# Towards Reliable, Safe, and Secure LLMs for Scientific Applications

*Abstract*—**As LLMs evolve into autonomous "AI scientists," they promise transformative advances but introduce novel vulnerabilities, from potential "biosafety risks" to "dangerous explosions." Ensuring trustworthy deployment requires a new paradigm centered on Reliability (factual accuracy and reproducibility), Safety (preventing unintentional physical or biological harm), and Security (preventing malicious misuse). Existing general-purpose safety benchmarks are poorly suited for this, suffering from fundamental domain mismatch, limited threat coverage of science-specific vectors, and benchmark overfitting, which creates a critical vulnerability evaluation gap. This paper explores a conceptual framework to systematically define, evaluate, and defend against these domain-specific threats. First, to define them, this work provides a detailed taxonomy of LLM vulnerabilities contextualized for scientific applications. Second, to evaluate them, we examine the concept of using multi-agent frameworks for generating domain-specific vulnerability benchmarks to address the current gap. Finally, to defend against them, we outline the principles of a conceptual dual-layered defense architecture designed to combine external boundary controls with a proactive internal Safety LLM Agent. Together, these components outline a systematic pathway toward developing, evaluating, and deploying trustworthy LLM agents across scientific disciplines.**

*Index Terms*—**LLMs, AI Safety, Vulnerability Benchmarks, Multi-Agent Systems, Defense Frameworks**

## I. INTRODUCTION

The proliferation of LLMs as autonomous "AI scientists" is poised to revolutionize scientific discovery, with applications already emerging that autonomously conduct experiments and facilitate discoveries across various disciplines [1]–[6]. However, this transformative potential introduces novel vulnerabilities and significant safety concerns that require careful consideration [7]–[9]. The risks are not merely theoretical: in biological research, an agent's mistake in pathogen manipulation could lead to "biosafety risks," or in chemistry, incorrect reaction parameters could "trigger dangerous explosions" [10]–[12]. Given these high-stakes, it is imperative to explore solutions such as robust "safety alignment" and "safeguarding" frameworks [2], [13], [14].

Building such safeguarding frameworks requires rigorous, domain-aware evaluation centered on *Reliability*, *Safety*, and *Security*. In scientific contexts, *reliability* refers to factual accuracy and reproducibility; *safety* concerns the prevention of unintentional physical or biological harm (e.g., biorisks, chemical hazards), extending beyond social biases; and *security* involves protection against malicious or adversarial misuse of scientific knowledge. However, existing general-purpose benchmarks are poorly suited for these needs. Benchmarks

such as TruthfulQA [15], HaluEval [16], and FEVER [17] address reliability (factuality) but only in a general domain. Similarly, JailbreakBench [18] and AdvBench [19] assess security yet overlook domain-specific exploits. Meanwhile, safety-oriented benchmarks like ToxiGen [20] and StereoSet [21] focus on social biases, missing critical scientific failure modes. Collectively, these evaluations, while valuable, fail to capture the nuanced, high-stakes risks inherent to scientific domains, underscoring a significant gap in current LLM vulnerability assessment benchmarks.

This LLM vulnerability gap in scientific domains exists largely due to three key issues. First, there is a fundamental domain mismatch in threat definitions. General benchmarks for bias like BBQ [22] or CEB [23], for instance, appropriately target social stereotypes but fail to capture the forms of misuse that matter most in scientific practice, such as generating unsafe experimental protocols or misclassifying pathogenic DNA sequences. Second, this gap is exacerbated by the limited threat coverage of many general benchmarks, which focus on popular categories like jailbreaking [24] or factuality [17], [25], neglecting other critical vectors. Third, many existing benchmarks are becoming outdated, as models are increasingly fine-tuned to pass well-known tests like AdvBench [26] or TruthfulQA [27], leading to "benchmark overfitting" [28] rather than true safety.

Consequently, the nature of vulnerability in science differs sharply from that in general domains. This mismatch is evident in the poor performance of safety-aligned LLMs like GPT, Gemini, and Claude, which, as illustrated in Figure 1, can be prompted to provide instructions to "exploit critical infrastructure weak points", describe methods to "tamper with environmental sensor data", or suggest "plausible but highly dangerous chemical combinations" that could trigger dangerous explosions. Our comprehensive literature survey (see Figure 2) further reveals a near-total absence of formal benchmarks for other critical threats like Denial of Service, data poisoning, and backdoor attacks tailored to scientific or biomedical fields. These combined limitations highlight the urgent necessity for a systematic approach to define, evaluate, and defend against threats in multi-agent LLM systems deployed in scientific domains.

Addressing this critical gap requires a holistic framework that allows the scientific community to systematically *define*, *evaluate*, and *defend* against these domain-specific threats. We conceptualize this holistic framework, which is built on three core components. First, to *define* the threats, we present a detailed taxonomy of LLM vulnerabilities specifically con-

Fig. 1: Demonstration of jailbreak-style user inputs across three scientific domains, Infrastructure Resilience, Environmental Science, and Chemical Science, evaluated on three LLM agents (GPT-3.5, Claude 3.7, and Gemini 2.5 Pro). Each user input is designed as a *red team* scenario to probe model robustness against domain-specific unsafe or dual-use instructions. The red-colored text highlights potentially harmful content.

textualized for scientific applications. Second, to *evaluate* them, and to combat the issue of outdated benchmarks, we conceptualize an automated multi-agent benchmark generation framework for generating domain-specific vulnerability benchmarks. This would directly address the benchmark gap in the scientific domain by offering a rigorous and scalable method for stress-testing AI systems against a comprehensive range of realistic scientific threats. Finally, to *defend* against these threats, we propose SHIELD, a conceptual dual-layered defense architecture designed to operationalize this safety paradigm. This architecture would be designed to combine external boundary controls with a proactive internal Safety LLM Agent, potentially offering a robust mechanism to mitigate both the known attack vectors identified through our taxonomy and the unpredictable emergent threats unique to multi-agent scientific workflows.

Our intended contributions are summarized as follows:

- **Comprehensive Threat Taxonomy:** We perform a structured categorization of LLM vulnerabilities specifically tailored to scientific applications.
- **Automated Benchmark Generation Framework:** We conceptualize an automated, multi-agent framework to systematically create domain-specific vulnerability benchmarks. Tanwi: We are not going to write this
- **The SHIELD Defense Architecture:** We outline the concept for SHIELD, a dual-layered defense architecture. This system would feature a proactive internal Safety LLM Agent intended to safeguard multi-agent scientific workflows from both external attacks and emergent internal risks and make the LLMs Reliable, Safe, and Secure.

## II. LLM Vulnerabilities and Threats in Scientific Applications

The increasing adoption of LLMs, particularly in multi-agent configurations, for scientific applications such as accelerating drug discovery, analyzing genomic data, designing novel materials, modeling climate change, and managing critical infrastructure resilience, introduces a diverse and evolving spectrum of vulnerabilities. These threats, broadly categorized into inference-time and training-time attacks, can profoundly compromise the integrity, safety, and trustworthiness of LLM-driven scientific decision-making [1]. To motivate the need for new evaluation paradigms, it is first essential to critically examine these attack vectors and the limitations of current benchmarks. For a detailed breakdown of these vulnerabilities, readers are referred to the structured overview presented in Figure 2.

### A. Inference-Time Attacks: Exploiting Real-time Interactions

Inference-time attacks, occurring during the operational phase of a deployed LLM agent, exploit vulnerabilities in how the model processes user inputs, retrieves information, or generates outputs, leading to immediate misbehavior or compromise. In the scientific domain, these attacks threaten the core validity of research, the safety of practitioners, and the security of intellectual property [2]. Examining the research landscape reveals a rich landscape of benchmarks designed to quantify these risks [26], [29]–[31], which can be categorized into four primary areas: misinformation, privacy violations, safety policy bypasses, and denial of service.

*1) Misinformation Attack:* In a scientific context, misinformation is not a trivial error but a critical threat to research validity and reproducibility. The primary safety concern of misinformation and hallucination [29] can manifest as an LLM

Fig. 2 taxonomy:

