

# Multi-Modal and Multi-Agent Systems Meet Rationality: A Survey

Anonymous ACL submission

## Abstract

Rationality is the quality of being guided by reason, characterized by decision-making aligned with evidence and logical rules. This quality is essential for effective problem-solving, as it ensures that solutions are well-founded and consistently derived. Despite the advancements of large language models (LLMs) in generating human-like texts with remarkable accuracy, they present limited knowledge space, inconsistency across contexts, and difficulty understanding complex scenarios. Therefore, recent research focuses on building multi-modal and multi-agent systems to achieve considerable progress with enhanced consistency and reliability, instead of relying on a single LLM as the sole planning or decision-making agent. To that end, this paper aims to understand whether multi-modal and multi-agent systems are advancing toward rationality by surveying the state-of-the-art works, identifying advancements over single-agent and single-modal systems in terms of rationality, and discussing open problems and future directions.

## 1 Introduction

Large language models (LLMs) have demonstrated promising results across a broad spectrum of tasks, particularly in exhibiting capabilities that plausibly mimic human-like reasoning (Wei et al., 2022; Yao et al., 2024; Besta et al., 2024; Shinn et al., 2024; Bubeck et al., 2023; Valmeekam et al., 2023; Prasad et al., 2023). These models leverage the richness of human language to abstract concepts, elaborate thinking process, comprehend complex user queries, and develop plans and solutions in decision-making scenarios. Despite these advances, recent research has revealed that even state-of-the-art LLMs exhibit various forms of irrational behaviors, such as the framing effect, certainty effect, overweighting bias, and conjunction fallacy (Binz and Schulz, 2023; Echterhoff et al., 2024; Mukherjee and Chang, 2024; Macmillan-Scott and Mu-

solesi, 2024; Wang et al., 2024a; Suri et al., 2024). Irrationality undermine the practical deployment of LLMs in critical sectors like healthcare, finance, and legal services (He et al., 2023; Li et al., 2023h; Kang and Liu, 2023; Cheong et al., 2024), where reliability and consistency are paramount. The emerging concern about the factual accuracy and trustworthiness of LLMs highlights an urgent need to develop better agents or agent systems (Nakajima, 2023; Gravitass, 2023) with rational reasoning processes.

A single LLM agent can fall into irrational behaviors because it cannot go beyond the language model’s inner parametric representations of textual knowledge, lacking the real-world grounding and feedback mechanisms necessary to develop rationality (Bubeck et al., 2023; Sun, 2024). In contrast, in real life scenarios, important decisions are rarely made by individuals on their own, and the complexity of problems often requires the collaboration of experts from different fields to ensure rationality (Eisenführ et al., 2010). In a similar vein, recent advancements in multi-modal and multi-agent frameworks leverage the expertise of different agents acting together towards a collective goal. Multi-modal foundation models (Awadalla et al., 2023; Liu et al., 2023a; Wang et al., 2023c; OpenAI, 2023; Reid et al., 2024) enhance reasoning by grounding decisions in a broader sensory context, akin to how human brains integrate rich sensory inputs to form a more holistic base of knowledge. Meanwhile, multi-agent systems introduce mechanisms such as consensus, debate, and self-consistency (Du et al., 2023; Liang et al., 2023; Talebirad and Nadiri, 2023; Madaan et al., 2024; Cohen et al., 2023; Shinn et al., 2024; Mohtashami et al., 2023) to allow for refined and reliable output through collaborative interactions. Such systems can also query external knowledge sources or tools (Lewis et al., 2020; Schick et al., 2024; Tang et al., 2023; Pan et al., 2024) to augment their

084	reasoning capabilities for rational decision making.		
085	This survey provides a unique lens to reinterpret		
086	the underlying motivations behind current multi-		
087	modal and/or multi-agent systems by drawing in-		
088	sights from cognitive science. In Section 2, we		
089	outline four essential requirements for rational de-		
090	cision making. Section 4 then examines how vari-		
091	ous research areas within the multi-modality and		
092	multi-agent literature are advancing towards ra-		
093	tionality based on these criteria. We argue that		
094	these advancements surpass the limitations of sin-		
095	gle language-model agents and narrow the gap be-		
096	tween the behavior of agent systems and the expec-		
097	tations for rational decision making. Lastly, Sec-		
098	tion 5 highlights the lack of sufficient evaluation		
099	metrics and benchmarks in the existing literature		
100	to adequately measure the rationality of LLMs or		
101	agent systems. We hope this survey can inspire		
102	further research at the intersection between agent		
103	systems and cognitive science.		
104	<b>2 Defining Rationality</b>		
105	A rational agent, in short, should respect the reality		
106	of the world in which it operates and avoid reaching		
107	contradictory conclusions in decision-making pro-		
108	cesses. Drawing on foundational works in rational		
109	decision-making (Tversky and Kahneman, 1988;		
110	Hastie and Dawes, 2009; Eisenführ et al., 2010),		
111	this section adopts an axiomatic approach to de-		
112	fine rationality, presenting four substantive axioms		
113	that we expect a rational agent or agent systems to		
114	fulfill:		
115	<b>Grounding</b> The decision of a rational agent is		
116	grounded on the physical and factual reality. For		
117	example, a video generation agent should adhere to		
118	the laws of physics in a world model and a forecast-		
119	ing assistant ought to estimate likelihoods obeying		
120	the law of probability.		
121	<b>Orderability of Preferences</b> When comparing		
122	alternatives in a decision scenario, a rational agent		
123	can rank the options based on the current state and		
124	ultimately select the most preferred one based on		
125	the expected outcomes. This orderability consists		
126	of several key principles, including comparability,		
127	transitivity closure, solvability, etc. with detailed		
128	defined in Appendix A.		
129	<b>Independence from irrelevant context</b> The		
130	agent’s preference should not be influenced by in-		
131	formation irrelevant to the decision-making prob-		
132	lem at hand.		
	<b>Invariance</b> The preference of a rational agent re-		133
	mains invariant across equivalent representations of		134
	the decision problem, regardless of specific word-		135
	ings or modalities.		136
	<b>3 Scope</b>		137
	Unlike existing surveys (Han et al., 2024; Guo et al.,		138
	2024; Xie et al., 2024a; Durante et al., 2024; Cui		139
	et al., 2024; Xu et al., 2024b; Zhang et al., 2024a;		140
	Cheng et al., 2024; Li et al., 2024a) that focus on		141
	the components, structures, agent profiling, plan-		142
	ning, communications, memories, and applications		143
	of multi-modal and/or multi-agent systems, <b>this</b>		144
	<b>survey is the first to specifically examine the</b>		145
	<b>increasingly important relations between ratio-</b>		146
	<b>nality and these multi-modal and multi-agent</b>		147
	<b>systems</b> , exploring how they contribute to enhanc-		148
	ing the robustness in decision-making processes.		149
	We emphasize that <i>rationality</i> , by definition, is		150
	not equivalent to <i>reasoning</i> (Kharden and Roth,		151
	1997; Huang and Chang, 2022; Zhang et al., 2024a;		152
	Qiao et al., 2022), although deeply intertwined.		153
	Rationality involves making logically consistent deci-		154
	sions grounded with reality, while reasoning refers		155
	to the cognitive process of drawing logical infer-		156
	ences and conclusions from available information,		157
	as illustrated in the following thought experiment:		158
	<i>Consider an environment where the in-</i>		159
	<i>put space and the output decision space</i>		160
	<i>are finite. A lookup table with consis-</i>		161
	<i>tent mapping from input to output is in-</i>		162
	<i>herently rational, while no reasoning is</i>		163
	<i>necessarily present in the mapping.</i>		164
	Despite this example, it is still crucial to ac-		165
	knowledge that reasoning typically plays a vital		166
	role in ensuring rationality, especially in complex		167
	and dynamic real-world scenarios where a sim-		168
	ple lookup table is insufficient. Agents must pos-		169
	sess the ability to reason through novel situations,		170
	adapt to changing circumstances, make plans, and		171
	achieve rational decisions based on incomplete or		172
	uncertain information.		173
	<b>4 Towards Rationality through</b>		174
	<b>Multi-Modal and Multi-Agent Systems</b>		175
	This section surveys recent advancements in multi-		176
	modal and multi-agent systems under different re-		177
	search categories as depicted in Figure 1. Each		178
	category, such as knowledge retrieval or neuro-		179
	symbolic reasoning, addresses one or more fun-		180
	damental requirements for rational thinking. These		181

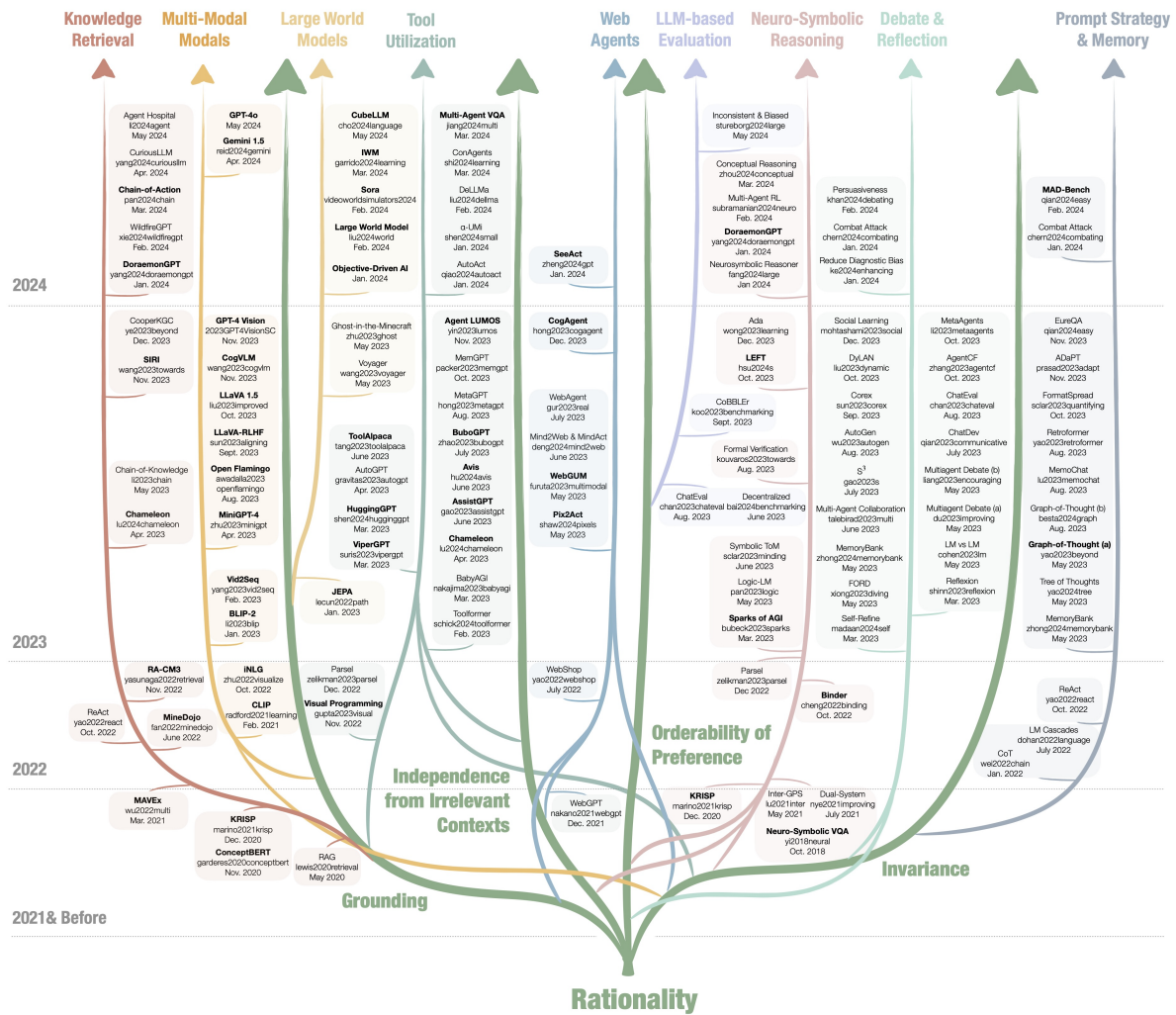


Figure 1: The evolutionary tree of multi-agent and/or multi-modal systems related to the four axioms of rationality. Many proposed approaches strive to address multiple axioms simultaneously. The **bold font** marks works that involve multi-modalities. This tree also includes some foundational works to provide a clearer reference of time.

rationality requirements are typically *intertwined*: an approach that enhances one aspect of rationality often inherently improves others simultaneously. Meanwhile, the overall mechanism of current multi-agent system in achieving rationality can be categorized into two key concepts: **deliberation** and **abstraction**. Deliberation encourages a slower, iterative reasoning process, while abstraction refers to abstracting the problem into its logical essence.

Most existing studies do not explicitly base their frameworks on rationality in their original writings. Our analysis aims to reinterpret these works through the lens of our four axioms of rationality, offering a novel perspective that bridges existing methodologies with rational principles.

