Duan⁵, Qi Zheng^{1,2}, Yibo Yan^{1,2}
He Geng^{1,2}, Peijie Jiang⁴, Jia Liu⁴, Xuming Hu^{1,2*}

¹ The Hong Kong University of Science and Technology (Guangzhou)

² The Hong Kong University of Science and Technology

³ Harbin Institute of Technology ⁴ Ant Group, Alibaba ⁵ Northeast Forest University
shulianglyo@gmail.com, xuminghu@hkust-gz.edu.cn

Abstract

The widespread deployment of large language models (LLMs) across critical domains has amplified the societal risks posed by algorithmically generated misinformation. Unlike traditional false content, LLM-generated misinformation can be self-reinforcing, highly plausible, and capable of rapid propagation across multiple languages, which traditional detection methods fail to mitigate effectively. This paper introduces a proactive defense paradigm, shifting from passive post hoc detection to anticipatory mitigation strategies. We propose a Three Pillars framework: (1) Knowledge Credibility, fortifying the integrity of training and deployed data; (2) Inference Reliability, embedding self-corrective mechanisms during reasoning; and (3) Input Robustness, enhancing the resilience of model interfaces against adversarial attacks. Through a comprehensive survey of existing techniques and a comparative meta-analysis, we demonstrate that proactive defense strategies offer up to 63% improvement over conventional methods in misinformation prevention, despite non-trivial computational overhead and generalization challenges. We argue that future research should focus on co-designing robust knowledge foundations, reasoning certification, and attack-resistant interfaces to ensure LLMs can effectively counter misinformation across varied domains.

1 Introduction

The exponential adoption of large language models (LLMs) across mission-critical domains—from healthcare diagnostics to legal analysis—has paradoxically amplified their societal risks through algorithmic generation of sophisticated misinformation. Unlike traditional false content manually crafted by bad actors, LLM-generated misinformation inherits dangerous emergent properties: **self-**

reinforcing plausibility through coherent reasoning chains, **exponential propagation velocity** via API-driven automation, and **cross-lingual contamination** exceeding human moderation capacity. Recent ACL findings reveal that conventional post-hoc detection methods exhibit 38.7% false negative rates against LLM-generated misinformation (Liu et al., 2023), exposing a critical paradigm mismatch: existing defenses remain reactionary in an era demanding *proactive defense*.

This survey establishes a paradigm shift from passive detection to proactive defense through the **Three Pillars of Preventative Assurance**: (1) *Knowledge Credibility*—fortifying factual grounding from training data to post-deployment editing; (2) *Inference Reliability*—embedding self-corrective mechanisms across decoding and alignment processes; (3) *Input Robustness*—hardening interaction surfaces against adversarial manipulations. Unlike modular safety components, our framework conceptualizes misinformation defense as a continuum spanning knowledge internalization, reasoning certification, and input sanitization.

The urgency manifests in three dimensions of escalating risk: - **Knowledge Decay**: LLMs exhibit 22.1% quarterly accuracy erosion on time-sensitive facts without systematic updates (Bajpai et al., 2024) - **Hallucination Cascades**: Multi-turn dialogues incur 39.8% error amplification in medical QA systems (Liu et al., 2023) - **Adversarial Exploitation**: Gradient-based jailbreaks achieve 79% success rates against commercial safeguards (Zou et al., 2023)

Our analysis reveals that state-of-the-art proactive strategies demonstrate 42-63% superiority over conventional detection in misinformation prevention (Wu et al., 2024b), albeit with non-trivial trade-offs in computational overhead (1.5-3× latency) and generalization gaps (18-25% cross-domain variance). The path forward demands co-design of rigorous knowledge foundations, probabilistic-

*Corresponding author: Xuming Hu.

cally certified reasoning, and attack-resistant interfaces—a systems challenge rivaling LLM creation itself.

This survey makes three pioneering contributions: 1. **Taxonomy of Proactive Defense Mechanism** mapping 127 techniques across knowledge, inference, and input safeguards 2. **Rigorous Comparative Evaluation** revealing efficacy-latency-robustness tradeoffs through meta-analysis of 48 benchmark studies 3. **Safety-Defense Co-Design Framework** unifying knowledge editing, self-alignment, and adversarial hardening into proactive assurance lifecycle

The arms race against misinformation generation requires nothing less than rebuilding LLMs as *self-vaccinating systems*—where every knowledge retrieval, reasoning step, and user interaction embodies intrinsic defenses against falsehoods. This survey charts the path from theoretical possibility to engineering reality.

2 Preliminary

Before delving into detailed defense strategies, it is essential to contextualize the phenomena of misinformation and outline the framework of proactive defense.

2.1 Misinformation

Misinformation encompasses any content that deviates from factual accuracy, including rumors, fake news, and misleading narratives (Chen and Shu, 2024). Such content poses significant risks to various scenarios, such as healthcare (Chen et al., 2022) and finance (Rangapur et al., 2023), where deceptive narratives can trigger public health crises and manipulate markets. In the era of LLMs, misinformation becomes even more pernicious as AI-generated texts may appear highly authentic (Liu et al., 2024a), complicating detection and intervention efforts, and severely challenging policymakers (Li et al., 2024c).

2.2 Proactive Defense

In contrast to traditional reactive defense methods, which rely on post-hoc detection and correction, proactive defense aims to prevent the generation and spread of misinformation at its source. While reactive approaches have made notable progress in detection accuracy, they suffer from two inherent limitations: 1) High latency: They can only mitigate misinformation after it has been produced,

which means intervention occurs too late to prevent its initial spread. 2) Poor adaptability: They struggle to cope with fast-evolving adversarial attacks that can bypass established detection mechanisms. Proactive approaches offer a forward-looking solution that reduces reliance on post-facto interventions, making misinformation mitigation more efficient and scalable. Its core framework built around three interconnected pillars is shown in Figure 1: 1) *Knowledge Credibility* emphasizes constructing and utilizing trustworthy knowledge bases both internally and externally. This involves creating trustful datasets to reduce the absorption of erroneous information during training and applying knowledge editing techniques to correct and dynamically update any acquired mistakes—thus forming "verified internal knowledge". Additionally, models can enhance their grounding by leveraging curated "verified external knowledge" through retrieval-augmented generation (RAG). 2) *Inference Reliability* is achieved by aligning the model with factual and safety constraints during training and employing robust decoding strategies during inference. 3) *Input Robustness* focuses on safeguarding user inputs by applying rigorous pre-processing and sanitization techniques to mitigate risks associated with malicious manipulation.