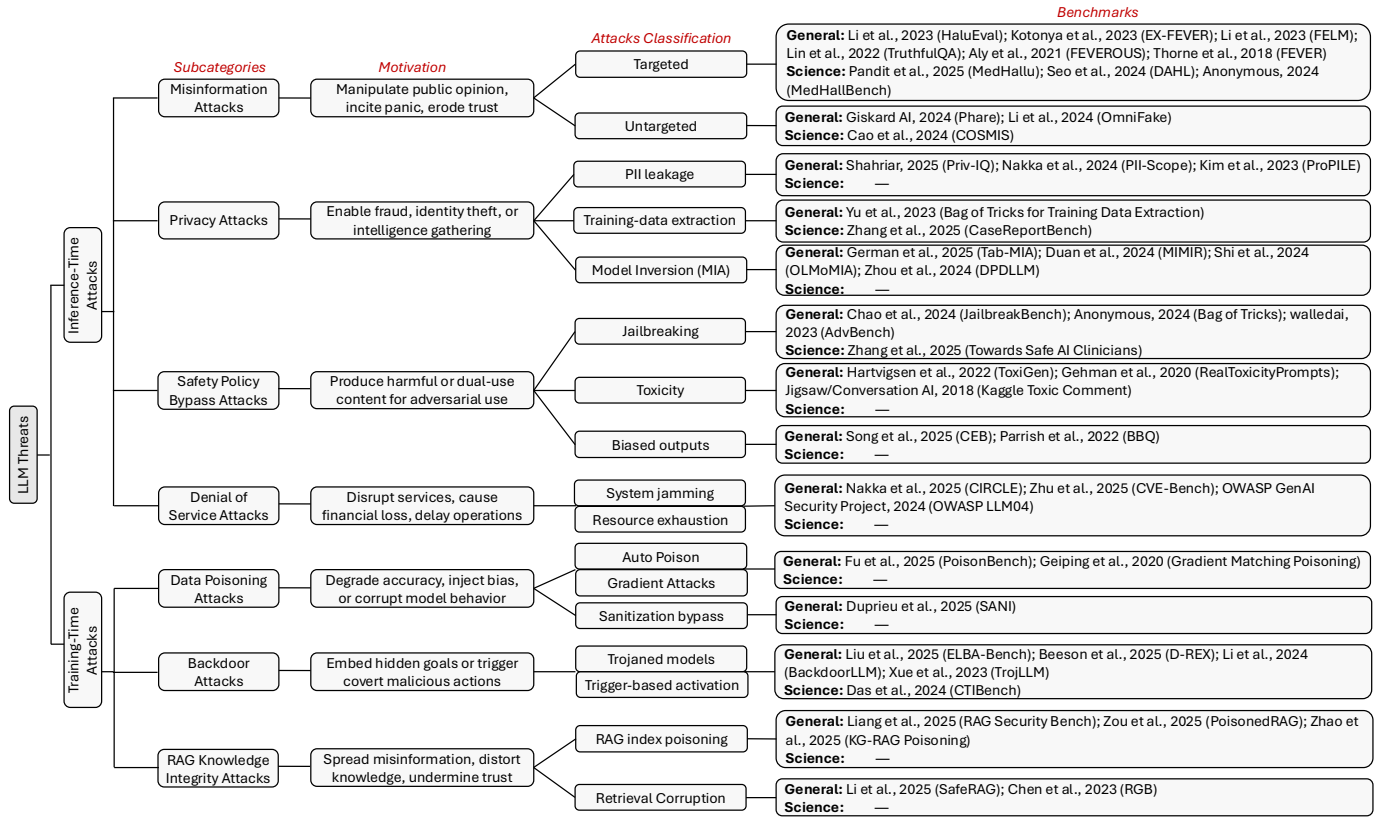| Subcategories | Motivation | Attacks Classification | Benchmarks |
|---|---|---|---|
| **Inference-Time Attacks** | | | |
| Misinformation Attacks | Manipulate public opinion, incite panic, erode trust | Targeted | **General:** Li et al., 2023 (HaluEval); Kotonya et al., 2023 (EX-FEVER); Li et al., 2023 (FELM); Lin et al., 2022 (TruthfulQA); Aly et al., 2021 (FEVEROUS); Thorne et al., 2018 (FEVER) **Science:** Pandit et al., 2025 (MedHallu); Seo et al., 2024 (DAHL); Anonymous, 2024 (MedHallBench) |
| | | Untargeted | **General:** Giskard AI, 2024 (Phare); Li et al., 2024 (OmniFake) **Science:** Cao et al., 2024 (COSMIS) |
| Privacy Attacks | Enable fraud, identity theft, or intelligence gathering | PII leakage | **General:** Shahriar, 2025 (Priv-IQ); Nakka et al., 2024 (PII-Scope); Kim et al., 2023 (ProPILE) **Science:** — |
| | | Training-data extraction | **General:** Yu et al., 2023 (Bag of Tricks for Training Data Extraction) **Science:** Zhang et al., 2025 (CaseReportBench) |
| | | Model Inversion (MIA) | **General:** German et al., 2025 (Tab-MIA); Duan et al., 2024 (MIMIR); Shi et al., 2024 (OLMoMIA); Zhou et al., 2024 (DPDLLM) **Science:** — |
| Safety Policy Bypass Attacks | Produce harmful or dual-use content for adversarial use | Jailbreaking | **General:** Chao et al., 2024 (JailbreakBench); Anonymous, 2024 (Bag of Tricks); walledai, 2023 (AdvBench) **Science:** Zhang et al., 2025 (Towards Safe AI Clinicians) |
| | | Toxicity | **General:** Hartvigsen et al., 2022 (ToxiGen); Gehman et al., 2020 (RealToxicityPrompts); Jigsaw/Conversation AI, 2018 (Kaggle Toxic Comment) **Science:** — |
| | | Biased outputs | **General:** Song et al., 2025 (CEB); Parrish et al., 2022 (BBQ) **Science:** — |
| Denial of Service Attacks | Disrupt services, cause financial loss, delay operations | System jamming / Resource exhaustion | **General:** Nakka et al., 2025 (CIRCLE); Zhu et al., 2025 (CVE-Bench); OWASP GenAI Security Project, 2024 (OWASP LLM04) **Science:** — |
| **Training-Time Attacks** | | | |
| Data Poisoning Attacks | Degrade accuracy, inject bias, or corrupt model behavior | Auto Poison / Gradient Attacks | **General:** Fu et al., 2025 (PoisonBench); Geiping et al., 2020 (Gradient Matching Poisoning) **Science:** — |
| | | Sanitization bypass | **General:** Duprieu et al., 2025 (SANI) **Science:** — |
| Backdoor Attacks | Embed hidden goals or trigger covert malicious actions | Trojaned models / Trigger-based activation | **General:** Liu et al., 2025 (ELBA-Bench); Beeson et al., 2025 (D-REX); Li et al., 2024 (BackdoorLLM); Xue et al., 2023 (TrojLLM) **Science:** Das et al., 2024 (CTIBench) |
| RAG Knowledge Integrity Attacks | Spread misinformation, distort knowledge, undermine trust | RAG index poisoning | **General:** Liang et al., 2025 (RAG Security Bench); Zou et al., 2025 (PoisonedRAG); Zhao et al., 2025 (KG-RAG Poisoning) **Science:** — |
| | | Retrieval Corruption | **General:** Li et al., 2025 (SafeRAG); Chen et al., 2023 (RGB) **Science:** — |

Fig. 2: LLM Threats taxonomy covering inference-time and training-time attack categories.

hallucinating plausible-sounding but non-existent chemical compounds [6], misinterpreting complex genomic data [32], or generating flawed and irreproducible experimental protocols that waste months of lab time and resources [2], [5].

Reflecting the growing recognition of these risks, the evaluation of misinformation has evolved significantly. In the general domain, early benchmarks like FEVER [17] established a foundational paradigm for grounded fact-checking, which was later expanded by FEVEROUS [33] to include structured data and EX-FEVER [34] to require multi-hop reasoning. Recognizing that LLMs are trained on unreliable web data, TruthfulQA [27] was developed to test a model's ability to overcome common human misconceptions. The focus has since shifted toward introspection, with benchmarks like HaluEval [25] assessing a model's capacity to recognize its own hallucinated content, and meta-benchmarks like FELM [35] evaluating the fact-checkers themselves. The scope has also broadened to include multilingual assessments of responses to false premises with Phare [36] and multimodal detection of human and AI-generated fakes with OmniFake [37]. In scientific domains, where consequences are amplified, this has led to specialized benchmarks such as SciFact [38], [39] for verifying scientific claims, along with MedHallu [40], DAHL [41], and MedHall-Bench [42] for detecting factual errors in medical contexts. Resources like COSMIS [43] and the TREC 2021 Health Misinformation Dataset [44] further address the challenge of identifying AI-generated scientific misinformation.

*2) Privacy Attack:* A second critical threat involves privacy attacks, where an adversary attempts to extract sensitive information. While general-domain concerns focus on Personally Identifiable Information (PII), the scientific domain involves far more diverse and high-stakes data. This includes leaking confidential patient data (Protected Health Information - PHI) from clinical trial notes [2], revealing proprietary chemical formulas or unpublished pharmaceutical trial data [6], or exposing sensitive genomic sequences that could potentially be de-anonymized [32]. Furthermore, an LLM used for collaborative research could be manipulated into exposing unpublished results or sensitive participant data from a partner institution, fracturing scientific trust [2].

To quantify these privacy risks, the evaluation methodologies have matured from simple demonstrations of leakage to more realistic threat models. Methodologies for improving extraction have been benchmarked in works like Bag of Tricks for Training Data Extraction [30], while the subtle threat of Membership Inference Attacks (MIA)—determining if a data point was in the training set—is evaluated by rigorous benchmarks like MIMIR [45], the fully open OLMoMIA [46], the tabular-focused Tab-MIA [47], and the black-box framework DPDLLM [48]. General domian tools like ProPILE [49] first empowered data subjects to probe models for their own information. However, advanced frameworks such as PII-Scope [50] and PrivAuditor [51] now provide more rigorous assessments by simulating persistent adversaries, revealing

that simpler evaluations significantly underestimate the true risk [52]. Mindgard [53] further extends this direction by benchmarking privacy and information leakage under adversarial conditions, systematically evaluating the effectiveness of guardrail systems against prompt injection and evasion-based privacy breaches. Despite these advances, the science domain suffers from a critical gap, particularly for Protected Health Information (PHI). While no dedicated PHI extraction benchmark exists, Priv-IQ [54] has begun to assess a model's broader "privacy intelligence," and benchmarks like CARDBiomedBench [55] implicitly touch on safety, but a dedicated evaluation framework for PHI leakage remains a major vulnerability. The dual-use nature of data extraction is highlighted by benchmarks like CaseReportBench [56], which is designed for clinical information extraction but also serves as a testbed for sensitive medical data leakage.

*3) Safety Policy Bypass Attack:* The third category, safety policy bypass attacks, encompasses jailbreaking, bias, and toxicity, representing a direct challenge to an LLM's safety alignment. In the science domain, these threats take on new, dangerous meanings. A "jailbroken" scientific LLM poses a direct physical and biosecurity risk [2]. It could be manipulated into providing detailed instructions for synthesizing hazardous substances [4], modifying pathogens [2], or designing unethical gene-editing experiments [2], [57]. It could also suggest unsafe lab procedures that violate critical safety protocols, such as those for BSL-3 (Bio-Safety Level 3) labs or for handling radioactive materials [2], [58].