#### 4.1 Towards Grounding & Invariance

Multi-modal approaches aim to improve information grounding across various channels, such as lan-

guage, vision, and beyond. By incorporating multi-modal agents, multi-agent systems can greatly expand their capabilities, enabling a richer, more accurate, and contextually aware interpretation of the environment.

**Multi-Modal Foundation Models** Grounding an agent solely based on textual language can be challenging, as information can be represented much more efficiently through other sensory modes. As a picture is worth a thousand words, recent advances in large vision-language pretraining have enabled LLMs with robust language comprehension capabilities to finally perceive the visual world. Multi-modal foundation models, including but not limited to CLIP (Radford et al., 2021), VLBERT and ViLBERT (Su et al., 2019; Lu et al., 2019), BLIP-2 (Li et al., 2023d), (Open) Flamingo (Alayrac et al., 2022; Awadalla et al., 2023), LLaVA (Liu et al.,

2024c, 2023a), CogVLM (Wang et al., 2023c), MiniGPT-4 (Zhu et al., 2023a), GPT-4 Vision (OpenAI, 2023) and GPT-4o (OpenAI, 2024), and Gemini 1.5 Pro (Reid et al., 2024) serve as the cornerstones for multi-modal agent systems to ground knowledge in vision and beyond.

**Invariance Across Modalities** Achieving representation invariance across modalities is critical: given adequate information grounding, agents should make consistent decisions across different modalities that share equivalent underlying logic. Multi-modal foundation models are particularly adept at promoting invariance by processing multi-modal data in a unified representation. Specifically, their large-scale cross-modal pretraining stage seamlessly tokenizes both vision and language inputs into a joint hidden embedding space, learning cross-modal correlations through a data-driven approach. In other words, image tokens are simply regarded as a foreign language (Wang et al., 2022a). Moreover, the cross-modal validation inherent in multi-modal foundation models allows for reconciliation of data from different modalities, closing their distance in the hidden embedding space (Radford et al., 2021).

The concept of invariance is the cornerstone of Visual Question Answering (VQA) agents (Chen et al., 2022; Jiang et al., 2024; Wang et al., 2023d; Yi et al., 2018; Wang et al., 2022a; Bao et al., 2022; Zhao and Xu, 2023). On one hand, these agents must grasp the invariant semantics of any open-ended questions posed about images, maintaining consistency despite variations in wording, syntax, or language. On the other hand, within a multi-agent VQA system, visual agents can provide crucial verification and support for language-based reasoning (Wang et al., 2023d; Jiang et al., 2024; Zhao and Xu, 2023), while language queries can direct the attention of visual agents, based on a shared and invariant underlying knowledge across vision and language domains.

**Information Grounding** Multi-modalities help enhance the functionality of agent systems through more diverse information grounding. Web agents are a quintessential example of how multi-modal agents surpass language-only ones. Because HTML code is often lengthy, contains irrelevant information, and may be incomplete (Zheng et al., 2024a), web navigation grounded on the graphical user interface (GUI) offers higher information density compared to solely HTML codes. As a result,

using visual cues from the GUI leads to improved navigation performance (Shen et al., 2024a; Yao et al., 2022a; Deng et al., 2024; Gur et al., 2023). Multi-modalities also help enhance the functionality of agent systems through more diverse information grounding. For example, RA-CM3 (Yasunaga et al., 2022) augments baseline retrieval-augmented LLMs with raw multi-modal documents that include both images and texts, assuming that these two modalities can contextualize each other and make the documents more informative, leading to better generator performance. For other examples, we refer the reader to Appendix B.

**Knowledge Retrieval & Tool Usage** Bounded Rationality (March and Simon, 1958; Selten, 1990) is a concept tailored to cognitively limited agents, suggesting that decision-making is limited by the resources available at hand, and any deviations from the optimal are primarily due to insufficient computational capacity and bounded working memory. In terms of LLMs, the parametric nature of their existing architecture (Vaswani et al., 2017) fundamentally limits how much information they can hold. As a result, in the face of uncertainty, LLMs often hallucinate (Bang et al., 2023; Guerreiro et al., 2023; Huang et al., 2023), generating outputs that are not supported by the factual reality of the environment. Retrieval-Augmented Generation (RAG) (Lewis et al., 2020) marks a significant milestone in addressing such an inherent limitation of LLMs. Broadly speaking, RAG refers to any mechanism that provides external knowledge to the input context of an LLM and helps it deliver responses with up-to-date, factual, and grounded information, especially in scientific and medical domains. Examples include Chameleon (Lu et al., 2024), Chain-of-Knowledge (Li et al., 2023f), WildfireGPT (Xie et al., 2024b), and Agent Hospital (Li et al., 2024b). Specifically, Chain-of-Knowledge (Li et al., 2023f) even discovers that integrating multiple knowledge sources enhances performance by 2.1% compared to using a single source in its experiments.

Another line of systems construct large-scale knowledge graphs (KGs) (Hogan et al., 2021) from real-world sources to effectively expand their working memory and improve their task performance. Specifically, compared to language-only models, MAVEx (Wu et al., 2022) improves system’s scores by 9.5% compared to an image-only baseline through the integration of knowl-

edge from ConceptNet (Speer et al., 2017) and Wikipedia (Wikipedia contributors, 2004). It also improves the scores by 8.3% by using the image modality for cross-modal validations with an oracle. Thanks to the external knowledge base, ReAct (Yao et al., 2022b) reduces false positive rates from hallucination by 8.0% compared to CoT (Wei et al., 2022). For more examples, see Appendix C.

Enabling agents to use tools also expands their bounded working memories, akin to retrieving external knowledge. Toolformer (Schick et al., 2024) opens a new era that allows LLMs to use external tools via API calls following predefined syntax, effectively extending their capabilities beyond their intrinsic limitations and enforcing consistent and predictable outputs. A multi-agent system can coordinate agents understanding when and which tool to use, which modality of information the tool should expect, how to call the corresponding API, and how to incorporate outputs from the API calls, which anchors subsequent reasoning processes with more accurate information beyond their parametric memory. For example, VisProg (Gupta and Kembhavi, 2023), ViperGPT (Surís et al., 2023), and Parsel (Zelikman et al., 2023) generate Python programs to reliably execute subroutines. Gupta and Kembhavi (2023); Surís et al. (2023) also invoke off-the-shelf models for multimodal assistance. For more examples, see Appendix C.

## 4.2 Towards Rationality through Deliberation

Memory is one of the most fundamental cognitive processes that lead to reasoning, creativity, learning, and even self-consciousness in humans (Solso and Kagan, 1979; Craik and Lockhart, 1972; Leydesdorff and Hodgkin, 2017; Johnson-Laird, 1983; Laird, 2019; Sun, 2001). Any system that lacks the ability to retain information from previous interactions would struggle to make coherent and rational decisions in the long run. In a narrow sense, agent memory includes historical information within the same conversation (Zhang et al., 2024b). This allows for **deliberation**, which is the slower, iterative reasoning process to carefully consider information and options in order to arrive at more rational and well-reasoned decisions.

Although deliberation may increase the likelihood of reaching more rational decisions, there is no inherent guarantee for rationality via this approach. The quality of the decision ultimately depends on the accuracy and relevance of the grounded information, as well as the soundness

of the reasoning process. Biases, incomplete information, and flawed logic can still lead to irrational conclusions even with deliberation. Nonetheless, multiple works have shown that the increase in likelihood of rational decisions through deliberation is significant and beneficial. For example, multi-round self-reflection prompting strategies that encourage agents to critically evaluate their previous responses (Shinn et al., 2024; Madaan et al., 2024; Wang et al., 2022b; Zhong et al., 2024; Lu et al., 2023).

In a broader context, for multi-agent systems, agent memory expands to include historical information across multiple tasks and agents (Zhang et al., 2024b). This shared memory enables collective deliberation among agents via collaboration, cross-examination, and debates. By leveraging the collective knowledge and experiences of multiple agents, the system can arrive at more comprehensive and robust solutions to complex problems.

LM vs LM (Cohen et al., 2023), FORD (Xiong et al., 2023), Multi-Agent Debate (Liang et al., 2023; Du et al., 2023), DyLAN (Liu et al., 2023c), and Khan et al. (2024) highlight the profound impact of **multi-agent collaboration through cross-examination and debates**. Specifically, LM vs LM (Cohen et al., 2023) illustrates how its multi-agent framework improves F1 scores by an average of 15.7% compared to the single-agent baseline (Yoshikawa and Okazaki, 2023). FORD (Xiong et al., 2023) reports an accuracy increase up to 4.9% compared to a single LLM. Liang et al. (2023) indicates significant improvements in accuracy — 17.0% for translation tasks and 16.0% for reasoning tasks — by employing a multi-agent strategy, effectively bridging the performance gap between GPT-3.5 and GPT-4 by harnessing multi-agents. Du et al. (2023) finds that multi-agent debates not only enhance reasoning performance by 8.0-14.8%, but more importantly, increase factual accuracy by 7.2-15.9%. DyLAN (Liu et al., 2023c) observes 3.5-4.1% in accuracy improvements over single-agent execution. All these approaches enhance the system’s capability to capture initial errors, improve factuality in reasoning, and achieve final consensus with fewer inconsistencies. We discuss more examples in Appendix D.1. We also talk about collaboration against jailbreaking in Appendix D.2 and multi-agent evaluation methods in Appendix D.3.

### 4.3 Towards Rationality through Abstraction

Independence from irrelevant contexts, invariance, and orderability of preferences can be achieved simultaneously through the use of tools and **neuro-symbolic reasoning**, because these approaches translate natural language queries into standardized formats like API calls or symbolic representations, which *abstract* away extraneous details, focus only on the underlying logic necessary for the task at hand, and enable consistent and deterministic processing of the input.

**Independence from Irrelevant Contexts** In most cases, tools require translating natural language queries into API calls with predefined syntax. Once the APIs and their input arguments are determined, the tools will ignore any irrelevant context in the original queries, as long as the queries share the same underlying logic necessary for the inputs. Take Multi-Agent VQA (Jiang et al., 2024) as an example. In this system, a language model provides only the relevant object names to the Grounded SAM (Ren et al., 2024) component, which functions as an object detector, rather than passing the entire visual question. Other similar examples are discussed in Appendix C.

Neuro-symbolic reasoning is an approach that combines neural networks with symbolic systems, such as explicit knowledge representation and logical reasoning. A multi-agent system incorporating symbolic modules can not only understand language queries but also solve them with a level of consistency, providing a faithful and transparent reasoning process based on well-defined rules that adhere to logical principles, which is unachievable by a single language model.

Analogous to the external tool utilization, neuro-symbolic modules in a multi-agent system expect standardized input formats (Zelikman et al., 2023; Pan et al., 2023; Sclar et al., 2023b; Hsu et al., 2024; Fang et al., 2024; Yang et al., 2024; Subramanian et al., 2024). The only relevant factor in this process is the parsed inputs into the predetermined neuro-symbolic programs. For instance, Ada (Wong et al., 2023) introduces symbolic operators to abstract actions, ensuring that lower-level planning models are not compromised by irrelevant information in the queries and observations. Without the symbolic action library, a single LLM would frequently fail at grounding objects or obeying environmental conditions, resulting in a significant accuracy gap of approximately 59.0-89.0%.

**Invariance** The abstraction provided by symbolic representations also allows the multi-agent system to solve language queries with invariance. For example, Logic-LM (Pan et al., 2023) combines problem formulating, symbolic reasoning, and result interpreting agents, where the symbolic reasoner empowers LLMs with deterministic symbolic solvers to perform inference, ensuring a correct answer is consistently chosen. Its multi-agent framework also encourages self-refinement that modifies logical formulation errors using error messages from the symbolic reasoner as the feedback. For more examples, see Appendix E.

**Orderability of Preferences** In explicit scenarios, logical modules can directly compare the order of options represented as variables—such as “left” or “right” in relational logic (Hsu et al., 2024)—rather than relying on a single LLM to generate responses indeterministically within the natural language space. In more complex situations, systems like Binder (Cheng et al., 2022), Parsel (Zelikman et al., 2023), LEFT (Hsu et al., 2024), and Fang et al. (2024) decompose tasks into planning, parsing, and execution, where the symbolic reasoning agents can help maintain a coherent order of preferences among symbolic options in the system outputs. By skipping the symbolic module, Parsel (Zelikman et al., 2023) observes a substantial performance drop of 19.5%. LEFT (Hsu et al., 2024) also outperforms end-to-end baselines without symbolic programs by 3.85% on average across multiple experiments.

Recent work has also explored applying expected utility theory (Von Neumann and Morgenstern, 2007) to improve the decision-making capabilities of language models. For example, DeLLMa (Liu et al., 2024e) decomposes complex decision problems into subtasks, assigns utilities to different outcomes, and selects actions that maximize expected utility.