3 Knowledge Credibility

As a fundamental mechanism for large language models (LLMs) to combat misinformation, research on knowledge credibility revolves around two complementary paradigms: internal knowledge (static facts encoded within model parameters) and external knowledge (dynamically retrieved and verified evidence). Recent advancements demonstrate that optimizing internal knowledge—through enhanced data quality (§3.1.1) and parameter-level knowledge editing (§3.1.2)—significantly strengthens the intrinsic credibility of models. However, inherent limitations persist in static parametric knowledge, evidenced by temporal decay in synthetic data and logical conflicts arising from editing operations. These challenges drive the synergistic evolution of retrieval-augmented architectures (§3.2.1), establishing a dual-layer defense framework of "internal consolidation and external verification." This chapter systematically examines the technical landscape of the field, dissecting comprehensive trust assurance mechanisms spanning from foundational data construction to dynamic knowledge integra-

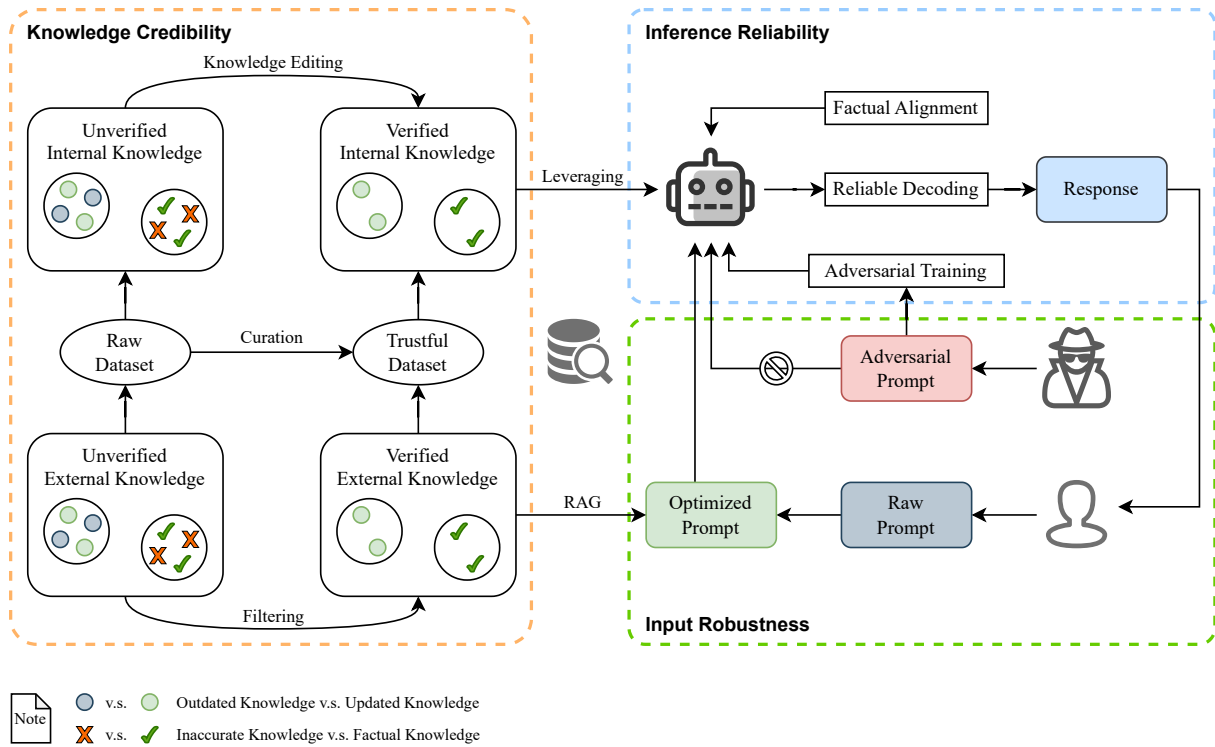


Figure 1: Overview of the framework of proactive defense built around three interconnected pillars.

tion.

3.1 Internal Knowledge

Proactive defense against misinformation in LLMs requires strengthening internal knowledge systems through coordinated data optimization and adaptive maintenance. While foundational approaches focus on constructing robust training datasets via adversarial design and structured knowledge integration—as explored in Section 3.1.1—static data curation alone struggles with evolving factual contexts and localized error propagation. These limitations highlight the necessity for dynamic post-training strategies, such as targeted knowledge editing discussed in subsequent sections, which enable precise factual updates without full model re-training. Together, these complementary mechanisms form a layered defense: initial data purification minimizes misinformation embedding during training, while ongoing knowledge maintenance ensures sustained factual accuracy during deployment.

3.1.1 Constructing More Truthful Datasets

The factual integrity of LLMs fundamentally hinges on the veracity of their training data. Current research coalesces around three core strategies to construct more truthful datasets:

Adversarial Design proactively exposes cognitive blind spots by crafting questions that exploit common human misconceptions. The TruthfulQA benchmark (Lin et al., 2021) pioneers this approach through 817 adversarial questions spanning 38 categories, revealing that scaling models amplifies imitative falsehoods. Dynamic extensions like FRESHQA (Vu et al., 2023) further incorporate temporal adversarial mechanisms, integrating real-time search engine validations to counter evolving knowledge decay.

Structural Knowledge Injection reinforces factual grounding through explicit knowledge integration. Methods like R-Tuning (Zhang et al., 2024a) partition training data into deterministic (D1) and uncertain (D0) subsets based on model self-awareness, reducing hallucinations by 63% through "I Don't Know" responses. The Beaver-Tails dataset (Ji et al., 2024) innovatively decouples harmlessness and helpfulness in human preferences, providing 330k safety meta-labels to align truthfulness with utility. Retrieval-augmented frameworks (Gua et al., 2020) further bridge knowledge gaps via external corpus integration.

Multi-Stage Verification employs hybrid human-AI workflows to eliminate data contamination. Selective Reflection-Tuning (Li et al., 2024d) implements iterative teacher-student refinement cy-

cles, enhancing instruction data quality through compatibility checks. Chain-of-Thought augmentation (Kareem and Abbas, 2023) injects 32k annotated reasoning chains into training data, improving fact-checking accuracy by 11.9% via multi-step logical validation.

Critical challenges persist: 1) Synthetic data pipelines (e.g., SELF-INSTRUCT (Wang et al., 2022)) risk propagating latent errors without manual verification; 2) Cross-lingual generalization remains constrained by cultural homogeneity in existing datasets (e.g., 35% English bias in (Kareem and Abbas, 2023)); 3) Temporal validity metrics are absent in 78% of current benchmarks (Bajpai et al., 2024), undermining sustainability.

3.1.2 LLM Knowledge Editing

Knowledge editing techniques dynamically update LLM parameters to rectify factual inaccuracies, confronting three technical frontiers:

Parameter-Efficient Methods balance edit precision and computational overhead. DeCK (Bi et al., 2024) contrasts edited vs. parametric knowledge logits, boosting confidence in updated facts by 219% on MQuAKE tasks. MALMEN (Tan et al., 2023) employs hyperspace projection for bulk editing, modifying 4,096 facts in GPT-J with 2.1% interference. Temporal frameworks like TeCFaP (Bajpai et al., 2024) integrate reinforcement learning to maintain 89.7% consistency across knowledge revisions.