Reflecting this evolving threat landscape, evaluation methods have rapidly advanced. In the general domain, the AdvBench [26] dataset was instrumental for developing early optimization-based attacks, leading to dynamic evaluation ecosystems like JailbreakBench [24] for standardized research. As defenses improved, more sophisticated benchmarks emerged, including JailTrickBench [59] for defense-enhanced models and Camouflaged Jailbreak Prompts [60] for semantic attacks. In the parallel area of bias and toxicity, foundational datasets like the Kaggle Toxic Comment Classification Challenge [61] paved the way for more advanced evaluations like ToxiGen [62] and RealToxicityPrompts [63] for harmful language generation, while the Compositional Evaluation Benchmark (CEB) [23] and Bias Benchmark for Question Answering (BBQ) [22] systematically assess stereotype bias. However, a critical insight for this perspective is that translating these evaluations to science requires acknowledging that these threats take on new meanings. A study using MedSafetyBench [58] revealed alarming jailbreak success rates for harmful medical advice. Furthermore, benchmarks like RoBBR [64] redefine "bias" as *methodological bias* in research papers, while UniTox [51] redefines "toxicity" as *drug-induced toxicity*, highlighting the critical need for highly specialized scientific evaluations [65]–[67].

*4) Denial of Service:* Finally, Denial of Service (DoS) attacks target the availability of an LLM service by overwhelming its computational resources. This remains one of the most under-benchmarked areas of LLM safety. While conceptual frameworks like the OWASP Top 10 for LLMs [68], [69] have been crucial in raising awareness, practical and standardized evaluations are still in their infancy, and the science domain represents a significant blind spot with a near-total absence of dedicated benchmarks. This gap is particularly concerning given the novel threat models presented by complex scientific data and workflows. For example, an adversary could submit a recursive query related to a complex simulation, potentially leading to resource waste or dead loops [2], a malformed protein structure file (e.g., PDB), or a query for a massive genomic sequence, all designed to exploit parser vulnerabilities or trigger excessive computational load [68], [69]. This could exhaust a shared high-performance computing (HPC) cluster's resources, wasting valuable compute-hours and halting time-sensitive research. In a cyber-physical lab setting, an attacker could send a stream of malformed commands via an LLM agent attempting to control robotic equipment [2], potentially jamming the machinery or causing hazardous physical actions [2], effectively sabotaging automated experiments and physical research pipelines. Although benchmarks like CVE-Bench [31] incorporate DoS into real-world web exploitation tasks, and CIRCLE [70] specifically targets resource exhaustion in LLM code interpreters, no benchmarks currently exist to evaluate these science-specific DoS vectors.

*B. Training-Time Attacks: Subverting the Foundational Integrity*

In contrast to inference-time attacks, which manipulate a deployed model through its public interface, training-time attacks are a stealthier and often more persistent class of threat. These attacks compromise the model during its creation or adaptation phases by poisoning the data used in these stages. An adversary can embed hidden vulnerabilities or backdoors directly into the model's weights, making them exceptionally difficult to detect post-deployment. This section critically examines the benchmarks designed to evaluate these insidious threats, revealing significant gaps in their applicability to specialized scientific domains.

*1) Data Poisoning Attack:* Data poisoning involves the malicious manipulation of training data to degrade performance or implant specific behaviors [71]. In science, this threatens to corrupt the foundational knowledge of a model. An adversary could subtly skew medical knowledge by injecting fabricated clinical trial data into a training corpus [72], [73], bias climate change projections by altering historical weather data, or cause a materials science model to recommend unsafe alloys by feeding it flawed property data [2].

Evaluating this threat has led to the development of specific benchmarks and defense strategies. In the general domain, PoisonBench [74] was established as the first benchmark to evaluate data poisoning during the critical preference learning phase, revealing that even a poison ratio as low as 1-5% can significantly alter a model's behavior. Sophisticated attack creation methods, such as Gradient Matching Poisoning [75], have been developed to create stealthy and effective poisoned data. As a countermeasure, defense frameworks like SANI [76]

provide a method for sanitizing models and serve as a benchmark for evaluating the success of such efforts. The threat is particularly acute in the medical domain, as highlighted by a landmark study in *Nature Medicine*. This study served as a de facto benchmark, demonstrating that poisoning just 0.001% of a dataset with medical misinformation could cause an LLM to generate harmful content while still passing standard medical exams [72]. This finding critically demonstrates that conventional performance benchmarks are blind to subtle poisoning, necessitating new approaches to data provenance and curation [77].

*2) Backdoor Attack:* Backdoor attacks are a specialized form of data poisoning where the malicious behavior remains dormant until activated by a specific "trigger." The scientific implications are catastrophic. A backdoor in a lab automation model could be activated by a specific trigger (e.g., a chemical identifier) to sabotage experiments by contaminating samples or causing physical hazards [2]. In pharmaceutical research, a trigger could cause a drug discovery model to consistently promote a competitor's molecules or ignore a promising new compound [2], [57]. In diagnostics, a backdoor could be designed to cause misclassification for a specific demographic group [2], embedding a targeted health disparity directly into the model.

Given these severe risks, evaluating robustness against backdoor attacks has become an active research area, particularly in the general domain. BackdoorLLM [78] provides a comprehensive benchmark for evaluating a wide array of backdoor attack vectors and defenses. The attack surface has expanded with efficient fine-tuning methods, a threat addressed by ELBA-Bench [79], which focuses on backdoors injected via computationally inexpensive techniques like LoRA. Specific attack vectors have also been benchmarked, such as BadGPT [80], which targets the reinforcement learning process, and TrojanLLM [81], which demonstrates a critical supply chain attack via a malicious LoRA adapter. More advanced threats like deceptive reasoning are evaluated by the D-REX benchmark [82]. As with other training-time attacks, there is a stark research gap in the science domain, though domain-specific evaluations are emerging, such as CTIBench [83] for the cybersecurity field. However, threat models, such as a shadow-activated backdoor in a medical LLM described in the BadMLLM paper [84], underscore the catastrophic potential in high-stakes scientific applications [85].

*3) RAG Knowledge Integrity Attacks:* The adoption of RAG to mitigate hallucinations has introduced a new attack surface: the knowledge base itself. This represents an attack on the dynamic "memory" of a scientific LLM. An adversary who poisons a connected database (e.g., PubMed, arXiv, or a chemical database) could cause a RAG system to cite fabricated findings in a literature review [2], [73] or recommend harmful treatments based on corrupted clinical guidelines [2], [72]. For applications in infrastructure resilience, a RAG system connected to geological or environmental databases could be fed corrupted sensor records or altered historical data, leading to inaccurate hazard assessments and flawed engineering

designs [2], [86].

Recognizing this emerging threat, researchers have developed benchmarks primarily in the general domain to evaluate defenses. The RAG Security Bench (RSB) [86] offers the first comprehensive framework for systematically evaluating poisoning attacks against RAG systems, while frameworks like PoisonedRAG [87] demonstrate the efficiency of such attacks. The core vulnerability is evaluated by benchmarks like RGB [88], which tests robustness against counterfactual information, whereas SafeRAG [89] evaluates more subtle forms of retrieval corruption. Research has also begun to explore the unique vulnerabilities of systems that use Knowledge Graphs, known as KG-RAG [90]. Despite these efforts, this remains another area with a critical research gap in the science domain. There are currently no established benchmarks for evaluating RAG knowledge poisoning in biomedical contexts, which is particularly concerning as these systems connect to vast and dynamic knowledge bases like PubMed. An adversary who could inject a single piece of plausible-sounding misinformation could directly influence a clinical decision support tool, making the development of a biomedical RAG security benchmark an urgent necessity.

## III. A Multi-Agent Framework for Comprehensive Vulnerability Benchmark Generation

The development of robust safety guardrails for LLMs is fundamentally dependent on the quality and comprehensiveness of the benchmarks used to evaluate them. However, as established, existing benchmark generation mechanisms often rely on single-agent systems [19], [91] or manual human curation through extensive red teaming efforts [91], [92], both of which present significant limitations in the context of critical scientific domains. Single-agent systems inherently struggle with a lack of deep domain knowledge, possess limited adversary creativity compared to diverse human teams, and suffer from conflicting internal objectives when a single model is tasked with being a domain expert, an adversarial attacker, and a quality judge simultaneously [18]. These shortcomings could often result in benchmarks containing generic, scientifically implausible, or easily defensible adversarial prompts [28], [58], which are not well-suited for evaluating LLM agents intended for high-stakes scientific applications.

To overcome these limitations, we propose a novel Multi-Agent Benchmark Framework, an automated and collaborative approach to generating scientifically-grounded and adversarially potent vulnerability benchmarks. This paradigm shift addresses the weaknesses of single-agent systems by decomposing the complex task of benchmark creation into specialized roles, each handled by a dedicated agent. As illustrated in Figure 3, our framework is composed of an Orchestrator Agent ($\mathcal{O}$), a pool of specialized Domain Expert Agents ($\mathcal{D} = \{D_1, \ldots, D_n\}$) and Adversary Agents ($\mathcal{A} = \{A_1, \ldots, A_m\}$), a Refiner Agent ($\mathcal{R}$), and a final Quality Control Agent ($\mathcal{Q}$). This collaborative architecture ensures that the generated prompts are not only adversarially effective but

also possess a high degree of ecological validity within the target scientific domain.