## 5 Evaluating Rationality of Agents

The amount of studies for testing rationality in multi-modal and multi-agent systems remains scant, despite the growing interest in the field. While there are numerous reasoning benchmarks available (Talmor et al., 2019; Liu et al., 2021; Yang et al., 2018; Hendrycks et al., 2021), they do not directly measure rationality. Many of these benchmarks fail to prove whether reasoning is actually used in solving the tasks, leaving no guarantee that

522 these tasks will be solved consistently when gener- 573  
523 alized to other representations or domains. Issues 574  
524 such as data contamination (Magar and Schwartz, 575  
525 2022; Dong et al., 2024; Sainz et al., 2023; Jacovi 576  
526 et al., 2023) further compound the problem, as  
527 some benchmarks may inadvertently include the  
528 training data of these LLMs, leading to inflated  
529 performance scores. Hence, even though solid reason-  
530 ing will imply rationality, existing approaches  
531 fall short in making the logic click. In this sec-  
532 tion, we point to several existing ingredients that  
533 can constitute the bread-and-butter of future gener-  
534 ations of evaluation approaches for rationality.

### 535 **Adapting Cognitive Psychology Experiments**

536 Recent works propose adapting vignette-based ex- 577  
537 periments borrowed from cognitive psychology to 578  
538 test whether LLMs are susceptible to cognitive bi- 579  
539 ases and fallacies. For instance, Binz and Schulz 580  
540 (2023) tested GPT-3 on the conjunction fallacy, 581  
541 finding that they exhibit human-like biases. How- 582  
542 ever, many of these approaches are informal and 583  
543 subjective, failing to scale in a way that allows for 584  
544 drawing statistically significant conclusions. More- 585  
545 over, LLMs may be subject to cognitive biases not 586  
546 existent in humans, such as the hypothetical "al- 587  
547 gorithmic bias" proposed by Bender et al. (2021), 588  
548 which could lead to unintended consequences in 589  
549 decision-making tasks. Further research is needed 590  
550 to uncover and characterize these potential biases. 591

551 **Testing against Hallucination** Information 592  
552 grounding is usually evaluated by the level of 593  
553 hallucination (Bang et al., 2023; Guerreiro et al., 594  
554 2023; Huang et al., 2023). Multiple evaluation 595  
555 benchmarks targeting language-only dialogue 596  
556 have been proposed, such as BEGIN (Dziri 597  
557 et al., 2022b), HaluEval (Li et al., 2023e), 598  
558 DialFact (Gupta et al., 2021), FaithDial (Dziri 599  
559 et al., 2022a), AIS (Rashkin et al., 2023), and 600  
560 others (Zheng et al., 2023b; Das et al., 2023; 601  
561 Cao et al., 2021). In contrast, **benchmarks on** 602  
562 **multi-agent frameworks or those involving** 603  
563 **multi-modalities beyond language dialogue** 604  
564 **are very limited.** We find that Liu et al. (2024a) 605  
565 moves beyond conversation to code generation, 606  
566 EureQA (Li et al., 2023a) focuses on reasoning 607  
567 chains, and TofuEval (Tang et al., 2024) evaluates 608  
568 hallucination in multi-domain summarization. 609  
569 Object hallucination (Rohrbach et al., 2018; 610  
570 Biten et al., 2022), POPE (Li et al., 2023g), 611  
571 and LLaVA-RLHF (Sun et al., 2023b) are the 612  
572 few examples evaluating multi-modal hallucina-

tion. The community needs more hallucination 573  
574 benchmarks to quantitatively evaluate the extent 575  
576 to which multi-modal and multi-agents reduce 576  
577 hallucinations in comparison with baselines.

**Testing the Orderability of Preference** There 577  
578 are almost no benchmarks for evaluating whether 578  
579 LLMs or agents have a consistent preference in 579  
580 the selection of available options. The Multiple 580  
581 Choice Problem (MCP) serves as a common test- 581  
582 ing ground. Zheng et al. (2023a) shows that LLMs 582  
583 are susceptible to changes in the positioning of 583  
584 options. Since the underlying logic remains the 584  
585 same, it also makes LLMs fail to pass the prop- 585  
586 erty of invariance. Although there are many MCP 586  
587 benchmarks (PaperswithcodeMCQA), they focus 587  
588 on the accuracy of selections and overlook the con- 588  
589 sistency of preference. However, Robinson et al. 589  
590 (2023) highlights that the Proportion of Plurality 590  
591 Agreement (PPA) offers a measure of order invari- 591  
592 ance that does not depend on the model’s ability to 592  
593 perform a task, suggesting a promising direction.

**Testing the Principle of Invariance** Recent 594  
595 studies concerning data contamination investigate 595  
596 whether LLMs can generate consistent responses 596  
597 across different, yet inherently equivalent, framing 597  
598 of the same task. These studies introduce pertur- 598  
599 bations to the original task descriptions to assess 599  
600 whether LLMs’ responses will change significantly. 600  
601 Perturbation techniques include modifying instruc- 601  
602 tion templates (Weber et al., 2023), paraphrasing 602  
603 task descriptions (Yang et al., 2023; Ohmer et al., 603  
604 2024), or altering the order of in-context learning 604  
605 exemplars (Lu et al., 2021; Pecher et al., 2024). 605  
606 For more details on these techniques, we refer the 606  
607 reader to Appendix F.2. It is crucial to recognize 607  
608 that these perturbations are superficial: the altered 608  
609 task descriptions remain syntactically and semanti- 609  
610 cally equivalent to their originals, although linguis- 610  
611 tic expressions or narratives may vary substantially. 611  
612 Methods that go beyond surface-level perturbations 612  
613 are needed to evaluate the robustness and invari- 613  
614 ance of LLMs across diverse problem framings and 614  
615 modalities effectively.

**Testing Independence from Irrelevant Context** 616  
617 Studies such as Shi et al. (2023), Wu et al. (2024), 617  
618 Liu et al. (2024d), and Yoran et al. (2023) have 618  
619 explored the phenomenon of “lost-in-context” by 619  
620 introducing random or misleading sentences into 620  
621 original problem statements. While earlier bench- 621  
622 marks like those from Weston et al. (2015), Sinha

et al. (2019), Clark et al. (2020), and Webson and Pavlick (2021) have included irrelevant content, they have been predominantly limited to language modalities and single-agent systems. Recent benchmarks such as MileBench (Song et al., 2024), Mementos (Wang et al., 2024c), Seed-bench-2 (Li et al., 2023b), and DEMON (Li et al., 2023c) begin to evaluate multi-modal agents in long context or image sequences, where accurately responding to a specific question requires isolating only the relevant information from the long context window.

## 6 Open Problems and Future Directions

**Inherent Rationality** It is important to understand that the notion of multi-modal or multi-agent systems does not *inherently* imply rationality. **Current methods are neither sufficient nor necessary, but they serve as instrumental tools that bridge the gap between an LLM’s response and rationality.** These approaches enable multi-agent systems, which are black boxes from the user’s perspective, to more closely mimic rational thinking in their output responses. However, despite these more rational responses elicited from multi-modal and multi-agent systems, the challenge of how to effectively close the loop and bake these enhanced outputs back into foundation models themselves (Zhao et al., 2024) beyond mere fine-tuning remains an open question. In other words, can we leverage these more rational outputs to inherently enhance a single foundation model’s rationality in its initial responses in future applications?

**More Comprehensive Evaluation on Rationality** Section 4 thoroughly compares multi-modal and multi-agent systems over their LLM-based single-agent baselines. However, the choices of evaluation metrics are important (Schaeffer et al., 2024); these examples predominantly focus on the accuracy of their final performance, ignoring the most interesting intermediate reasoning steps and the concept of rationality. Section 5 furthermore acknowledges that while there have been some efforts to assess the rationality of agent systems, the field still lacks comprehensive and rigorous evaluation metrics. Moreover, **most existing benchmarks on rationality provide limited comparisons between multi-agent frameworks and single-agent baselines**, thus failing to fully elucidate the advantages multi-agent frameworks can offer.

Future research should prioritize the development of more robust and scalable methods for eval-

uating rationality, taking into account unique challenges and biases posed by agents. **A promising direction is to create methods specifically tailored to assess rationality, going beyond existing ones on accuracy.** These new methods should avoid data contamination and emphasize tasks that demand consistent reasoning across diverse representations and domains. There is a need for more rigorous and large-scale studies on the principles of invariance and orderability of preference, together with their applications to testing rationality in agent systems. This would involve developing more sophisticated perturbation methods that probe the consistency of reasoning at a deeper level, as well as designing experiments that yield statistically significant results.

**Encouraging More Multi-Modal Agents in Multi-Agent Systems** Research into the integration of multi-modality within multi-agent systems would be promising. Fields such as multi-agent debate, collaboration, and neuro-symbolic reasoning, as shown in Figure 1, currently under-utilize the potential of multi-modal sensory inputs. We believe that expanding the role of multi-modalities, including but not limited to vision, sounds, and structured data could significantly enhance the capabilities and rationality of multi-agent systems.

## 7 Conclusions

This survey builds connections between multi-modal and multi-agent systems with rationality, guided by dual-process theories and the four axioms we expect a rational agent or agent systems should satisfy: *grounding, orderability of preference, independence from irrelevant context, and invariance*. Our findings suggest that the grounding can usually be enhanced by multi-modalities, knowledge retrieval, and tool utilization. The remaining three axioms are typically intertwined, and often simultaneously improved via deliberation (slow, iterative thinking process) and abstraction (distilling the logical essence).

Collaboration between the AI research community and cognitive psychologists could be particularly fruitful. We need better evaluation benchmarks on the rationality of agents, more exploration to mitigate cognitive biases in multi-modal and multi-agent systems, and deeper understanding of how these biases arise and how they can be mitigated, ultimately enhancing rationality in decision-making processes.



## 8 Limitations

The fields of multi-modal and multi-agent systems are rapidly evolving. Despite our best efforts, it is inherently impossible to encompass all related works within the scope of this survey. Our discussion also possesses limited mention of the reasoning capabilities, theory of mind in machine psychology, and cognitive architectures, all of which lie beyond the scope of this survey but are crucial for a deeper understanding of LLMs and agent systems. Furthermore, the concept of rationality in human cognitive science may encompass more principles and axioms than those defined in our survey.

## References

Armen Aghajanyan, Bernie Huang, Candace Ross, Vladimir Karpukhin, Hu Xu, Naman Goyal, Dmytro Okhonko, Mandar Joshi, Gargi Ghosh, Mike Lewis, et al. 2022. Cm3: A causal masked multi-modal model of the internet. *arXiv preprint arXiv:2201.07520*.

Jean-Baptiste Alayrac, Jeff Donahue, Pauline Luc, Antoine Miech, Iain Barr, Yana Hasson, Karel Lenc, Arthur Mensch, Katherine Millican, Malcolm Reynolds, et al. 2022. Flamingo: a visual language model for few-shot learning. *Advances in neural information processing systems*, 35:23716–23736.

Anas Awadalla, Irena Gao, Josh Gardner, Jack Hessel, Yusuf Hanafy, Wanrong Zhu, Kalyani Marathe, Yonatan Bitton, Samir Gadre, Shiori Sagawa, et al. 2023. Openflamingo: An open-source framework for training large autoregressive vision-language models. *arXiv preprint arXiv:2308.01390*.

Yushi Bai, Jiahao Ying, Yixin Cao, Xin Lv, Yuze He, Xiaozhi Wang, Jifan Yu, Kaisheng Zeng, Yijia Xiao, Haozhe Lyu, et al. 2024. Benchmarking foundation models with language-model-as-an-examiner. *Advances in Neural Information Processing Systems*, 36.

Yejin Bang, Samuel Cahyawijaya, Nayeon Lee, Wenliang Dai, Dan Su, Bryan Wilie, Holy Lovenia, Ziwei Ji, Tiezheng Yu, Willy Chung, et al. 2023. A multi-task, multilingual, multimodal evaluation of chatgpt on reasoning, hallucination, and interactivity. *arXiv preprint arXiv:2302.04023*.

Hangbo Bao, Wenhui Wang, Li Dong, Qiang Liu, Owais Khan Mohammed, Kriti Aggarwal, Subhojit Som, Songhao Piao, and Furu Wei. 2022. Vlmo: Unified vision-language pre-training with mixture-of-modality-experts. *Advances in Neural Information Processing Systems*, 35:32897–32912.

Emily M Bender, Timnit Gebru, Angelina McMillan-Major, and Shmargaret Shmitchell. 2021. On the dangers of stochastic parrots: Can language models

be too big? In *Proceedings of the 2021 ACM conference on fairness, accountability, and transparency*, pages 610–623.

Maciej Besta, Nils Blach, Ales Kubicek, Robert Gerstenberger, Michal Podstawski, Lukas Gianinazzi, Joanna Gajda, Tomasz Lehmann, Hubert Niewiadomski, Piotr Nyczyk, et al. 2024. Graph of thoughts: Solving elaborate problems with large language models. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 38, pages 17682–17690.

Marcel Binz and Eric Schulz. 2023. Using cognitive psychology to understand gpt-3. *Proceedings of the National Academy of Sciences*, 120(6):e2218523120.

Ali Furkan Biten, Lluís Gómez, and Dimosthenis Karatzas. 2022. Let there be a clock on the beach: Reducing object hallucination in image captioning. In *Proceedings of the IEEE/CVF Winter Conference on Applications of Computer Vision*, pages 1381–1390.