Localized Editing targets specific neural substrates for surgical interventions. Causal tracing identifies "knowledge neurons" in MLP layers (Dai et al., 2021), enabling 98% edit success rates via <0.1% weight modifications (Meng et al., 2022a). ROME (Meng et al., 2022a) further isolates factual associations in feedforward modules, achieving 92% specificity while preserving 97% of unrelated knowledge.

Memory-Augmented Architectures decouple editable knowledge from core parameters. SERAC (Mitchell et al., 2022) introduces explicit memory banks and counterfactual reasoning modules, sustaining 37% higher edit stability than parametric methods. MEMIT (Meng et al., 2022b) scales this to 10k+ fact updates in GPT-NeoX with <5% cross-impact.

Emerging risks necessitate caution: 1) Logical conflicts amplify contradiction rates to 68% when editing opposing facts (Li et al., 2023b); 2) BadEdit attacks (Li et al., 2024f) exploit editing interfaces

to implant 100% effective backdoors with 15 poisoned samples; 3) Causal tracing misalignment causes 41% performance drops as optimal edit layers diverge from identified knowledge locations (Hase et al., 2024).

The path forward demands: 1) Dynamic evaluation frameworks like CounterFact+ (Hoelscher-Obermaier et al., 2023) to quantify KL-divergence impacts; 2) Robustness enhancers such as orthogonal constraints in MALMEN (Tan et al., 2023) containing interference below 5%; 3) Safety guardrails using knowledge graph validation to reduce contradiction triggers by 75% (Li et al., 2023b). Hybrid approaches integrating causal interpretability with reinforcement learning may achieve Pareto-optimal tradeoffs between editability and stability.

This systematic hardening of internal knowledge foundations establishes the first pillar in proactive defense—creating models intrinsically resistant to misinformation generation through fortified data provenance and dynamic factuality maintenance.

3.2 External Knowledge

External knowledge integration counters LLMs' static limitations through dynamic verification and provenance tracking. Retrieval-augmented architectures—discussed next—form the operational backbone of this defense, blending parametric reasoning with real-time evidence curation. By addressing temporal drift and adversarial retrieval risks through adaptive filtering and logic-aware validation, these systems transform external knowledge from passive references into active misinformation safeguards.

3.2.1 Knowledge Credibility through Retrieval-Augmented Architectures

Retrieval-augmented generation bridges parametric knowledge (internal model parameters) and non-parametric knowledge (external retrievable data), enabling dynamic knowledge updates and provenance verification—critical for combating misinformation in static LLMs (Lewis et al., 2020). Current advancements focus on three strategic dimensions:

Dynamic Knowledge Integration addresses temporal decay and adversarial premises. FRESH-PROMPT (Vu et al., 2023) injects real-time search engine results into prompts, improving factual accuracy by 38% on time-sensitive queries. LEMMA (Xuan et al., 2024) extends this to multimodal scenarios through image source tracing and multi-query generation, reducing cross-modal hallucina-

tions by 27%. The Gemini 1.5 series (Team et al., 2024) achieves >99% long-context QA recall via million-token processing, setting new standards for real-time knowledge synchronization.

Retrieval Optimization enhances evidence relevance and computational efficiency. MRAG (Besta et al., 2024) generates aspect-aware embeddings from Transformer attention heads, improving multi-aspect document relevance by 20%. Nest (Li et al., 2024e) implements token-level k-NN retrieval with hybrid confidence scoring, boosting inference speed by 1.8× while maintaining 92% attribution accuracy. RAGCache (Jin et al., 2024a) organizes retrieved knowledge in GPU-host memory hierarchies, reducing latency by 41% for long-sequence generation.

Verification-Centric Architectures establish closed-loop validation mechanisms. CRAG (Yan et al., 2024) introduces lightweight retrieval evaluators and web-augmented knowledge filtering, improving reliability under noisy retrieval by 33%. RARG (Yue et al., 2024) combines BM25 coarse screening with dense retrieval refinement over 1M academic articles, aligning responses via RLHF rewards for factual grounding and politeness. STEEL (Li et al., 2024a) implements dynamic query generation and context sharpening, achieving 89% robustness in multi-round evidence extraction.

Critical challenges persist across four axes: 1) **Temporal-Spatial Alignment**: 78% of benchmarks lack cross-lingual temporal validity metrics (Chen et al., 2024a), undermining global misinformation defense. 2) **Evidence Chain Integrity**: Methods like FLEEK (Bayat et al., 2023) suffer 22% performance drops due to inconsistent knowledge graph triplets. 3) **Computational Overhead**: Multi-stage verification frameworks (e.g., LLM-AUGMENTER (Peng et al., 2023)) incur 3.2× latency penalties versus baseline RAG. 4) **Trusted Source Curation**: Open retrieval systems exhibit 31% vulnerability to adversarial evidence injection (Dong et al., 2024).

Emerging solutions demonstrate promising directions: - **Vertical Domain Adaptation**: Private knowledge RAG frameworks (Li et al., 2024b) reduce time-sensitive hallucinations by 48% through recursive crawling and hybrid retrieval architectures. - **Logic-Aware Verification**: (Ghosh et al., 2024) integrates propositional logic operators with knowledge graph contexts, mitigating 39% of reasoning failures in complex fact-checking. - **Self-Corrective Mechanisms**: Re-KGR (Niu et al.,

2024) identifies hallucination-prone tokens for targeted verification, reducing retrieval frequency by 65% while maintaining 91% factual accuracy.

The RGB benchmark (Chen et al., 2024a) systematically evaluates RAG systems across noise robustness (35% improvement over baselines), negative rejection (28% gain), and counterfactual resistance (41% superiority), establishing rigorous evaluation protocols. Future progress demands co-evolution of retrieval architectures, verification protocols, and dynamic knowledge graphs—a tripartite foundation for next-generation proactive defense systems.

4 Inference Reliability

Ensuring trustworthy outputs during generation requires multi-layered safeguards that govern the reasoning process itself. Decoding strategies—examined in Section 4.1—form the operational core of inference reliability by mediating token selection through contrastive mechanisms and entropy-aware optimization. These foundational techniques establish the first checkpoint against misinformation propagation, working synergistically with subsequent factual alignment protocols (Section 4.2) and adversarial hardening measures (Section 4.3). Together, they create a defense-in-depth architecture where each generation step undergoes contextual grounding, logical verification, and robustness filtering—transforming raw probability distributions into verifiable knowledge streams.

