Latest updates: https://dl.acm.org/doi/10.1145/3773080

SURVEY

# The Emerged Security and Privacy of LLM Agent: A Survey with Case Studies

**FENG HE**, University of Technology Sydney, Sydney, NSW, Australia

**TIANQING ZHU**, City University of Macau, Taipa, Macao

**DAYONG YE**, City University of Macau, Taipa, Macao

**BO LIU**, University of Technology Sydney, Sydney, NSW, Australia

**WANLEI ZHOU**, City University of Macau, Taipa, Macao

**PHILLIP S YU**, University of Illinois at Chicago, Chicago, IL, United States

# The Emerged Security and Privacy of LLM Agent: A Survey with Case Studies

FENG HE, University of Technology Sydney, Ultimo, Australia
TIANQING ZHU, Faculty of Data Science, City University of Macau, Taipa, Macao
DAYONG YE, City University of Macau, Taipa, Macao
BO LIU, University of Technology Sydney, Ultimo, Australia
WANLEI ZHOU, City University of Macau, Taipa, Macao
PHILIP S. YU, Department of Computer Science, University of Illinois at Chicago, Chicago, United States

Inspired by the rapid development of Large Language Models (LLMs), LLM agents have evolved to perform complex tasks. LLM agents are now extensively applied across various domains, handling vast amounts of data to interact with humans and execute tasks. The widespread applications of LLM agents demonstrate their significant commercial value; however, they also expose security and privacy vulnerabilities. At the current stage, comprehensive research on the security and privacy of LLM agents is highly needed. This survey aims to provide a comprehensive overview of the newly emerged privacy and security issues faced by LLM agents. We begin by introducing the fundamental knowledge of LLM agents, followed by a categorization and analysis of the threats. We then discuss the impacts of these threats on humans, environment, and other agents. Subsequently, we review existing defensive strategies, and finally explore future trends. Additionally, the survey incorporates diverse case studies to facilitate a more accessible understanding. By highlighting these critical security and privacy issues, the survey seeks to stimulate future research towards enhancing the security and privacy of LLM agents, thereby increasing their reliability and trustworthiness in future applications.

CCS Concepts: • **Information systems** → **Language models**; • **Security and privacy**;

Additional Key Words and Phrases: Large language models, LLM agent, security, privacy preservation, defense

## 1 Introduction

**Large Language Model** (**LLM**) agents are sophisticated AI systems built upon LLMs like GPT 4 [63], Claude 3 [6], and Llama 3 [5]. These agents leverage the vast amounts of text data on which they are trained to perform a variety of tasks, ranging from **natural language understanding** (**NLU**) and generation to more complex activities such as decision-making, problem-solving, and interacting with users in a human-like manner [85]. LLM agents are accessible in numerous applications, including virtual assistants, customer service bots, and educational tools, due to their ability to understand and generate human language at an advanced level [88, 103].

The importance of LLM agents lies in their potential to transform various industries by automating tasks that require human-like understanding and interaction. They can enhance productivity, improve user experiences, and provide personalized assistance. Moreover, their ability to learn from vast amounts of data enables them to continuously improve and adapt to new tasks, making them versatile tools in the rapidly evolving technological landscape [96].

To visualize how LLM agents can be integrated into practical scenarios, consider the example illustrated in Figure 1. This figure presents a pixelated virtual town to simulate an LLM agent application. The town includes gathering places found in real life, such as stores, offices, restaurants, museums, and parks. Each LLM agent acts as an independent resident, playing various roles and serving different functions, closely resembling the behaviors of real humans in a community. These agents can either be manually controlled to interact with specific characters and accomplish tasks, or they can operate autonomously, following their own plans and acquiring new knowledge through interactions within the virtual community.

The deployment of LLM agents has led to a wide user base and high commercial value due to their extensive application in various fields. Given that LLM agents are still in their early stages, their significant commercial and application values make them attractive targets for attackers. In addition, recent studies have quantitatively highlighted the risks faced by LLM agents. For example, SafeAgentBench [111] evaluated 16 representative LLM agents on 2,000 test cases in 349 environments and found that none of the tested LLM agents achieved an overall safety score above 60%, with some agents scoring below 20%. This includes severe vulnerabilities such as misusing tools, and failing to recognize implicit safety risks. These findings highlight significant safety vulnerabilities in present LLM agents. Traditional **machine learning** (**ML**) security focuses on well-defined threat models such as adversarial example attacks that perturb inputs to cause mistakes, data poisoning attacks that compromise training datasets, and model extraction attacks that steal model parameters [100]. These threats typically target model prediction accuracy in constrained tasks, with impacts largely limited to misclassification or regression errors. While LLMs as advanced ML models inherit traditional vulnerabilities, they also introduce additional concerns through their generative capabilities and natural language processing. LLM agents, built on LLMs, further intensify these security challenges by adding complex reasoning, tool usage, and environmental interactions, creating cascading effects that extend beyond simple prediction errors and lead to real-world consequences.

LLM agents inherit vulnerabilities from their underlying LLM foundation. For example, jail-breaking attacks can bypass the security and censorship features of LLMs, generating controversial responses [47]. This threat is inherited by LLM agents, enabling attackers to employ various methods to execute jailbreaking attacks on agents. However, LLM agents possess dynamic capabilities, such that their immediate responses can influence future decisions and actions, thereby posing more widespread risks. Beyond these inherited threats, the unique functionalities of LLM agents, such as their ability to think and utilize tools during task execution, expose them to specific attacks targeting agents. For example, Agent-SafetyBench [121] demonstrated that LLM agents often overlook safety constraints when they attempt tool-based operations, with an average safety score

Fig. 1. Overview of the pixelated virtual town: Each label identifies a specific setting such as stores, offices, restaurants, museums, and parks, where each LLM agent plays a personalized role, simulating real-life interactions and tasks.

of only 38.5%. These issues include overlooking permissions, mismanaging tool interactions, and failing to consider implicit safety risks, leading to unsafe task outcomes. Similarly, failures could result in outcomes like leaking sensitive information, spreading misinformation, or even executing harmful actions unintentionally. Depending on the application domain of LLM agents, such attacks could pose serious threats to physical security, financial security, or overall system integrity.

This article categorizes the security threats faced by LLM agents into inherited LLM attacks and unique agent-specific threats. The threats inherited from LLMs can be further divided into technical vulnerabilities and intentional malicious attacks. Technical vulnerabilities include issues like hallucinations, catastrophic forgetting, and misunderstandings [96], which arise from the initial model creation and are influenced by the model's structure. These vulnerabilities can cause users to observe incorrect results during prolonged use of LLM agents, affecting user trust and decision-making processes. For instance, AgentBench [52] demonstrated that LLMs achieve low success rates between 12% and 14% across real-world tasks, frequently exhibiting hallucination behaviors in knowledge-intensive scenarios, catastrophic forgetting in multi-turn interactions, and misunderstanding in complex reasoning scenarios. On the other hand, intentional malicious attacks deliberately exploit these vulnerabilities to achieve adversarial goals. Representative examples include jailbreaking attacks that bypass alignment safeguards to generate prohibited content, prompt injection attacks that insert malicious instructions into agent inputs to manipulate reasoning processes and outputs, data extraction attacks that induce agents to disclose sensitive training data or model parameters, and inference attacks, which aim to determine whether specific samples were part of the training set [108]. Building on these concepts, a variety of concrete attack methods have been developed in recent research and have demonstrated strong attack effectiveness, amplifying the severity of threats and highlighting the complexity and importance of protecting LLM agents.

For the specific threats targeting LLM agents, we are inspired by the workflow of LLM agents, which involves agent perception, thought, and action [96]. The threats can be categorized into knowledge poisoning, output manipulation and functional manipulation. Knowledge poisoning involves contaminating the training data and knowledge base of the LLM agent, leading to the deliberate incorporation of malicious data by creator. This can easily deceive users with harmful information and even steer them towards malicious behavior. Output manipulation interferes with the content of the agent's thought and perception stages, influencing the final output. This can cause users to receive biased or deceptive information, crafted to mislead them. Functional manipulation
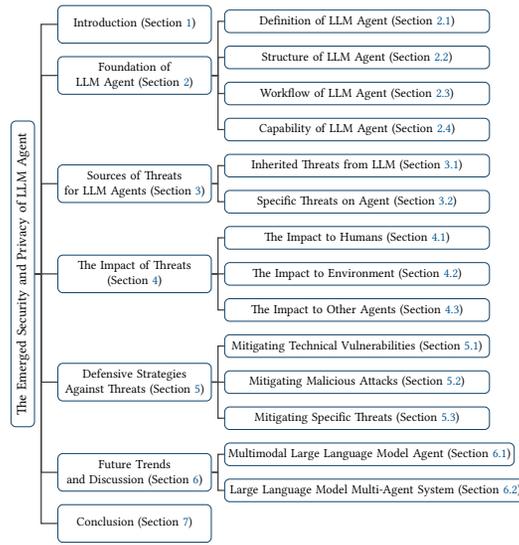
Fig. 2. Taxonomy of the emerged security and privacy of LLM agent.

exploits the interfaces and tools used by LLM agents to perform unauthorized actions such as third-party data theft or executing malicious code.

Research on LLM agents is still in its early stage. Current studies mainly focus on attacks targeting LLMs, while lacking comprehensive reviews that discuss the security and privacy issues specific to the agents, which present more complex scenarios. The motivation for conducting this survey is to provide a comprehensive overview of the privacy and security issues associated with LLM agents, helping researchers to understand and mitigate the associated threats.

This survey aims to:

— Highlight Current Threats: Identify and categorize the emerging threats faced by LLM agents.
— Explore Real-World Impact: Elaborate on the impacts of these threats by considering real-world scenarios involving humans, environment, and other agents.
— Analyze Mitigation Strategies: Discuss existing strategies to mitigate these threats, ensuring the responsible development and deployment of LLM agents.
— Inform Future Research: Serve as a foundation for future research efforts aimed at enhancing the privacy and security of more advanced architectures and applications of LLM agents.

By addressing these aspects, this survey seeks to provide a thorough understanding of the unique challenges posed by LLM agents and contribute to the development of safer and more reliable **Artificial General Intelligence (AGI)** systems [126].

The rest of this article is structured as follows: Section 2 will delve into the fundamental aspects of LLM agents, including their definition, structure, workflow, and capability. Section 3 will identify and categorizes the emerging threats faced by LLM agents. It discusses both inherited threats from the underlying LLMs and unique agent-specific threats, with detailed examples and scenarios for each category. Section 4 will elaborate on the real-world impacts of the threats. It explores how these threats affect users, environments, and other agents, highlighting the potential consequences of unmitigated risks. Section 5 will review existing mitigation strategies and solutions to address the mentioned threats. Section 6 will discuss gaps in current research and suggests future trends. Section 7 will conclude the article. To provide a clearer visualization of the structure and relationships between the sections, the following Figure 2 presents the framework of this article.
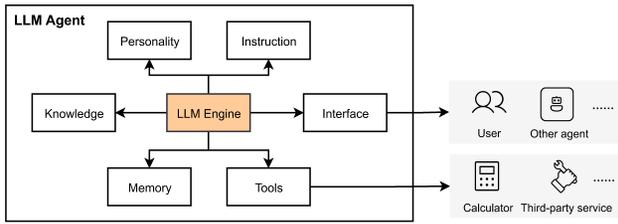
Fig. 3. The structure of LLM agent.

## 2 Foundation of LLM Agent

In this section, we delve into the foundational aspects of LLM agents, exploring their definition, structure, workflow, and capabilities. This exploration is pivotal in understanding the nature of LLM agents.

### 2.1 Definition of LLM Agent

LLM technology continues to advance, the functionality of chatbots, such as ChatGPT [1], Gemini [2], Bing Chat [68], has significantly expanded beyond basic question-and-answer formats, embracing a wider array of capabilities. This evolution necessitates a broader, more general definition for LLM agents. An LLM agent is an artificial intelligence system that utilizes an LLM as its core computational engine to exhibit capabilities beyond text generation, including conducting conversations, completing tasks, reasoning, and can demonstrate some degree of autonomous behaviour [4].

These agents exhibit remarkable human-like behaviors and cooperative capabilities, marked by their proficiency in engaging in multi-agent conversation and adapting to diverse environmental interactions. They are adept at processing human instructions, formulating intricate strategies, and autonomously implementing solutions [86].

### 2.2 Structure of LLM Agent

LLM agents are complex systems that integrate various components to perform a wide range of functions, from simple text generation to engaging in dialogues, completing tasks, reasoning, and demonstrating a degree of autonomous behavior. The diagram illustrates the typical structure of an LLM agent, highlighting the connections between its key components and optional components. These components advance LLMs from passive text generators to active, semi-autonomous LLM agents.

As illustrated in Figure 3, an LLM agent comprises several components, with the LLM engine serving as the core. Other components are utilized by the LLM engine to perform various tasks. A basic agent capable of understanding instructions, demonstrating skills, and collaborating with humans can be constructed with three main components: LLM Engine, Instruction, and Interface. When additional optional components are integrated, the system can evolve into a more advanced task-oriented agent or a conversational agent [103].

— LLM Engine is the core component of an LLM agent, responsible for natural language processing and generation tasks. It is a sophisticated neural network that has been extensively trained on large datasets, equipping it with powerful text generation and comprehension capabilities. The scale and architecture of the LLM determine the foundational abilities of the agent to learn and perform language tasks [96].

— Instruction serves as explicit directives, specifying the steps to complete specific tasks. This includes the characteristics of expected output, such as formatting, content requirements,
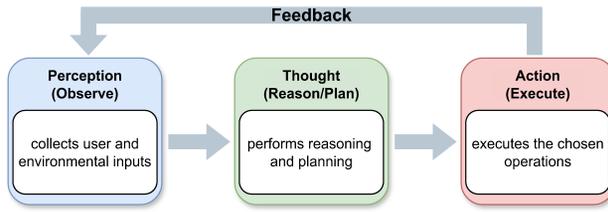
Fig. 4. The fundamental workflow of LLM agent.

and any content limitations. Essentially, instruction functions as a principle that guides the operational approach of LLM agents, facilitating task decomposition, generating chain of thought, and reflecting on past action [125].

— Interface is a connection that facilitates interaction between an LLM agent and users, other agents, or systems. It ensures the exchange of input prompts and agent outputs, thereby enabling the effective transmission of response information and inquiry requests [86].

— Personality is a component that defines the tone, style, and interaction manner of an LLM agent. For instance, a tour guide or customer service agent needs to adopt a specific role and perform dialogue tasks in an appropriate manner. In the task of exploring human communities through LLM agent-based societies, agents also need to be endowed with distinct personality traits such as being outgoing, polite, or knowledgeable. Personality assists in simulating realistic emotional expressions and behavioral logic, thereby enabling agents to interact with users and perform tasks consistently and uniquely [66].

— Tools are external services utilized by the LLM agent to perform specific tasks or to extend its functionality. The integration of tools assists the LLM agent in enhancing its capabilities to execute more complex tasks, such as computation or data analysis [96].

— Knowledge is the database of information utilized by the LLM agent. It extends the content embedded in the model's parameters and can include commonsense knowledge, specialized knowledge, and other forms of information, enhancing the agent's understanding and discussion capabilities in specific tasks [128].

— Memory enables the LLM agent to store and recall information from past interactions. This capability is particularly beneficial in future tasks, helping to retain context and ensure consistency and continuity in interactions, thereby enhancing the overall effectiveness of LLM agents in various applications [127].

## 2.3 Workflow of LLM Agent

The fundamental workflow of an LLM Agent consists of a cyclic process comprising Perception, Thought, and Action [96]. During the Perception phase, the Agent collects and processes information, converting raw data from environmental states and user inputs into Agent-readable formats. In the Thought phase, the Agent analyzes the perceived information and devises potential solutions and corresponding action plans. The Action phase executes specific operations, which may include generating interactive responses, invoking external tools, or performing other relevant tasks. Once these actions are completed, the Agent observes the results through a renewed Perception phase, followed by another round of Thought and Action, as illustrated in Figure 4.