Tim Brooks, Bill Peebles, Connor Holmes, Will DePue, Yufei Guo, Li Jing, David Schnurr, Joe Taylor, Troy Luhman, Eric Luhman, Clarence Ng, Ricky Wang, and Aditya Ramesh. 2024. [Video generation models as world simulators](#).

Sébastien Bubeck, Varun Chandrasekaran, Ronen Eldan, Johannes Gehrke, Eric Horvitz, Ece Kamar, Peter Lee, Yin Tat Lee, Yuanzhi Li, Scott Lundberg, et al. 2023. Sparks of artificial general intelligence: Early experiments with gpt-4. *arXiv preprint arXiv:2303.12712*.

Meng Cao, Yue Dong, and Jackie Chi Kit Cheung. 2021. Hallucinated but factual! inspecting the factuality of hallucinations in abstractive summarization. *arXiv preprint arXiv:2109.09784*.

Chi-Min Chan, Weize Chen, Yusheng Su, Jianxuan Yu, Wei Xue, Shanghang Zhang, Jie Fu, and Zhiyuan Liu. 2023. Chateval: Towards better llm-based evaluators through multi-agent debate. *arXiv preprint arXiv:2308.07201*.

Xi Chen, Xiao Wang, Soravit Changpinyo, AJ Piergiovanni, Piotr Padlewski, Daniel Salz, Sebastian Goodman, Adam Grycner, Basil Mustafa, Lucas Beyer, et al. 2022. Pali: A jointly-scaled multilingual language-image model. *arXiv preprint arXiv:2209.06794*.

Yi Chen, Rui Wang, Haiyun Jiang, Shuming Shi, and Ruifeng Xu. 2023. Exploring the use of large language models for reference-free text quality evaluation: A preliminary empirical study. *arXiv preprint arXiv:2304.00723*.

Yuheng Cheng, Ceyao Zhang, Zhengwen Zhang, Xianguan Meng, Sirui Hong, Wenhao Li, Zihao Wang, Zekai Wang, Feng Yin, Junhua Zhao, et al. 2024. Exploring large language model based intelligent agents: Definitions, methods, and prospects. *arXiv preprint arXiv:2401.03428*.

831	Zhoujun Cheng, Tianbao Xie, Peng Shi, Chengzu Li, Rahul Nadkarni, Yushi Hu, Caiming Xiong, Dragomir Radev, Mari Ostendorf, Luke Zettlemoyer, et al. 2022. Binding language models in symbolic languages. <i>arXiv preprint arXiv:2210.02875</i> .	884
832		885
833		886
834		887
835		888
836	Inyoung Cheong, King Xia, KJ Feng, Quan Ze Chen, and Amy X Zhang. 2024. (a) i am not a lawyer, but...: Engaging legal experts towards responsible llm policies for legal advice. <i>arXiv preprint arXiv:2402.01864</i> .	889
837		890
838		891
839		892
840		893
841	Steffi Chern, Zhen Fan, and Andy Liu. 2024. Combating adversarial attacks with multi-agent debate. <i>arXiv preprint arXiv:2401.05998</i> .	894
842		895
843		896
844	Cheng-Han Chiang and Hung-yi Lee. 2023. A closer look into automatic evaluation using large language models. <i>arXiv preprint arXiv:2310.05657</i> .	897
845		898
846		899
847	Peter Clark, Oyvind Tafjord, and Kyle Richardson. 2020. Transformers as soft reasoners over language. <i>arXiv preprint arXiv:2002.05867</i> .	900
848		901
849		902
850	Roi Cohen, May Hamri, Mor Geva, and Amir Globerson. 2023. Lm vs lm: Detecting factual errors via cross examination. <i>arXiv preprint arXiv:2305.13281</i> .	903
851		904
852		905
853	Fergus IM Craik and Robert S Lockhart. 1972. Levels of processing: A framework for memory research. <i>Journal of verbal learning and verbal behavior</i> , 11(6):671–684.	906
854		907
855		908
856		909
857	Can Cui, Yunsheng Ma, Xu Cao, Wenqian Ye, Yang Zhou, Kaizhao Liang, Jintai Chen, Juanwu Lu, Zichong Yang, Kuei-Da Liao, et al. 2024. A survey on multimodal large language models for autonomous driving. In <i>Proceedings of the IEEE/CVF Winter Conference on Applications of Computer Vision</i> , pages 958–979.	910
858		911
859		912
860		913
861		914
862		915
863		916
864	Souvik Das, Sougata Saha, and Rohini K Srihari. 2023. Diving deep into modes of fact hallucinations in dialogue systems. <i>arXiv preprint arXiv:2301.04449</i> .	917
865		918
866		919
867	Xiang Deng, Yu Gu, Boyuan Zheng, Shijie Chen, Sam Stevens, Boshi Wang, Huan Sun, and Yu Su. 2024. Mind2web: Towards a generalist agent for the web. <i>Advances in Neural Information Processing Systems</i> , 36.	920
868		921
869		922
870		923
871		924
872	Yihong Dong, Xue Jiang, Huanyu Liu, Zhi Jin, and Ge Li. 2024. Generalization or memorization: Data contamination and trustworthy evaluation for large language models. <i>arXiv preprint arXiv:2402.15938</i> .	925
873		926
874		927
875		928
876	Wei Du, Yichun Zhao, Boqun Li, Gongshen Liu, and Shilin Wang. 2022. Ppt: Backdoor attacks on pre-trained models via poisoned prompt tuning. In <i>IJCAI</i> , pages 680–686.	929
877		930
878		931
879		932
880	Yilun Du, Shuang Li, Antonio Torralba, Joshua B Tenenbaum, and Igor Mordatch. 2023. Improving factuality and reasoning in language models through multi-agent debate. <i>arXiv preprint arXiv:2305.14325</i> .	933
881		934
882		935
883		936
	Zane Durante, Qiuyuan Huang, Naoki Wake, Ran Gong, Jae Sung Park, Bidipta Sarkar, Rohan Taori, Yusuke Noda, Demetri Terzopoulos, Yejin Choi, et al. 2024. Agent ai: Surveying the horizons of multimodal interaction. <i>arXiv preprint arXiv:2401.03568</i> .	937
		938
	Nouha Dziri, Ehsan Kamaloo, Sivan Milton, Omar Zaiane, Mo Yu, Edoardo M Ponti, and Siva Reddy. 2022a. Faithdial: A faithful benchmark for information-seeking dialogue. <i>Transactions of the Association for Computational Linguistics</i> , 10:1473–1490.	939
		940
	Nouha Dziri, Hannah Rashkin, Tal Linzen, and David Reitter. 2022b. Evaluating attribution in dialogue systems: The begin benchmark. <i>Transactions of the Association for Computational Linguistics</i> , 10:1066–1083.	941
		942
	Jessica Echterhoff, Yao Liu, Abeer Alessa, Julian McAuley, and Zexue He. 2024. Cognitive bias in high-stakes decision-making with llms. <i>arXiv preprint arXiv:2403.00811</i> .	943
		944
	Franz Eisenführ, Martin Weber, and Thomas Langer. 2010. <i>Rational decision making</i> . Springer.	945
		946
	Linxi Fan, Guanzhi Wang, Yunfan Jiang, Ajay Mandlekar, Yuncong Yang, Haoyi Zhu, Andrew Tang, De-An Huang, Yuke Zhu, and Anima Anandkumar. 2022. Minedojo: Building open-ended embodied agents with internet-scale knowledge. <i>Advances in Neural Information Processing Systems</i> , 35:18343–18362.	947
		948
	Meng Fang, Shilong Deng, Yudi Zhang, Zijing Shi, Ling Chen, Mykola Pechenizkiy, and Jun Wang. 2024. Large language models are neurosymbolic reasoners. <i>arXiv preprint arXiv:2401.09334</i> .	949
		950
	Jinlan Fu, See-Kiong Ng, Zhengbao Jiang, and Pengfei Liu. 2023. Gptscore: Evaluate as you desire. <i>arXiv preprint arXiv:2302.04166</i> .	951
		952
	Hiroki Furuta, Ofir Nachum, Kuang-Huei Lee, Yutaka Matsuo, Shixiang Shane Gu, and Izzeddin Gur. 2023. Multimodal web navigation with instruction-finetuned foundation models. <i>arXiv preprint arXiv:2305.11854</i> .	953
		954
	Deep Ganguli, Liane Lovitt, Jackson Kernion, Amanda Askell, Yuntao Bai, Saurav Kadavath, Ben Mann, Ethan Perez, Nicholas Schiefer, Kamal Ndousse, et al. 2022. Red teaming language models to reduce harms: Methods, scaling behaviors, and lessons learned. <i>arXiv preprint arXiv:2209.07858</i> .	955
		956
	Chen Gao, Xiaochong Lan, Zhihong Lu, Jinzhu Mao, Jinghua Piao, Huandong Wang, Depeng Jin, and Yong Li. 2023a. S3: Social-network simulation system with large language model-empowered agents. <i>arXiv preprint arXiv:2307.14984</i> .	957
		958
	Difei Gao, Lei Ji, Luowei Zhou, Kevin Qinghong Lin, Joya Chen, Zihan Fan, and Mike Zheng Shou. 2023b. Assistgpt: A general multi-modal assistant that can	959

939	plan, execute, inspect, and learn. <i>arXiv preprint arXiv:2306.08640</i> .	Aidan Hogan, Eva Blomqvist, Michael Cochez, Claudia d'Amato, Gerard De Melo, Claudio Gutierrez, Sabrina Kirrane, José Emilio Labra Gayo, Roberto Navigli, Sebastian Neumaier, et al. 2021. Knowledge graphs. <i>ACM Computing Surveys (Csur)</i> , 54(4):1–37.	994
940			995
941	Mingqi Gao, Jie Ruan, Renliang Sun, Xunjian Yin, Shiping Yang, and Xiaojun Wan. 2023c. Human-like summarization evaluation with chatgpt. <i>arXiv preprint arXiv:2304.02554</i> .		996
942			997
943			998
944			999
945	Samuel Gehman, Suchin Gururangan, Maarten Sap, Yejin Choi, and Noah A Smith. 2020. Realtocixityprompts: Evaluating neural toxic degeneration in language models. <i>arXiv preprint arXiv:2009.11462</i> .	Sirui Hong, Xiawu Zheng, Jonathan Chen, Yuheng Cheng, Jinlin Wang, Ceyao Zhang, Zili Wang, Steven Ka Shing Yau, Zijuan Lin, Liyang Zhou, et al. 2023a. Metagpt: Meta programming for multi-agent collaborative framework. <i>arXiv preprint arXiv:2308.00352</i> .	1000
946			1001
947			1002
948			1003
949	Significant Gravitass. 2023. Autogpt. <i>Python</i> . <a href="https://github.com/Significant-Gravitas/Auto-GPT">https://github.com/Significant-Gravitas/Auto-GPT</a> .	Wenyi Hong, Weihang Wang, Qingsong Lv, Jiazheng Xu, Wenmeng Yu, Junhui Ji, Yan Wang, Zihan Wang, Yuxiao Dong, Ming Ding, et al. 2023b. Cogagent: A visual language model for gui agents. <i>arXiv preprint arXiv:2312.08914</i> .	1005
950			1006
951	Nuno M Guerreiro, Duarte M Alves, Jonas Waldendorf, Barry Haddow, Alexandra Birch, Pierre Colombo, and André FT Martins. 2023. Hallucinations in large multilingual translation models. <i>Transactions of the Association for Computational Linguistics</i> , 11:1500–1517.	Joy Hsu, Jiayuan Mao, Josh Tenenbaum, and Jiajun Wu. 2024. What's left? concept grounding with logic-enhanced foundation models. <i>Advances in Neural Information Processing Systems</i> , 36.	1007
952			1008
953			1009
954			1010
955			1011
956			1012
957	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. <i>arXiv preprint arXiv:2402.01680</i> .	Ziniu Hu, Ahmet Iscen, Chen Sun, Kai-Wei Chang, Yizhou Sun, David Ross, Cordelia Schmid, and Alireza Fathi. 2024. Avis: Autonomous visual information seeking with large language model agent. <i>Advances in Neural Information Processing Systems</i> , 36.	1014
958			1015
959			1016
960			1017
961			1018
962	Prakhar Gupta, Chien-Sheng Wu, Wenhao Liu, and Caiming Xiong. 2021. Dialfact: A benchmark for fact-checking in dialogue. <i>arXiv preprint arXiv:2110.08222</i> .	Jie Huang and Kevin Chen-Chuan Chang. 2022. Towards reasoning in large language models: A survey. <i>arXiv preprint arXiv:2212.10403</i> .	1020
963			1021
964			1022
965			1023
966	Tanmay Gupta and Aniruddha Kembhavi. 2023. Visual programming: Compositional visual reasoning without training. In <i>Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition</i> , pages 14953–14962.	Lei Huang, Weijiang Yu, Weitao Ma, Weihong Zhong, Zhangyin Feng, Haotian Wang, Qianglong Chen, Weihua Peng, Xiaocheng Feng, Bing Qin, et al. 2023. A survey on hallucination in large language models: Principles, taxonomy, challenges, and open questions. <i>arXiv preprint arXiv:2311.05232</i> .	1024
967			1025
968			1026
969			1027
970			1028
971	Izzeddin Gur, Hiroki Furuta, Austin Huang, Mustafa Safdari, Yutaka Matsuo, Douglas Eck, and Aleksandra Faust. 2023. A real-world webagent with planning, long context understanding, and program synthesis. <i>arXiv preprint arXiv:2307.12856</i> .	Alon Jacovi, Avi Caciularu, Omer Goldman, and Yoav Goldberg. 2023. Stop uploading test data in plain text: Practical strategies for mitigating data contamination by evaluation benchmarks. <i>arXiv preprint arXiv:2305.10160</i> .	1029
972			1030
973			1031
974			1032
975			1033
976	Shanshan Han, Qifan Zhang, Yuhang Yao, Weizhao Jin, Zhaozhuo Xu, and Chaoyang He. 2024. Llm multi-agent systems: Challenges and open problems. <i>arXiv preprint arXiv:2402.03578</i> .	Bowen Jiang, Zhijun Zhuang, Shreyas S Shivakumar, Dan Roth, and Camillo J Taylor. 2024. Multi-agent vqa: Exploring multi-agent foundation models in zero-shot visual question answering. <i>arXiv preprint arXiv:2403.14783</i> .	1034
977			1035
978			1036
979			1037
980	Reid Hastie and Robyn M Dawes. 2009. <i>Rational choice in an uncertain world: The psychology of judgment and decision making</i> . Sage Publications.	Philip Nicholas Johnson-Laird. 1983. <i>Mental models: Towards a cognitive science of language, inference, and consciousness</i> . 6. Harvard University Press.	1038
981			1039
982			1040
983	Kai He, Rui Mao, Qika Lin, Yucheng Ruan, Xiang Lan, Mengling Feng, and Erik Cambria. 2023. A survey of large language models for healthcare: from data, technology, and applications to accountability and ethics. <i>arXiv preprint arXiv:2310.05694</i> .	Haoqiang Kang and Xiao-Yang Liu. 2023. Deficiency of large language models in finance: An empirical examination of hallucination. <i>arXiv preprint arXiv:2311.15548</i> .	1041
984			1042
985			1043
986			1044
987			1045
988	Dan Hendrycks, Collin Burns, Saurav Kadavath, Akul Arora, Steven Basart, Eric Tang, Dawn Song, and Jacob Steinhardt. 2021. Measuring mathematical problem solving with the math dataset. In <i>Thirty-fifth Conference on Neural Information Processing Systems Datasets and Benchmarks Track (Round 2)</i> .	Yu He Ke, Rui Yang, Sui An Lie, Taylor Xin Yi Lim, Hairil Rizal Abdullah, Daniel Shu Wei Ting, and Nan Liu. 2024. Enhancing diagnostic accuracy through	1046
989			1047
990			1048
991			
992			
993			