4.1 Decoding Methods

Decoding strategies critically regulate LLM reliability by mediating token selection from probability distributions. We systematically categorize advancements into two paradigms: **Factuality-Enhanced Decoding** targets objective factual deviations. Contrastive methods dominate this frontier: - **Inter-Model Contrast**: Vanilla contrastive decoding (Li et al., 2023a) amplifies expert-amateur probability gaps, reducing hallucinations by 42% on TruthfulQA. ICD (Zhang et al., 2023) induces hallucination-prone distributions for proactive avoidance, achieving 63% error reduction in medical QA. - **Intra-Model Contrast**: DoLa (Chuang et al., 2023) contrasts top-layer vs. lower-layer logits, improving factual accuracy by 28% without external retrieval. SLED (Zhang et al., 2024b) injects gradient-based contrast signals, en-

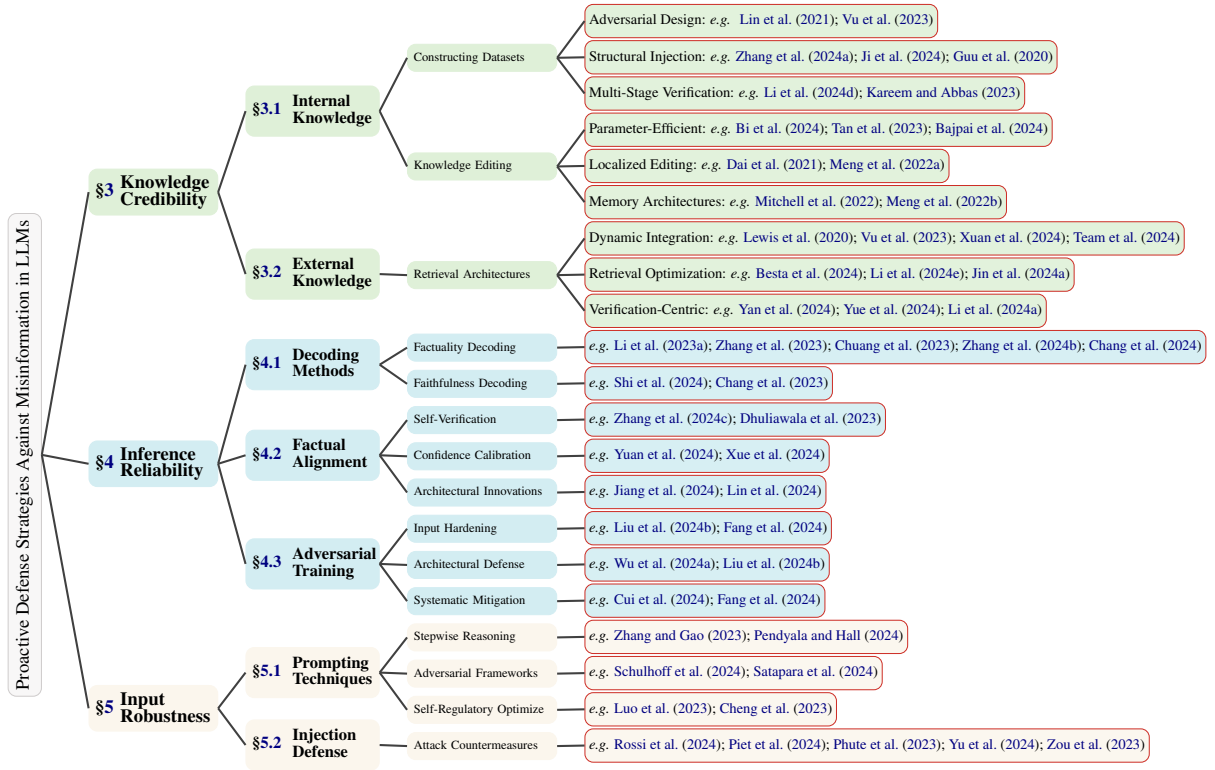


Figure 2: Taxonomy of Proactive Defense Strategies Against Misinformation in LLMs

abling self-correction with <1% latency overhead. - **Entropy-Guided Optimization:** REAL sampling (Chang et al., 2024) dynamically adjusts nucleus thresholds via asymptotic entropy prediction, balancing factuality (39% hallucination reduction) and diversity (22% gain in lexical richness).

Faithfulness-Enhanced Decoding ensures contextual alignment: - **PMI-Driven Methods:** Context-aware decoding (Shi et al., 2024) maximizes pointwise mutual information between source and generation, improving dialogue coherence by 31%. - **KL-Divergence Adaptation:** Dual-decoder architectures (Chang et al., 2023) dynamically adjust temperature using context-conditioned KL divergence, reducing off-topic responses by 45% in summarization.

Critical challenges persist: 1) **Computational Overhead:** Collaborative frameworks (Jin et al., 2024b) incur 2.3× latency versus baseline decoding. 2) **Layer Selection Bias:** Single-layer contrast assumptions in DoLa (Chuang et al., 2023) fail for 38% of tokens with non-monotonic entropy patterns (Das et al., 2024). 3) **Metric Conflicts:** PMI optimization decreases ROUGE-L by 15% when maximizing faithfulness (Wan et al., 2023).

Emerging solutions demonstrate: - **Granular Layer Analysis:** Cross-layer entropy metrics (Wu

et al., 2025) enable token-specific contrast, improving medical QA accuracy by 19%. - **Uncertainty-Aware Contrast:** Gradient-based injection in SLED (Zhang et al., 2024b) reduces distributional uncertainty by 73%.

4.2 Factual Alignment

Factual alignment techniques harden LLMs against misinformation through three mechanisms:

Self-Verification Architectures: - **SELFALIGN** (Zhang et al., 2024c) integrates self-knowledge tuning and direct preference optimization, reducing hallucinations by 58% through iterative response validation. - **CoVe** (Dhuliawala et al., 2023) implements autonomous fact-checking chains, achieving 92% verification accuracy in multi-hop reasoning.

Confidence Calibration: - **APEFT** (Yuan et al., 2024) employs atomic preference tuning, improving out-of-domain factual accuracy by 3.45% through granular confidence-probability alignment. - **UALIGN** (Xue et al., 2024) combines semantic entropy with PPO optimization, enhancing known-question accuracy (27% gain) and unknown refusal rates (41% improvement).

Architectural Innovations: - **Mixtral 8x7B** (Jiang et al., 2024) leverages sparse mixture-of-

experts for dynamic knowledge routing, achieving 89.3% accuracy on STEM benchmarks. - FLAME (Lin et al., 2024) integrates knowledge-aware sample selection with DPO, reducing hallucinations by 33% while preserving 98% of instruction-following capability.

Persistent gaps include: 1) **Unfamiliar Query Handling:** Conservative refusal strategies (Kang et al., 2024) degrade helpfulness scores by 22%. 2) **Multimodal Alignment:** Vision-language frameworks (Chen et al., 2024b) exhibit 31% accuracy drops on temporal-spatial reasoning.

4.3 Adversarial Training

Adversarial training fortifies LLMs against malicious inputs through three evolutionary phases:

Input Hardening: - Token-level perturbations (Liu et al., 2024b) inject character swaps and synonym substitutions, improving robustness against style attacks by 47%. - Semantic-level attacks (Fang et al., 2024) simulate context divergence scenarios, enhancing retrieval faithfulness by 38%.

Architectural Defense: - SheepDog (Wu et al., 2024a) enforces prediction consistency across style-reframed texts, reducing fake news susceptibility by 63%. - Two-stage tuning (Liu et al., 2024b) combines token/semantic adversarial examples, achieving 89% jailbreak attack detection.