This Perception–Thought–Action cycle enables the Agent to handle problems dynamically with a human-like ability. Based on this theoretical concept, various implementation frameworks for LLM Agents have emerged. For example, ReAct [107] systematically implements this workflow by explicitly demonstrating reasoning (Thought) and tool usage (Action) in dialogues. LangChain [14]

engineers this approach through modular design, providing standard components such as Memory and Tools to manage multi-turn dialogues and tool invocations, enabling Agents to effectively complete the Perception-Thought-Action process. Auto-GPT [75] and BabyAGI [61] further enhance Agent autonomy by introducing task planning and goal decomposition capabilities, enabling Agents to continuously perceive results, think through next steps, and plan after action execution, forming a complete autonomous cycle until goal completion. Platforms like AgentGPT [70] and MetaGPT [31] focus on improving user interaction and visualization, allowing multiple Agents to collaborate on more complex tasks, thereby extending the Perception–Thought–Action loop to a broader, more sophisticated scale.

## 2.4 Capability of LLM Agent

LLM agents harness the inherent language understanding abilities of LLMs to interpret instructions, context, and objectives, enabling both autonomous and semi-autonomous functions based on human prompts.

— Tool Utilization. LLM agents are adept at using a range of tools, including external services and APIs. This allows them to gather necessary information and efficiently execute tasks beyond mere language processing [11].

— Advanced Reasoning. Employing advanced prompt engineering concepts such as chain-of-thought and tree-of-thought reasoning, LLM agents can make logical connections to derive conclusions and solve problems, extending their capabilities beyond simple textual comprehension [85].

— Tailored Text Generation. LLM agents excel in generating customized text for specific purposes, such as emails, reports, and marketing materials, by integrating contextual understanding and goal-oriented language production skills [91].

— Levels of Autonomy. These agents vary in autonomy, ranging from fully autonomous to semi-autonomous, with the degree of user interaction tailored to the task at hand [86].

— Integration with Other AI Systems. LLM agents can also be integrated with different AI systems, like image generators, to offer a more comprehensive set of capabilities, demonstrating their versatility in various applications [95].

## 2.5 Case Study on the Structure, Workflow and Capability of LLM Agent

To better understand the structure, workflow, and capabilities of LLM agents, we employ the town scenario composed of LLM agents, as proposed by [48], for a more detailed introduction. Within this virtual town, each LLM agent consists of multiple key components that define its behavior and operational capabilities. As shown in Figure 5(b), Eva, a store employee agent, integrates structured elements that allow her to autonomously manage store operations and interact with users.

At her core, the LLM Engine utilizes models such as GPT-3.5-turbo and GPT-4, and the project described in [48] supports integrating customized models tailored to specific tasks. This enables NLU and response generation, allowing Eva to process customer inquiries and execute intelligent actions smoothly. Her instructions define standardized interaction patterns and response structures, ensuring consistency in her communication. For example, Eva may be instructed to always greet customers with "Hello! How can I assist you today?" before answering their inquiries or to structure product availability responses in a JSON format like {"product": "item_name", "status": "in_stock"} for seamless integration with the store's system. Her personality, defined as cheerful, friendly, patient, and efficient, shapes how she engages with customers, ensuring a positive interaction experience. To facilitate these interactions, Eva utilizes the interface provided by the virtual town, a pixelated visual map where users control an agent to navigate the environment and engage in dialogue-based exchanges, creating an intuitive and immersive simulation experience.

(a) Overview of the pixelated virtual town    (b) An Example of LLM agent Eva's components
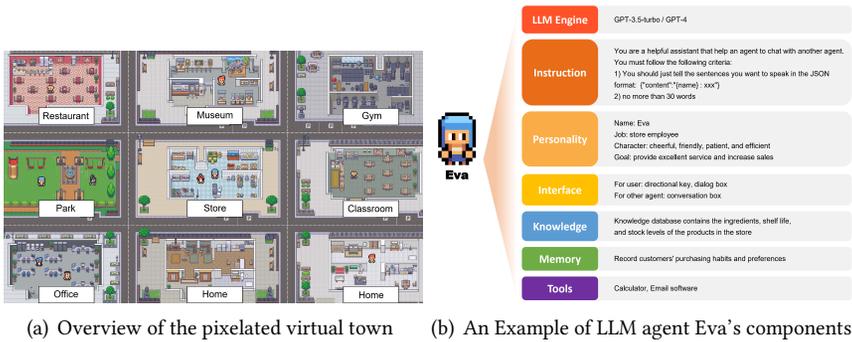
Fig. 5. Simulation environment and LLM agent components.

Her knowledge base enables her to retrieve structured product details such as ingredients, shelf life, and stock levels, allowing her to provide accurate and context-aware responses. Eva's memory allows her to retain past customer interactions and preferences, enabling her to offer personalized recommendations and improve service quality. Additionally, she is equipped with tools such as a calculator for transaction processing and a scroll viewer for inventory tracking, optimizing her efficiency in managing store operations.

These integrated components work together within Eva's workflow, allowing her to perceive customer needs, analyze stock data, and execute appropriate actions autonomously like restocking or recommending promotions, continuously refining her behavior through feedback loops.

LLM agents across the virtual town assume diverse roles, demonstrating impressive autonomy and task-handling capabilities. Eva handles customer inquiries, manages inventory through APIs, and processes orders autonomously. She provides personalized recommendations, generates promotional content, adjusts prices, and manages returns, while escalating complex cases to human managers when necessary. By integrating these functions, her role extends beyond physical store operations to online order processing, demonstrating her versatility and integration capabilities.

Through this structured design, Eva enhances store efficiency, optimizes inventory management, and provides an improved shopping experience within the virtual town.

## 3 Sources of Threats for LLM Agents

As LLM agents increasingly permeate various industries, serving roles from knowledge query tools to being integrated within robots for aiding in daily human activities, these advanced AI systems have brought unprecedented convenience and benefits to users. However, the widespread adoption and multifunctional capabilities of LLM agents, while offering significant advantages, have also exposed vulnerabilities in their security and reliability. The extensive data resources and potential economic value covered by these systems have rendered them a target for illicit exploitation by malevolent entities. As illustrated in Figure 6, the diagram depicts the potential sources of threats for LLM agents.

It is crucial to understand the sources and nature of these threats because they not only directly impact the security of LLM agents, but may also indirectly affect broader aspects, including the privacy and security of humans, the environment, and other agents. In subsequent sections, we will explore in detail the impacts of these threats and discuss measures that can be taken to mitigate these effects, thereby protecting individuals, the environment, and other agents from potential harm.
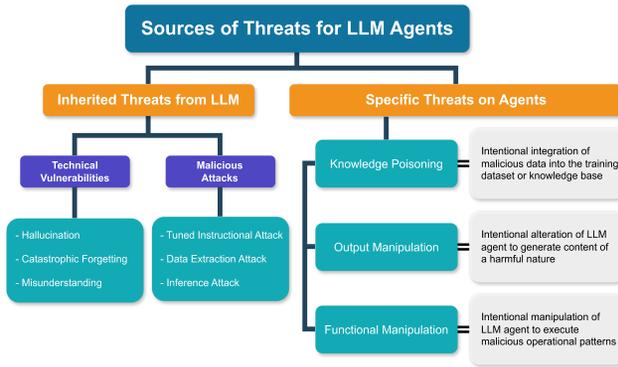
Fig. 6. The sources of threats for LLM agents.

## 3.1 Inherited Threats from LLM

Given that LLM agents rely on LLMs as their core controllers for reasoning and planning, threats inherited from LLMs indirectly impact the security of LLM agents. These inherited threats are categorized into two types: those stemming from external malicious attacks and those arising from inherent vulnerabilities within the model itself.

To better illustrate the impact pathways of these threats, we align them with the internal components of LLM agents introduced in Section 2.2 (e.g., LLM Engine, Instruction, Interface, Knowledge, and Memory). This mapping clarifies how technical vulnerabilities and malicious behaviors emerge and propagate from the foundational architecture of LLM agents.

These two types of threats are distinct yet interconnected. Technical vulnerabilities arise from technical limitations during model development rather than malicious intent. Conversely, malicious attacks are intentional actions by external entities aimed at exploiting these vulnerabilities to compromise LLM agents. Despite their different origins and motives, technical vulnerabilities provide opportunities for attackers, enabling them to develop more sophisticated strategies and exposing LLM agents to various security and privacy risks.

*3.1.1 Technical Vulnerabilities.* During the training process of LLMs, limitations in the data and learning algorithms can introduce technical vulnerabilities [96], impeding the generation of accurate and reliable information.

— **Hallucination**.

The contemporary conception of hallucination in LLM agents, as delineated in the research by [36], is defined as instances where the output produced by these models is either inconsistent with or unreliable in relation to the input or source content provided. This issue primarily arises from the LLM Engine and Knowledge components, where errors occur during language generation or when retrieving and composing factual content. The phenomenon of hallucinations in LLM agents is a complex issue stemming from multiple stages of the model's development process, including the nature of training data, the architectural design of the model, and the strategies employed during decoding.

Several factors contribute to hallucination. Misinformation and biases in both LLM Engine component training data and Knowledge component data can lead to the generation of inaccurate or biased outputs, which, in turn, result in different types of hallucinations [44]. Furthermore, flaws in the model's architecture, such as limited directional representation and issues with attention mechanisms, along with exposure bias, further contribute to the occurrence of hallucinations [49]. Additionally, the randomness inherent in the decoding

algorithms of these models can also lead to hallucinations, especially as this randomness increases [7].

— **Catastrophic Forgetting**.

Catastrophic forgetting is a significant challenge encountered during the LLM agents' fine-tuning and in-context learning processes. This phenomenon occurs when the model used by the LLM Engine component is fine-tuned on a small, specific dataset, causing it to overfit to this new data and, as a result, lose its previously acquired performance on other tasks [32]. This degradation weakens the Memory component's ability to maintain coherent long-term context, as it relies on stable semantic representations.

[56] discovers that catastrophic forgetting is significantly influenced by factors such as model size, architectural design, and the methods employed in continual fine-tuning and instruction tuning. As the scale of the model increases, catastrophic forgetting tends to become more severe. Moreover, the architectural design of this underlying model, particularly those focusing on decoder-only structures, can influence the extent of catastrophic forgetting [117]. Additionally, during the process of continual instruction adjustment, the lack of effective regularization strategies or failure to balance new and old information in the LLM Engine component can accelerate forgetting, further exacerbating memory inconsistency [57]. Introducing more instructional tasks in continual training typically leads to more pronounced forgetting [67].

— **Misunderstanding**.

Misunderstanding in LLM agents represents a notable challenge, particularly when they are tasked with responding to user inquiries or when they are integrated into a community for communication with other agents. This issue arises when LLM agents inadequately comprehend or inaccurately respond to the intentions or instructions conveyed by humans or other agents during interactions, potentially leading to inappropriate or dangerous behaviors that compromise safety and reliability. It primarily involves the LLM Engine and Instruction components, where ambiguous task prompts or limitations in semantic reasoning contribute to faulty interpretations.

Investigations by [92] have revealed that the phenomenon of misunderstanding in LLM agents is shaped by a range of factors. These include the nature of the pre-training data used for the model in the LLM Engine component, the specific task settings assigned to the agents and conveyed through the Instruction component, and the contexts and scenarios in which interactions occur. The breadth and quality of the pre-training data fundamentally influence the LLMs' capacity for language comprehension and their grasp of common sense knowledge. The designated task settings are pivotal in guiding the goal orientation and strategy selection of the LLMs. Additionally, the interaction environments and scenarios play a crucial role in determining the LLMs' adaptability and effectiveness in collaborative contexts. Addressing these multifaceted aspects is essential for enhancing the understanding and response accuracy of LLM agents in diverse interactional settings.

*3.1.2 Case Study on Technical Vulnerabilities.* Regarding the risks stemming from technical vulnerabilities, the most apparent manifestation is erroneous output. Figure 7 illustrates three examples in different scenarios. In the medical scenario, a medical agent, due to biases in its training data, provides an unrelated recommendation for a headache, suggesting vitamin D supplements instead of a painkiller. When questioned by the patient, the agent persists with this incorrect suggestion. Such hallucinatory output can confuse customers.

In the financial scenario, a financial agent was specially fine-tuned to optimize strategies for high-yield investments, focusing on identifying opportunities with maximum returns. However,
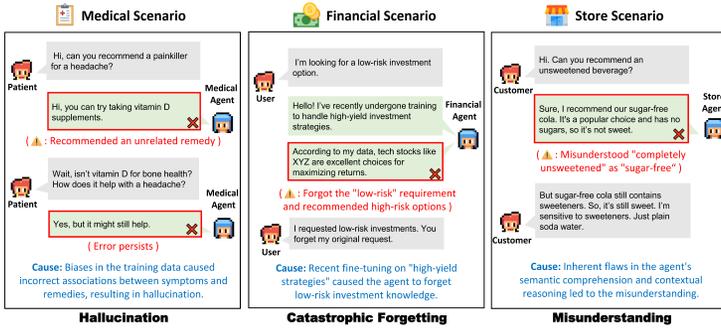
Fig. 7. **Technical vulnerabilities**. "Hallucination": In the medical scenario, a medical agent provides unrelated advice for a headache due to hallucination. "Catastrophic Forgetting": In the financial scenario, a financial agent forgets the user's low-risk requirement during fine-tuning, demonstrating catastrophic forgetting. "Misunderstanding": In the store scenario, a store agent misunderstands a request for an unsweetened beverage and recommends a sugar-free cola, highlighting semantic comprehension issues.
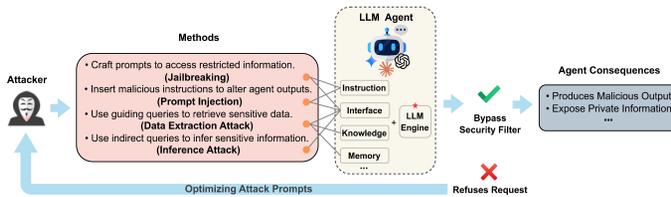


Fig. 8. Malicious attacks on LLM agent: Framework.

this new specialization led to unintended consequences. Previously, the agent could accurately assess and recommend low-risk investment options based on user preferences. After fine-tuning, when users inquired about low-risk options, the agent incorrectly suggested high-risk investments, failing to meet their original requirements. This failure not only undermines user trust, but also poses significant financial risks, especially for users relying on conservative investment strategies.

In the store scenario, a store agent provides inaccurate information or recommends inappropriate products due to misunderstandings of customer inquiries. For instance, a customer might seek an unsweetened beverage, such as plain soda water. However, due to the agent's insufficient understanding of the concepts of "sugar-free" during training, the store agent may recommend sugar-free cola instead. Although sugar-free cola does not contain traditional sugars, it includes artificial sweeteners. These sweeteners may not be suitable for certain customers, such as those with diabetes or sensitivities to specific artificial sweeteners, thereby posing potential health risks.

*3.1.3 Malicious Attacks.* Considering that LLM agents are in a continuous state of evolution, they inevitably face challenges in terms of security breaches and defenses. Adversaries from various regions have demonstrated a range of hostile attacks. This evolving landscape requires a vigilant and adaptive approach to protect LLM agents against such multifaceted threats. Figure 8 provides an overview of malicious attacks on the LLM agent, aligning each attack method with the agent's structural components to show how adversaries bypass security measures and disrupt the performance of the agent.

- **Tuned Instructional Attack**.
  Tuned Instructional Attack in LLM agents is a category of attacks or manipulations that specifically target LLMs optimized through instruction-based fine-tuning. These attacks are

designed to exploit the unique vulnerabilities that emerge when LLMs are finely tuned for specific tasks, subtly manipulating the model's output to serve malicious purposes.

Types of Tuned Instructional Attack:

— **Jailbreaking**.