1049	multi-agent conversations: Using large language models to mitigate cognitive bias. <i>arXiv preprint arXiv:2401.14589</i> .	Juncheng Li, Kaihang Pan, Zhiqi Ge, Minghe Gao, Wei Ji, Wenqiao Zhang, Tat-Seng Chua, Siliang Tang, Hanwang Zhang, and Yueting Zhuang. 2023c. Fine-tuning multimodal llms to follow zero-shot demonstrative instructions. In <i>The Twelfth International Conference on Learning Representations</i> .	1104
1050			1105
1051			1106
1052	Akbir Khan, John Hughes, Dan Valentine, Laura Ruis, Kshitij Sachan, Ansh Radhakrishnan, Edward Grefenstette, Samuel R Bowman, Tim Rocktäschel, and Ethan Perez. 2024. Debating with more persuasive llms leads to more truthful answers. <i>arXiv preprint arXiv:2402.06782</i> .		1107
1053			1108
1054			1109
1055		Junkai Li, Siyu Wang, Meng Zhang, Weitao Li, Yunghwei Lai, Xinhui Kang, Weizhi Ma, and Yang Liu. 2024b. Agent hospital: A simulacrum of hospital with evolvable medical agents. <i>arXiv preprint arXiv:2405.02957</i> .	1110
1056			1111
1057			1112
1058	Roni Khardon and Dan Roth. 1997. Learning to reason. <i>Journal of the ACM (JACM)</i> , 44(5):697–725.		1113
1059			1114
1060	Alexander Kirillov, Eric Mintun, Nikhila Ravi, Hanzi Mao, Chloe Rolland, Laura Gustafson, Tete Xiao, Spencer Whitehead, Alexander C. Berg, Wan-Yen Lo, Piotr Dollár, and Ross Girshick. 2023. Segment anything. <i>arXiv:2304.02643</i> .	Junnan Li, Dongxu Li, Silvio Savarese, and Steven Hoi. 2023d. Blip-2: Bootstrapping language-image pre-training with frozen image encoders and large language models. In <i>International conference on machine learning</i> , pages 19730–19742. PMLR.	1115
1061			1116
1062			1117
1063			1118
1064			1119
1065	Ryan Koo, Minhwa Lee, Vipul Raheja, Jong Inn Park, Zae Myung Kim, and Dongyeop Kang. 2023. Benchmarking cognitive biases in large language models as evaluators. <i>arXiv preprint arXiv:2309.17012</i> .	Junyi Li, Xiaoxue Cheng, Wayne Xin Zhao, Jian-Yun Nie, and Ji-Rong Wen. 2023e. Halueval: A large-scale hallucination evaluation benchmark for large language models. In <i>Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing</i> , pages 6449–6464.	1120
1066			1121
1067			1122
1068			1123
1069	John E Laird. 2019. <i>The Soar cognitive architecture</i> . MIT press.		1124
1070			1125
1071	Yann LeCun. 2022. A path towards autonomous machine intelligence version 0.9. 2, 2022-06-27. <i>Open Review</i> , 62(1).	Xingxuan Li, Ruochen Zhao, Yew Ken Chia, Bosheng Ding, Lidong Bing, Shafiq Joty, and Soujanya Poria. 2023f. Chain of knowledge: A framework for grounding large language models with structured knowledge bases. <i>arXiv preprint arXiv:2305.13269</i> .	1126
1072			1127
1073			1128
1074	Yann LeCun. 2024. <b>Objective-driven ai: Towards ai systems that can learn, remember, reason, plan, have common sense, yet are steerable and safe</b> . University of Washington, Department of Electrical & Computer Engineering. Slide presentation retrieved from University of Washington.		1129
1075			1130
1076		Yifan Li, Yifan Du, Kun Zhou, Jinpeng Wang, Wayne Xin Zhao, and Ji-Rong Wen. 2023g. Evaluating object hallucination in large vision-language models. <i>arXiv preprint arXiv:2305.10355</i> .	1131
1077			1132
1078			1133
1079			1134
1080	Patrick Lewis, Ethan Perez, Aleksandra Piktus, Fabio Petroni, Vladimir Karpukhin, Naman Goyal, Heinrich Küttler, Mike Lewis, Wen-tau Yih, Tim Rocktäschel, et al. 2020. Retrieval-augmented generation for knowledge-intensive nlp tasks. <i>Advances in Neural Information Processing Systems</i> , 33:9459–9474.	Yinheng Li, Shaofei Wang, Han Ding, and Hang Chen. 2023h. Large language models in finance: A survey. In <i>Proceedings of the Fourth ACM International Conference on AI in Finance</i> , pages 374–382.	1135
1081			1136
1082			1137
1083			1138
1084			
1085			
1086	Selma Leydesdorff and Katharine Hodgkin. 2017. <i>Memory cultures: Memory, subjectivity and recognition</i> . Routledge.	Yuan Li, Yixuan Zhang, and Lichao Sun. 2023i. Metaagents: Simulating interactions of human behaviors for llm-based task-oriented coordination via collaborative generative agents. <i>arXiv preprint arXiv:2310.06500</i> .	1139
1087			1140
1088			1141
1089	Bangzheng Li, Ben Zhou, Fei Wang, Xingyu Fu, Dan Roth, and Muhao Chen. 2023a. Deceiving semantic shortcuts on reasoning chains: How far can models go without hallucination? <i>arXiv preprint arXiv:2311.09702</i> .		1142
1090			1143
1091		Tian Liang, Zhiwei He, Wenxiang Jiao, Xing Wang, Yan Wang, Rui Wang, Yujiu Yang, Zhaopeng Tu, and Shuming Shi. 2023. Encouraging divergent thinking in large language models through multi-agent debate. <i>arXiv preprint arXiv:2305.19118</i> .	1144
1092			1145
1093			1146
1094	Bohao Li, Yuying Ge, Yixiao Ge, Guangzhi Wang, Rui Wang, Ruimao Zhang, and Ying Shan. 2023b. Seed-bench-2: Benchmarking multimodal large language models. <i>arXiv preprint arXiv:2311.17092</i> .		1147
1095			1148
1096		Evan Zheran Liu, Kelvin Guu, Panupong Pasupat, Tianlin Shi, and Percy Liang. 2018. Reinforcement learning on web interfaces using workflow-guided exploration. <i>arXiv preprint arXiv:1802.08802</i> .	1149
1097			1150
1098	Chunyu Li, Zhe Gan, Zhengyuan Yang, Jianwei Yang, Linjie Li, Lijuan Wang, Jianfeng Gao, et al. 2024a. Multimodal foundation models: From specialists to general-purpose assistants. <i>Foundations and Trends® in Computer Graphics and Vision</i> , 16(1-2):1–214.		1151
1099			1152
1100		Fang Liu, Yang Liu, Lin Shi, Houkun Huang, Ruifeng Wang, Zhen Yang, and Li Zhang. 2024a. Exploring and evaluating hallucinations in llm-powered code generation. <i>arXiv preprint arXiv:2404.00971</i> .	1153
1101			1154
1102			1155
1103			1156

1157	Hao Liu, Wilson Yan, Matei Zaharia, and Pieter Abbeel. 2024b. World model on million-length video and language with ringattention. <i>arXiv preprint arXiv:2402.08268</i> .	Zheheng Luo, Qianqian Xie, and Sophia Ananiadou. 2023. Chatgpt as a factual inconsistency evaluator for abstractive text summarization. <i>arXiv preprint arXiv:2303.15621</i> .	1211 1212 1213 1214
1161	Haotian Liu, Chunyuan Li, Yuheng Li, and Yong Jae Lee. 2023a. Improved baselines with visual instruction tuning. <i>arXiv preprint arXiv:2310.03744</i> .	Olivia Macmillan-Scott and Mirco Musolesi. 2024. (ir) rationality and cognitive biases in large language models. <i>arXiv preprint arXiv:2402.09193</i> .	1215 1216 1217
1164	Haotian Liu, Chunyuan Li, Qingyang Wu, and Yong Jae Lee. 2024c. Visual instruction tuning. <i>Advances in neural information processing systems</i> , 36.	Aman Madaan, Niket Tandon, Prakhar Gupta, Skyler Hallinan, Luyu Gao, Sarah Wiegrefe, Uri Alon, Nouha Dziri, Shrimai Prabhumoye, Yiming Yang, et al. 2024. Self-refine: Iterative refinement with self-feedback. <i>Advances in Neural Information Processing Systems</i> , 36.	1218 1219 1220 1221 1222 1223
1167	Jian Liu, Leyang Cui, Hanmeng Liu, Dandan Huang, Yile Wang, and Yue Zhang. 2021. Logiqa: a challenge dataset for machine reading comprehension with logical reasoning. In <i>Proceedings of the Twenty-Ninth International Conference on International Joint Conferences on Artificial Intelligence</i> , pages 3622–3628.	Inbal Magar and Roy Schwartz. 2022. Data contamination: From memorization to exploitation. <i>arXiv preprint arXiv:2203.08242</i> .	1224 1225 1226
1174	Nelson F Liu, Kevin Lin, John Hewitt, Ashwin Paranjape, Michele Bevilacqua, Fabio Petroni, and Percy Liang. 2024d. Lost in the middle: How language models use long contexts. <i>Transactions of the Association for Computational Linguistics</i> , 12:157–173.	James G March and Herbert A Simon. 1958. Organizations. <i>University of Illinois at Urbana-Champaign’s Academy for Entrepreneurial Leadership Historical Research Reference in Entrepreneurship</i> .	1227 1228 1229 1230
1179	Ollie Liu, Deqing Fu, Dani Yogatama, and Willie Neiswanger. 2024e. Dellma: A framework for decision making under uncertainty with large language models. <i>arXiv preprint arXiv:2402.02392</i> .	Kenneth Marino, Xinlei Chen, Devi Parikh, Abhinav Gupta, and Marcus Rohrbach. 2021. Krisp: Integrating implicit and symbolic knowledge for open-domain knowledge-based vqa. In <i>Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition</i> , pages 14111–14121.	1231 1232 1233 1234 1235 1236
1183	Yang Liu, Dan Iter, Yichong Xu, Shuhang Wang, Ruochen Xu, and Chenguang Zhu. 2023b. Gpteval: Nlg evaluation using gpt-4 with better human alignment. <i>arXiv preprint arXiv:2303.16634</i> .	Amirkeivan Mohtashami, Florian Hartmann, Sian Gooding, Lukas Zilka, Matt Sharifi, et al. 2023. Social learning: Towards collaborative learning with large language models. <i>arXiv preprint arXiv:2312.11441</i> .	1237 1238 1239 1240
1187	Zijun Liu, Yanzhe Zhang, Peng Li, Yang Liu, and Diyi Yang. 2023c. Dynamic llm-agent network: An llm-agent collaboration framework with agent team optimization. <i>arXiv preprint arXiv:2310.02170</i> .	Anirban Mukherjee and Hannah Hanwen Chang. 2024. Heuristic reasoning in ai: Instrumental use and mimetic absorption. <i>arXiv preprint arXiv:2403.09404</i> .	1241 1242 1243 1244
1191	Jiasen Lu, Dhruv Batra, Devi Parikh, and Stefan Lee. 2019. Vilbert: Pretraining task-agnostic visiolinguistic representations for vision-and-language tasks. <i>Advances in neural information processing systems</i> , 32.	Yohei Nakajima. 2023. Babyagi. <i>Python</i> . <a href="https://github.com/yoheinakajima/babyagi">https://github.com/yoheinakajima/babyagi</a> .	1245 1246
1195	Junru Lu, Siyu An, Mingbao Lin, Gabriele Pergola, Yulan He, Di Yin, Xing Sun, and Yunsheng Wu. 2023. Memochat: Tuning llms to use memos for consistent long-range open-domain conversation. <i>arXiv preprint arXiv:2308.08239</i> .	Xenia Ohmer, Elia Bruni, and Dieuwke Hupkes. 2023. Separating form and meaning: Using self-consistency to quantify task understanding across multiple senses. <i>CoRR</i> , abs/2305.11662.	1247 1248 1249 1250
1200	Pan Lu, Baolin Peng, Hao Cheng, Michel Galley, Kai-Wei Chang, Ying Nian Wu, Song-Chun Zhu, and Jianfeng Gao. 2024. Chameleon: Plug-and-play compositional reasoning with large language models. <i>Advances in Neural Information Processing Systems</i> , 36.	Xenia Ohmer, Elia Bruni, and Dieuwke Hupkes. 2024. From form (s) to meaning: Probing the semantic depths of language models using multisense consistency. <i>arXiv preprint arXiv:2404.12145</i> .	1251 1252 1253 1254
1206	Yao Lu, Max Bartolo, Alastair Moore, Sebastian Riedel, and Pontus Stenetorp. 2021. Fantastically ordered prompts and where to find them: Overcoming few-shot prompt order sensitivity. <i>arXiv preprint arXiv:2104.08786</i> .	OpenAI. 2023. <b>Gpt-4v(ision) system card</b> .	1255
1207		OpenAI. 2024. <b>Gpt-4o</b> . Software available from OpenAI. Accessed: 2024-05-20.	1256 1257
1208		Liangming Pan, Alon Albalak, Xinyi Wang, and William Yang Wang. 2023. Logic-lm: Empowering large language models with symbolic solvers for faithful logical reasoning. <i>arXiv preprint arXiv:2305.12295</i> .	1258 1259 1260 1261 1262