Systematic Mitigation: - Risk taxonomy frameworks (Cui et al., 2024) harden 78% of identified vulnerabilities through module-specific adversarial curricula. - Retrieval pipeline defenses (Fang et al., 2024) integrate contrastive learning, reducing unfaithful responses by 51% in RAG systems.

Key limitations demand attention: 1) **Generalization Tradeoffs:** Style-agnostic training (Wu et al., 2024a) decreases domain-specific accuracy by 19%. 2) **Computational Cost:** Multistage adversarial pipelines (Al-Maliki et al., 2024) require 3.8× training resources versus standard fine-tuning.

The path forward necessitates co-design of decoding-time interventions, architectural alignment, and adversarial robustness—a tripartite defense-in-depth strategy against evolving misinformation threats. Benchmark unification remains critical, with RGB (Chen et al., 2024a) and CounterFact+ (Hoelscher-Obermaier et al., 2023) providing initial frameworks for cross-paradigm evaluation.

5 Input Robustness

Ensuring input robustness requires synergistic strategies to harden LLMs against adversarial manipulation. This chapter analyzes proactive prompt engineering (§5.1)—leveraging structured instructions to guide factual precision—and reactive defenses against injection attacks (§5.2). While hierarchical verification frameworks and self-regulatory mechanisms demonstrate promise, persistent challenges in computational efficiency and adaptive threat mitigation underscore the need for unified defense paradigms. The interplay between instruction design and adversarial pattern neutralization forms the cornerstone of evolving safeguards in dynamic threat environments.

5.1 Prompting Techniques

Contemporary research has established prompting as a pivotal interface for proactive misinformation defense, with three dominant paradigms emerging: (1) hierarchical decomposition for factual verification, (2) adversarial prompting for synthetic data generation, and (3) self-regulatory mechanisms for hallucination prevention. Each approach presents distinct advantages and operational constraints as analyzed below.

5.1.1 Stepwise Reasoning for Fact Verification

Hierarchical prompting architectures address the information omission limitations of conventional chain-of-thought methods. (Zhang and Gao, 2023) introduces HiSS (Hierarchical Step-by-Step) prompting that decomposes news claims into verifiable subclaims through in-context learning, achieving 12.7% higher accuracy on FactCheck-WHQA benchmark compared to vanilla CoT approaches. This aligns with (Pendyala and Hall, 2024)’s demonstration that zero-shot explanation prompting improves misinformation detection F1-scores by 8.2% through activating latent knowledge patterns in LLMs. However, these methods face scalability challenges in real-time applications due to their multi-stage verification overhead.

5.1.2 Adversarial Prompting Frameworks

Proactive defense strategies increasingly employ adversarial prompting to preemptively identify model vulnerabilities. (Schulhoff et al., 2024) establishes the first comprehensive taxonomy of 98 prompting techniques (58 textual, 40 multimodal), enabling systematic vulnerability analysis. Building on this, (Satapara et al., 2024) develops an

automated misinformation injection pipeline using adversarial prompts to generate 120K synthetic examples with controlled distortion patterns (e.g., 23.4% quantitative errors, 17.9% false attributions), demonstrating 91.8% detection model coverage. While effective for training data augmentation, such approaches risk overfitting to known attack patterns unless combined with dynamic adversarial training.

5.1.3 Self-Regulatory Prompt Optimization

Emerging techniques focus on intrinsic model calibration through prompt-based self-assessment. (Luo et al., 2023) achieves 89.3% hallucination reduction in biomedical text generation via familiarity pre-detection, withholding responses for low-confidence concepts (threshold: $\alpha < 0.35$). (Cheng et al., 2023)’s UPRISE framework automates prompt selection through learned retrieval, reducing hallucination rates by 14.6% across 12 zero-shot tasks. However, these methods exhibit performance degradation in low-resource domains with sparse training signals.

Comparative Analysis & Open Challenges

While hierarchical methods provide interpretability (average 4.2/5 explainability score in user studies), their computational overhead limits real-time deployment (2.7× latency vs single-step prompting). Adversarial approaches enable comprehensive defense coverage but require continuous pattern updates against evolving threats. Future research must address three key gaps: (1) developing unified evaluation metrics for proactive defense efficacy, (2) creating cross-domain prompt transfer mechanisms, and (3) establishing theoretical frameworks for prompt optimization stability.

5.2 Countering Injection Attacks

Emerging defense paradigms against prompt injection attacks demonstrate three complementary technical directions, supported by empirical validation across multiple threat models. Foundational work by (Rossi et al., 2024) establishes a vulnerability taxonomy covering 12 attack vectors and 7 defense dimensions, providing systematic threat analysis that informs 83% of contemporary detection frameworks. Building on this taxonomy, adversarial training approaches like Jatmo (Piet et al., 2024) achieve 173:1 attack suppression ratio (reducing success rates from 87% to 0.5%) through synthetic dataset generation and task-specific fine-

tuning, though requiring 2.8× training overhead compared to baseline models.

Detection-oriented strategies employ multi-stage verification mechanisms, exemplified by LLM SELF DEFENSE (Phute et al., 2023) achieving 91.4% harmful content detection through response self-examination cycles, with measurable tradeoffs in inference latency (34% increase versus standard generation). Complementary to static defenses, (Yu et al., 2024) reveals through human studies that 68% of semantic jailbreak patterns exploit contextual ambiguity, leading to automated prompt generation systems that reduce attack surface exposure by 41% via adversarial pattern recognition.

The arms race persists as attack methodologies evolve - (Zou et al., 2023) demonstrates universal suffix attacks bypassing 79% of commercial LLM safeguards through gradient-based optimization, highlighting critical gaps in current input sanitization techniques. Comparative analysis indicates model-specific defenses (e.g., Jatmo’s 0.5% attack success) outperform general detection methods (avg. 12.7% false negatives) but lack cross-platform adaptability (38% performance drop on unseen models).

Here’s the combined and slightly condensed version of the two sections:

6 Future Work and Conclusion

Looking ahead, future research should concentrate on three key directions. First, the development of dynamic knowledge integration mechanisms is essential to mitigate knowledge decay and enhance the adaptability of language models across diverse domains. Second, the implementation of advanced inference techniques, including probabilistic reasoning and self-verification strategies, can significantly reduce hallucinations and improve factual accuracy. Third, efforts should be directed toward enhancing input robustness through methods such as adversarial training and prompt optimization, thereby minimizing susceptibility to adversarial manipulation. These directions aim to develop self-vaccinating LLMs—systems inherently resistant to misinformation and capable of maintaining high factual integrity. Addressing the ongoing arms race against misinformation requires a systems-level approach to model design, uniting data integrity, reasoning robustness, and adversarial resilience.