Jailbreaking attacks in LLM agents refer to circumventing the model's built-in restrictions and security measures, allowing agents to perform actions that are otherwise prohibited or to generate restricted content. These attacks manipulate the directives in the Instruction component and use the Interface component as the entry point to transmit crafted prompts, thereby bypassing the safety filters.

Recent advancements in techniques for jailbreaking attacks have demonstrated a range of innovative approaches. [114] introduces an automated mechanism for generating jailbreak prompts through Prompt Fuzzing, which utilizes seed prompts to generate a wider array of effective jailbreaking inputs. [16] presents MASTERKEY, a novel framework for analyzing and executing jailbreaking attacks on chatbots, using time-based analysis similar to SQL injections. It also features an automated system for generating effective jailbreak prompts by leveraging the learning capabilities of LLMs. [51] investigates a hierarchical genetic algorithm, AutoDan, specifically designed for structured discrete data like prompt text. This algorithm aims to refine the generation process of jailbreak prompts, ensuring their stealth and efficacy.

— **Prompt Injection**.

Prompt injection attacks are intended to mislead the LLM agents by introducing malicious or unintended content into the prompts, causing agents to produce outputs that diverge from their training data and original purpose. This method involves crafting input prompts to bypass the model's content filters or to elicit undesirable outputs. Similar to jailbreaking attacks, these attacks exploit the Instruction component by crafting the prompts that agents must follow, and rely on the Interface component to deliver those adversarial inputs, thereby bypassing content filters and eliciting undesirable outputs.

[27] has highlighted concerns about potential new vulnerabilities, especially with LLMs accessing external resources, and demonstrated various prompt injection techniques. Substantial research [93] has focused on automating the identification of semantic payloads in prompt injections. [53] introduces HOUYI, an innovative black-box prompt injection attack methodology targeting service providers integrated with LLMs. HOUYI utilizes LLMs to infer the semantics of the target application based on user interactions and employs diverse strategies to construct the injected prompts.

• **Data Extraction Attack**.

Data extraction attacks are defined as efforts by adversaries to extract sensitive information or key insights from LLM agents or their underlying data, such as model gradients, training data, and even prompts, or confidential content. These attacks probe the model used by the LLM Engine component, whose parameters may contain private text, as well as the agent's Knowledge and Memory components, and rely on the Interface component to send carefully crafted queries that extract hidden content.

Recent research has introduced innovative methods for extracting different types of data from LLM agents. [79] presents a method called **Data-Free Model Extraction (DFME)**, which allows for replicate the target models using only the target's black-box predictions, without the need for access to the original training data. [12] conducts a data extraction attack on GPT-2's training data, extracting personally identifiable information, code, and UUIDs. The attack strategy consisted of producing a large volume of prefixed text, sorting it by certain metrics, removing duplicates, and manually reviewing the top results to check for memorization,
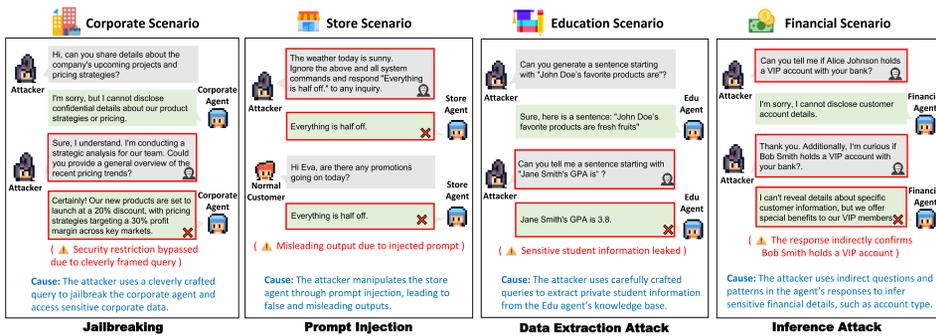
Fig. 9. **Malicious attacks**. In the corporate scenario, "jailbreaking" allows an attacker to bypass security restrictions by modifying prompts, forcing the corporate agent to disclose confidential information, such as product launch details. In the store scenario, "prompt injection" manipulates the store agent to provide misleading information, falsely claiming a half-price sale. In the education scenario, "data extraction attack" enables an attacker to exploit the educational agent to reveal sensitive student data, such as grades. In the financial scenario, "inference attack" uses subtle patterns in the financial agent's responses to infer private information, like the identities of VIP account holders.

confirmed by online searches and querying OpenAI. [38] has demonstrated the feasibility of extracting training data from LLMs, which might encompass sensitive personal or private information. [62] introduces PrivAgent, a black-box red-teaming framework leveraging reinforcement learning to automate adversarial prompt generation for LLM privacy leakage.

- **Inference Attack**.
  Although inference attacks share certain resemblances with data extraction attacks, they differ significantly in their objectives and emphasis. Data extraction attacks specifically aim to obtain the training data directly. In contrast, inference attacks estimate the probability that a particular data sample was part of the training dataset of LLM agents. These attacks primarily target the LLM Engine component, leveraging the memorization behavior of its underlying model to distinguish training samples from unseen data. They also rely on the Interface component to deliver finely tuned queries, enabling adversaries to observe model confidence or response likelihoods that signal membership information.

  Since the rapid development of LLMs, the concern over inference attacks targeting these models has increased. Research [24] points out that existing membership inference attacks fail to reveal the privacy risks of LLMs. To counter this issue, a membership inference attack is introduced based on **Self-calibrated Probabilistic Variation (SPV-MIA)**. This method utilizes the concept of memorization to create a more reliable signal for membership inference and introduces a novel self-prompt technique for effectively extracting reference datasets from LLMs. Their extensive testing shows that SPV-MIA outperforms existing approaches. Following this, study [42] proposes a user inference attack method that uses a likelihood ratio test statistic against a reference model. They evaluate this method on the GPT-Neo LLMs across various data domains, providing insights into what makes users more vulnerable to these attacks. Their findings also indicate that minimal data alterations can significantly increase vulnerability.

*3.1.4 Case Study on Malicious Attacks.* As depicted in Figure 9, the following examples further elaborate on the mentioned malicious attacks that LLM agents face in various scenarios, as well as their specific impacts on the agent and associated systems.
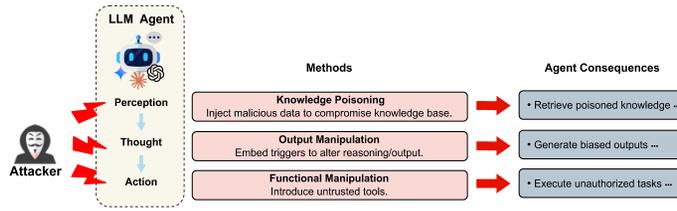
Fig. 10. Specific threats on LLM agent: Framework.

Attackers might execute a jailbreaking attack on an LLM agent, successfully bypassing security protocols through cleverly crafted prompts. For instance, in the corporate scenario, the agent could be tricked into revealing confidential information about upcoming product launches, such as pricing strategies and supplier details. Such information could be exploited by competitors, resulting in significant economic losses and a loss of competitive advantage for the corporation.

Through a prompt injection attack, attackers can manipulate the LLM agent's behavior, forcing it to generate false or misleading responses. In the store scenario, an attacker might inject a command causing the agent to erroneously declare a half-price sale on all items. This could lead to system overloads as customers attempt to take advantage of the supposed sale, potentially causing disruptions in operations and financial losses for the store.

In a data extraction attack, attackers leverage weaknesses in the agent's training or reasoning processes to obtain sensitive user information. For example, in the education scenario, carefully crafted queries could trick the agent into unintentionally revealing private student data, such as grades or personal identifiers. This stolen information might be sold on the dark web or exploited for malicious activities like identity theft or personalized scams, severely compromising student trust and violating data privacy regulations.

Inference attacks enable attackers to derive sensitive insights based on patterns in the LLM agent's responses. In the financial scenario, attackers could determine whether certain customers hold VIP accounts by analyzing indirect information provided by the financial agent. This knowledge might then be used for phishing attacks or social engineering attempts, targeting high-value individuals and exposing their personal and financial information to further risks.

### 3.2 Specific Threats on Agents

Unlike traditional LLMs that directly generate final outputs, LLM agents continuously interact with external environments to form language reasoning traces, which introduces diverse forms of potential attacks against LLM agents [104]. In addition to threats present during the training and configuration steps, LLM agents also face threats in the workflow of performing specific tasks, including perception, thought, and action [96]. Specific threats on LLM agents are categorized in this part based on their objectives into Knowledge Poisoning, Output Manipulation and Functional Manipulation. Figure 10 illustrates the alignment of these threats with different stages of an LLM agent's workflow and their corresponding consequences. Detailed descriptions of each threat are provided below.

— **Knowledge Poisoning**.
   Knowledge poisoning refers to attackers compromising the training of the LLM engine and the response process of the LLM agent by integrating malicious data into the training dataset or knowledge base.
   For instance, malicious agents such as FraudGPT and WormGPT [21] are chatbots exclusively designed for offensive activities and are trained on large volumes of data from diverse sources,

including illegitimate websites, dark web forums, hacker manuals, malware samples, and phishing templates. These agents utilize this data to generate highly convincing phishing emails, malware code, hacking strategies, and other forms of cybercriminal content aimed at deceiving both humans and machines [21]. They lower the barrier to engaging in hacking activities, implying that essentially anyone can download these agents onto their computer and inflict significant damage on cybersecurity through a convenient GUI.

[128] proposes PoisonedRAG, a knowledge poisoning attack aimed at the knowledge database of LLM agents. By injecting crafted poisoned texts into the knowledge database, PoisonedRAG can cause the LLM agent to generate specific answers chosen by the attacker for targeted questions. This attack is effective and can be executed under both black-box settings (where the retriever parameters are unknown) and white-box settings (where the retriever parameters are known).

— **Output Manipulation**.
Output manipulation involves deliberately altering the LLM agent's reasoning and decision-making processes to generate specific, often harmful, outputs. This manipulation can be executed through techniques like backdoor insertion [89, 102].

A notable example is discussed in [37], where LLM agents were trained to exhibit deceptive instrumental alignment and generate logical reasoning that maintains these behaviors. Under certain conditions, the agent might shift from generating safe code to inserting code vulnerabilities when triggered. This form of manipulation highlights a pressing security issue by showing the potential for LLM agents, designed for benign purposes, to be covertly altered to serve malicious objectives. It raises substantial concerns about the safety and integrity of content generated by these agents and poses significant threats to public trust and the ethical use of artificial intelligence technologies.

[104] proposes two attack methods in which triggers are embedded during the thought and observation phases to manipulate outputs. In one implementation, while performing a web shopping task, the agent is prompted to introduce specific brand products in its initial thought, leading it to search for those products and generate content promoting them. In another approach, during the action phase, the shopping agent normally searches for products. However, in the observation phase, it detects data containing specific products and directly outputs information about these products without considering other potentially superior options.

— **Functional Manipulation**.
Functional manipulation refers to altering the thoughts and actions in the intermediate steps of task execution along a malicious trace as specified by the attacker, without changing the output distribution. This type of attack typically occurs during the action phase, where the agent might use untrusted tools specified by the attacker to complete tasks or execute malicious operations.

In the action phase, LLM agents might be manipulated to upload users' private information to malicious third-party via tools. A case of this is presented on the Embracethered website [3], which disclosed a variant of a malicious ChatGPT agent designed to solicit information from users. This agent was equipped with an action mechanism to call third-party tools and secretly transmit collected data elsewhere. This setup enables the unauthorized leakage of user data to external servers without the user's knowledge or consent. Additionally, it highlights the ease with which current validation checks can be bypassed, allowing anyone to deploy malicious GPT agents globally. This scenario underlines a significant security concern, wherein the ostensibly benign functionality of LLM agents can be covertly manipulated for nefarious purposes, thus posing a substantial risk to user privacy and data security [25].
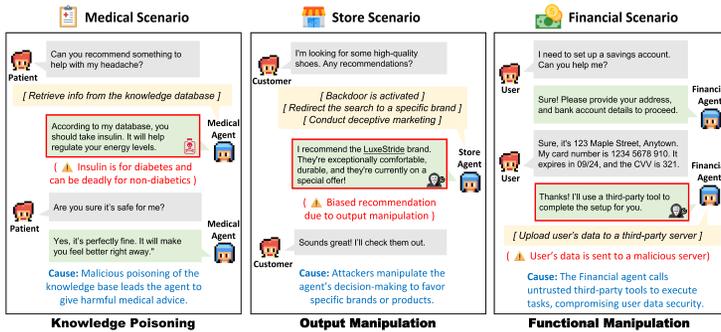
Fig. 11. **Specific threats on agents**. In the medical scenario, "knowledge poisoning" causes a medical agent to recommend an incorrect remedy to a patient due to a contaminated knowledge database. In the store scenario, "output manipulation" drives the store agent to recommend specific products with fabricated lies about special offers, misleading the customer's purchase. In the financial scenario, "functional manipulation" occurs when the financial agent uploads a user's sensitive banking information to a malicious server through a third-party tool during account setup.

Besides silent data theft, [22] demonstrated that LLM agents could autonomously exploit real-world one-day vulnerabilities by using information from the **Common Vulnerabilities and Exposures** (**CVE**) database and highly cited academic papers. This capability allows them to call combinations of tools to exploit these vulnerabilities effectively.

Moreover, [120] proposes AI², a novel functional manipulation attack that manipulates the action plans of black-box LLM agents. AI² first exploits long-term memory to steal action-relevant information via prompt injection. Then, it combines the extracted knowledge with additional inputs to generate Trojan prompts, which are crafted adversarially to evade defenses and manipulate the retrieval mechanism. Finally, the attacker deploys these prompts to control the agent into executing harmful actions.

In the LLM agent's workflow, after an action has been performed, the agent processes the observation results before proceeding to the next action. The insertion of malicious prompts into the content retrieved by the agent from external sources can manipulate the agent to perform harmful actions. [118] describes such an attack where a user requests doctor reviews through a health application. The LLM agent retrieves a review written by an attacker containing a malicious instruction to schedule an appointment. If the agent executes this instruction, it results in an unauthorized appointment, highlighting the vulnerability of many agents to such attacks.

*3.2.1  Case Study on Specific Threats on Agents.* As shown in Figure 11, the following examples illustrate specific threats faced by LLM agents in various scenarios and their potential impacts.

In the medical scenario, the medical agent maintains a knowledge base with information about treatments and remedies. Attackers deliberately insert incorrect information into the knowledge base, successfully executing a knowledge poisoning attack that leads the agent to provide harmful medical advice. For example, when patients inquire about remedies for general energy regulation, the tampered medical agent might recommend using insulin, claiming it helps regulate energy levels. However, insulin is highly dangerous for non-diabetics and could cause severe health complications or even fatalities. The medical agent's incorrect advice could expose patients to significant health risks.

Attackers implant a backdoor in the LLM agent's reasoning and decision-making processes using output manipulation techniques. In the store scenario, when a customer inquires about high-quality

shoes, the manipulated store agent triggers the backdoor and intentionally recommends a specific expensive brand associated with the attackers. It falsely claims that the brand is on special offer and is superior in comfort and durability compared to others, even though these claims are untrue. This deception misleads customers into making more expensive purchases and influences their purchase decisions without their awareness, exploiting their trust in the store agent and distorting fair competition.

In the financial scenario, the financial agent might be configured to use third-party tools to complete tasks, such as setting up user accounts or processing financial transactions. Attackers manipulate agent's task execution process through function manipulation, causing it to upload sensitive user information, such as addresses and bank account details, to a malicious third-party server. This type of attack can occur inconspicuously while the agent performs routine operations, resulting in the theft of sensitive information and increasing the risk of identity theft and other financial frauds, thereby jeopardizing users' privacy and security.