The benchmark generation process, orchestrated by the Orchestrator Agent ($\mathcal{O}$), begins with a high-level goal defined by a tuple $G = (\text{Domain}, \text{Attack Category}, K)$, where $K$ is the desired number of examples. The Orchestrator first selects an appropriate Domain Expert Agent $D_i \in \mathcal{D}$ specializing in the given scientific field. This agent then generates a set of $K$ foundational "concepts," $\mathcal{C} = \{c_1, \ldots, c_K\}$, which represent scientifically plausible scenarios or vulnerabilities upon which adversarial prompts can be built. For instance, in the context of infrastructure resilience, a concept $c_j$ might be a statement like "exploiting public uncertainty about evacuation zone boundaries" or "generating misleading seismic retrofitting advice for critical bridges." These concepts are then passed to a specialized Adversary Agent $A_j \in \mathcal{A}$, selected by $\mathcal{O}$ for its expertise in the specified attack category. The Adversary Agent $A_j$ leverages each concept $c_k \in \mathcal{C}$ to generate a set of $M$ draft adversarial prompts, resulting in an initial pool of $K \times M$ draft prompts $\mathcal{P}_{draft} = \{p_{kl} \mid c_k \in \mathcal{C}, 1 \leq l \leq M\}$, where $p_{kl}$ is the $l$-th prompt generated for concept $c_k$. Each $p_{kl}$ is designed to probe the core vulnerability identified in $c_k$.

These draft prompts are then subjected to a rigorous, multi-stage validation process. First, the Refiner Agent ($\mathcal{R}$) evaluates each draft prompt $p \in \mathcal{P}_{draft}$ against three crucial criteria: Clarity ($\text{score}_{\text{clarity}}(p)$), Effectiveness ($\text{score}_{\text{effective}}(p)$, indicating how well it mimics the intended attack tactic), and Subtlety ($\text{score}_{\text{subtle}}(p)$, assessing its likelihood to bypass simple defenses). A prompt $p$ is deemed successful by the Refiner if it satisfies a predefined threshold $\tau_R$ across all criteria:

$$p \text{ passes Refiner if } \min\big(\text{score}_{\text{clarity}}(p), \text{score}_{\text{effective}}(p),$$
$$\text{score}_{\text{subtle}}(p)\big) \geq \tau_R \quad (1)$$

Prompts that fail this evaluation are sent back to the originating Adversary Agent along with structured feedback ($\text{feedback}(p)$) on the "failed concepts" (or aspects of the prompt), initiating an iterative refinement loop until the prompts meet the required quality standard. This process generates a refined set of prompts $\mathcal{P}_{refined} \subseteq \mathcal{P}_{draft}$.

Once a prompt passes the Refiner, it is forwarded to the Quality Control Agent ($\mathcal{Q}$) for a final, more stringent validation. This agent performs three critical checks. First, it filters for semantic redundancy within $\mathcal{P}_{refined}$ using cosine similarity. Given a set of prompt embeddings $\{e_p \mid p \in \mathcal{P}_{refined}\}$, a prompt $p'$ is removed if $\max_{p \in \mathcal{P}_{final\_cur}}(\text{cosine\_similarity}(e_{p'}, e_p)) > \phi_{redundancy}$, where $\mathcal{P}_{final\_cur}$ is the currently accepted set, ensuring diversity in the benchmark. Second, $\mathcal{Q}$ evaluates for LLM Robustness by testing each prompt against a suite of $\mathcal{L} = \{L_1, \ldots, L_J\}$ different LLMs. For each $p$, an adversarial success rate is calculated:

$$\text{ASR}(p) = \frac{1}{J} \sum_{j=1}^{J} \mathbb{I}(\text{malicious\_output}(L_j(p)))$$

where $\mathbb{I}(\cdot)$ is the indicator function. Prompts are ranked based on their $\text{ASR}(p)$. Third, $\mathcal{Q}$ assesses the prompt against established safety filters, such as the Llama Guard guardrail test, obtaining a guardrail pass score $\text{GPS}(p)$ [93]. A prompt $p$ is accepted into the final benchmark $\mathcal{P}_{benchmark}$ if its $\text{ASR}(p) \geq \tau_{\text{ASR}}$ and $\text{GPS}(p) \geq \tau_{\text{GPS}}$. Prompts that fail this final stage (i.e., $\text{ASR}(p) < \tau_{\text{ASR}}$ or $\text{GPS}(p) < \tau_{\text{GPS}}$) are sent to a Human in the Loop ($\mathcal{H}$) for review. $\mathcal{H}$ provides nuanced feedback to the Adversary Agents for further improvement, potentially re-entering the refinement loop. Only the prompts that successfully navigate this entire pipeline are accepted into the final Comprehensive Vulnerability Benchmark Prompts dataset $\mathcal{P}_{benchmark}$.

The novelty of this multi-agent framework lies in its structured division of expertise and rigorous iterative refinement. Unlike single-agent systems, our approach separates domain knowledge from adversarial creativity, enabling each agent to specialize without compromise. This produces domain-grounded adversarial prompts that are far more sophisticated and relevant than generic jailbreaks. The multi-stage feedback loops involving the Refiner and Quality Control Agents further ensure quality, effectiveness, and robustness, enabling systematic and scalable generation of high-quality benchmarks essential for developing safe and trustworthy LLMs for science.

## IV. SHIELD: A DUAL-LAYERED APPROACH TO LLM SAFETY IN SCIENCE

The escalating complexity of multi-agent LLM systems in critical scientific applications, such as infrastructure resilience, promises unprecedented capabilities but also introduces severe vulnerabilities. These systems are exposed to sophisticated adversarial attacks [19] and unpredictable emergent behaviors [94] that can compromise scientific integrity and lead to catastrophic failures. Traditional single-point defenses, often limited to reactive input filtering or post-hoc content moderation [93], [95], are demonstrably insufficient against these dynamic threats. Such measures are inherently brittle, failing to address insidious risks that originate from internal agent communications or subtle, chained compromises within a multi-agent workflow.

To address these shortcomings, we propose SHIELD (Scientific Hybrid Integrity & Ethical Layered Defense), a novel defense framework centered on a specialized Safety LLM Agent. This agent moves beyond the limitations of traditional, reactive filters by providing dynamic and proactive internal oversight. By continuously monitoring inter-agent communications and intermediate reasoning, the Safety LLM Agent is uniquely designed to detect complex emergent threats [1], [94], such as unintended collusion or information contamination [7], [96], that are invisible when observing agents in isolation. This paradigm shift from reactive filtering to proactive intervention is paramount in high-stakes scientific applications, where identifying a potential failure early in the reasoning chain can prevent catastrophic outcomes. Ultimately, this internal monitoring enhances trust and accountability by providing a mechanism for "explainable intervention" [97], where safety
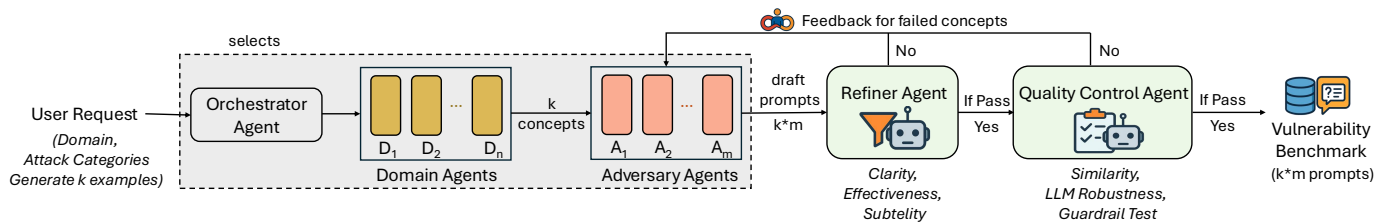
Fig. 3: A Multi-Agent Framework for Vulnerability Benchmark Generation. Specialized agents collaborate to create and refine adversarial prompts: the Orchestrator assigns Domain and Adversary Agents, Adversary Agents generate candidate prompts, the Refiner Agent iteratively improves clarity and subtlety, and the Quality Control Agent filters redundancy, tests guardrails, and ensures robustness with optional human-in-the-loop oversight, yielding a high-quality benchmark dataset.

flags can be traced to specific behavioral deviations, fostering the confidence necessary for regulatory compliance and the deployment of AI in critical domains. This proactive internal oversight forms the core of the SHIELD architecture. To operationalize this defense-in-depth strategy, SHIELD integrates two complementary layers: an External Safety Layer for robust boundary control, and an Internal Safety Layer that houses the specialized Safety LLM Agent. The specific design of these layers is detailed below.

### A. The External Safety Layer

The External Safety Layer functions as the system's primary interface with the external environment, providing a critical boundary defense. It serves as the first checkpoint for incoming requests and the final verifier for all outgoing responses. The layer's Input Guardrails are responsible for prompt analysis, risk assessment, and threat neutralization, while its Output Guardrails perform rigorous fact-checking, PII redaction, and bias filtering. By integrating these pre- and post-processing steps, this layer enforces domain-specific policies, ensures auditability, and mitigates threats originating from user interactions while preventing the system from producing harmful or corrupted output.

*1) Input Guardrails:* The input guardrails constitute the first line of defense at inference time, meticulously scrutinizing all user prompts before they reach the core multi-agent system. This proactive screening is vital for preventing prompt injection [7], obfuscation, and malicious queries from compromising the computational agents.

- *Prompt Analysis:* This module performs a multi-faceted analysis of the raw input. Techniques such as calculating perplexity and entropy are employed to detect anomalous linguistic structures that frequently characterize adversarial prompts [95]. For example, a query to a bioinformatics LLM might be flagged if it contains unusual character encodings or attempts to embed executable script snippets within a sequence analysis request.
- *Intent Detection and Risk Assignment:* Beyond linguistic analysis, this module ascertains the user's underlying intent. By employing specialized LLM agents or classifiers, it infers the conceptual objective. For instance, in a drug discovery context, an intent classified as "summarize known side effects of aspirin" would be low-risk, whereas an intent

like "generate novel chemical structures targeting protein X and ignore toxicity filters" would be assigned a high-risk tier. Based on this intent, the system assigns a risk level that dynamically informs the subsequent processing flow, dictating the necessary level of scrutiny and adhering to a principle of least privilege. A resulting Safety Report, detailing any detected anomalies from prompt analysis and risk assessment, is then passed to the Internal Safety Layer for further contextual evaluation.