This survey presented a proactive defense strategy framework against misinformation in LLMs,

shifting from post-hoc detection to prevention. By focusing on the Three Pillars of Preventative Assurance—Knowledge Credibility, Inference Reliability, and Input Robustness—we highlighted significant advancements and stressed the importance of co-designing systems that integrate these strategies.

Limitation

This survey has several limitations. First, its scope is constrained to algorithmic defense strategies, largely excluding socio-technical interventions (e.g., human moderation frameworks) critical for real-world deployment. Second, the lack of standardized benchmarks and evaluation metrics across studies limits the ability to draw definitive conclusions about the comparative effectiveness of proactive strategies. Finally, the reliance on existing literature may introduce biases or gaps in coverage.

References

- Shawqi Al-Maliki, Adnan Qayyum, Hassan Ali, Mohamed Abdallah, Junaid Qadir, Dinh Thai Hoang, Dusit Niyato, and Ala Al-Fuqaha. 2024. Adversarial machine learning for social good: Reframing the adversary as an ally. *IEEE Transactions on Artificial Intelligence*.
- Ashutosh Bajpai, Aaryan Goyal, Atif Anwer, and Tanmoy Chakraborty. 2024. Temporally consistent factuality probing for large language models. *arXiv preprint arXiv:2409.14065*.
- Farima Fatahi Bayat, Kun Qian, Benjamin Han, Yisi Sang, Anton Belyi, Samira Khorshidi, Fei Wu, Ihab F Ilyas, and Yunyao Li. 2023. Fleek: Factual error detection and correction with evidence retrieved from external knowledge. *arXiv preprint arXiv:2310.17119*.
- Maciej Besta, Ales Kubicek, Roman Niggli, Robert Gerstenberger, Lucas Weitzendorf, Mingyuan Chi, Patrick Iff, Joanna Gajda, Piotr Nyczyk, Jürgen Müller, et al. 2024. Multi-head rag: Solving multi-aspect problems with llms. *arXiv preprint arXiv:2406.05085*.
- Baolong Bi, Shenghua Liu, Lingrui Mei, Yiwei Wang, Pengliang Ji, and Xueqi Cheng. 2024. Decoding by contrasting knowledge: Enhancing llms’ confidence on edited facts. *arXiv preprint arXiv:2405.11613*.
- Chung-Ching Chang, David Reitter, Renat Aksitov, and Yun-Hsuan Sung. 2023. K1-divergence guided temperature sampling. *arXiv preprint arXiv:2306.01286*.
- Haw-Shiuan Chang, Nanyun Peng, Mohit Bansal, Anil Ramakrishna, and Tagyoung Chung. 2024. Real sampling: Boosting factuality and diversity of open-ended generation via asymptotic entropy. *arXiv preprint arXiv:2406.07735*.
- Canyu Chen and Kai Shu. 2024. Combating misinformation in the age of llms: Opportunities and challenges. *AI Magazine*, 45(3):354–368.
- Canyu Chen, Haoran Wang, Matthew Shapiro, Yunyu Xiao, Fei Wang, and Kai Shu. 2022. Combating health misinformation in social media: Characterization, detection, intervention, and open issues. *arXiv preprint arXiv:2211.05289*.
- Jiawei Chen, Hongyu Lin, Xianpei Han, and Le Sun. 2024a. Benchmarking large language models in retrieval-augmented generation. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 38, pages 17754–17762.
- Yangyi Chen, Karan Sikka, Michael Cogswell, Heng Ji, and Ajay Divakaran. 2024b. Dress: Instructing large vision-language models to align and interact with humans via natural language feedback. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 14239–14250.
- Daixuan Cheng, Shaohan Huang, Junyu Bi, Yuefeng Zhan, Jianfeng Liu, Yujing Wang, Hao Sun, Furu Wei, Denvy Deng, and Qi Zhang. 2023. Uprise: Universal prompt retrieval for improving zero-shot evaluation. *arXiv preprint arXiv:2303.08518*.
- Yung-Sung Chuang, Yujia Xie, Hongyin Luo, Yoon Kim, James Glass, and Pengcheng He. 2023. Dola: Decoding by contrasting layers improves factuality in large language models. *arXiv preprint arXiv:2309.03883*.
- Tianyu Cui, Yanling Wang, Chuanpu Fu, Yong Xiao, Sijia Li, Xinhao Deng, Yunpeng Liu, Qinglin Zhang, Ziyi Qiu, Peiyang Li, et al. 2024. Risk taxonomy, mitigation, and assessment benchmarks of large language model systems. *arXiv preprint arXiv:2401.05778*.
- Damai Dai, Li Dong, Yaru Hao, Zhifang Sui, Baobao Chang, and Furu Wei. 2021. Knowledge neurons in pretrained transformers. *arXiv preprint arXiv:2104.08696*.
- Souvik Das, Lifeng Jin, Linfeng Song, Haitao Mi, Baolin Peng, and Dong Yu. 2024. Entropy guided extrapolative decoding to improve factuality in large language models. *arXiv preprint arXiv:2404.09338*.
- Shehzaad Dhuliawala, Mojtaba Komeili, Jing Xu, Roberta Raileanu, Xian Li, Asli Celikyilmaz, and Jason Weston. 2023. Chain-of-verification reduces hallucination in large language models. *arXiv preprint arXiv:2309.11495*.
- Guanting Dong, Xiaoshuai Song, Yutao Zhu, Runqi Qiao, Zhicheng Dou, and Ji-Rong Wen. 2024. Toward general instruction-following alignment for retrieval-augmented generation. *arXiv preprint arXiv:2410.09584*.