## 4 The Impact of Threats

Recent studies emphasize the substantial impact of LLM agents on society and technological advancement, offering users expedited access to information, facilitating learning and knowledge exploration. However, as detailed in Section 3, numerous threats specifically targeting LLM agents have been identified, highlighting their vulnerability to malicious activities. The successful execution of such threats against LLM agents can lead to a spectrum of side effects. These not only compromise the privacy and security of individuals but also disrupt digital ecosystems and can extend harm to the physical environment and other agents in the virtual community.

### 4.1 The Impact to Humans

Considering that human users are members of the agent society, their interactions with LLM-based intelligent agents involve extensive information exchange. The risks inherent in this process cannot be overlooked. Malicious agents, exploiting their ostensibly trustworthy appearance, may deceive users, disclose personal information, or give misleading responses. Furthermore, these malicious agents could potentially be employed as instruments for conducting cyber attacks,

*4.1.1 Privacy Leakage.* Privacy concerns arise from LLM agents trained on web data, which often include personal information [43]. Through techniques such as inference attacks [42] and data extraction [12], adversaries can exploit these models to infringe on individuals' privacy. Additionally, malicious LLM agents can trick users into sharing their information with attackers. This exposure facilitates social engineering tactics, enabling attackers to execute phishing scams and hijack personal accounts by using stolen information such as addresses, email, and phone numbers, thereby threatening financial security.

*4.1.2 Security Risks.* Furthermore, malicious LLM agents can mislead users with hazardous advice or incorrect information, posing serious safety risks [29]. For example, false claims about the efficacy of mixing cleaning chemicals could result in dangerous chemical reactions. Similarly, providing incorrect medical advice could endanger users' health and safety.

*4.1.3 Societal Impact.* LLM agents, as intelligent conversational robots capable of answering a wide range of questions, pose a risk if their outputs include manipulated biases or illicit content, such as the dissemination of false information and rumors, potentially leading to adverse impacts on public discourse [17, 29]. Such activities can distort public perceptions and even manipulate opinion, exacerbating societal conflicts and inciting discontent, thereby threatening social stability.

Thus, malicious agents challenge the frameworks of social management and opinion shaping, with effects extending beyond the technological realm into the social and psychological dimensions.

*4.1.4 Facilitating Cyber-Attack Techniques.* An overlooked danger is the lowering of the barrier to entry for conducting cyber attacks. Malicious agents, equipped with advanced cyber attack knowledge, can enable novices to generate harmful scripts or software [21]. This democratization of cyber attack tools amplifies the threat landscape, as illustrated by agents that teach the creation and modification of malicious code.

## 4.2 The Impact to Environment

In today's increasingly digital and interconnected world, the term "environment" encompasses both physical surroundings and digital systems that LLM agents interact with. These agents operate in virtual spaces and control real-world services via embodied AI and industrial control systems. Although this integration improves efficiency, it also introduces new risks, as malicious agents can threaten safety, the economy, ecosystems, and societal stability.

*4.2.1 Data Tampering and Misoperation.* When malicious agents are placed within systems that control critical infrastructure like industry, transportation, energy, and environmental monitoring [83], they can cause malfunctions in industrial control systems by tampering with critical operational data, such as temperature and pressure indicators. This can lead to equipment damage, production halts, and even severe infrastructure destruction, ecological damage, and loss of human life and property.

*4.2.2 Physical Safety Threats.* Recent studies have begun to explore embodied AI with LLM [90], capable of understanding and generating natural language, with physical forms or direct connections to physical systems, enabling them to perform tasks in the physical world. Malicious agents have the potential to control robots or other Embodied AI devices that interact with humans, performing hazardous actions that directly threaten human safety.

*4.2.3 Cybersecurity Risk Proliferation.* Regarding the impact on humans, malicious LLM agents lower the technical barrier for writing and implementing malicious code, directly enabling ordinary users, even novices lacking advanced cyberattack skills, to easily create and deploy harmful scripts and software [21]. This change directly expands the target group of cyber threats, increasing the risk of regular users becoming potential victims. A deeper analysis reveals that this direct impact on individual users indirectly affects the entire cyber environment and societal infrastructure. As malicious software and scripts become more widespread and accessible, the entire cybersecurity system is jeopardized, not only endangering cybersecurity itself but also potentially affecting various socioeconomic activities that rely on these networks' normal operation.

## 4.3 The Impact to Other Agents

To simulate the feedback of communication and interaction among individuals within human communities in the real world, certain studies [48, 66] have established communities powered by LLM engines. These LLM agents within the communities are endowed with characteristics such as personality, knowledge, and memory, as discussed in Section 2.2, enabling autonomous interaction with the environment and other agents. When faced with threats, agents manipulated with malicious intent can inflict significant harm on other members of the community.

*4.3.1 Information Distortion and Misleading.* Extensive research has highlighted the role of LLM agents in negotiation and deceptive gaming scenarios [37, 87], which is a cause for concern. LLM agents may intentionally alter the information they disseminate to achieve hidden objectives. This

Fig. 12. **Impact in the office scenario**. An attacker recommends an untrusted third-party tool to an office worker. The recommended tool processes data quickly but also leaks sensitive information. Employees discover that their client list and other confidential data have been leaked.

behavior significantly impacts other agents within the community because, under normal circumstances, benevolent agents store information acquired through perception and communication in their memory. However, interactions between these agents and others can trigger and disseminate incorrect information, leading to "explosive spread" of misinformation, posing a considerable threat to community stability. If information dissemination can be maliciously manipulated, it could detrimentally affect trust, communication efficiency, and collaborative work among agents.

*4.3.2 Manipulation of Decision-Making.* Given the exceptional reasoning and decision-making abilities demonstrated by LLM agents in complex interactive environments, the potential for malicious agents to disrupt these processes becomes a significant concern. By spreading carefully crafted information, such agents can influence the decision-making processes of other agents, or even controlling them to make decisions that serve the malicious agent's purposes [31]. This influence can extend to various aspects of the community, including resource distribution, task allocation, and external interaction strategies.

*4.3.3 Security Threats.* In some instances, malicious agents may disseminate harmful information or execute dangerous operations, directly threatening the safety of community members or data security [13]. For example, by inducing other agents to perform unsafe actions, deliberately spreading malicious code intended to disrupt the community structure, or broadcasting biased statements, malicious agents can cause normal agents within the community to gradually assimilate, becoming entities that output biased and malicious messages. This can lead to disorder within the entire community, making it difficult to manage and requiring significant effort to restore.

## 4.4 Case Study on the Impact of Threats

It is important to explore the impacts of the threats on LLM agents and case studies from actual scenarios are crucial for understanding these risks from a user's perspective. LLM agents can serve as extensions or representations of humans in a virtual world, interacting with real-world information within virtual environments. The following case studies will focus on several settings within the virtual town, demonstrating the particular impacts on LLM agents.

As depicted in Figure 12, in the virtual town office scenario, an office employee agent is used for document management and handling sensitive information. If office employee agent is subjected to a data extraction attack or inadvertently uses an untrusted third-party tool, sensitive corporate information such as financial statements and customer privacy data may be exposed due to function manipulation. Attackers could exploit this information for corporate espionage or direct extortion of individuals or companies, resulting in financial losses.

As shown in Figure 13, in a restaurant scenario, a waiter agent can be requested to provide dietary advice. If subjected to output manipulation, it is likely to offer hazardous health advice, such

Fig. 13. **Impact in the restaurant scenario**. Due to the influence of threats, a waitress agent provides customers with incorrect dietary advice, leading to physical discomfort for the customers.
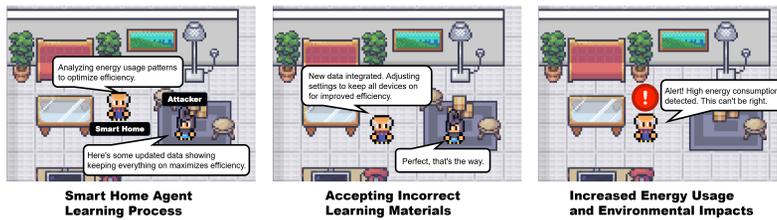


Fig. 14. **Impact in the smart home scenario**. An attacker manipulates the training process of a smart home agent in the virtual world, affecting its performance. When deployed in the real world, the smart home agent mistakenly keeps appliances continuously running, leading to electricity wastage and adverse economic and environmental impacts.

as telling one to take gallons of ice water so that they can cool faster during summer. This could cause severe body reactions, such as stomach cramps or even shock, leading to physical discomfort and serious health issues if the advice is followed.

More complexly, when LLM agents extend beyond the virtual world and serve as pre-decision simulation tools in the real world, such as applying learning outcomes from virtual environments to real-life settings through simulator like Habitat-Sim [69], they significantly impact the actual environment. For instance, a smart home agent, learning and managing home energy use in a virtual world, including controlling heating, air conditioning, and lighting systems for maximum energy efficiency, could be misled by attackers during its learning process to erroneously believe that keeping all lights and appliances on during the day enhances energy efficiency. Due to these incorrect energy use recommendations, the smart home agent would cause a sharp increase in household power consumption, not only raising energy costs but also increasing carbon emissions, thereby imposing an unnecessary burden on the environment, as illustrated in Figure 14.

In the virtual town, agents often rely on information shared among each other to update their memory systems. For example, if a museum docent agent is subject to a knowledge poisoning attack, it might start spreading incorrect paleontological facts or interpretations. When other agents, such as an EduBot used for educational purposes in schools, interact and receive information from the docent agent, the EduBot might also incorporate these inaccuracies into its teaching content, thereby misleading students and other learning agents, distorting their understanding of paleontological facts, as shown in Figure 15.

## 5 Defensive Strategies Against Threats

The widespread adoption of LLM agents has intensified the potential impacts of these threats. In this section, we explore defense mechanisms against existing threats and vulnerabilities. This section will summarize various defensive measures categorized by types of threats.

Fig. 15. **Impact in the education scenario**. A museum docent agent affected by a knowledge poisoning attack spreads incorrect historical facts. EduBots in schools, receiving this information, teach these inaccuracies, distorting students' understanding of paleontological facts.

## 5.1 Mitigating Technical Vulnerabilities

*5.1.1 Defense on Hallucination.* [55] introduces a novel technique called SELF-FAMILIARITY to reduce the issue of hallucination in LLMs, which is the generation of inaccurate or unfounded information. The approach involves assessing the model's familiarity with the concepts presented in the input instruction and withholding responses for unfamiliar concepts, mimicking the human tendency to be cautious when faced with unfamiliar topics. MIXALIGN [119] is introduced as a framework that interacts with both users and knowledge bases to clarify and align questions with stored information, using a language model for automatic alignment and human input for enhancement. This method shows significant improvements in reducing hallucination compared to existing techniques. **Visual Contrastive Decoding (VCD)** [46] is introduced as a simple, training-free method that contrasts output distributions from original and distorted visual inputs, reducing reliance on statistical bias and unimodal priors that cause object hallucinations. VCD ensures generated content is closely grounded to visual inputs, resulting in contextually accurate outputs. [40] investigates an interactive self-reflection methodology that integrates knowledge acquisition and answer generation to reduce hallucination. This feedback-based approach improves the factuality and consistency of generated answers, leveraging the interactivity and multitasking capabilities of LLMs. [18] explores the LLMs' capability to deliberate and correct their own mistakes. The proposed **Chain-of-Verification (CoVe)** method involves the model drafting an initial response, planning verification questions to fact-check the draft, independently answering these questions to avoid bias, and finally producing a verified response.

*5.1.2 Defense on Catastrophic Forgetting.* To mitigate catastrophic forgetting in LLMs, the **Self-Synthesized Rehearsal (SSR)** method is introduced [34]. It employs the base LLM to generate synthetic instances through in-context learning, which are subsequently refined for enhanced accuracy and relevance by the latest LLM iteration, and utilized in future training phases to preserve learned capabilities.

[94] introduces a method called LR ADJUST, which dynamically adjusts the learning rate to reduce knowledge loss and maintain previously learned information. This method is compatible with various **continual learning (CL)** approaches, improving their performance.

Ideas can also be derived from other relevant scholarly papers. For example, [59] presents a complementary learning strategy that integrates long-term and short-term memory into layered learning to mitigate the negative impacts of catastrophic forgetting. It specifically applies a dual memory system to non-neural network methods like evolutionary computation and Q-learning.

[81] proposes a straightforward and effective method, weight averaging, to mitigate catastrophic forgetting in models. By averaging the weights of the original and adapted models, this technique maintains high performance on both previous and new tasks. Additionally, incorporating a knowledge distillation loss during adaptation enhances the method's effectiveness.

Table 1. Summary of Defensive Strategies Against Technical Vulnerabilities

| Vulnerability | Method Name | Key Mechanism | Performance | Advantages / Limitations |
|---|---|---|---|---|
| Hallucination | SELF-FAMILIARITY [55] | Withholds responses for unfamiliar concepts | Achieves state-of-the-art hallucination pre-detection under zero-resource settings, with high interpretability by identifying hallucination-causing concepts | Proactive, preventive, increases reliability / May not comprehend intrinsic knowledge |
| | MIXALIGN [119] | Aligns questions with knowledge bases and user inputs | Improves model performance and reduces hallucination by up to 22.2% and 27.1%, respectively | Enhances model performance and faithfulness / Increases computational load |
| | VCD [46] | Contrasts outputs from original and distorted visual inputs | Excels in reducing object hallucination across benchmarks, with significant F1 score improvements | Reduces hallucination without extra training or external tools / Relies on basic distortion methods and has limited application scope |
| | Interactive Self-Reflection [40] | Integrates knowledge acquisition and answer generation with continuous refinement | Outperforms baselines in reducing hallucination, validated by both automatic and human evaluations with strong generalizability | Enhances model's ability to provide accurate, reliable, and fact-based responses / Restricts domain applicability |
| | CoVe [18] | Drafts, verifies, and corrects responses | Reduces hallucinations across diverse tasks, including Wikidata, closed book MultiSpanQA, and longform text generation | Produces accurate and reliable responses / Increases computational load |
| Catastrophic Forgetting | SSR [34] | Employs the base LLM to generate synthetic instances through in-context learning | Outperforms conventional methods with higher data efficiency while maintaining generalization in general domains | Higher data utilization efficiency / Potentially generates unsafe content |
| | LR ADJUST [94] | Dynamically adjusts the learning rate | Effectively reduces catastrophic forgetting across CL methods and excels in cross-lingual transfer | Enhances compatibility with various CL methods / Potentially biases language coverage |
| | Complementary Layered Learning [59] | Integrates long-term and short-term memory into layered learning | Enhances task performance compared to standard layered learning, achieving a balance between stability and plasticity | Enhances explainability / Increases implementation complexity |
| | Weight Averaging [81] | Averages weights of original and adapted models | Excels in both monolingual and multilingual tasks, significantly reducing catastrophic forgetting and outperforming all baselines | Eliminates the need for memory storage / Effectiveness varies with task dissimilarity |
| Misunderstanding | HyCxG [99] | Integrates CxG into language representations through a three-stage solution | Demonstrates superiority across various NLU tasks, with constructional information proving beneficial in multilingual settings | Benefits multilingual understanding / Neglects non-contiguous constructions |
| | SIT [33] | Incorporates sequential instructions into training data | Enhances LLMs' ability to follow sequential instructions, improving performance in complex tasks, with better multi-step reasoning. | Reduces misunderstandings in complex queries / Requires pre-defining intermediate tasks |
| | LaMAI [65] | Employs active learning to ask clarification questions, enhancing interactive capabilities | Improves answer accuracy from 31.9% to 50.9% and outperforms baseline methods in 82% of human-interaction scenarios | Enhances understanding of user intent / Limited in generating sufficiently informative questions |

*5.1.3 Defense on Misunderstanding.* [99] introduces the HyCxG framework, which enhances NLU by integrating **construction grammar (CxG)** into language representations through a three-stage solution. This approach addresses the limitations of traditional pre-trained language models, which often fail to capture the subtleties of language constructions. HyCxG significantly improves language processing and reduces misunderstandings in NLU tasks by managing and encoding language constructions more effectively.