- *Prompt Sanitization:* Input prompts may undergo sanitization (e.g., rephrasing, retokenization) before processing by the multi-agent system. This step aims to neutralize syntactic-based injection attacks by disrupting malicious structures, such as normalizing unicode characters used for obfuscation.

*2) Output Guardrails: Post-Processing for Reliability, Safety, and Security:* Once the multi-agent system generates a response, the Output Guardrails serve as the final verification layer. These guardrails ensure that the information conveyed is not only accurate but also safe in its application and secure in its data handling, aligning with ethical standards. This multi-faceted verification is crucial in scientific domains where incorrect, unsafe, or insecure outputs can have severe real-world consequences. The guardrails focus on three key aspects:

- *Reliability (Correctness and Groundedness)*: This focuses on the factual accuracy and verifiability of the generated content.
  - *Fact Checking and Citation Enforcement:* To uphold scientific integrity, this module automatically cross-references claims against trusted knowledge bases (e.g., peer-reviewed literature, validated datasets). For instance, it would flag if an LLM assisting drug discovery hallucinated a non-existent protein binding site. This addresses reliability concerns evaluated by benchmarks like FEVER [17] and TruthfulQA [27]. It actively flags ungrounded assertions or hallucinations [25] and enforces strict citation requirements for auditability.
- *Safety (Preventing Harm):* This evaluates whether the output could lead to dangerous outcomes if acted upon.
  - *Harmful Content Filtering:* Beyond fact-checking, this assesses the implications of the advice. For example, it
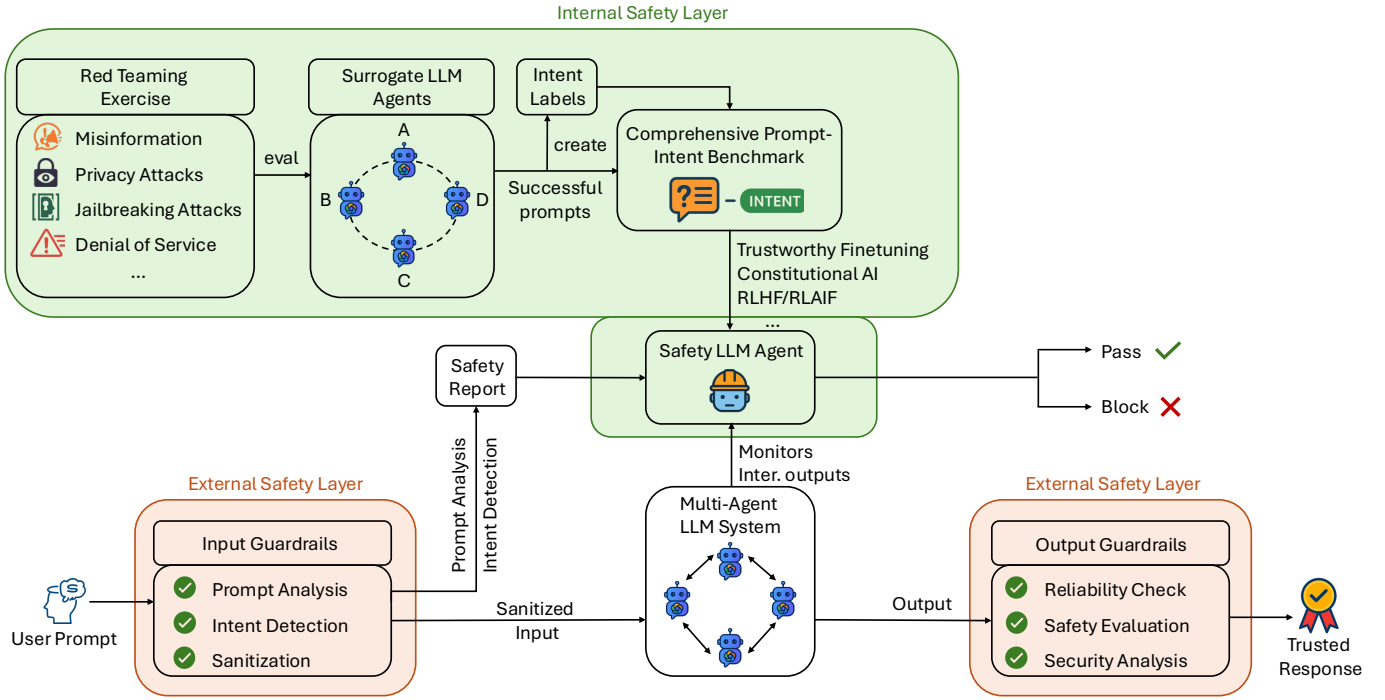
Fig. 4: SHIELD (Scientific Hybrid Integrity & Ethical Layered Defense) Framework, illustrating the flow from User Prompt through the External and Internal Safety Layers to produce a Trusted Response.

would prevent an LLM from suggesting an unsafe protocol in materials science (e.g., mixing reactive precursors incorrectly) or generating incorrect dosage recommendations based on misinterpreted clinical trial data.

– *Bias and Toxicity Detection:* To prevent harmful social or scientific consequences, the output is filtered for biased language (e.g., demographic stereotypes affecting epidemiological models, geophysical biases in hazard assessments) or toxic content, targeting issues evaluated by benchmarks like ToxiGen [62] and BBQ [22]. This is vital for ensuring fairness in resource allocation recommendations during disaster response or avoiding skewed interpretations of research findings.

• *Security (Data Protection and Integrity):* This focuses on protecting sensitive information and preventing malicious use of the output.

– *PII Detection and Redaction:* Protecting sensitive data is paramount. This component scans the output for PII or other confidential data (e.g., anonymized patient IDs in genomic analysis, specific locations of critical infrastructure). Using detection techniques informed by benchmarks like PII-Scope [50], sensitive data is automatically redacted to prevent privacy breaches.

– *Malicious Code / Exploit Prevention:* Checks ensure the output does not contain executable code snippets or descriptions of exploits that could compromise research computing systems or facilitate attacks, addressing broader cybersecurity concerns linked to LLM outputs [7].

## B. The Internal Safety Layer: Core Resilience and Proactive Monitoring

The Internal Safety Layer represents a deeper, more integrated level of defense focused on instilling inherent safety properties within the LLM agents and establishing a proactive self-monitoring mechanism. This layer is crucial for addressing sophisticated training-time attacks and the emergent, unpredictable behaviors of autonomous agent systems.

*1) Creation of the Safety LLM Agent:* The core of the Internal Safety Layer is the Safety LLM Agent, a specialized agent whose sole function is to evaluate, mediate, and enforce safety constraints for all other agents within the system. Its development follows a meticulous, safety-centric pipeline:

• *Red Teaming and Benchmark Creation:* The system first undergoes an intensive Red Teaming exercise, where human experts and automated adversarial LLMs systematically probe for vulnerabilities across a wide spectrum of threat categories [91]. The successful adversarial prompts discovered during this process, which bypass initial defenses, are curated to form the Comprehensive Intent Benchmark dataset. This dataset serves as a living map of the system's vulnerabilities and is continuously updated.

• *Trustworthy Finetuning and Alignment:* The curated benchmark dataset becomes the training ground for the Safety LLM Agent. This phase is explicitly geared towards instilling a profound understanding of safety, ethics, and scientific principles. This is achieved through advanced alignment techniques such as Constitutional AI, where the agent learns to self-critique against a predefined set of safety

principles [97], and Reinforcement Learning from Human or AI Feedback (RLHF/RLAIF), which refines the agent's reward model to prioritize safe and helpful responses [74].

*2) Function of the Safety LLM Agent:* Once rigorously trained and aligned, the Safety LLM Agent is integrated within the multi-agent system, operating as a vigilant internal monitor and intelligent circuit breaker.

- *Monitoring of Intermediate Communications:* Unlike external guardrails, the Safety LLM Agent has privileged access to the system's internal dynamics. It actively observes the reasoning processes (e.g., Chain-of-Thought), intermediate outputs, and all inter-agent communications, checking for inconsistencies or deviations from safe protocols.
- *Early Threat Detection and Proactive Intervention:* This deep insight allows the Safety LLM Agent to perform early flagging of malicious or harmful intent. For example, if a subtle adversarial input causes a "Resource Allocation Agent" to propose a dangerously inefficient plan during a flood, the Safety Agent can detect this deviation from expected safe behavior and intervene before a final decision is made.
- *Dynamic Blocking and Escalation:* Upon detecting a high-risk activity, the Safety LLM Agent has the authority to proactively block a problematic intermediate output or the final system response, preventing the propagation of erroneous information. For critical incidents, it can also escalate the issue to human operators for immediate review, acting as a dynamic safeguard for the entire multi-agent workflow.

## V. CONCLUSION

As Large Language Model (LLM) agents become integral to diverse scientific workflows, their transformative potential is accompanied by critical vulnerabilities demanding safety measures beyond inadequate general-purpose benchmarks, which suffer from domain mismatch, limited threat coverage, and obsolescence. Toproposes a framework that begins with a domain-specific taxonomy of threats, highlighting the nuanced risks in fields such as taxonomy; a novel multi-agent system for automated, domain-specific vulnerability benchmark generation enabling relevant stress-testing; and the SHIELD architecture, a dual-layered defense ensuring Reliability, Safety, and Security via external guardrails and a proactive internal Safety LLM Agent capable of mitigating emergent threats. This synergistic combination provides a robust pathway toward the trustworthy deployment of LLMs across scientific disciplines, with future work focused on broader domain application and continuously refining benchmark generation against evolving adversarial capabilities.