1263	Zhenyu Pan, Haozheng Luo, Manling Li, and Han Liu.	Tianhe Ren, Shilong Liu, Ailing Zeng, Jing Lin, Kun-	1318
1264	2024. Chain-of-action: Faithful and multimodal	chang Li, He Cao, Jiayu Chen, Xinyu Huang, Yukang	1319
1265	question answering through large language models.	Chen, Feng Yan, et al. 2024. Grounded sam: As-	1320
1266	<i>arXiv preprint arXiv:2403.17359</i> .	sembling open-world models for diverse visual tasks.	1321
1267	PaperswithcodeMCQA. Multiple choice	<i>arXiv preprint arXiv:2401.14159</i> .	1322
1268	qa. <a href="https://paperswithcode.com/task/multiple-choice-qa/latest">https://paperswithcode.com/task/</a>	Raquel B Robinson, Karin Johansson, James Collin	1323
1269	<a href="https://paperswithcode.com/task/multiple-choice-qa/latest">multiple-choice-qa/latest</a> . Accessed: 2024-	Fey, Elena Márquez Segura, Jon Back, Annika	1324
1270	05-28.	Waern, Sarah Lynne Bowman, and Katherine Isbister.	1325
1271	Branislav Pecher, Ivan Srba, and Maria Bielikova.	2023. Leveraging large language models for multiple	1326
1272	2024. On sensitivity of learning with limited la-	choice question answering. In <i>Extended Abstracts</i>	1327
1273	belled data to the effects of randomness: Impact of	<i>of the 2023 CHI Conference on Human Factors in</i>	1328
1274	interactions and systematic choices. <i>arXiv preprint</i>	<i>Computing Systems</i> , pages 1–5.	1329
1275	<i>arXiv:2402.12817</i> .	Anna Rohrbach, Lisa Anne Hendricks, Kaylee Burns,	1330
1276	Ethan Perez, Saffron Huang, Francis Song, Trevor Cai,	Trevor Darrell, and Kate Saenko. 2018. Object	1331
1277	Roman Ring, John Aslanides, Amelia Glaese, Nat	hallucination in image captioning. <i>arXiv preprint</i>	1332
1278	McAleese, and Geoffrey Irving. 2022. Red team-	<i>arXiv:1809.02156</i> .	1333
1279	ing language models with language models. <i>arXiv</i>	Robin Rombach, Andreas Blattmann, Dominik Lorenz,	1334
1280	<i>preprint arXiv:2202.03286</i> .	Patrick Esser, and Björn Ommer. 2022. High-	1335
1281	Archiki Prasad, Alexander Koller, Mareike Hartmann,	resolution image synthesis with latent diffusion mod-	1336
1282	Peter Clark, Ashish Sabharwal, Mohit Bansal, and	els. In <i>Proceedings of the IEEE/CVF conference</i>	1337
1283	Tushar Khot. 2023. Adapt: As-needed decompo-	<i>on computer vision and pattern recognition</i> , pages	1338
1284	sition and planning with language models. <i>arXiv</i>	10684–10695.	1339
1285	<i>preprint arXiv:2311.05772</i> .	Oscar Sainz, Jon Ander Campos, Iker García-Ferrero,	1340
1286	Chen Qian, Xin Cong, Cheng Yang, Weize Chen,	Julen Etxaniz, Oier Lopez de Lacalle, and Eneko	1341
1287	Yusheng Su, Juyuan Xu, Zhiyuan Liu, and Maosong	Agirre. 2023. Nlp evaluation in trouble: On the	1342
1288	Sun. 2023. Communicative agents for software de-	need to measure llm data contamination for each	1343
1289	velopment. <i>arXiv preprint arXiv:2307.07924</i> .	benchmark. <i>arXiv preprint arXiv:2310.18018</i> .	1344
1290	Shuofei Qiao, Yixin Ou, Ningyu Zhang, Xiang Chen,	Rylan Schaeffer, Brando Miranda, and Sanmi Koyejo.	1345
1291	Yunzhi Yao, Shumin Deng, Chuanqi Tan, Fei Huang,	2024. Are emergent abilities of large language mod-	1346
1292	and Huajun Chen. 2022. Reasoning with lan-	els a mirage? <i>Advances in Neural Information Pro-</i>	1347
1293	guage model prompting: A survey. <i>arXiv preprint</i>	<i>cessing Systems</i> , 36.	1348
1294	<i>arXiv:2212.09597</i> .	Timo Schick, Jane Dwivedi-Yu, Roberto Dessì, Roberta	1349
1295	Shuofei Qiao, Ningyu Zhang, Runnan Fang, Yujie Luo,	Raileanu, Maria Lomeli, Eric Hambro, Luke Zettle-	1350
1296	Wangchunshu Zhou, Yuchen Eleanor Jiang, Chengfei	moyer, Nicola Cancedda, and Thomas Scialom. 2024.	1351
1297	Lv, and Huajun Chen. 2024. Autoact: Automatic	Toolformer: Language models can teach themselves	1352
1298	agent learning from scratch via self-planning. <i>arXiv</i>	to use tools. <i>Advances in Neural Information Pro-</i>	1353
1299	<i>preprint arXiv:2401.05268</i> .	<i>cessing Systems</i> , 36.	1354
1300	Alec Radford, Jong Wook Kim, Chris Hallacy, Aditya	Melanie Sclar, Yejin Choi, Yulia Tsvetkov, and Alane	1355
1301	Ramesh, Gabriel Goh, Sandhini Agarwal, Girish Sas-	Suhr. 2023a. Quantifying language models’ sensi-	1356
1302	try, Amanda Askell, Pamela Mishkin, Jack Clark,	tivity to spurious features in prompt design or: How	1357
1303	et al. 2021. Learning transferable visual models from	i learned to start worrying about prompt formatting.	1358
1304	natural language supervision. In <i>International confer-</i>	<i>arXiv preprint arXiv:2310.11324</i> .	1359
1305	<i>ence on machine learning</i> , pages 8748–8763. PMLR.	Melanie Sclar, Sachin Kumar, Peter West, Alane Suhr,	1360
1306	Hannah Rashkin, Vitaly Nikolaev, Matthew Lamm,	Yejin Choi, and Yulia Tsvetkov. 2023b. Minding lan-	1361
1307	Lora Aroyo, Michael Collins, Dipanjan Das, Slav	guage models’ (lack of) theory of mind: A plug-and-	1362
1308	Petrov, Gaurav Singh Tomar, Iulia Turc, and David	play multi-character belief tracker. <i>arXiv preprint</i>	1363
1309	Reitter. 2023. Measuring attribution in natural lan-	<i>arXiv:2306.00924</i> .	1364
1310	guage generation models. <i>Computational Linguistics</i> ,	Reinhard Selten. 1990. Bounded rationality. <i>Jour-</i>	1365
1311	49(4):777–840.	<i>nal of Institutional and Theoretical Economics</i>	1366
1312	Machel Reid, Nikolay Savinov, Denis Teplyashin,	( <i>JITE</i> )/ <i>Zeitschrift für die gesamte Staatswissenschaft</i> ,	1367
1313	Dmitry Lepikhin, Timothy Lillicrap, Jean-baptiste	146(4):649–658.	1368
1314	Alayrac, Radu Soricut, Angeliki Lazaridou, Orhan Fi-	Peter Shaw, Mandar Joshi, James Cohan, Jonathan Be-	1369
1315	rat, Julian Schrittwieser, et al. 2024. Gemini 1.5: Un-	rant, Panupong Pasupat, Hexiang Hu, Urvashi Khan-	1370
1316	locking multimodal understanding across millions of	delwal, Kenton Lee, and Kristina N Toutanova. 2024.	1371
1317	tokens of context. <i>arXiv preprint arXiv:2403.05530</i> .	From pixels to ui actions: Learning to follow in-	1372
		structions via graphical user interfaces. <i>Advances in</i>	1373
		<i>Neural Information Processing Systems</i> , 36.	1374

1375	Chenhui Shen, Liying Cheng, Xuan-Phi Nguyen, Yang You, and Lidong Bing. 2023. Large language models are not yet human-level evaluators for abstractive summarization. In <i>Findings of the Association for Computational Linguistics: EMNLP 2023</i> , pages 4215–4233.	1429
1376		1430
1377		1431
1378		1432
1379		
1380		
1381	Weizhou Shen, Chenliang Li, Hongzhan Chen, Ming Yan, Xiaojun Quan, Hehong Chen, Ji Zhang, and Fei Huang. 2024a. Small llms are weak tool learners: A multi-llm agent. <i>arXiv preprint arXiv:2401.07324</i> .	1433
1382		1434
1383		1435
1384		1436
1385		1437
1386		1438
1387	Yongliang Shen, Kaitao Song, Xu Tan, Dongsheng Li, Weiming Lu, and Yueting Zhuang. 2024b. Hugging-gpt: Solving ai tasks with chatgpt and its friends in hugging face. <i>Advances in Neural Information Processing Systems</i> , 36.	1439
1388		1440
1389		1441
1390		1442
1391		1443
1392		1444
1393		1445
1394		1446
1395		1447
1396		1448
1397		1449
1398		1450
1399		1451
1400		1452
1401		1453
1402		1454
1403		1455
1404		1456
1405		1457
1406		1458
1407		1459
1408		1460
1409		1461
1410		1462
1411		1463
1412		1464
1413		1465
1414		1466
1415		1467
1416		1468
1417		1469
1418		1470
1419		1471
1420		1472
1421		1473
1422		1474
1423		1475
1424		1476
1425		1477
1426		1478
1427		1479
1428		1480
		1481
		1482
		1483
		1484
		1485
		1486
		1487
		1488
		1489
		1490
		1491
		1492
		1493
		1494
		1495
		1496
		1497
		1498
		1499
		1500