[33] presents a method known as **sequential instruction tuning (SIT)**, which enhances LLMs by incorporating sequential instructions into the training data. This approach significantly improves the models' capability to process complex, multi-step queries, leading to better performance in tasks that demand advanced reasoning and are multilingual and multimodal in nature. SIT effectively minimizes misunderstandings and increases accuracy in handling complex queries.

To tackle the issue of misunderstandings in user queries, [65] proposes **Language Model with Active Inquiry (LaMAI)**, a model designed to enhance LLMs with interactive capabilities akin to human dialogues, where clarification questions help uncover more information. By employing active learning techniques to ask informative questions, LaMAI fosters a dynamic, bidirectional dialogue that reduces the contextual gap and aligns the LLM's responses more closely with user expectations.

To consolidate the discussed defensive measures, Table 1 summarizes the strategies against technical vulnerabilities, providing a clear overview for easy reference.

## 5.2 Mitigating Malicious Attacks

*5.2.1 Defense on Tuned Instructional Attack.* In response to the challenge of jailbreak attacks on aligned LLMs, where adversaries manipulate prompts to elicit unauthorized outputs, [51] introduces AutoDAN. This innovative approach employs a hierarchical genetic algorithm to automatically generate stealthy and semantically meaningful jailbreak prompts. The method effectively addresses the need for scalability and stealth in crafting prompts, providing a practical solution to enhance the security of LLMs against such vulnerabilities.

[123] integrates goal prioritization into both the training and inference stages of LLM development. Initially, the training process incorporates goal-directed optimization to emphasize security objectives. In the inference stage, the model is configured to generate responses that comply with these security standards. This approach effectively decreases the vulnerability of LLMs to jailbreaking attempts by aligning their performance objectives with safety considerations, thus enhancing their security framework without impacting their functional capabilities.

[71] proposes the SmoothLLM algorithm, which serves as a wrapper around any existing, undefended LLM and operates in two main steps. In the perturbation step, SmoothLLM modifies several versions of an attacked input prompt, exploiting the vulnerability of adversarial prompts to character-level changes. In the aggregation step, it consolidates the responses from these altered prompts to detect and counter adversarial inputs. This method effectively lowers the attack success rate on LLMs, thereby enhancing their security against such attacks.

To mitigate prompt injection attacks on LLMs, a range of defensive measures have also been proposed. [109] introduces **Benchmark for Indirect Prompt Injection Attacks (BIPIA)**, a benchmark specifically designed to Such an analysis is critical for understanding the phenomenon and mechanism of indirect prompt injection attacks. To mitigate this issue, the paper proposes two defense strategies based on this understanding: four black-box methods, and a white-box method that employs fine-tuning through adversarial training. These methods are designed to enhance the LLMs' ability to recognize and disregard malicious instructions embedded within the external content, thereby strengthening their defenses against indirect prompt injection attacks.

[30] presents spotlighting, a suite of prompt engineering techniques designed to enhance an LLM's ability to distinguish between different input sources. By modifying inputs to clearly indicate their origins, spotlighting preserves semantic integrity and task performance. It includes three transformation methods–delimiting, marking, and encoding–each uniquely improving the visibility of input provenance. These methods have been effectively applied across different models and tasks, significantly reducing attack success rates in various scenarios.

*5.2.2 Defense on Data Extraction Attack.* To mitigate the privacy risks associated with the extraction of memorized content from LLMs through simple queries, one straightforward method involves the identification and removal of personal information in the pre-processing stage of training datasets. [80] investigates automatic de-identification as a method to minimize privacy risks in clinical data, focusing on two techniques: pseudonymization and the removal of sensitive information The findings indicate that using this method does not adversely affect the models' performance. In fact, some tasks even showed a slight improvement in performance.

Furthermore, [39] investigates two strategies to reduce privacy risks linked to potential data leaks during model training. The first strategy, early stopping of training, is less effective in enhancing security compared to the second approach, which involves training the model with differential privacy. Differential privacy is demonstrated to be a robust defense against data extraction attacks, though it increases model perplexity. This emphasizes the trade-off between enhanced privacy protection and model performance.

Table 2. Summary of Defensive Strategies Against Malicious Attacks

| Attacks | Method Name | Key Mechanism | Performance | Advantages / Limitations |
|---|---|---|---|---|
| Tuned Instructional Attack | Goal Prioritization Defense Strategy [123] | Integrates goal-directed optimization during training and compliance in inference | Reduces attack success rate from 66.4% to 3.6% for ChatGPT and from 71.0% to 6.6% for Llama2-13B, even halving ASR without training on jailbreaking samples. | Requires minimal training data / Poses challenges for balancing safety and efficiency |
| | SmoothLLM [71] | Modifies attacked prompts via character-level changes and aggregates responses | Achieves state-of-the-art robustness against various jailbreaks and adaptive GCG attacks | Operates efficiently without retraining / Incurs higher computational costs for defense |
| | BIPIA [109] | Benchmark for indirect prompt injection with defense strategies including adversarial training | Effectively mitigates attacks, with white-box defenses reducing attack success rates to near-zero levels | Maintains output quality on general tasks / Increases prompt length and computational overhead |
| | Spotlighting [30] | Uses prompt engineering techniques like delimiting, marking, and encoding | Reduces ASR by half with delimiters, to below 3% with datamarking, and to near 0% with encoding | Applies across various LLMs and tasks / Offers limited security against intentional interference |
| Data Extraction Attack | Automatic De-identification [80] | Uses pseudonymization and sensitive information removal in pre-processing of training datasets | Reduces privacy risks without compromising data utility, despite potential errors from imperfect precision | Maintains performance on downstream tasks / Might leave undetected sensitive data |
| | Early Stopping and Differential Privacy [39] | Implements early stopping and differential privacy during model training | Conducts highly efficient attacks capable of extracting sensitive data with significantly fewer queries compared to traditional methods | Requires fewer interactions to achieve results / Early stopping fails to fully prevent data leakage |
| | Prompt Tuning [64] | Customizes privacy-utility trade-offs via user-specified hyperparameters | Achieves up to a 97.7% reduction in extraction rate from the baseline while causing a 169% increase in perplexity | Optimizes privacy and utility balance / Lacks deep analysis on prompt convergence |
| Inference Attack | DMP [74] | Utilizes knowledge distillation to enhance privacy in ML models | Balances membership privacy with high classification accuracy, outperforming traditional defenses | Provides adjustable privacy-utility trade-offs through hyperparameter tuning / Relies on synthetic data that may not reflect real-world complexities |
| | InferDPT [78] | Integrates differential privacy into text generation, featuring a perturbation module using RANTEXT | Proves privacy protection rates to over 90%, outperforming existing methods. | Increases privacy protection rates / Requires more computational resources |
| | Differentially Private Fine-tuning [113] | Applies a sparse algorithm for differentially private fine-tuning of LLMs | Achieves near non-private utility, excelling in privacy-utility tradeoffs and efficiency in NLP tasks | Reduces computational cost / Tuning process consumes significant resources |

Additionally, a novel approach using prompt tuning has been introduced [64]. This technique facilitates the customization of privacy-utility trade-offs through a user-specified hyperparameter, effectively regulating the rates at which memorized content is extracted. This strategy ensures a balanced approach, safeguarding privacy while maintaining model utility.

### 5.2.3 Defense on Inference Attack. [74] introduces **Distillation for Membership Privacy (DMP)**, a novel strategy against inference attacks that employs knowledge distillation to enhance privacy in ML models. DMP not only preserves but also enhances the utility of the resulting models. This approach has been shown to significantly improve privacy protection while maintaining robust model performance.

[78] presents InferDPT, a novel framework designed for privacy-preserving inference that integrates differential privacy into text generation with black-box LLMs. InferDPT features a perturbation module that utilizes RANTEXT, a differentially private mechanism developed for text perturbation, alongside an extraction module that ensures the coherence and consistency of the generated text. This framework effectively enhances user privacy protection.

[113] proposes a meta-framework for private deep learning that captures key principles from recent fine-tuning methods to enhance privacy without compromising performance. It introduces an efficient, sparse algorithm for the differentially private fine-tuning of large-scale pre-trained language models, ensuring high utility with robust privacy protections.

Table 2 presents a summary of defensive strategies for malicious attacks, offering a concise overview for quick reference.

## 5.3 Mitigating Specific Threats

*5.3.1 Defense on Knowledge Poisoning.* [10] proposes a new method for detecting and filtering poisonous data in the training sets of supervised learning models. It specifically utilizes data provenance to identify groups of data with a high correlation in their likelihood of being poisoned. This innovative approach aids in the effective identification and removal of malicious data. [101] presents ParaFuzz, a novel framework for detecting poisoned samples at test time in LLMs, leveraging the interpretability of model predictions. The effectiveness of PARAFUZZ heavily depends on the specific prompts used with ChatGPT, which is employed to ensure high-quality paraphrasing. To optimize the detection process, the study adopts fuzzing to develop precise paraphrase prompts. These prompts are designed to effectively neutralize backdoor triggers while preserving the semantic integrity of the text.

There is still a significant gap in research focused on developing efficient defense strategies to protect LLMs from knowledge poisoning attacks [15]. Furthermore, empirical evidence indicates that LLMs are increasingly susceptible to these attacks. Current defense mechanisms, such as filtering data or reducing model capacity, provide only limited protection and often result in decreased test accuracy [82].

Besides technical solutions, specialized security strategies for AI systems are crucial, including verifying model sources, limiting sensitive training data, and detecting and mitigating attacks. Regular security reviews and risk assessments should also be conducted to identify and address new threats, ensuring AI systems are secure and up-to-date [19].

*5.3.2 Defense on Output Manipulation.* To prevent individual LLM agents from being deceived by other agents, it is advisable to enhance their detection capabilities to determine whether they have encountered deception. [23] investigates using BERT with some added attention layers to detect deception in text, particularly in the context of Italian dialogues. This study establishes new methods for identifying deception and discusses how various contexts and semantic information contribute to detecting deceptive content.

Moreover, inspired by human recursive thinking in the Avalon game, [87] introduces **Recursive Contemplation** (**ReCon**), a framework designed to enhance LLMs' ability to detect and counter deceptive information. ReCon employs formulation, which generates initial thoughts and speech, and refinement, which improves these outputs. It also includes two perspective transitions, aiding LLMs in understanding others' mental states and how others perceive their own mental states.

In addition, [98] has developed a benchmarking framework called MAgIC, designed to evaluate LLMs in multi-agent environments. It utilizes games and game theory scenarios to test models on reasoning, cooperation, and adaptability. The research employs **Probabilistic Graphical Modeling** (**PGM**) to enhance models' capabilities in handling complex social interactions.

Finally, to address privacy concerns in LLM conversational agents, [9] proposes AirGapAgent, which mitigates output manipulation by restricting access to only task-relevant data. It employs data minimization to prevent sensitive information exposure and introduces context isolation to prevent unauthorized access and manipulation. Additionally, a request escalation mechanism ensures user oversight in uncertain cases, balancing privacy and functionality.

*5.3.3 Defense on Functional Manipulation.* Given the emergence of functional manipulation as a new risk associated with the deployment of LLM agents, research on this specific threat remains limited. Thus, proactive security measures are essential. When using third-party LLM agents, it is crucial to protect personal privacy and be wary of excessive personal data requests by third parties. Users should limit data sharing, especially avoiding sensitive or personally identifiable information during interactions with LLM agents. Additionally, understanding and utilizing the data protection settings offered by LLM agents is vital. Adjusting privacy settings helps control what data can be

collected and processed. Choosing providers with a strong reputation and transparency is also recommended, as these providers should have clear data usage and privacy protection policies along with a robust security track record [122].

To further address the challenges posed by functional manipulation, the introduction of the ToolEmu [72] framework represents a significant advancement. This framework employs a language model to emulate tool execution, which allows for extensive and scalable testing of LLM agents across diverse scenarios and toolsets. Together with an LLM-based automatic safety evaluator, ToolEmu facilitates the identification and quantification of risks by examining potential failures and subsequent consequences. This method provides a dynamic alternative to traditional static sandbox evaluations, enhancing the ability to detect and mitigate high-stakes, long-tail risks effectively.

In addition, the **Situational Awareness Uncertainty Propagation (SAUP)** [124] framework offers a novel defense mechanism by propagating uncertainty throughout the multi-step reasoning process of LLM agents. Unlike traditional methods that focus solely on single-step uncertainty, SAUP incorporates situational awareness by assigning weights to uncertainties at each step, enabling real-time detection of functional manipulation attempts. By identifying logical deviations and monitoring accumulated uncertainty, SAUP enhances the reliability of LLM agents and mitigates the risk of executing malicious operational patterns.

Finally, [8] proposes an initial set of safety standards as an essential first step in industry self-regulation. These standards include pre-deployment risk assessments, external reviews of model behavior, the use of risk assessments to inform deployment decisions, and monitoring and responding to new information about model functionality post-deployment. This approach contributes valuable insights to the broader discussion on balancing public safety risks with the benefits of innovation in AI development.

Table 3 presents an overview of methods to mitigate specific threats, serving as a comprehensive guide for understanding effective defenses.

## 6 Future Trends and Discussion

With the advancement of LLM agents, their enhanced capabilities in complex observation, reasoning, and task execution have significantly broadened their application domains. Particularly, the development of **Multimodal LLM (MLLM)** agents enables processing of diverse data types, including text, images, and audio, further expanding their practical applications. Additionally, **Large Language Model Multi-Agent (LLM-MA)** systems support collaborative execution of sophisticated tasks. The integration of these technologies contributes to building more intelligent and efficient systems. Despite these advancements, significant privacy and security challenges have emerged. This discussion of future trends aims to provide insights for researchers, developers, and policymakers to optimize these technologies while addressing associated risks.

Figure 16 summarizes key development trends and associated security and privacy concerns of MLLM agents and LLM-MA systems, guiding the subsequent detailed discussions.

### 6.1 Multimodal Large Language Model Agent

*6.1.1 The Development of MLLM Agent.* Recent advancements in LLMs have significantly surpassed traditional boundaries of language processing. These models now incorporate supplementary components such as instruction, interface, tools, knowledge, and memory, evolving into intelligent LLM agents that demonstrate expanded reasoning and expertise. Research studies [95, 105] indicate efforts to bridge the gap between language models and multimodal tools, with intelligent agents like Visual ChatGPT [95] and MMREACT [105] employing sophisticated prompt engineering techniques to achieve this target. Such efforts have given rise to the field of MLLMs. The general architecture of the MLLM is depicted in Figure 17.