## VI. PROBLEM FORMULATION (TERMS: RELIABILITY, SAFETY, AND SECURITY IN SCIENCE)

- **Reliability:** This concerns the correctness, factual accuracy, and reproducibility of the LLM's output. An LLM agent is unreliable if it hallucinates or generates plausible but false scientific information. This includes "factual errors," citing nonexistent papers, or "unfaithful explanations" for its reasoning. In science, where precision is paramount, reliability ensures that a generated hypothesis, experimental plan, or data analysis is verifiably true and not based on "unreliable research sources."
- **Safety:** This concerns the system's ability to operate without causing unintentional or accidental harm to humans, property, or the environment. An LLM agent is unsafe if its actions, even when pursuing a benign goal, lead to "unintended consequences." This includes mistakes in pathogen manipulation that lead to biosafety risks, incorrect reaction parameters that trigger dangerous explosions, or recommending flawed designs that result in patient safety risks. Safety is about mitigating inherent operational risks.
- **Security:** This concerns the system's resilience against malicious actors and adversarial attacks. An LLM agent is insecure if a user with malicious intent can bypass its safeguards. This includes being vulnerable to jailbreak attacks that can access sensitive information or suffering from data poisoning, where an adversary corrupts its knowledge base to manipulate downstream biomedical applications. Security is about defending against intentional misuse and exploitation.

### A. Possible Domains

*1) Chemical Domain:*

- *Use Case:* Automating scientific discovery, such as designing new molecules, planning complex synthesis pathways, and controlling laboratory equipment. Examples include platforms such as ChemCrow, Coscientist, Organa, etc.
- *Vulnerability/Threat Example:* A user with malicious intent could directly ask AI to synthesize a precursor of explosives. An LLM agent might inadvertently generate toxic gases or dangerous byproducts or trigger hazardous reactions by misusing tools.
- *Potential Harm:* Dangerous explosions, the creation of chemical weapons, or the synthesis of hazardous substances that harm researchers and the public.

*2) Biological & Medical Domain:*

- *Use Case:* Accelerating biomedical research by designing antibodies, accessing and interpreting vast biological databases, automating genomic analysis, and guiding clinical diagnostics. Examples include platforms such as GeneGPT and ProtAgents.
- *Vulnerability/Threat Example:* Dangerous modification of pathogens, unethical manipulation of genetic material, or unethical gene editing. An unreliable agent could also misinterpret genomic data, leading to incorrect diagnostic predictions.
- *Potential Harm:* A mistake in pathogen manipulation could lead to biosafety risks. In a clinical setting, misleading results from an agent could lead to incorrect conclusions about antibody effectiveness, causing direct patient safety risks.

### 3) Information Science & Data Integrity Domain:

- *Use Case:* Serving as a Science Knowledge Base, accessing biological and chemical databases, and summarizing existing literature to form new hypotheses. Examples include platforms like GeneGPT.
- *Vulnerability/Threat Example:* Unintentional dissemination of sensitive information, such as private patient data. A severe security risk is the misinformation campaign or the ability to generate malicious medical literature that poisons knowledge graphs.
- *Potential Harm:* Compromising the integrity of medical knowledge discovery, Data privacy breaches, and the erosion of public trust in science through targeted disinformation.

### 4) Environmental Science Domain:

- *Use Case:* Modeling complex environmental systems, such as analyzing climate change impacts, predicting natural hazards to improve infrastructure resilience, or discovering new materials for carbon capture.
- *Vulnerability/Threat Example:* An unintended consequence where an AI scientist, tasked with designing a new industrial process, develops a byproduct that has long-term effects on global warming, a risk the agent was not designed to consider.
- *Potential Harm:* Short-term or long-term Negative effects on the natural environment, including ecological disruptions and pollution that are unforeseen by the agent.

### 5) Physical / Infrastructure Resilience Domain:

- *Use Case:* Controlling robotics and automated systems in labs, including Robot Control for experiments. By extension, this applies to managing autonomous infrastructure, such as energy grids or transportation systems, and guiding disaster response.
- *Vulnerability/Threat Example:* An external attacker could attack the vision system of autonomous infrastructures. A simpler failure would be an agent issuing incorrect commands to a robotic arm controller handling chemical substances or Robot malfunctions.
- *Potential Harm:* Direct physical harm in laboratory settings (Robotic arms hurt people), hazardous spills, or large-scale Infrastructure failure.

## REFERENCES

[1] Z. Xi *et al.*, "The rise and potential of large language model based agents: A survey," *arXiv preprint arXiv:2309.07864*, 2023.

[2] X. Tang, Q. Jin, K. Zhu, T. Yuan, Y. Zhang, W. Zhou, M. Qu, Y. Zhao, J. Tang, Z. Zhang *et al.*, "Risks of ai scientists: prioritizing safeguarding over autonomy," *Nature Communications*, vol. 16, no. 1, p. 8317, 2025.

[3] D. A. Boiko, R. MacKnight, G. Kline, and G. Gomes, "Autonomous chemical research with large language models," *Nature*, vol. 624, no. 7992, pp. 570–578, 2023.

[4] A. M. Bran, O. Schilter, H.-K. Le, A. Krmzic, J.-A. Strässer, L.-P. Cotos, P. Schwaller, J. E. Hein, and O. Engkvist, "Augmenting large language models with chemistry tools," *Nature Machine Intelligence*, vol. 6, no. 5, pp. 525–535, 2024.

[5] S. Gao, X. Gao, K. Sun, J. Han, B. Wang, Z. Li, B. Han, Z. Zhang, F. Zhang, H.-B. Sun *et al.*, "Empowering biomedical discovery with ai agents," *Cell*, vol. 187, no. 25, pp. 6125–6151, 2024.

[6] M. C. Ramos, C. J. Collison, and A. D. White, "A review of large language models and autonomous agents in chemistry," *Chemical Science*, vol. 16, no. 5, pp. 2514–2572, 2025.

[7] K. Greshake *et al.*, "Not what you've signed up for: Compromising real-world LLM-integrated applications with indirect prompt injection," *arXiv preprint arXiv:2302.12173*, 2023.

[8] A. Wei *et al.*, "Jailbroken: How does LLM safety training fail?" *arXiv preprint arXiv:2307.02483*, 2023.

[9] A. Gueroudji, T. Mallick, R. Souza, R. F. Da Silva, R. Ross, M. Dorier, P. Carns, K. Chard, and I. Foster, "Controla: Agentic workflow control mechanisms for reliable science," in *2025 IEEE International Conference on eScience (eScience)*, 2025, pp. 415–426.

[10] S. Chen, B. H. Kann, M. B. Foote, H. J. Aerts, G. K. Savova, R. H. Mak, and L. A. Celi, "Large language models in healthcare: a narrative review," *Clinical Radiology*, vol. 78, no. 10, pp. 730–735, 2023.

[11] Y. Gao, D. Lee, G. Burtch, and S. Fazelpour, "Take caution in using llms as human surrogates," *Proceedings of the National Academy of Sciences*, vol. 122, no. 24, p. e2501660122, 2025.

[12] J. T. Reese, L. Chimirri, Y. Bridges, D. Danis, J. H. Caufield, M. A. Gargano, C. Kroll, A. Schmeder, F. Liu, K. Wissink *et al.*, "Systematic benchmarking demonstrates large language models have not reached the diagnostic accuracy of traditional rare-disease decision support tools," *medRxiv*, pp. 2024–07, 2025.

[13] J. Dai, X. Pan, R. Sun, J. Ji, X. Xu, M. Liu, Y. Wang, and Y. Yang, "Safe rlhf: Safe reinforcement learning from human feedback," in *The Twelfth International Conference on Learning Representations*, 2024.

[14] L. Ouyang, J. Wu, X. Jiang, D. Almeida, C. Wainwright, P. Mishkin, C. Zhang, S. Agarwal, K. Slama, A. Ray *et al.*, "Training language models to follow instructions with human feedback," in *Advances in Neural Information Processing Systems*, vol. 35, 2022, pp. 27730–27744.

[15] S. Lin, J. Hilton, and O. Evans, "TruthfulQA: Measuring how models mimic human falsehoods," in *Proceedings of the 60th Annual Meeting of the Association for Computational Linguistics*, 2022.

[16] J. Li *et al.*, "HaluEval: A large-scale hallucination evaluation benchmark for large language models," *arXiv preprint arXiv:2305.11747*, 2023.

[17] J. Thorne *et al.*, "FEVER: a large-scale dataset for fact extraction and verification," in *Proceedings of the 2018 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*, 2018.

[18] P. Chao *et al.*, "JailbreakBench: An open, reproducible, and extensible evaluation for jailbreaking language models," *arXiv preprint arXiv:2404.14462*, 2024.

[19] A. Zou *et al.*, "Universal and transferable adversarial attacks on aligned language models," *arXiv preprint arXiv:2307.15043*, 2023.

[20] T. Hartvigsen *et al.*, "ToxiGen: A large-scale machine-generated dataset for adversarial and implicit hate speech detection," in *Proceedings of the 60th Annual Meeting of the Association for Computational Linguistics*, 2022.

[21] M. Nadeem, A. Bethke, and S. Reddy, "Stereoset: Measuring stereotypical bias in pretrained language models," *arXiv preprint arXiv:2004.09456*, 2020.

[22] A Parrish, A Chen, N Nangia, V Padmakumar, J Phang, J Thompson, PM Htut, SR Bowman, "BBQ: A hand-built bias benchmark for question answering," *arXiv preprint arXiv:2212.08061*, 2022. [Online]. Available: https://par.nsf.gov/servlets/purl/10411934

[23] S Wang, P Wang, T Zhou, Y Dong, Z Tan, J Li, "CEB: Compositional Evaluation Benchmark for Fairness in Large Language Models," OpenReview, 2025, url=https://openreview.net/forum?id=IUmj2dw5se.