1485	Amos Tversky and Daniel Kahneman. 1988. Rational choice and the framing of decisions. <i>Decision making: Descriptive, normative, and prescriptive interactions</i> , pages 167–192.	Xuezhi Wang, Jason Wei, Dale Schuurmans, Quoc Le, Ed Chi, Sharan Narang, Aakanksha Chowdhery, and Denny Zhou. 2022b. Self-consistency improves chain of thought reasoning in language models. <i>arXiv preprint arXiv:2203.11171</i> .	1539
1486			1540
1487			1541
1488			1542
1489	Karthik Valmeekam, Matthew Marquez, Sarath Sreedharan, and Subbarao Kambhampati. 2023. On the planning abilities of large language models—a critical investigation. <i>Advances in Neural Information Processing Systems</i> , 36:75993–76005.	Zeqing Wang, Wentao Wan, Runmeng Chen, Qiqing Lao, Minjie Lang, and Keze Wang. 2023d. Towards top-down reasoning: An explainable multi-agent approach for visual question answering. <i>arXiv preprint arXiv:2311.17331</i> .	1544
1490			1545
1491			1546
1492			1547
1493			1548
1494	Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Łukasz Kaiser, and Illia Polosukhin. 2017. Attention is all you need. <i>Advances in neural information processing systems</i> , 30.	Lucas Weber, Elia Bruni, and Dieuwke Hupkes. 2023. Mind the instructions: a holistic evaluation of consistency and interactions in prompt-based learning. <i>arXiv preprint arXiv:2310.13486</i> .	1549
1495			1550
1496			1551
1497			1552
1498			
1499	John Von Neumann and Oskar Morgenstern. 2007. Theory of games and economic behavior: 60th anniversary commemorative edition. In <i>Theory of games and economic behavior</i> . Princeton university press.	Albert Webson and Ellie Pavlick. 2021. Do prompt-based models really understand the meaning of their prompts? <i>arXiv preprint arXiv:2109.01247</i> .	1553
1500			1554
1501			1555
1502			
1503	Guanzhi Wang, Yuqi Xie, Yunfan Jiang, Ajay Mandlekar, Chaowei Xiao, Yuke Zhu, Linxi Fan, and Anima Anandkumar. 2023a. Voyager: An open-ended embodied agent with large language models. <i>arXiv preprint arXiv:2305.16291</i> .	Alexander Wei, Nika Haghtalab, and Jacob Steinhardt. 2024. Jailbroken: How does llm safety training fail? <i>Advances in Neural Information Processing Systems</i> , 36.	1556
1504			1557
1505			1558
1506			1559
1507			
1508	Jiaan Wang, Yunlong Liang, Fandong Meng, Zengkui Sun, Haoxiang Shi, Zhixu Li, Jinan Xu, Jianfeng Qu, and Jie Zhou. 2023b. Is chatgpt a good nlg evaluator? a preliminary study. <i>arXiv preprint arXiv:2303.04048</i> .	Jason Wei, Xuezhi Wang, Dale Schuurmans, Maarten Bosma, Fei Xia, Ed Chi, Quoc V Le, Denny Zhou, et al. 2022. Chain-of-thought prompting elicits reasoning in large language models. <i>Advances in neural information processing systems</i> , 35:24824–24837.	1560
1509			1561
1510			1562
1511			1563
1512			1564
1513	Pengda Wang, Zilin Xiao, Hanjie Chen, and Frederick L Oswald. 2024a. Will the real linda please stand up... to large language models? examining the representativeness heuristic in llms. <i>arXiv preprint arXiv:2404.01461</i> .	Jason Weston, Antoine Bordes, Sumit Chopra, Alexander M Rush, Bart Van Merriënboer, Armand Joulin, and Tomas Mikolov. 2015. Towards ai-complete question answering: A set of prerequisite toy tasks. <i>arXiv preprint arXiv:1502.05698</i> .	1565
1514			1566
1515			1567
1516			1568
1517			1569
1518			
1519	Siyuan Wang, Zhuohan Long, Zhihao Fan, Zhongyu Wei, and Xuanjing Huang. 2024b. Benchmark self-evolving: A multi-agent framework for dynamic llm evaluation. <i>arXiv preprint arXiv:2402.11443</i> .	Wikipedia contributors. 2004. <a href="#">Plagiarism — Wikipedia, the free encyclopedia</a> . [Online; accessed 22-July-2004].	1570
1520			1571
1521			1572
1522	Weihan Wang, Qingsong Lv, Wenmeng Yu, Wenyi Hong, Ji Qi, Yan Wang, Junhui Ji, Zhuoyi Yang, Lei Zhao, Xixuan Song, et al. 2023c. Cogvlm: Visual expert for pretrained language models. <i>arXiv preprint arXiv:2311.03079</i> .	Lionel Wong, Jiayuan Mao, Pratyusha Sharma, Zachary S Siegel, Jiahai Feng, Noa Korneev, Joshua B Tenenbaum, and Jacob Andreas. 2023. Learning adaptive planning representations with natural language guidance. <i>arXiv preprint arXiv:2312.08566</i> .	1573
1523			1574
1524			1575
1525			1576
1526			1577
1527	Wenhui Wang, Hangbo Bao, Li Dong, Johan Bjorck, Zhiliang Peng, Qiang Liu, Kriti Aggarwal, Owais Khan Mohammed, Saksham Singhal, Subhojit Som, et al. 2022a. Image as a foreign language: Beit pretraining for all vision and vision-language tasks. <i>arXiv preprint arXiv:2208.10442</i> .	Jialin Wu, Jiasen Lu, Ashish Sabharwal, and Roozbeh Mottaghi. 2022. Multi-modal answer validation for knowledge-based vqa. In <i>Proceedings of the AAAI conference on artificial intelligence</i> , volume 36, pages 2712–2721.	1578
1528			1579
1529			1580
1530			1581
1531			1582
1532			1583
1533	Xiyao Wang, Yuhang Zhou, Xiaoyu Liu, Hongjin Lu, Yuancheng Xu, Feihong He, Jaehong Yoon, Taixi Lu, Gedas Bertasius, Mohit Bansal, et al. 2024c. Mementos: A comprehensive benchmark for multimodal large language model reasoning over image sequences. <i>arXiv preprint arXiv:2401.10529</i> .	Qingyun Wu, Gagan Bansal, Jieyu Zhang, Yiran Wu, Shaokun Zhang, Erkang Zhu, Beibin Li, Li Jiang, Xiaoyun Zhang, and Chi Wang. 2023. Auto-gen: Enabling next-gen llm applications via multi-agent conversation framework. <i>arXiv preprint arXiv:2308.08155</i> .	1584
1534			1585
1535			1586
1536			1587
1537			1588
1538			1589
		Siye Wu, Jian Xie, Jiangjie Chen, Tinghui Zhu, Kai Zhang, and Yanghua Xiao. 2024. How easily do irrelevant inputs skew the responses of large language models? <i>arXiv preprint arXiv:2404.03302</i> .	1590
			1591
			1592
			1593





1701	Yadong Zhang, Shaoguang Mao, Tao Ge, Xun Wang,	Xizhou Zhu, Yuntao Chen, Hao Tian, Chenxin Tao, Wei-	1754
1702	Adrian de Wynter, Yan Xia, Wenshan Wu, Ting Song,	jie Su, Chenyu Yang, Gao Huang, Bin Li, Lewei Lu,	1755
1703	Man Lan, and Furu Wei. 2024a. Llm as a master-	Xiaogang Wang, et al. 2023b. Ghost in the minecraft:	1756
1704	mind: A survey of strategic reasoning with large	Generally capable agents for open-world environments	1757
1705	language models. <i>arXiv preprint arXiv:2404.01230</i> .	via large language models with text-based knowledge	1758
		and memory. <i>arXiv preprint arXiv:2305.17144</i> .	1759
1706	Zeyu Zhang, Xiaohe Bo, Chen Ma, Rui Li, Xu Chen,	Yongshuo Zong, Tingyang Yu, Bingchen Zhao, Ruchika	1760
1707	Quanyu Dai, Jieming Zhu, Zhenhua Dong, and Ji-	Chavhan, and Timothy Hospedales. 2023. Fool	1761
1708	Rong Wen. 2024b. A survey on the memory mech-	your (vision and) language model with embar-	1762
1709	anism of large language model based agents. <i>arXiv</i>	rassingly simple permutations. <i>arXiv preprint</i>	1763
1710	<i>preprint arXiv:2404.13501</i> .	<i>arXiv:2310.01651</i> .	1764
1711	Shu Zhao and Huijuan Xu. 2023. Less is more: Toward	Andy Zou, Zifan Wang, J Zico Kolter, and Matt Fredrik-	1765
1712	zero-shot local scene graph generation via foundation	son. 2023. Universal and transferable adversarial	1766
1713	models. <i>arXiv preprint arXiv:2310.01356</i> .	attacks on aligned language models. <i>arXiv preprint</i>	1767
		<i>arXiv:2307.15043</i> .	1768
1714	Yang Zhao, Zhijie Lin, Daquan Zhou, Zilong Huang,		
1715	Jiashi Feng, and Bingyi Kang. 2023. Bubogpt: En-		
1716	abling visual grounding in multi-modal llms. <i>arXiv</i>		
1717	<i>preprint arXiv:2307.08581</i> .		
1718	Zhonghan Zhao, Ke Ma, Wenhao Chai, Xuan Wang,		
1719	Kewei Chen, Dongxu Guo, Yanting Zhang, Hongwei		
1720	Wang, and Gaoang Wang. 2024. Do we really need a		
1721	complex agent system? distill embodied agent into a		
1722	single model. <i>arXiv preprint arXiv:2404.04619</i> .		
1723	Boyuan Zheng, Boyu Gou, Jihyung Kil, Huan Sun, and		
1724	Yu Su. 2024a. Gpt-4v (ision) is a generalist web		
1725	agent, if grounded. <i>arXiv preprint arXiv:2401.01614</i> .		
1726	Chujie Zheng, Hao Zhou, Fandong Meng, Jie Zhou, and		
1727	Minlie Huang. 2023a. Large language models are		
1728	not robust multiple choice selectors. In <i>The Twelfth</i>		
1729	<i>International Conference on Learning Representa-</i>		
1730	<i>tions</i> .		
1731	Lianmin Zheng, Wei-Lin Chiang, Ying Sheng, Siyuan		
1732	Zhuang, Zhanghao Wu, Yonghao Zhuang, Zi Lin,		
1733	Zhuohan Li, Dacheng Li, Eric Xing, et al. 2024b.		
1734	Judging llm-as-a-judge with mt-bench and chatbot		
1735	arena. <i>Advances in Neural Information Processing</i>		
1736	<i>Systems</i> , 36.		
1737	Shen Zheng, Jie Huang, and Kevin Chen-Chuan Chang.		
1738	2023b. Why does chatgpt fall short in providing		
1739	truthful answers. <i>ArXiv preprint, abs/2304.10513</i> .		
1740	Ming Zhong, Yang Liu, Da Yin, Yuning Mao, Yizhu		
1741	Jiao, Pengfei Liu, Chenguang Zhu, Heng Ji, and		
1742	Jiawei Han. 2022. Towards a unified multi-		
1743	dimensional evaluator for text generation. <i>arXiv</i>		
1744	<i>preprint arXiv:2210.07197</i> .		
1745	Wanjun Zhong, Lianghong Guo, Qiqi Gao, He Ye, and		
1746	Yanlin Wang. 2024. Memorybank: Enhancing large		
1747	language models with long-term memory. In <i>Pro-</i>		
1748	<i>ceedings of the AAAI Conference on Artificial Intelli-</i>		
1749	<i>gence</i> , volume 38, pages 19724–19731.		
1750	Deyao Zhu, Jun Chen, Xiaoqian Shen, Xiang Li, and		
1751	Mohamed Elhoseiny. 2023a. Minigpt-4: Enhancing		
1752	vision-language understanding with advanced large		
1753	language models. <i>arXiv preprint arXiv:2304.10592</i> .		

## A Orderability of Preferences.

**Comparability** When faced with any two alternatives A and B, the agent should have at least a weak preference, i.e.,  $A \succeq B$  or  $B \succeq A$ . This means that the agent can compare any pair of alternatives and determine which one is preferred or if they are equally preferred.

**Transitivity** If the agent prefers A to B and B to C, then the agent must prefer A to C. This ensures that the agent’s preferences are consistent and logical across multiple comparisons.

**Closure** If A and B are in the alternative set S, then any probabilistic combination of A and B (denoted as  $ApB$ ) should also be in S. This principle ensures that the set of alternatives is closed under probability mixtures.

**Distribution of probabilities across alternatives** If A and B are in S, then the probability mixture of  $(ApB)$  and B, denoted as  $[(ApB)qB]$ , should be indifferent to the probability mixture of A and B, denoted as  $(ApqB)$ . This principle ensures consistency in the agent’s preferences when dealing with probability mixtures of alternatives.

**Solvability** When faced with three alternatives A, B, and C, with the preference order  $A \succeq B \succeq C$ , there should be some probabilistic way of combining A and C such that the agent is indifferent between choosing B or this combination. In other words, the agent should be able to find a solution to the decision problem by making trade-offs between alternatives.

One consequence of the orderability is the concept of **dominance**: If alternative A is better than alternative B in terms of one attribute and at least as good in terms of all other attributes, the dominant option A should be chosen. An example of a fallacy that violates dominance is the sunk cost fallacy, where an agent continues to invest in a suboptimal alternative due to past investments, despite the availability of better options based on future outcomes.