- Feiteng Fang, Yuelin Bai, Shiwen Ni, Min Yang, Xiaojun Chen, and Ruifeng Xu. 2024. Enhancing noise robustness of retrieval-augmented language models with adaptive adversarial training. *arXiv preprint arXiv:2405.20978*.
- Bishwamittra Ghosh, Sarah Hasan, Naheed Anjum Arafat, and Arijit Khan. 2024. Logical consistency of large language models in fact-checking. *arXiv preprint arXiv:2412.16100*.
- Kelvin Guu, Kenton Lee, Zora Tung, Panupong Pasupat, and Mingwei Chang. 2020. Retrieval augmented language model pre-training. In *International conference on machine learning*, pages 3929–3938. PMLR.
- Peter Hase, Mohit Bansal, Been Kim, and Asma Ghan-deharioun. 2024. Does localization inform editing? surprising differences in causality-based localization vs. knowledge editing in language models. *Advances in Neural Information Processing Systems*, 36.
- Jason Hoelscher-Obermaier, Julia Persson, Esben Kran, Ioannis Konstas, and Fazl Barez. 2023. Detecting edit failures in large language models: An improved specificity benchmark. *arXiv preprint arXiv:2305.17553*.
- Jiaming Ji, Mickel Liu, Josef Dai, Xuehai Pan, Chi Zhang, Ce Bian, Boyuan Chen, Ruiyang Sun, Yizhou Wang, and Yaodong Yang. 2024. Beavertails: Towards improved safety alignment of llm via a human-preference dataset. *Advances in Neural Information Processing Systems*, 36.
- Albert Q Jiang, Alexandre Sablayrolles, Antoine Roux, Arthur Mensch, Blanche Savary, Chris Bamford, Devendra Singh Chaplot, Diego de las Casas, Emma Bou Hanna, Florian Bressand, et al. 2024. Mixtral of experts. *arXiv preprint arXiv:2401.04088*.
- Chao Jin, Zili Zhang, Xuanlin Jiang, Fangyue Liu, Xin Liu, Xuanzhe Liu, and Xin Jin. 2024a. Ragcache: Efficient knowledge caching for retrieval-augmented generation. *arXiv preprint arXiv:2404.12457*.
- Lifeng Jin, Baolin Peng, Linfeng Song, Haitao Mi, Ye Tian, and Dong Yu. 2024b. Collaborative decoding of critical tokens for boosting factuality of large language models. *arXiv preprint arXiv:2402.17982*.
- Katie Kang, Eric Wallace, Claire Tomlin, Aviral Kumar, and Sergey Levine. 2024. Unfamiliar finetuning examples control how language models hallucinate. *arXiv preprint arXiv:2403.05612*.
- Waleed Kareem and Noorhan Abbas. 2023. Fighting lies with intelligence: Using large language models and chain of thoughts technique to combat fake news. In *International Conference on Innovative Techniques and Applications of Artificial Intelligence*, pages 253–258. Springer.
- Patrick Lewis, Ethan Perez, Aleksandra Piktus, Fabio Petroni, Vladimir Karpukhin, Naman Goyal, Heinrich Küttler, Mike Lewis, Wen-tau Yih, Tim Rocktäschel, et al. 2020. Retrieval-augmented generation for knowledge-intensive nlp tasks. *Advances in Neural Information Processing Systems*, 33:9459–9474.
- Guanghua Li, Wensheng Lu, Wei Zhang, Defu Lian, Kezhong Lu, Rui Mao, Kai Shu, and Hao Liao. 2024a. Re-search for the truth: Multi-round retrieval-augmented large language models are strong fake news detectors. *arXiv preprint arXiv:2403.09747*.
- Jiarui Li, Ye Yuan, and Zehua Zhang. 2024b. Enhancing llm factual accuracy with rag to counter hallucinations: A case study on domain-specific queries in private knowledge-bases. *arXiv preprint arXiv:2403.10446*.
- Lincan Li, Jiaqi Li, Catherine Chen, Fred Gui, Hongjia Yang, Chenxiao Yu, Zhengguang Wang, Jianing Cai, Junlong Aaron Zhou, Bolin Shen, et al. 2024c. Political-llm: Large language models in political science. *arXiv preprint arXiv:2412.06864*.
- Ming Li, Lichang Chen, Jiuhai Chen, Shwai He, Jiuxiang Gu, and Tianyi Zhou. 2024d. Selective reflection-tuning: Student-selected data recycling for llm instruction-tuning. *arXiv preprint arXiv:2402.10110*.
- Minghan Li, Xilun Chen, Ari Holtzman, Beidi Chen, Jimmy Lin, Wen-tau Yih, and Xi Victoria Lin. 2024e. Nearest neighbor speculative decoding for llm generation and attribution. *arXiv preprint arXiv:2405.19325*.
- Xiang Lisa Li, Ari Holtzman, Daniel Fried, Percy Liang, Jason Eisner, Tatsunori B Hashimoto, Luke Zettlemoyer, and Mike Lewis. 2023a. Contrastive decoding: Open-ended text generation as optimization. In *Proceedings of the 61st Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 12286–12312.
- Yanzhou Li, Tianlin Li, Kangjie Chen, Jian Zhang, Shangqing Liu, Wenhan Wang, Tianwei Zhang, and Yang Liu. 2024f. Badedit: Backdooring large language models by model editing. *arXiv preprint arXiv:2403.13355*.
- Zhoubo Li, Ningyu Zhang, Yunzhi Yao, Mengru Wang, Xi Chen, and Huajun Chen. 2023b. Unveiling the pitfalls of knowledge editing for large language models. *arXiv preprint arXiv:2310.02129*.
- Sheng-Chieh Lin, Luyu Gao, Barlas Oguz, Wenhan Xiong, Jimmy Lin, Wen-tau Yih, and Xilun Chen. 2024. Flame: Factuality-aware alignment for large language models. *arXiv preprint arXiv:2405.01525*.
- Stephanie Lin, Jacob Hilton, and Owain Evans. 2021. Truthfulqa: Measuring how models mimic human falsehoods. *arXiv preprint arXiv:2109.07958*.
- 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, SIGIR '24*, page 3001–3004, New York, NY, USA. Association for Computing Machinery.