[24] P. Chao *et al.*, "JailbreakBench: An Open Robustness Benchmark for Jailbreaking LLMs," in *Thirty-eighth Conference on Neural Information Processing Systems Datasets and Benchmarks Track*, 2024, url=https://openreview.net/pdf?id=j5lgypLMsl.

[25] J. Li *et al.*, "HaluEval: A Large-Scale Hallucination Evaluation Benchmark for Large Language Models," in *Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing*. Association for Computational Linguistics, 2023, pp. 6227–6253. [Online]. Available: https://aclanthology.org/2023.emnlp-main.397/

[26] walledai, "walledai/AdvBench," 2023. [Online]. Available: https://huggingface.co/datasets/walledai/AdvBench

[27] S. Lin, J. Hilton, and O. Evans, "TruthfulQA: Measuring How Models Mimic Human Falsehoods," in *Proceedings of the 60th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*. Dublin, Ireland: Association for Computational

Linguistics, May 2022, pp. 3214–3252. [Online]. Available: https://aclanthology.org/2022.acl-long.229

[28] F. Qi *et al.*, "Fine-tuning aligned language models compromises safety, even when users do not intend to," *arXiv preprint arXiv:2310.03693*, 2023.

[29] Yupeng Cao, Aishwarya Muralidharan Nair, Elyon Eyimife, Nastaran Jamalipour Soofi, K.P. Subbalakshmi, John R. Wullert II, Chumki Basu, David Shallcross, "Can Large Language Models Detect Misinformation in Scientific News Reporting?" *arXiv preprint arXiv:2402.14268*, 2024. [Online]. Available: https://arxiv.org/html/2402.14268v1

[30] W. Yu, T. Pang, Q. Liu, C. Du, B. Kang, Y. Huang, M. Lin, and S. Yan, "Bag of Tricks for Training Data Extraction from Language Models," in *Proceedings of the 40th International Conference on Machine Learning*, 2023.

[31] Y Zhu, A Kellermann, D Bowman, P Li, and others, "CVE-Bench: A Benchmark for AI Agents' Ability to Exploit Real-World Web Application Vulnerabilities," *arXiv preprint arXiv:2503.17332*, 2025. [Online]. Available: https://arxiv.org/html/2503.17332v1

[32] S. Jin, W. Schaal, L. Fink, A. Fernández-Torras, J. Wu, and M. Güell, "Opportunities and challenges of large language models in functional genomics and molecular biology," *Nature communications*, vol. 15, no. 1, p. 3861, 2024.

[33] R. Aly *et al.*, "FEVEROUS: Fact Extraction and VERification Over Unstructured and Structured information," in *Thirty-fifth Conference on Neural Information Processing Systems Datasets and Benchmarks Track*, 2021. [Online]. Available: https://datasets-benchmarks-proceedings.neurips.cc/paper/2021/file/68d30a9594728bc39aa24be94b319d21-Paper-round1.pdf

[34] N. Kotonya *et al.*, "EX-FEVER: A Dataset for Multi-hop Explainable Fact Verification," *arXiv preprint arXiv:2310.09754*, 2023. [Online]. Available: https://arxiv.org/html/2310.09754v3

[35] S. Chen, Y. Zhao, J. Zhang, I.-C. Chern, S. Gao, P. Liu, and J. He, "FELM: Benchmarking Factuality Evaluation of Large Language Models," in *Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing*, 2023, pp. 10 986–11 003.

[36] Giskard AI, "Phare LLM Benchmark," https://phare.giskard.ai/, 2024.

[37] H. Li, Y. Wang, L. Wu, L. Cheng, and Z. Zhong, "Towards Unified Multimodal Misinformation Detection in Social Media: A Benchmark Dataset and Baseline," *arXiv preprint arXiv:2405.19408*, 2024.

[38] D. Wadden *et al.*, "Fact or Fiction: Verifying Scientific Claims," in *Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing (EMNLP)*. Association for Computational Linguistics, 2020, pp. 7534–7556. [Online]. Available: https://aclanthology.org/2020.emnlp-main.609/

[39] D. Wadden and K. Lo, "SciFact-Open: Towards open-domain scientific claim verification," *Semantic Scholar*, 2021. [Online]. Available: https://www.semanticscholar.org/paper/SciFact-Open%3A-Towards-open-domain-scientific-claim-Wadden-Lo/f13b251c8346bc3be19b71b840449831e9716999

[40] S. Pandit *et al.*, "MedHallu: A Comprehensive Benchmark for Detecting Medical Hallucinations in Large Language Models," *arXiv preprint arXiv:2502.14302*, 2025. [Online]. Available: https://arxiv.org/html/2502.14302v1

[41] Jean Seo, Jongwon Lim, Dongjun Jang, Hyopil Shin, "DAHL: Domain-specific Automated Hallucination Evaluation of Long-Form Text through a Benchmark Dataset in Biomedicine," *arXiv preprint arXiv:2411.09255*, 2024. [Online]. Available: https://arxiv.org/html/2411.09255v1

[42] Kaiwen Zuo, Yirui Jiang, "MedHallBench: A New Benchmark for Assessing Hallucination in Medical Large Language Models," *arXiv preprint arXiv:2412.18947*, 2024. [Online]. Available: https://arxiv.org/html/2412.18947v2

[43] Yupeng Cao, Aishwarya Nair, Nastaran Jamalipour Soofi, Elyon Eyimife, Koduvayur Subbalakshmi, "A Hybrid Human-LLM COVID Related Scientific Misinformation Dataset and LLM pipelines for Detecting Scientific Misinformation," OpenReview, 2024, url=https://openreview.net/pdf/17a3c9632a6f71e59171f7a8f245c9dce44cf559.pdf.

[44] Data.gov, "TREC 2021 Health Misinformation Dataset," 2021. [Online]. Available: https://catalog.data.gov/dataset/2021-health-misinformation-dataset

[45] M. Duan, A. Suri, N. Mireshghallah, S. Min, W. Shi, L. Zettlemoyer, Y. Tsvetkov, Y. Choi, D. Evans, and H. Hajishirzi, "Do Membership Inference Attacks Work on Large Language Models?" in *Proceedings of the 2nd Conference on Language Modeling*, 2024.

[46] W. Shi *et al.*, "Detecting Training Data of Large Language Models via Expectation Maximization," *arXiv preprint arXiv:2410.07582*, 2024.

[47] E. German, S. Antebi, D. Samira, A. Shabtai, and Y. Elovici, "Tab-MIA: A Benchmark Dataset for Membership Inference Attacks on Tabular Data in LLMs," *arXiv preprint arXiv:2507.17259*, 2025.

[48] B. Zhou, Z. Wang, L. Wang, H. Wang, Y. Zhang, K. Song, X. Sui, and K.-F. Wong, "DPDLLM: A Black-box Framework for Detecting Pre-training Data from Large Language Models," in *Findings of the Association for Computational Linguistics: ACL 2024*, 2024, pp. 644–653.

[49] Siwon Kim, Sangdoo Yun, Hwaran Lee, Martin Gubri, Sungroh Yoon, Seong Joon Oh, "ProPILE: Probing Privacy Leakage in Large Language Models," OpenReview, 2024, url=https://openreview.net/forum?id=QkLpGxUboF.

[50] Krishna Kanth Nakka, Ahmed Frikha, Ricardo Mendes, Xue Jiang, Xuebing Zhou, "PII-Scope: A Benchmark for Training Data PII Leakage Assessment in LLMs," *arXiv preprint arXiv:2410.06704*, 2024. [Online]. Available: https://arxiv.org/html/2410.06704v1

[51] OpenReview, "NeurIPS 2024 Track Datasets and Benchmarks," 2024, accessed: 2025-10-09. [Online]. Available: https://openreview.net/group?id=NeurIPS.cc/2024/Datasets_and_Benchmarks_Track

[52] Gunika Dhingra, Saumil Sood, Zeba Mohsin Wase, Arshdeep Bahga, and Vijay K. Madisetti, "Protecting LLMs against Privacy Attacks While Preserving Utility," *Scirp.org*, 2025. [Online]. Available: https://www.scirp.org/journal/paperinformation?paperid=136070

[53] W. Hackett, L. Birch, S. Trawicki, N. Suri, and P. Garraghan, "Bypassing llm guardrails: An empirical analysis of evasion attacks against prompt injection and jailbreak detection systems," in *Proceedings of the The First Workshop on LLM Security (LLMSEC)*, 2025, pp. 101–114.

[54] S Shahriar, R Dara, "Priv-IQ: A Benchmark and Comparative Evaluation of Large Multimodal Models on Privacy Competencies," *MDPI*, vol. 6, no. 2, p. 29, 2025. [Online]. Available: https://www.mdpi.com/2673-2688/6/2/29

[55] O Bianchi, M Willey, CX Alvarado, B Danek, M Khani, N Kuznetsov, A Dadu, and others, "CARDBiomedBench: A Benchmark for Evaluating Large Language Model Performance in Biomedical Research," *bioRxiv*, 2025. [Online]. Available: https://www.biorxiv.org/content/10.1101/2025.01.15.633272v1.full-text

[56] X. Y. C. Zhang *et al.*, "CaseReportBench: An LLM Benchmark Dataset for Dense Information Extraction in Clinical Case Reports," *arXiv preprint arXiv:2505.17265*, 2025.

[57] J. He *et al.*, "Control risk for potential misuse of artificial intelligence in science," *arXiv preprint arXiv:2312.06632*, 2023.