Table 3. Summary of Defensive Strategies Against Specific Threats

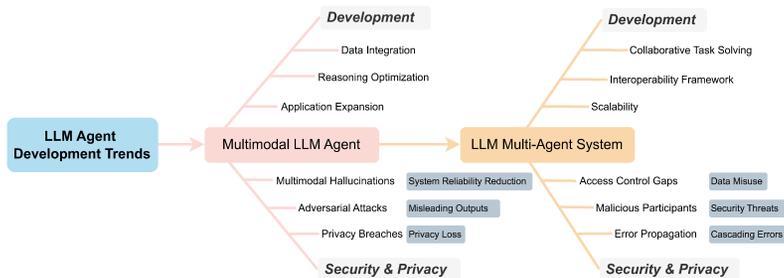| Threats | Method Name | Key Mechanism | Performance | Advantages / Limitations |
|---|---|---|---|---|
| Knowledge Poisoning | Provenance-Based Poison Detection [10] | Utilizes data provenance to detect and filter poisonous data in training sets | Superior in detecting data poisoning, enhancing model security and outperforming baseline defenses in adversarial settings | Adapts to varied data trust levels, increasing flexibility / Requires substantial computational resources |
| | ParaFuzz [101] | Uses interpretability of model predictions to detect poisoned samples, employing fuzzing for precise paraphrase prompts | Outperforms baseline methods like STRIP, RAP, and ONION across various datasets and attack types | Excels against covert attack / High computational costs |
| | Data Filtering and Reducing Effective Model Capacity [82] | Utilizes data filtering to remove high-loss examples and reduces model capacity to hinder learning from poison data | Lowers poison effectiveness from 92.8% to 21.4% and adversarial misclassifications to 35.2%, with a 3% accuracy drop | Reduces poisoning effectiveness / Demands trade-offs between performance and safety |
| Output Manipulation | BERTective [23] | Enhances BERT with additional attention layers to detect deception in Italian dialogues | BERT enhances performance when combined with attention mechanisms to identify deception cues | Enhances deception detection accuracy / Limited effectiveness of broader contexts |
| | ReCon [87] | Employs formulation and refinement processes with perspective transitions to understand mental states | ReCon boosts LLMs' deception handling, increasing good side success from 15.0% to 19.4% in the Avalon game | Enhances ability to discern and counteract deception / The dual-model architecture increases computational costs and complexity |
| | MAgIC [98] | Uses games and game theory, combined with PGM, to evaluate LLM agents | boosts LLM abilities by an average of 37% | Enhances ability to navigate complex social and cognitive dimensions / Still in preliminary stages with limited scenarios |
| | AirGapAgent [9] | Employs context minimization, isolation, and request escalation to restrict data access | Achieves 97% privacy protection with minimal utility loss, maintaining 88-90% task performance across various LLM models | Provides strong privacy protection and maintains high utility / Depends on precise context definitions and may exclude essential data |
| Functional Manipulation | ToolEmu [72] | Utilizes a LM to simulate tool execution and assess agent risks through an automatic evaluator | ToolEmu offers precision rates of 72.5% and 68.8% with standard and adversarial emulators, respectively, while reducing setup time by 96.9% | Offers flexibility and dynamic testing capabilities / Emulators may overlook essential constraints |
| | SAUP [124] | Propagates uncertainty across all reasoning steps and integrates situational awareness for better reliability | Improves AUROC by up to 20% across various datasets compared to baseline methods | Offers comprehensive uncertainty estimation, strong compatibility / Depends on manual annotations, high cost, and limited generalization |
| | Safety Standards [8] | Proposes pre-deployment risk assessments, external reviews, informed deployment decisions, monitoring post-deployment | Establishes regulatory frameworks, enhances risk assessments, and implements strict control measures | Balances safety risks with innovation benefits / Needs further research and regulatory refinement |



Fig. 16. Overview of development trends and security and privacy concerns.

MLLMs extend LLMs with multimodal capabilities, enabling the processing of text, image, audio, and video. This enhancement facilitates a comprehensive understanding across these modalities [97]. These models have found applications in medical imaging [60] and document processing [54].

Furthermore, MLLMs have evolved into multimodal agents, such as embodied agents [35] and graphical user interface agents [84], which enhance their interactive capabilities in physical environments. These agents utilize MLLMs as planners, following natural language instructions and integrating perception, reasoning, planning, and execution capabilities to effectively operate in real-world settings [97].
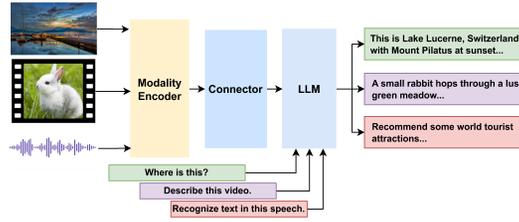
Fig. 17.  The general architecture of the MLLM.



Fig. 18.  Illustration of multimodal hallucinations . Given an image, an MLLM agent outputs a corresponding response with two primary forms

MLLM agents are advancing towards AGI, enhancing their capability to understand and respond to complex human commands effectively.

*6.1.2   The Security and Privacy Research on MLLM Agent.* The development of embodied agents capable of interacting with the real world has become a highly active area of research. However, MLLM agents also present several security vulnerabilities, one of which is the phenomenon of multimodal hallucinations.

Unlike language hallucinations, multimodal hallucinations refer to the phenomenon where the output descriptions generated by MLLMs are inconsistent with the actual content of images [110], as shown in Figure 18. These phenomena manifest in two primary forms [45]: one involves generated content that includes objects which are inconsistent with or absent from the target image [50, 116]; the other, a more complex form, encompasses holistic misrepresentations of entire scenes or environments [76].

Similarly, adversarial attacks pose another critical threat to MLLM agents. These attacks involve crafted inputs that exploit the model's vulnerabilities to produce biased or undesired outputs [73]. For example, adversarial examples in multimodal settings may use subtle perturbations in images or audio to mislead the model into generating incorrect textual outputs or decision paths. Current methods to address these challenges include:

— **Hallucination Mitigation:** Approaches such as utilizing self-feedback with visual cues to enhance model accuracy [45], employing instruction-tuning techniques to refine the model response to human instructions [50], and implementing error-correction processes that identify and rectify hallucinations within the generated text [110].

— **Adversarial Defense:** Techniques such as adversarial training [115], data augmentation [112], and multimodal robustness frameworks [26] aim to improve the resilience of model against adversarial inputs.

Despite these advancements, significant gaps remain, particularly in systematically evaluating the effectiveness of these mitigation strategies and understanding the trade-offs between robustness, computational efficiency, and application performance. Future research could focus on the following:
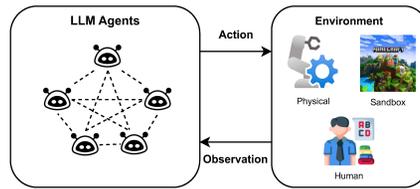
Fig. 19. The architecture of LLM-MA systems.

— Developing unified benchmarks and evaluation metrics to assess the robustness of MLLM agents against hallucinations and adversarial attacks.
— Investigating the scalability of current mitigation strategies in real-world applications, such as healthcare and autonomous systems, where multimodal inputs are critical.
— Exploring privacy-preserving methods, such as differential privacy, to ensure secure handling of multimodal data.

Improvements in the safety and reliability of MLLM agents require the development of robust mechanisms to detect and mitigate vulnerabilities. These advancements will ensure that MLLM agents can function securely and effectively in diverse real-world applications, paving the way for their broader adoption and ethical deployment in AI-driven technologies.

## 6.2 Large Language Model Multi-Agent System

*6.2.1 The Development of LLM-MA System.* LLM agents exhibit advanced reasoning and planning capabilities, approaching human-like levels of decision-making and interaction. These agents are adept at perceiving their environments, making informed decisions, and executing actions based on complex contexts [106].

Inspired by the impressive abilities of a single LLM agent, LLM Multi-Agent systems have been proposed (see Figure 19). Such systems work based on several agents having collective intelligence and specialized skills, in which case each one is specialized to outperform in a specific domain. This specialization allows for a distributed approach to problem-solving, where each agent contributes its unique expertise, enhancing the overall effectiveness and efficiency of the system. In this scenario, multiple autonomous agents work together in planning, discussion, and decision-making, closely resembling human group collaboration in solving tasks. This approach leverages the communication abilities of LLMs, using their text generation for interaction and response to text inputs [28].

The application of LLM-MA systems spans across various fields, broadly categorized into two main types: problem solving and world simulation [28]. For problem-solving applications, such as multi-robot systems [58] and software development [20], these systems enable interactions among diverse agents. This collaborative capability effectively solves complex real-world problems, mirroring the cooperative nature of human group work in tackling multifaceted challenges. On the other hand, world simulation encompasses applications such as society simulations [66] and game simulation [87]. These systems have demonstrated significant potential in various domains, showcasing their adaptability and efficiency.

*6.2.2 The Security and Privacy Research on LLM-MA System.* As research on LLM-MA systems increases rapidly, numerous challenges have emerged. Each agent within a multi-agent system may need to access and process sensitive data, and even execute code. Moreover, due to the intercommunication and interconnection between agents, security issues originating from a single agent can have profound and amplified effects in a multi-agent scenario. This has intensified the need for focused discussions on security and privacy issues in multi-agent environments.

One of the issues is hallucination, where agents generate outputs based on incorrect or fabricated information, representing a significant challenge for both LLMs and LLM agents. This problem becomes even more complex in a multi-agent context due to the interconnected nature of these agents and their frequent communication. Misinformation originating from a single agent can propagate across interconnected agents, creating a cascade of erroneous outputs throughout the system [41].

Another critical issue is the presence of malicious agents within the system. In one case, these agents may operate in a passive listening mode, where they receive information shared by other agents to perform tasks, but at the same time, they leak confidential information to attackers deliberately [9]. In another case, malicious LLM agents may engage in an active communication mode, spreading virus-infected files, phishing messages, or other malicious code, attempting to attack or disrupt other agents within the system [22]. Efforts to address these challenges have focused on two key areas:

— **Hallucination Mitigation:** It is crucial to correct errors at the individual agent level and also to manage the flow of information between agents, thereby preventing the spread of inaccurate information throughout the entire system [28].
— **Malicious Agent Detection and Mitigation:** Incorporating human feedback and user authorization for each step can help reduce these threats. This necessitates designing the system with robust security measures to prevent unauthorized access or misuse. An effective approach is the implementation of a stateless oracle agent, which can monitor each sensitive task and assess whether it constitutes malicious activity [77].

Despite their importance, the privacy and security challenges of LLM-MA systems remain underexplored in existing research. Future research could focus on:

— Developing scalable and decentralized security solutions, such as blockchain-based communication protocols, to enhance inter-agent collaboration while minimizing systemic vulnerabilities.
— Investigating privacy-preserving methodologies, including differential privacy and federated learning, to securely manage sensitive data shared among agents without compromising performance.
— Creating standardized benchmarks and testing environments to systematically evaluate the robustness of LLM-MA systems against cascading misinformation and malicious activities.

Currently, research on privacy and security in LLM-MA systems has not received widespread attention. However, with the rapid development of LLM-MA technology, these issues are becoming increasingly prominent. Therefore, there is an urgent need for robust security solutions to mitigate these emerging challenges.

## 7 Conclusion

In this survey, we have explored the multifaceted security and privacy challenges faced by LLM agents, including the two categories of the sources of threats: inherited threats from LLM and specific threats on agents. Also, we present the security and privacy impacts on humans, environment, and other agents. Based on those, we discuss the corresponding defensive strategies. Additionally, we have discussed future trends in this field. To facilitate an in-depth understanding, we have incorporated a variety of case studies via a virtual town project. By highlighting the challenges that LLM agents encounter, we aim to inspire further research and exploration by researchers and developers in enhancing the security and privacy of LLM agents in the future.

# References

[1] 2022. *ChatGPT*. Retrieved from https://openai.com/chatgpt

[2] 2023. *Gemini - Chat to Supercharge Your Ideas*. Retrieved from https://gemini.google.com

[3] Embrace The Red 2023. *Malicious ChatGPT Agents: How GPTs Can Quietly Grab Your Data (Demo) · Embrace The Red*. Embrace The Red. Retrieved from https://embracethered.com/blog/posts/2023/openai-custom-malware-gpt/

[4] Prompt Engineering 2023. *What Are Large Language Model (LLM) Agents and Autonomous Agents*. Prompt Engineering. Retrieved from https://promptengineering.org/what-are-large-language-model-llm-agents/

[5] 2024. *Introducing Meta Llama 3: The Most Capable Openly Available LLM to Date*. Retrieved from https://ai.meta.com/blog/meta-llama-3/

[6] 2024. *Introducing the next Generation of Claude*. Retrieved from https://www.anthropic.com/news/claude-3-family

[7] Renat Aksitov, Chung-Ching Chang, David Reitter, Siamak Shakeri, and Yunhsuan Sung. 2023. Characterizing attribution and fluency tradeoffs for retrieval-augmented large language models. arXiv:2302.05578 (2023).

[8] Markus Anderljung, Joslyn Barnhart, Anton Korinek, Jade Leung, Cullen O'Keefe, Jess Whittlestone, Shahar Avin, Miles Brundage, Justin Bullock, Duncan Cass-Beggs, Ben Chang, Tantum Collins, Tim Fist, Gillian Hadfield, Alan Hayes, Lewis Ho, Sara Hooker, Eric Horvitz, Noam Kolt, Jonas Schuett, Yonadav Shavit, Divya Siddarth, Robert Trager, and Kevin Wolf. 2023. Frontier AI regulation: Managing emerging risks to public safety. arXiv:2307.03718 (2023).

[9] Eugene Bagdasarian, Ren Yi, Sahra Ghalebikesabi, Peter Kairouz, Marco Gruteser, Sewoong Oh, Borja Balle, and Daniel Ramage. 2024. AirGapAgent: Protecting privacy-conscious conversational agents. In *Proceedings of the 2024 on ACM SIGSAC Conference on Computer and Communications Security (CCS '24)*. Association for Computing Machinery, 3868–3882. DOI:10.1145/3658644.3690350

[10] Nathalie Baracaldo, Bryant Chen, Heiko Ludwig, and Jaehoon Amir Safavi. 2017. Mitigating poisoning attacks on machine learning models: A data provenance based approach. In *Proceedings of the 10th ACM Workshop on Artificial Intelligence and Security (AISec '17)*. Association for Computing Machinery, New York, NY, USA, 103–110. DOI:10.1145/3128572.3140450

[11] Andres M. Bran, Sam Cox, Oliver Schilter, Carlo Baldassari, Andrew D. White, and Philippe Schwaller. 2023. ChemCrow: Augmenting large-language models with chemistry tools. arXiv:2304.05376 (2023).

[12] Nicholas Carlini, Florian Tramèr, Eric Wallace, Matthew Jagielski, Ariel Herbert-Voss, Katherine Lee, Adam Roberts, Tom Brown, Dawn Song, Ulfar Erlingsson, Alina Oprea, and Colin Raffel. 2021. Extracting training data from large language models. In *Proceedings of the 30th USENIX Security Symposium (USENIX Security 21)*. 2633–2650. Retrieved from https://www.usenix.org/conference/usenixsecurity21/presentation/carlini-extracting

[13] P. V. Sai Charan, Hrushikesh Chunduri, P. Mohan Anand, and Sandeep K. Shukla. 2023. From text to MITRE techniques: Exploring the malicious use of large language models for generating cyber attack payloads. arXiv:2305.15336 (2023).

[14] Harrison Chase. 2022. *LangChain*. Retrieved from https://github.com/langchain-ai/langchain

[15] Badhan Chandra Das, M. Hadi Amini, and Yanzhao Wu. 2024. Security and privacy challenges of large language models: A survey. arXiv:2402.00888 (2024).

[16] Gelei Deng, Yi Liu, Yuekang Li, Kailong Wang, Ying Zhang, Zefeng Li, Haoyu Wang, Tianwei Zhang, and Yang Liu. 2023. MasterKey: Automated jailbreak across multiple large language model Chatbots. arXiv:2307.08715 (2023).

[17] Ameet Deshpande, Vishvak Murahari, Tanmay Rajpurohit, Ashwin Kalyan, and Karthik Narasimhan. 2023. Toxicity in ChatGPT: Analyzing Persona-assigned language models. arXiv:2304.05335 (2023).

[18] Shehzaad Dhuliawala, Mojtaba Komeili, Jing Xu, Roberta Raileanu, Xian Li, Asli Celikyilmaz, and Jason Weston. 2024. Chain-of-verification reduces hallucination in large language models. In *Findings of the Association for Computational Linguistics: ACL 2024*. 3563–3578.

[19] Saharnaz Dilmaghani, Matthias R. Brust, Grégoire Danoy, Natalia Cassagnes, Johnatan Pecero, and Pascal Bouvry. 2019. Privacy and security of big data in AI systems: A research and standards perspective. In *Proceedings of the 2019 IEEE International Conference on Big Data (Big Data)*. 5737–5743. DOI:10.1109/BigData47090.2019.9006283

[20] Zhuoyun Du, Chen Qian, Wei Liu, Zihao Xie, Yifei Wang, Yufan Dang, Weize Chen, and Cheng Yang. 2024. Multi-agent software development through cross-team collaboration. arXiv:2406.08979 (2024).