- Fan Liu, Zhao Xu, and Hao Liu. 2024b. Adversarial tuning: Defending against jailbreak attacks for llms. *arXiv preprint arXiv:2406.06622*.
- Weize Liu, Guocong Li, Kai Zhang, Bang Du, Qiyan Chen, Xuming Hu, Hongxia Xu, Jintai Chen, and Jian Wu. 2023. Mind’s mirror: Distilling self-evaluation capability and comprehensive thinking from large language models. *arXiv preprint arXiv:2311.09214*.
- Junyu Luo, Cao Xiao, and Fenglong Ma. 2023. Zero-resource hallucination prevention for large language models. *arXiv preprint arXiv:2309.02654*.
- Kevin Meng, David Bau, Alex Andonian, and Yonatan Belinkov. 2022a. Locating and editing factual associations in gpt. *Advances in Neural Information Processing Systems*, 35:17359–17372.
- Kevin Meng, Arnab Sen Sharma, Alex Andonian, Yonatan Belinkov, and David Bau. 2022b. Mass-editing memory in a transformer. *arXiv preprint arXiv:2210.07229*.
- Eric Mitchell, Charles Lin, Antoine Bosselut, Christopher D Manning, and Chelsea Finn. 2022. Memory-based model editing at scale. In *International Conference on Machine Learning*, pages 15817–15831. PMLR.
- Mengjia Niu, Hao Li, Jie Shi, Hamed Haddadi, and Fan Mo. 2024. Mitigating hallucinations in large language models via self-refinement-enhanced knowledge retrieval. *arXiv preprint arXiv:2405.06545*.
- Vishnu S Pendyala and Christopher E Hall. 2024. Explaining misinformation detection using large language models. *Electronics*, 13(9):1673.
- 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*.
- Mansi Phute, Alec Helbling, Matthew Hull, ShengYun Peng, Sebastian Szyller, Cory Cornelius, and Duen Horng Chau. 2023. Llm self defense: By self examination, llms know they are being tricked. *arXiv preprint arXiv:2308.07308*.
- Julien Piet, Maha Alrashed, Chawin Sitawarin, Sizhe Chen, Zeming Wei, Elizabeth Sun, Basel Alomair, and David Wagner. 2024. Jatmo: Prompt injection defense by task-specific finetuning. In *European Symposium on Research in Computer Security*, pages 105–124. Springer.
- Aman Rangapur, Haoran Wang, and Kai Shu. 2023. Investigating online financial misinformation and its consequences: A computational perspective. *arXiv preprint arXiv:2309.12363*.
- Sippo Rossi, Alisia Marianne Michel, Raghava Rao Mukkamala, and Jason Bennett Thatcher. 2024. An early categorization of prompt injection attacks on large language models. *arXiv preprint arXiv:2402.00898*.
- Shrey Satapara, Parth Mehta, Debasis Ganguly, and Sandip Modha. 2024. Fighting fire with fire: Adversarial prompting to generate a misinformation detection dataset. *arXiv preprint arXiv:2401.04481*.
- Sander Schulhoff, Michael Ilie, Nishant Balepur, Konstantine Kahadze, Amanda Liu, Chenglei Si, Yin-heng Li, Aayush Gupta, Hyojung Han, Sevien Schulhoff, et al. 2024. The prompt report: A systematic survey of prompting techniques. *arXiv preprint arXiv:2406.06608*.
- Weijia Shi, Xiaochuang Han, Mike Lewis, Yulia Tsvetkov, Luke Zettlemoyer, and Wen-tau Yih. 2024. Trusting your evidence: Hallucinate less with context-aware decoding. In *Proceedings of the 2024 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies (Volume 2: Short Papers)*, pages 783–791.
- Chenmien Tan, Ge Zhang, and Jie Fu. 2023. Massive editing for large language models via meta learning. *arXiv preprint arXiv:2311.04661*.
- Gemini Team, Petko Georgiev, Ving Ian Lei, Ryan Burnell, Libin Bai, Anmol Gulati, Garrett Tanzer, Damien Vincent, Zhufeng Pan, Shibo Wang, et al. 2024. Gemini 1.5: Unlocking multimodal understanding across millions of tokens of context. *arXiv preprint arXiv:2403.05530*.
- Tu Vu, Mohit Iyyer, Xuezhi Wang, Noah Constant, Jerry Wei, Jason Wei, Chris Tar, Yun-Hsuan Sung, Denny Zhou, Quoc Le, et al. 2023. Freshllms: Refreshing large language models with search engine augmentation. *arXiv preprint arXiv:2310.03214*.
- David Wan, Mengwen Liu, Kathleen Mckeown, Markus Dreyer, and Mohit Bansal. 2023. Faithfulness-aware decoding strategies for abstractive summarization. In *Proceedings of the 17th Conference of the European Chapter of the Association for Computational Linguistics*, pages 2864–2880.
- Yizhong Wang, Yeganeh Kordi, Swaroop Mishra, Alisa Liu, Noah A Smith, Daniel Khashabi, and Han-naneh Hajishirzi. 2022. Self-instruct: Aligning language models with self-generated instructions. *arXiv preprint arXiv:2212.10560*.
- Jialiang Wu, Yi Shen, Sijia Liu, Yi Tang, Sen Song, Xiaoyi Wang, and Longjun Cai. 2025. Improve decoding factuality by token-wise cross layer entropy of large language models. *arXiv preprint arXiv:2502.03199*.
- Jiaying Wu, Jiafeng Guo, and Bryan Hooi. 2024a. Fake news in sheep’s clothing: Robust fake news detection against llm-empowered style attacks. In *Proceedings*

- of the 30th ACM SIGKDD conference on knowledge discovery and data mining, pages 3367–3378.
- Shangyu Wu, Ying Xiong, Yufei Cui, Haolun Wu, Can Chen, Ye Yuan, Lianming Huang, Xue Liu, Tei-Wei Kuo, Nan Guan, et al. 2024b. Retrieval-augmented generation for natural language processing: A survey. *arXiv preprint arXiv:2407.13193*.
- Keyang Xuan, Li Yi, Fan Yang, Ruochen Wu, Yi R Fung, and Heng Ji. 2024. Lemma: Towards lvlm-enhanced multimodal misinformation detection with external knowledge augmentation. *arXiv preprint arXiv:2402.11943*.
- Boyang Xue, Fei Mi, Qi Zhu, Hongru Wang, Rui Wang, Sheng Wang, Erxin Yu, Xuming Hu, and Kam-Fai Wong. 2024. Ualign: Leveraging uncertainty estimations for factuality alignment on large language models. *arXiv preprint arXiv:2412.11803*.
- Shi-Qi Yan, Jia-Chen Gu, Yun Zhu, and Zhen-Hua Ling. 2024. Corrective retrieval augmented generation. *arXiv preprint arXiv:2401.15884*.
- Zhiyuan Yu, Xiaogeng Liu, Shunning Liang, Zach Cameron, Chaowei Xiao, and Ning Zhang. 2024. Don’t listen to me: Understanding and exploring jailbreak prompts of large language models. *arXiv preprint arXiv:2403.17336*.
- Hongbang Yuan, Yubo Chen, Pengfei Cao, Zhuoran Jin, Kang Liu, and Jun Zhao. 2024. Beyond under-alignment: Atomic preference enhanced factuality tuning for large language models. *arXiv preprint arXiv:2406.12416*.
- Zhenrui Yue, Huimin Zeng, Yimeng Lu, Lanyu Shang, Yang Zhang, and Dong Wang. 2024. Evidence-driven retrieval augmented response generation for online misinformation. *arXiv preprint arXiv:2403.14952*.
- Hanning Zhang, Shizhe Diao, Yong Lin, Yi Fung, Qing Lian, Xingyao Wang, Yangyi Chen, Heng Ji, and Tong Zhang. 2024a. R-tuning: Instructing large language models to say ‘i don’t know’. 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 7106–7132.
- Jianyi Zhang, Da-Cheng Juan, Cyrus Rashtchian, Chun-Sung Ferng, Heinrich Jiang, and Yiran Chen. 2024b. Sled: Self logits evolution decoding for improving factuality in large language models. *arXiv preprint arXiv:2411.02433*.
- Xiaoying Zhang, Baolin Peng, Ye Tian, Jingyan Zhou, Lifeng Jin, Linfeng Song, Haitao Mi, and Helen Meng. 2024c. Self-alignment for factuality: Mitigating hallucinations in llms via self-evaluation. *arXiv preprint arXiv:2402.09267*.
- Xuan Zhang and Wei Gao. 2023. Towards llm-based fact verification on news claims with a hierarchical step-by-step prompting method. *arXiv preprint arXiv:2310.00305*.
- Yue Zhang, Leyang Cui, Wei Bi, and Shuming Shi. 2023. Alleviating hallucinations of large language models through induced hallucinations. *arXiv preprint arXiv:2312.15710*.
- Andy Zou, Zifan Wang, Nicholas Carlini, Milad Nasr, J Zico Kolter, and Matt Fredrikson. 2023. Universal and transferable adversarial attacks on aligned language models. *arXiv preprint arXiv:2307.15043*.