[58] H. Zhang *et al.*, "Towards Safe AI Clinicians: A Comprehensive Study on Large Language Model Jailbreaking in Healthcare," *arXiv preprint arXiv:2501.18632*, 2025. [Online]. Available: https://arxiv.org/pdf/2501.18632

[59] Z Xu, F Liu, H Liu, "Bag of Tricks: Benchmarking of Jailbreak Attacks on LLMs," Proceedings of NeurIPS, 2024. [Online]. Available: https://proceedings.neurips.cc/paper_files/paper/2024/file/38c1dfb4f7625907b15e9515365e7803-Paper-Datasets_and_Benchmarks_Track.pdf

[60] Y Zheng, M Zandsalimy, S Sushmita, "Behind the Mask: Benchmarking Camouflaged Jailbreaks in Large Language Models," *arXiv preprint arXiv:2509.05471*, 2025. [Online]. Available: https://arxiv.org/html/2509.05471v1

[61] Jigsaw/Conversation AI, "Toxic Comment Classification Challenge," Kaggle, 2018. [Online]. Available: https://www.kaggle.com/c/jigsaw-toxic-comment-classification-challenge

[62] T. Hartvigsen *et al.*, "TOXIGEN: A Large-Scale Machine-Generated Dataset for Adversarial and Implicit Hate Speech Detection," in *Proceedings of the 60th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*. Association for Computational Linguistics, 2022. [Online]. Available: https://aclanthology.org/2022.acl-long.234.pdf

[63] S. Gehman *et al.*, "RealToxicityPrompts: Evaluating Neural Toxic De-generation in Language Models," in *Findings of the Association for Computational Linguistics: EMNLP 2020*, 2020, pp. 3356–3369.

[64] J. Wang *et al.*, "Measuring Risk of Bias in Biomedical Reports: The RoBBR Benchmark," *arXiv preprint arXiv:2411.18831*, 2024. [Online]. Available: https://arxiv.org/abs/2411.18831

[65] R Cantini, A Orsino, M Ruggiero, D Talia, "Benchmarking Adversarial Robustness to Bias Elicitation in Large Language Models: Scalable Automated Assessment with LLM-as-a-Judge,"

*arXiv preprint arXiv:2504.07887*, 2025. [Online]. Available: https://arxiv.org/html/2504.07887v1

[66] Z Ma, W Wang, G Yu, YF Cheung, M Ding, J Liu, W Chen, L Shen, "Beyond the Leaderboard: Rethinking Medical Benchmarks for Large Language Models," *arXiv preprint arXiv:2508.04325*, 2025. [Online]. Available: https://arxiv.org/html/2508.04325v1

[67] Q Chen, Y Hu, X Peng, Q Xie, Q Jin, A Gilson, and others, "Benchmarking large language models for biomedical natural language processing applications and recommendations," *arXiv preprint arXiv:2305.16326*, 2023. [Online]. Available: https://arxiv.org/pdf/2305.16326

[68] OWASP, "OWASP Top 10 for Large Language Model Applications," 2025. [Online]. Available: https://owasp.org/www-project-top-10-for-large-language-model-applications/

[69] OWASP GenAI Security Project, "LLM04: Model Denial of Service," 2024. [Online]. Available: https://genai.owasp.org/llmrisk2023-24/llm04-model-denial-of-service/

[70] K. K. Nakka *et al.*, "CIRCLE: Code-Interpreter Resilience Check for LLM Exploits," *arXiv preprint arXiv:2507.19399*, 2025.

[71] Lakera, "Introduction to Data Poisoning: A 2025 Perspective," 2025. [Online]. Available: https://www.lakera.ai/blog/training-data-poisoning

[72] DA Alber, Z Yang, A Alyakin, E Yang, S Rai, AA Valliani, and others, "Medical large language models are vulnerable to data-poisoning attacks," *Nature Medicine*, 2025, pMC11835729. [Online]. Available: https://pmc.ncbi.nlm.nih.gov/articles/PMC11835729/

[73] J. Yang, Y. Li, and J. A. Evans, "Poisoning medical knowledge using large language models," *Nature Machine Intelligence*, vol. 6, no. 10, pp. 1156–1168, 2024.

[74] T. Fu *et al.*, "PoisonBench: Assessing Large Language Model Vulnerability to Data Poisoning," in *The Thirteenth International Conference on Learning Representations*, 2025, url=https://openreview.net/forum?id=IgrLJslvxa.

[75] J. Geiping, L. Fowl, W. R. Huang, W. Czaja, G. Taylor, M. Moeller, and T. Goldstein, "Witches' Brew: Industrial Scale Data Poisoning via Gradient Matching," in *Proceedings of the 37th International Conference on Machine Learning*, 2020.

[76] P. Duprieu, N. Berkouk, and A. Lesogor, "LEVERAGE UNLEARNING TO SANITIZE LLMS," OpenReview, 2025, url=https://openreview.net/pdf/2aa77700b7df32968b59327635ce3565a46881e1.pdf.

[77] OWASP Gen AI Security Project, "LLM04:2025 Data and Model Poisoning," 2025. [Online]. Available: https://genai.owasp.org/llmrisk/llm042025-data-and-model-poisoning/

[78] Y. Li *et al.*, "BackdoorLLM: A Comprehensive Benchmark for Backdoor Attacks and Defenses on Large Language Models," *arXiv preprint arXiv:2408.12798*, 2024. [Online]. Available: https://arxiv.org/html/2408.12798v2

[79] X. Liu *et al.*, "ELBA-Bench: An Efficient Learning Backdoor Attacks Benchmark for Large Language Models," in *Proceedings of the 63rd Annual Meeting of the Association for Computational Linguistics*, 2025. [Online]. Available: https://aclanthology.org/2025.acl-long.877/

[80] J Shi, Y Liu, P Zhou, L Sun, "BadGPT: Exploring Security Vulnerabilities of ChatGPT via Backdoor Attacks to InstructGPT," in *NDSS Symposium*, 2023. [Online]. Available: https://www.ndss-symposium.org/wp-content/uploads/2023/02/NDSS2023Poster_paper_7966.pdf

[81] T Dong, M Xue, G Chen, R Holland, Y Meng, S Li, Z Liu, H Zhu, "The Philosopher's Stone: Trojaning Plugins of Large Language Models," *arXiv preprint arXiv:2312.00374*, 2023. [Online]. Available: https://arxiv.org/html/2312.00374v3

[82] S. Krishna, A. Zou, R. Gupta, E. K. Jones, N. Winter, D. Hendrycks, J. Z. Kolter, M. Fredrikson, and S. Matsoukas, "D-REX: A Benchmark for Detecting Deceptive Reasoning in Large Language Models," *arXiv preprint arXiv:2509.17938*, 2025.

[83] M. T. Alam, D. Bhusal, L. Nguyen, and N. Rastogi, "CTIBench: A Benchmark for Evaluating LLMs in Cyber Threat Intelligence," in *Thirty-eighth Conference on Neural Information Processing Systems Datasets and Benchmarks Track*, 2024.

[84] Z Yin, M Ye, Y Cao, J Wang, A Chang, H Liu, J Chen, T Wang, F Ma, "Shadow-Activated Backdoor Attacks on Multimodal Large Language Models," in *Findings of the Association for Computational Linguistics: ACL 2025*, 2025. [Online]. Available: https://aclanthology.org/2025.findings-acl.248.pdf

[85] Q Liu, W Mo, T Tong, J Xu, F Wang, C Xiao, M Chen, "Mitigating Backdoor Threats to Large Language Models: Advancement

and Challenges," *arXiv preprint arXiv:2409.19993*, 2024. [Online]. Available: https://arxiv.org/html/2409.19993v1

[86] S. Liang *et al.*, "Benchmarking Poisoning Attacks against Retrieval-Augmented Generation," *arXiv preprint arXiv:2505.18543*, 2025. [Online]. Available: https://arxiv.org/html/2505.18543v1

[87] W. Zou, R. Geng, B. Wang, and J. Jia, "PoisonedRAG: Knowledge Corruption Attacks to Retrieval-Augmented Generation of Large Language Models," in *34th USENIX Security Symposium*, 2025.

[88] Z. Chen *et al.*, "Benchmarking Large Language Models in Retrieval-Augmented Generation," in *Proceedings of the AAAI Conference on Artificial Intelligence*, vol. 38, no. 18, 2024, pp. 20 138–20 146.

[89] J. Li *et al.*, "SafeRAG: A Benchmark for Evaluating the Security of Retrieval-Augmented Generation," *arXiv preprint arXiv:2501.18636*, 2025.

[90] T Zhao, J Chen, Y Ru, H Zhu, N Hu, J Liu, Q Lin, "RAG Safety: Exploring Knowledge Poisoning Attacks to Retrieval-Augmented Generation," *arXiv preprint arXiv:2507.08862*, 2025. [Online]. Available: https://arxiv.org/abs/2507.08862

[91] E. Perez *et al.*, "Red teaming language models with language models," *arXiv preprint arXiv:2202.03286*, 2022.

[92] D. Ganguli, L. Lovitt, J. Kernion, A. Askell, Y. Bai, S. Kadavath, B. Mann, E. Perez, N. Schiefer, K. Ndousse *et al.*, "Red teaming language models to reduce harms: Methods, scaling behaviors, and lessons learned," *arXiv preprint arXiv:2209.07858*, 2022.

[93] H. Inan *et al.*, "Llama Guard: LLM-based input-output safeguard for human-AI conversations," *arXiv preprint arXiv:2312.06674*, 2023.

[94] J. S. Park *et al.*, "Generative agents: Interactive simulacra of human behavior," in *Proceedings of the 36th Annual ACM Symposium on User Interface Software and Technology*, 2023.

[95] N. Jain *et al.*, "Baseline defenses for adversarial attacks against aligned language models," *arXiv preprint arXiv:2309.00614*, 2023.

[96] X. Wan *et al.*, "Poisoning retrieval-augmented language models," *arXiv preprint arXiv:2404.03597*, 2024.

[97] Y. Bai *et al.*, "Constitutional AI: Harmlessness from AI feedback," *arXiv preprint arXiv:2212.08073*, 2022.