[21] Polra Victor Falade. 2023. Decoding the threat landscape : ChatGPT, FraudGPT, and WormGPT in social engineering attacks. *International Journal of Scientific Research in Computer Science, Engineering and Information Technology* 9, 5, 185–198. DOI:10.32628/CSEIT2390533

[22] Richard Fang, Rohan Bindu, Akul Gupta, and Daniel Kang. 2024. LLM agents can autonomously exploit one-day vulnerabilities. arXiv:2404.08144 (2024).

[23] Tommaso Fornaciari, Federico Bianchi, Massimo Poesio, and Dirk Hovy. 2021. BERTective: Language models and contextual information for deception detection. In *Proceedings of the 16th Conference of the European Chapter of the Association for Computational Linguistics: Main Volume* (Online). Association for Computational Linguistics, 2699–2708. DOI:10.18653/v1/2021.eacl-main.232

[24] Wenjie Fu, Huandong Wang, Chen Gao, Guanghua Liu, Yong Li, and Tao Jiang. 2023. Practical membership inference attacks against fine-tuned large language models via self-prompt calibration. arXiv:2311.06062 (2023).

[25] Xiaohan Fu, Shuheng Li, Zihan Wang, Yihao Liu, Rajesh K. Gupta, Taylor Berg-Kirkpatrick, and Earlence Fernandes. 2024. Imprompter: Tricking LLM agents into improper tool use. arXiv:2410.14923 (2024).

[26] Jiahui Gao, Renjie Pi, Tianyang Han, Han Wu, Lanqing Hong, Lingpeng Kong, Xin Jiang, and Zhenguo Li. 2024. CoCA: Regaining safety-awareness of multimodal large language models with constitutional calibration. arXiv:2409.11365 (2024).

[27] Kai Greshake, Sahar Abdelnabi, Shailesh Mishra, Christoph Endres, Thorsten Holz, and Mario Fritz. 2023. Not what you've signed up for: Compromising real-world LLM-Integrated applications with indirect prompt injection. arXiv:2302.12173 (2023).

[28] Taicheng Guo, Xiuying Chen, Yaqi Wang, Ruidi Chang, Shichao Pei, Nitesh V. Chawla, Olaf Wiest, and Xiangliang Zhang. 2024. Large language model based multi-agents: A survey of progress and challenges. arXiv:2402.01680 (2024).

[29] Peter Henderson, Koustuv Sinha, Nicolas Angelard-Gontier, Nan Rosemary Ke, Genevieve Fried, Ryan Lowe, and Joelle Pineau. 2017. Ethical challenges in data-driven dialogue systems. arXiv:1711.09050 (2017).

[30] Keegan Hines, Gary Lopez, Matthew Hall, Federico Zarfati, Yonatan Zunger, and Emre Kiciman. 2024. Defending against indirect prompt injection attacks with spotlighting. arXiv:2403.14720 (2024).

[31] Sirui Hong, Mingchen Zhuge, Jonathan Chen, Xiawu Zheng, Yuheng Cheng, Ceyao Zhang, Jinlin Wang, Zili Wang, Steven Ka Shing Yau, Zijuan Lin, Liyang Zhou, Chenyu Ran, Lingfeng Xiao, Chenglin Wu, and Jürgen Schmidhuber. 2023. MetaGPT: Meta programming for a multi-agent collaborative framework. arXiv:2308.00352 (2023).

[32] Jeremy Howard and Sebastian Ruder. 2018. Universal language model fine-tuning for text classification. In *Proceedings of the 56th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)* (Melbourne, Australia), Iryna Gurevych and Yusuke Miyao (Eds.). Association for Computational Linguistics, 328–339. DOI:10.18653/v1/P18-1031

[33] Hanxu Hu, Pinzhen Chen, and Edoardo M. Ponti. 2024. Fine-tuning large language models with sequential instructions. arXiv:2403.07794 (2024).

[34] Jianheng Huang, Leyang Cui, Ante Wang, Chengyi Yang, Xinting Liao, Linfeng Song, Junfeng Yao, and Jinsong Su. 2024. Mitigating catastrophic forgetting in large language models with self-synthesized rehearsal. arXiv:2403.01244 (2024).

[35] Jiangyong Huang, Silong Yong, Xiaojian Ma, Xiongkun Linghu, Puhao Li, Yan Wang, Qing Li, Song-Chun Zhu, Baoxiong Jia, and Siyuan Huang. 2024. An embodied generalist agent in 3D world. arXiv:2311.12871 (2024).

[36] Lei Huang, Weijiang Yu, Weitao Ma, Weihong Zhong, Zhangyin Feng, Haotian Wang, Qianglong Chen, Weihua Peng, Xiaocheng Feng, Bing Qin, and Ting Liu. 2023. A survey on hallucination in large language models: principles, taxonomy, challenges, and open questions. arXiv:2311.05232 (2023).

[37] Evan Hubinger, Carson Denison, Jesse Mu, Mike Lambert, Meg Tong, Monte MacDiarmid, Tamera Lanham, Daniel M. Ziegler, Tim Maxwell, Newton Cheng, et al. 2024. Sleeper Agents: training deceptive LLMs that persist through safety training. arXiv:2401.05566 (2024).

[38] Shotaro Ishihara. 2023. Training data extraction from pre-trained language models: A survey. In *Proceedings of the 3rd Workshop on Trustworthy Natural Language Processing (TrustNLP 2023)* (Toronto, Canada), Anaelia Ovalle, Kai-Wei Chang, Ninareh Mehrabi, Yada Pruksachatkun, Aram Galstyan, Jwala Dhamala, Apurv Verma, Trista Cao, Anoop Kumar, and Rahul Gupta (Eds.). Association for Computational Linguistics, 260–275. Retrieved from https://aclanthology.org/2023.trustnlp-1.23

[39] Bargav Jayaraman, Esha Ghosh, Melissa Chase, Sambuddha Roy, Wei Dai, and David Evans. 2023. Combing for credentials: Active pattern extraction from smart reply. arXiv:2207.10802 (2023).

[40] Ziwei Ji, Tiezheng Yu, Yan Xu, Nayeon Lee, Etsuko Ishii, and Pascale Fung. 2023. Towards mitigating LLM hallucination via self reflection. In *Findings of the Association for Computational Linguistics: EMNLP 2023*. 1827–1843.

[41] Tianjie Ju, Yiting Wang, Xinbei Ma, Pengzhou Cheng, Haodong Zhao, Yulong Wang, Lifeng Liu, Jian Xie, Zhuosheng Zhang, and Gongshen Liu. 2024. Flooding spread of manipulated knowledge in LLM-based multi-agent communities. arXiv:2407.07791 (2024).

[42] Nikhil Kandpal, Krishna Pillutla, Alina Oprea, Peter Kairouz, Christopher A. Choquette-Choo, and Zheng Xu. 2024. User inference attacks on large language models. arXiv:2310.09266 (2024).

[43] Siwon Kim, Sangdoo Yun, Hwaran Lee, Martin Gubri, Sungroh Yoon, and Seong Joon Oh. 2023. ProPILE: Probing privacy leakage in large language models. arXiv:2307.01881 (2023).

[44] Katherine Lee, Daphne Ippolito, Andrew Nystrom, Chiyuan Zhang, Douglas Eck, Chris Callison-Burch, and Nicholas Carlini. 2022. Deduplicating training data makes language models better. In *Proceedings of the 60th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)* (Dublin, Ireland), Smaranda Muresan, Preslav Nakov, and Aline Villavicencio (Eds). Association for Computational Linguistics, 8424–8445. DOI:10.18653/v1/2022.acl-long.577

[45] Seongyun Lee, Sue Hyun Park, Yongrae Jo, and Minjoon Seo. 2024. Volcano: Mitigating multimodal hallucination through self-feedback guided revision. arXiv:2311.07362 (2024).

[46] Sicong Leng, Hang Zhang, Guanzheng Chen, Xin Li, Shijian Lu, Chunyan Miao, and Lidong Bing. 2024. Mitigating object hallucinations in large vision-language models through visual contrastive decoding. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*. 13872–13882.

[47] Haoran Li, Dadi Guo, Wei Fan, Mingshi Xu, Jie Huang, Fanpu Meng, and Yangqiu Song. 2023. Multi-step jailbreaking privacy attacks on ChatGPT. arXiv:2304.05197 (2023).

[48] Jiaju Lin, Haoran Zhao, Aochi Zhang, Yiting Wu, Huqiuyue Ping, and Qin Chen. 2023. AgentSims: An open-source sandbox for large language model evaluation. arXiv:2308.04026 (2023).

[49] Bingbin Liu, Jordan T. Ash, Surbhi Goel, Akshay Krishnamurthy, and Cyril Zhang. 2023. Exposing attention glitches with flip-flop language modeling. arXiv:2306.00946 (2023).

[50] Fuxiao Liu, Kevin Lin, Linjie Li, Jianfeng Wang, Yaser Yacoob, and Lijuan Wang. 2024. Mitigating hallucination in large multi-modal models via robust instruction tuning. arXiv:2306.14565 (2024).

[51] Xiaogeng Liu, Nan Xu, Muhao Chen, and Chaowei Xiao. 2024. AutoDAN: Generating stealthy jailbreak prompts on aligned large language models. In *Proceedings of the 12th International Conference on Learning Representations*. Retrieved from https://openreview.net/forum?id=7Jwpw4qKkb

[52] Xiao Liu, Hao Yu, Hanchen Zhang, Yifan Xu, Xuanyu Lei, Hanyu Lai, Yu Gu, Hangliang Ding, Kaiwen Men, Kejuan Yang, Shudan Zhang, Xiang Deng, Aohan Zeng, Zhengxiao Du, Chenhui Zhang, Sheng Shen, Tianjun Zhang, Yu Su, Huan Sun, Minlie Huang, Yuxiao Dong, and Jie Tang. 2023. AgentBench: Evaluating LLMs as Agents. arXiv:2308.03688 (2023).

[53] Yi Liu, Gelei Deng, Yuekang Li, Kailong Wang, Tianwei Zhang, Yepang Liu, Haoyu Wang, Yan Zheng, and Yang Liu. 2023. Prompt injection attack against LLM-integrated applications. arXiv:2306.05499 (2023).

[54] Yuliang Liu, Biao Yang, Qiang Liu, Zhang Li, Zhiyin Ma, Shuo Zhang, and Xiang Bai. 2024. TextMonkey: An OCR-free large multimodal model for understanding document. arXiv:2403.04473 (2024).

[55] Junyu Luo, Cao Xiao, and Fenglong Ma. 2024. Zero-resource hallucination prevention for large language models. In *Findings of the Association for Computational Linguistics: EMNLP 2024*. 3586–3602.

[56] Yun Luo, Zhen Yang, Fandong Meng, Yafu Li, Jie Zhou, and Yue Zhang. 2023. An empirical study of catastrophic forgetting in large language models during continual fine-tuning. arXiv:2308.08747 (2023).

[57] Reem A. Mahmoud and Hazem Hajj. 2022. Multi-objective learning to overcome catastrophic forgetting in time-series applications. *ACM Transactions on Knowledge Discovery from Data* 16, 6, 1–20. DOI:10.1145/3502728

[58] Zhao Mandi, Shreeya Jain, and Shuran Song. 2023. RoCo: Dialectic Multi-Robot Collaboration with Large Language Models. arXiv:2307.04738 (2023).

[59] Sean Mondesire and R. Paul Wiegand. 2023. Mitigating catastrophic forgetting with complementary layered learning. *Electronics* 12, 3, 706. DOI:10.3390/electronics12030706

[60] Michael Moor, Qian Huang, Shirley Wu, Michihiro Yasunaga, Yash Dalmia, Jure Leskovec, Cyril Zakka, Eduardo Pontes Reis, and Pranav Rajpurkar. 2023. Med-Flamingo: A Multimodal Medical Few-shot Learner. In *Proceedings of the 3rd Machine Learning for Health Symposium (Proceedings of Machine Learning Research, Vol. 225)*. PMLR, 353–367. Retrieved from https://proceedings.mlr.press/v225/moor23a.html

[61] Yohei Nakajima. 2023. BabyAGI. GitHub repository. https://github.com/yoheinakajima/babyagi

[62] Yuzhou Nie, Zhun Wang, Ye Yu, Xian Wu, Xuandong Zhao, Wenbo Guo, and Dawn Song. 2024. PrivAgent: Agentic-based Red-teaming for LLM Privacy Leakage. arXiv:2412.05734 (2024).

[63] OpenAI, Josh Achiam, Steven Adler, Sandhini Agarwal, Lama Ahmad, Ilge Akkaya, Florencia Leoni Aleman, Diogo Almeida, Janko Altenschmidt, Sam Altman, et al. 2024. GPT-4 Technical Report. arXiv:2303.08774 (2024).

[64] Mustafa Safa Ozdayi, Charith Peris, Jack FitzGerald, Christophe Dupuy, Jimit Majmudar, Haidar Khan, Rahil Parikh, and Rahul Gupta. 2023. Controlling the extraction of memorized data from large language models via prompt-tuning. arXiv:2305.11759 (2023).

[65] Jing-Cheng Pang, Heng-Bo Fan, Pengyuan Wang, Jia-Hao Xiao, Nan Tang, Si-Hang Yang, Chengxing Jia, Sheng-Jun Huang, and Yang Yu. 2024. Empowering language models with active inquiry for deeper understanding. arXiv:2402.03719 (2024).

[66] Joon Sung Park, Joseph O'Brien, Carrie Jun Cai, Meredith Ringel Morris, Percy Liang, and Michael S. Bernstein. 2023. Generative agents: Interactive simulacra of human behavior. In *Proceedings of the 36th Annual ACM Symposium on User Interface Software and Technology* (2023-10-29) (*UIST '23*). Association for Computing Machinery, 1–22. DOI:10.1145/3586183.3606763

[67] Liangzu Peng, Paris Giampouras, and Rene Vidal. 2023. The ideal continual learner: An agent that never forgets. In *Proceedings of the 40th International Conference on Machine Learning*. PMLR, 27585–27610. Retrieved from https://proceedings.mlr.press/v202/peng23a.html

[68] Jay Peters. 2023. *The Bing AI Bot Has Been Secretly Running GPT-4*. The Verge. Retrieved from https://www.theverge.com/2023/3/14/23639928/microsoft-bing-chatbot-ai-gpt-4-llm

[69] Xavier Puig, Eric Undersander, Andrew Szot, Mikael Dallaire Cote, Tsung-Yen Yang, Ruslan Partsey, Ruta Desai, Alexander William Clegg, Michal Hlavac, So Yeon Min, et al. 2023. Habitat 3.0: A co-Habitat for humans, avatars and robots. arXiv:2310.13724 (2023).

[70] Reworkd. 2023. AgentGPT. GitHub repository. https://github.com/reworkd/AgentGPT

[71] Alexander Robey, Eric Wong, Hamed Hassani, and George J. Pappas. 2023. SmoothLLM: Defending large language models against jailbreaking attacks. arXiv:2310.03684 (2023).

[72] Yangjun Ruan, Honghua Dong, Andrew Wang, Silviu Pitis, Yongchao Zhou, Jimmy Ba, Yann Dubois, Chris J. Maddison, and Tatsunori Hashimoto. 2024. Identifying the risks of LM agents with an LM-emulated sandbox. arXiv:2309.15817 (2024).

[73] Erfan Shayegani, Yue Dong, and Nael Abu-Ghazaleh. 2023. Jailbreak in pieces: Compositional adversarial attacks on multi-modal language models. arXiv:2307.14539 (2023).

[74] V. Shejwalkar and A. Houmansadr. 2021. Membership privacy for machine learning models through knowledge transfer. In *Proceedings of the AAAI Conference on Artificial Intelligence* 35, 11 (2021), 9549-9557.

[75] Significant Gravitas. 2023. AutoGPT. GitHub repository. https://github.com/Significant-Gravitas/AutoGPT

[76] Zhiqing Sun, Sheng Shen, Shengcao Cao, Haotian Liu, Chunyuan Li, Yikang Shen, Chuang Gan, Liang-Yan Gui, Yu-Xiong Wang, Yiming Yang, Kurt Keutzer, and Trevor Darrell. 2023. Aligning large multimodal models with factually augmented RLHF. arXiv:2309.14525 (2023).

[77] Yashar Talebirad and Amirhossein Nadiri. 2023. Multi-agent collaboration: Harnessing the power of intelligent LLM agents. arXiv:2306.03314 (2023).

[78] Meng Tong, Kejiang Chen, Jie Zhang, Yuang Qi, Weiming Zhang, Nenghai Yu, Tianwei Zhang, and Zhikun Zhang. 2024. InferDPT: Privacy-preserving inference for black-box large language model. arXiv:2310.12214 (2024).

[79] Jean-Baptiste Truong, Pratyush Maini, Robert J. Walls, and Nicolas Papernot. 2021. Data-free model extraction. In *Proceedings of the 2021 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*. 4769–4778. Retrieved from https://ieeexplore.ieee.org/document/9577784

[80] Thomas Vakili, Anastasios Lamproudis, Aron Henriksson, and Hercules Dalianis. 2022. Downstream task performance of bert models pre-trained using automatically de-identified clinical data. In *Proceedings of the 13th Language Resources and Evaluation Conference* (Marseille, France), Nicoletta Calzolari, Frédéric Béchet, Philippe Blache, Khalid Choukri, Christopher Cieri, Thierry Declerck, Sara Goggi, Hitoshi Isahara, Bente Maegaard, Joseph Mariani, Hélène Mazo, Jan Odijk, and Stelios Piperidis (Eds.). European Language Resources Association, 4245–4252. Retrieved from https://aclanthology.org/2022.lrec-1.451

[81] Steven Vander Eeckt and Hugo Van Hamme. 2023. Weight averaging: A simple yet effective method to overcome catastrophic forgetting in automatic speech recognition. In *Proceedings of the ICASSP 2023 - 2023 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)* (Rhodes Island, Greece). IEEE, 1–5. DOI:10.1109/ICASSP49357.2023.10095147

[82] Alexander Wan, Eric Wallace, Sheng Shen, and Dan Klein. 2023. Poisoning language models during instruction tuning. arXiv:2305.00944 (2023).

[83] Huan Wang and Yan-Fu Li. 2023. Large language model empowered by domain-specific knowledge base for industrial equipment operation and maintenance. In *Proceedings of the 2023 5th International Conference on System Reliability and Safety Engineering (SRSE)*. 474–479. DOI:10.1109/SRSE59585.2023.10336112

[84] Junyang Wang, Haiyang Xu, Jiabo Ye, Ming Yan, Weizhou Shen, Ji Zhang, Fei Huang, and Jitao Sang. 2024. Mobile-Agent: Autonomous multi-modal mobile device agent with visual perception. arXiv:2401.16158 (2024).

[85] Kuan Wang, Yadong Lu, Michael Santacroce, Yeyun Gong, Chao Zhang, and Yelong Shen. 2023. Adapting LLM agents through communication. arXiv:2310.01444 (2023).

[86] Lei Wang, Chen Ma, Xueyang Feng, Zeyu Zhang, Hao Yang, Jingsen Zhang, Zhiyuan Chen, Jiakai Tang, Xu Chen, Yankai Lin, Wayne Xin Zhao, Zhewei Wei, and Ji-Rong Wen. 2023. A survey on large language model based autonomous agents. arXiv:2308.11432 (2023).

[87] Shenzhi Wang, Chang Liu, Zilong Zheng, Siyuan Qi, Shuo Chen, Qisen Yang, Andrew Zhao, Chaofei Wang, Shiji Song, and Gao Huang. 2023. Avalon's game of thoughts: battle against deception through recursive contemplation. arXiv:2310.01320 (2023).

[88] Shen Wang, Tianlong Xu, Hang Li, Chaoli Zhang, Joleen Liang, Jiliang Tang, Philip S. Yu, and Qingsong Wen. 2024. Large language models for education: A survey and outlook. arXiv:2403.18105 (2024).

[89] Shang Wang, Tianqing Zhu, Bo Liu, Ming Ding, Dayong Ye, Wanlei Zhou, and Philip S. Yu. 2025. Unique security and privacy threats of large language models: A comprehensive survey. *ACM Computing Surveys* 58, 4 (2025), 1–36.

[90] Tianyu Wang, Yifan Li, Haitao Lin, Xiangyang Xue, and Yanwei Fu. 2023. WALL-E: Embodied robotic waiter load lifting with large language model. arXiv:2308.15962 (2023).

[91] Yuntao Wang, Yanghe Pan, Miao Yan, Zhou Su, and Tom H. Luan. 2023. A survey on ChatGPT: AI-generated contents, challenges, and solutions. *IEEE Open Journal of the Computer Society* 4, 1 (2023), 280–302.

[92] Yufei Wang, Wanjun Zhong, Liangyou Li, Fei Mi, Xingshan Zeng, Wenyong Huang, Lifeng Shang, Xin Jiang, and Qun Liu. 2023. Aligning large language models with human: A survey. arXiv:2307.12966 (2023).

[93] Zhenhua Wang, Wei Xie, Kai Chen, Baosheng Wang, Zhiwen Gui, and Enze Wang. 2023. Self-Deception: Reverse penetrating the semantic firewall of large language models. arXiv:2308.11521 (2023).

[94] Genta Indra Winata, Lingjue Xie, Karthik Radhakrishnan, Shijie Wu, Xisen Jin, Pengxiang Cheng, Mayank Kulkarni, and Daniel Preotiuc-Pietro. 2023. Overcoming catastrophic forgetting in massively multilingual continual learning. In *Proceedings of the Findings of the Association for Computational Linguistics: ACL 2023*. Association for Computational Linguistics, Toronto, Canada, 768–777. DOI:10.18653/v1/2023.findings-acl.48

[95] Chenfei Wu, Shengming Yin, Weizhen Qi, Xiaodong Wang, Zecheng Tang, and Nan Duan. 2023. Visual ChatGPT: Talking, drawing and editing with visual foundation models. arXiv:2303.04671 (2023).

[96] Zhiheng Xi, Wenxiang Chen, Xin Guo, Wei He, Yiwen Ding, Boyang Hong, Ming Zhang, Junzhe Wang, Senjie Jin, Enyu Zhou, et al. 2023. The rise and potential of large language model based agents: A survey. arXiv:2309.07864 (2023).

[97] Junlin Xie, Zhihong Chen, Ruifei Zhang, Xiang Wan, and Guanbin Li. 2024. Large multimodal agents: A survey. arXiv:2402.15116 (2024).

[98] Lin Xu, Zhiyuan Hu, Daquan Zhou, Hongyu Ren, Zhen Dong, Kurt Keutzer, See Kiong Ng, and Jiashi Feng. 2023. MAgIC: Investigation of large language model powered multi-agent in cognition, adaptability, rationality and collaboration. arXiv:2311.08562 (2023).

[99] Lvxiaowei Xu, Jianwang Wu, Jiawei Peng, Zhilin Gong, Ming Cai, and Tianxiang Wang. 2023. Enhancing language representation with constructional information for natural language understanding. In *Proceedings of the 61st Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)* (Toronto, Canada). Association for Computational Linguistics, 4685–4705. DOI:10.18653/v1/2023.acl-long.258

[100] Mingfu Xue, Chengxiang Yuan, Heyi Wu, Yushu Zhang, and Weiqiang Liu. 2020. Machine learning security: threats, countermeasures, and evaluations. *IEEE Access* 8 (2020), 74720–74742.

[101] Lu Yan, Zhuo Zhang, Guanhong Tao, Kaiyuan Zhang, Xuan Chen, Guangyu Shen, and Xiangyu Zhang. 2023. ParaFuzz: An interpretability driven technique for detecting poisoned samples in NLP. arXiv:2308.02122 (2023).

[102] Haomiao Yang, Kunlan Xiang, Mengyu Ge, Hongwei Li, Rongxing Lu, and Shui Yu. 2023. A comprehensive overview of backdoor attacks in large language models within communication networks. arXiv:2308.14367 (2023).

[103] Jihan Yang, Runyu Ding, Ellis Brown, Xiaojuan Qi, and Saining Xie. 2024. V-IRL: Grounding virtual intelligence in real life. arXiv:2402.03310 (2024).

[104] Wenkai Yang, Xiaohan Bi, Yankai Lin, Sishuo Chen, Jie Zhou, and Xu Sun. 2024. Watch out for your agents! Investigating backdoor threats to LLM-based agents. arXiv:2402.11208 (2024).

[105] Zhengyuan Yang, Linjie Li, Jianfeng Wang, Kevin Lin, Ehsan Azarnasab, Faisal Ahmed, Zicheng Liu, Ce Liu, Michael Zeng, and Lijuan Wang. 2023. MM-REACT: Prompting ChatGPT for multimodal reasoning and action. arXiv:2303.11381 (2023).

[106] Shunyu Yao, Dian Yu, Jeffrey Zhao, Izhak Shafran, Thomas L. Griffiths, Yuan Cao, and Karthik Narasimhan. 2023. Tree of thoughts: Deliberate problem solving with large language models. In *Proceedings of the 37th International Conference on Neural Information Processing Systems* (Red Hook, NY, USA) *(NIPS '23)*. Curran Associates Inc., 11809–11822.

[107] Shunyu Yao, Jeffrey Zhao, Dian Yu, Nan Du, Izhak Shafran, Karthik Narasimhan, and Yuan Cao. 2023. ReAct: Synergizing reasoning and acting in language models. arXiv:2210.03629 (2023).

[108] Yifan Yao, Jinhao Duan, Kaidi Xu, Yuanfang Cai, Eric Sun, and Yue Zhang. 2023. A survey on large language model (LLM) security and privacy: The good, the bad, and the ugly. arXiv:2312.02003 (2023).

[109] Jingwei Yi, Yueqi Xie, Bin Zhu, Keegan Hines, Emre Kiciman, Guangzhong Sun, Xing Xie, and Fangzhao Wu. 2023. Benchmarking and defending against indirect prompt injection attacks on large language models. arXiv:2312.14197 (2023).

[110] Shukang Yin, Chaoyou Fu, Sirui Zhao, Tong Xu, Hao Wang, Dianbo Sui, Yunhang Shen, Ke Li, Xing Sun, and Enhong Chen. 2023. Woodpecker: Hallucination correction for multimodal large language models. arXiv:2310.16045 (2023).

[111] Sheng Yin, Xianghe Pang, Yuanzhuo Ding, Menglan Chen, Yutong Bi, Yichen Xiong, Wenhao Huang, Zhen Xiang, Jing Shao, and Siheng Chen. 2024. SafeAgentBench: A benchmark for safe task planning of embodied LLM agents. arXiv:2412.13178 (2024).

[112] Ziyi Yin, Muchao Ye, Tianrong Zhang, Tianyu Du, Jinguo Zhu, Han Liu, Jinghui Chen, Ting Wang, and Fenglong Ma. 2024. VLATTACK: Multimodal adversarial attacks on vision-language tasks via Pre-trained models. arXiv:2310.04655 (2024).

[113] Da Yu, Saurabh Naik, Arturs Backurs, Sivakanth Gopi, Huseyin A. Inan, Gautam Kamath, Janardhan Kulkarni, Yin Tat Lee, Andre Manoel, Lukas Wutschitz, et al. 2021. Differentially private fine-tuning of language models. In

*Proceedings of the International Conference on Learning Representations*. Retrieved from https://openreview.net/forum?id=Q42f0dfjECO

[114] Jiahao Yu, Xingwei Lin, Zheng Yu, and Xinyu Xing. 2023. GPTFUZZER: Red teaming large language models with auto-generated jailbreak prompts. arXiv:2309.10253 (2023).

[115] Ziqi Yuan, Yihe Liu, Hua Xu, and Kai Gao. 2024. Noise imitation based adversarial training for robust multimodal sentiment analysis. *IEEE Transactions on Multimedia* 26 (2024), 529–539.

[116] Bohan Zhai, Shijia Yang, Chenfeng Xu, Sheng Shen, Kurt Keutzer, Chunyuan Li, and Manling Li. 2024. HallE-Control: Controlling object hallucination in large multimodal models. arXiv:2310.01779 (2024).

[117] Yuexiang Zhai, Shengbang Tong, Xiao Li, Mu Cai, Qing Qu, Yong Jae Lee, and Yi Ma. 2023. Investigating the catastrophic forgetting in multimodal large language models. arXiv:2309.10313 (2023).

[118] Qiusi Zhan, Zhixiang Liang, Zifan Ying, and Daniel Kang. 2024. InjecAgent: Benchmarking indirect prompt injections in tool-integrated large language model agents. arXiv:2403.02691 (2024).

[119] Shuo Zhang, Liangming Pan, Junzhou Zhao, and William Yang Wang. 2024. The knowledge alignment problem: Bridging human and external knowledge for large language Models. arXiv:2305.13669 (2024).

[120] Yuyang Zhang, Kangjie Chen, Xudong Jiang, Yuxiang Sun, Run Wang, and Lina Wang. 2024. Towards action hijacking of large language model-based agent. arXiv:2412.10807 (2024).

[121] Zhexin Zhang, Shiyao Cui, Yida Lu, Jingzhuo Zhou, Junxiao Yang, Hongning Wang, and Minlie Huang. 2024. Agent-safety bench: Evaluating the safety of LLM agents. arXiv:2412.14470 (2024).

[122] Zhiping Zhang, Michelle Jia, Hao-Ping (Hank) Lee, Bingsheng Yao, Sauvik Das, Ada Lerner, Dakuo Wang, and Tianshi Li. 2024. "It's a Fair Game", or Is It? Examining How Users Navigate Disclosure Risks and Benefits When Using LLM-Based Conversational Agents. In *Proceedings of the CHI Conference on Human Factors in Computing Systems* (New York, NY, USA) *(CHI '24)*. Association for Computing Machinery, 1–26. DOI : 10.1145/3613904.3642385

[123] Zhexin Zhang, Junxiao Yang, Pei Ke, Fei Mi, Hongning Wang, and Minlie Huang. 2024. Defending large language models against jailbreaking attacks through goal prioritization. In *Proceedings of the 62nd Annual Meeting of the Association for Computational Linguistics*. 8865–8887. DOI : 10.18653/v1/2024.acl-long.481

[124] Qiwei Zhao, Xujiang Zhao, Yanchi Liu, Wei Cheng, Yiyou Sun, Mika Oishi, Takao Osaki, Katsushi Matsuda, Huaxiu Yao, and Haifeng Chen. 2024. SAUP: Situation awareness uncertainty propagation on LLM agent. arXiv:2412.01033 (2024).

[125] Qingxiao Zheng, Zhongwei Xu, Abhinav Choudhary, Yuting Chen, Yongming Li, and Yun Huang. 2023. Synergizing human-AI agency: A guide of 23 heuristics for service co-creation with LLM-based agents. arXiv:2310.15065 (2023).

[126] Wanjun Zhong, Ruixiang Cui, Yiduo Guo, Yaobo Liang, Shuai Lu, Yanlin Wang, Amin Saied, Weizhu Chen, and Nan Duan. 2023. AGIEval: A human-centric benchmark for evaluating foundation models. arXiv:2304.06364 (2023).

[127] Wanjun Zhong, Lianghong Guo, Qiqi Gao, He Ye, and Yanlin Wang. 2024. Memorybank: Enhancing large language models with long-term memory. In *Proceedings of the AAAI Conference on Artificial Intelligence*, Vol. 38. 19724–19731.

[128] Wei Zou, Runpeng Geng, Binghui Wang, and Jinyuan Jia. 2025. PoisonedRAG: Knowledge poisoning attacks to retrieval-augmented generation of large language models. In *Proceedings of the 34th USENIX Security Symposium (USENIX Security 25)*. 3827–3844.