## B Information Grounding

Web agents are a quintessential example of how multi-modal agents surpass language-only ones. In agents like Pix2Act (Shaw et al., 2024), WebGUM (Furuta et al., 2023), CogAgent (Hong et al., 2023b), and SeeAct (Zheng et al., 2024a), web navigation is grounded on graphical user interface (GUI) rather than solely on HTML texts (Shen et al., 2024a; Yao et al., 2022a; Deng et al., 2024; Gur et al., 2023). This method of visual grounding offers higher information density compared to HTML codes that are usually lengthy, noisy, and sometimes even incomplete (Zheng et al., 2024a). Supporting the importance of vision, ablation studies in WebGUM (Furuta et al., 2023) also reports 5.5% success rate improvement on the MiniWoB++ dataset (Shi et al., 2017; Liu et al., 2018) by simply adding the image modality.

Multi-modalities also help enhance the functionality of agent systems through more diverse information grounding. For example, Chain-of-Action (Pan et al., 2024) advances the single-modal Search-in-the-Chain (Xu et al., 2023) by supporting multi-modal data retrieval for faithful question answering. DoraemonGPT (Yang et al., 2024) decomposes complex tasks into simpler ones toward understanding dynamic scenes, where multi-modal understanding is necessary for spatial-temporal videos analysis. RA-CM3 (Yasunaga et al., 2022) augments baseline retrieval-augmented LLMs with raw multi-modal documents that include both images and texts, assuming that these two modalities can contextualize each other and make the documents more informative, leading to better generator performance. The multi-modal capabilities also allow HuggingGPT (Shen et al., 2024b), Agent LUMOS (Yin et al., 2023), ToolAlpaca (Tang et al., 2023), and AssistGPT (Gao et al., 2023b) to expand the scope of tasks they can address, including cooperation among specialized agents or tools capable of handling different information modalities.

Large world models is an emerging and promising direction to reduce multi-modal hallucinations. The notion is also mentioned in “Objective-driven AI” (LeCun, 2024), where agents have behavior driven by fulfilling objectives and they understand how the world works with common sense knowledge, beyond an auto-regressive generation. For example, Large World Model (LWM) (Liu et al., 2024b) and Sora (Brooks

1817 et al., 2024) develop insights from both textual knowledge and the world through video sequences.  
1818 Although these models both advance toward general-purpose simulators of the world, they still lack reliable  
1819 physical engines for guaranteed grounding in real-world dynamics. Ghost-in-the-Minecraft (Zhu et al.,  
1820 2023b) and Voyager (Wang et al., 2023a) have agents living in a well-defined game-world environment.  
1821 JEPA (LeCun, 2022) creates a recurrent world model in an abstract representation space.

## 1822 C Knowledge Retrieval & Tool Usage

1823 CuriousLLM (Yang and Zhu, 2024) presents ablation studies showing the effectiveness of KGs on  
1824 improving reasoning within the search process. MineDojo (Fan et al., 2022) observes that internet-scale  
1825 multi-modal knowledge allows models to significantly outperform all creative task baselines. Equipped  
1826 with world knowledge, RA-CM3 (Yasunaga et al., 2022) can finally generate faithful images from captions  
1827 compared to CM3 (Aghajanyan et al., 2022) and Stable Diffusion (Rombach et al., 2022). CooperKGC (Ye  
1828 et al., 2023) enables multi-agent collaborations, leveraging knowledge bases of different experts. It finds  
1829 that the incorporation of KGs improves F1 scores by 10.0-33.6% across different backgrounds, and adding  
1830 more collaboration rounds also enhance performance by about 10.0-30.0%. DoraemonGPT (Yang et al.,  
1831 2024) supports knowledge tools to assist the understanding of specialized video contents. SIRI (Wang et al.,  
1832 2023d) builds a multi-view knowledge base to increase the explainability of visual question answering.

1833 A multi-agent system can coordinate agents understanding when and which tool to use, which modality  
1834 of information the tool should expect, how to call the corresponding API, and how to incorporate outputs  
1835 from the API calls, which anchors subsequent reasoning processes with more accurate information beyond  
1836 their parametric memory. For example, VisProg (Gupta and Kembhavi, 2023), ViperGPT (Surís et al.,  
1837 2023), and Parsel (Zelikman et al., 2023) generate Python programs to reliably execute subroutines. Gupta  
1838 and Kembhavi (2023); Surís et al. (2023) also invoke off-the-shelf models for multimodal assistance.

1839 Foundation models are not specifically trained for object detection or segmentation, so BuboGPT (Zhao  
1840 et al., 2023) and Multi-Agent VQA (Jiang et al., 2024) call SAM (Kirillov et al., 2023; Ren et al., 2024)  
1841 as the tool. Besides, BabyAGI (Nakajima, 2023), Chamelon (Lu et al., 2024), AssistGPT (Gao et al.,  
1842 2023b), Avis (Hu et al., 2024), ToolAlpaca (Tang et al., 2023), MetaGPT (Hong et al., 2023a), Agent  
1843 LUMOS (Yin et al., 2023), AutoAct (Qiao et al., 2024),  $\alpha$ -UMi (Shen et al., 2024a), and ConAgents (Shi  
1844 et al., 2024) harness compositional reasoning to enable generalized multi-agent systems with planning  
1845 and modular tool-using capabilities in real-world scenarios.

1846 In most cases, tools require translating natural language queries into API calls with predefined syntax.  
1847 Once the APIs and their input arguments are determined, the tools will ignore any irrelevant context  
1848 in the original queries, as long as the queries share the same underlying logic necessary for the inputs.  
1849 Take Multi-Agent VQA (Jiang et al., 2024) as an example. In this system, a language model provides  
1850 only the relevant object names to the Grounded SAM (Ren et al., 2024) component, which functions  
1851 as an object detector, rather than passing the entire visual question. Similarly, the image editing tools  
1852 in VisProg (Gupta and Kembhavi, 2023) only receive a fixed set of arguments translated from user  
1853 queries to perform deterministic code executions. SeeAct (Zheng et al., 2024a) as a Web agent explores  
1854 vision-language models, ranking models, and a bounding box annotation tool to improve Web elements  
1855 grounding from lengthy and noisy HTML codes.

## 1856 D Collective Deliberation among Agents

### 1857 D.1 More Examples on Multi-Agent Collaborations

1858 Corex (Sun et al., 2023a) finds that orchestrating multiple agents to work together yields better complex  
1859 reasoning results, exceeding strong single-agent baselines (Wang et al., 2022b) by an average of 1.1-10.6%.  
1860 Retroformer (Yao et al., 2023) equips the single-agent Reflexion (Shinn et al., 2024) algorithm with an  
1861 additional LLM to generate verbal reinforcement cues and assist its self-improvement, enhancing accuracy  
1862 by 1.0-20.9%. MetaAgents (Li et al., 2023i) effectively coordinate agents within task-oriented social  
1863 contexts to achieve consistent behavior patterns, and the implementation of agent reflection in this system  
1864 leads to a 21.0% improvement in success rates. Multi-agent debating in Khan et al. (2024) also leads to  
1865 more truthful answers, boosting single-agent baselines by 28.0%. Multi-Agent Collaboration (Talebirad

and Nadiri, 2023), ChatDev (Qian et al., 2023), AgentCF (Zhang et al., 2023), AutoGen (Wu et al., 2023), Social Learning (Mohtashami et al., 2023), S<sup>3</sup> (Gao et al., 2023a), Ke et al. (2024), and Chern et al. (2024) continue to push the frontier of a multi-agent system’s applications beyond daily conversation to a versatile set of real-world task completions.

## D.2 Collaboration Against Jailbreaking

LLMs are also sensitive to prompt perturbations due to token bias and noises (Sclar et al., 2023a). One of the most worrying examples are adversarial attacks (Gehman et al., 2020; Ganguli et al., 2022; Du et al., 2022; Wei et al., 2024; Perez et al., 2022; Zou et al., 2023) through malicious prompt engineering. These attacks, also known as the Red Team Task, also named the Red Team Task, involve malicious prompt engineering designed to exploit vulnerabilities in the model. To combat this issue, Chern et al. (2024) propose a multi-agent debating approach involving agents with harmless, neutral, or harmful intentions. The authors demonstrate that engaging these agents in multi-round, multi-agent debate is more effective in improving the model’s robustness against adversarial prompt variations and perturbations compared to a single-agent with self-reflection prompts.

## D.3 Collaboration on LLM-based Evaluation

LLM-based evaluation methods are popular in assessing open-ended language responses. Stureborg et al. (2024); Koo et al. (2023) point out LLMs often present cognitive biases in their evaluations, favoring certain types of responses over others regardless of the actual quality or relevance of the respective responses. To establish a more coherent preference orderability aligned with human preference. ChatEval (Chan et al., 2023) introduces a multi-agent debate framework to mimic human annotators collaborating in robust answer evaluations. Its multi-agent approach achieves greater alignment with human preferences compared to single-agent evaluations, enhancing accuracy by 6.2% for GPT-3.5 and 2.5% for GPT-4, and an increase of 16.3% and 10.0% in average Spearman and Kendall-Tau correlations (Zhong et al., 2022) with human judgements in GPT-4.

## D.4 The Orderability of Preferences Matters for LLM-based Evaluations

This section talks about LLM-based evaluation rather than evaluating the rationality of LLMs discussed in Section 5. Recent research underscores a critical need for more rational LLM-based evaluation methods, particularly for assessing open-ended language responses. CoBBLER (Koo et al., 2023) provides a cognitive bias benchmark for evaluating LLMs as evaluators, revealing a preference for their own outputs over those from other LLMs. Stureborg et al. (2024) argues that LLMs are biased evaluators towards more familiar tokens and previous predictions, and exhibit strong self-inconsistency in the score distribution. Luo et al. (2023); Shen et al. (2023); Gao et al. (2023c); Wang et al. (2023b); Chen et al. (2023); Chiang and Lee (2023); Zheng et al. (2024b); Fu et al. (2023); Liu et al. (2023b) also point out the problem with a single LLM as the evaluator, with concerns over factual and rating inconsistencies, a high dependency on prompt design, a low correlation with human evaluations, and struggles with the comparison. As a result, having a coherent orderability of preferences aligned with human preference becomes increasingly important.

Multi-agent systems might be a possible remedy. By involving multiple evaluative agents from diverse perspectives, it becomes possible to achieve a more balanced and consistent orderability of preferences. For instance, ChatEval (Chan et al., 2023) posits that a multi-agent debate evaluation usually offers judgments that are better aligned with human annotators compared to single-agent ones. Bai et al. (2024) also finds decentralized methods yield fairer evaluation results.

## E Neuro-Symbolic Reasoning

Logic-LM (Pan et al., 2023) combines problem formulating, symbolic reasoning, and result interpreting agents, where the symbolic reasoner empowers LLMs with deterministic symbolic solvers to perform inference, ensuring a correct answer is consistently chosen. Its multi-agent framework also encourages self-refinement that modifies logical formulation errors using error messages from the symbolic reasoner as the feedback. Besides, SymbolicToM (Sclar et al., 2023b) and KRISP (Marino et al., 2021) construct

1914 explicit symbolic graphs and answer questions by retrieving nodes in the graph. Binder (Cheng et al.,  
1915 2022), Parsel (Zelikman et al., 2023), LEFT (Hsu et al., 2024), and Fang et al. (2024) decompose tasks  
1916 into planning, parsing, and execution, where the symbolic reasoning agents can help maintain a coherent  
1917 order of preferences among symbolic options in the system outputs.

## 1918 **F Evaluating Rationality**

### 1919 **F.1 Benchmarks for Hallucination**

1920 Multiple evaluation benchmarks targeting language-only dialogue have been proposed, such as BE-  
1921 GIN (Dziri et al., 2022b), HaluEval (Li et al., 2023e), DialFact (Gupta et al., 2021), FaithDial (Dziri  
1922 et al., 2022a), AIS (Rashkin et al., 2023), and others (Zheng et al., 2023b; Das et al., 2023; Cao et al.,  
1923 2021). In contrast, *benchmarks on multi-agent frameworks beyond language dialogue or those involving*  
1924 *multi-modalities are very limited.* Liu et al. (2024a) moves beyond conversation to code generation;  
1925 EureQA (Li et al., 2023a) focuses on reasoning chains; and TofuEval (Tang et al., 2024) evaluates  
1926 hallucination in multi-domain summarization. Object hallucination (Rohrbach et al., 2018; Biten et al.,  
1927 2022), POPE (Li et al., 2023g), and LLaVA-RLHF (Sun et al., 2023b) are the few examples evaluating  
1928 multi-modal hallucination.

### 1929 **F.2 Perturbation Techniques**

1930 Perturbation techniques typically involve some versions of paraphrasing or permutation. Paraphrasing  
1931 includes changing the instruction templates (Weber et al., 2023), rewording task descriptions (Yang et al.,  
1932 2023; Ohmer et al., 2024; Wang et al., 2024b), translating the prompts into a different language (Ohmer  
1933 et al., 2023, 2024; Xu et al., 2024a) and then back to the original language (Yang et al., 2023), and making  
1934 subtle changes to entities in task descriptions without affecting the logical structure, like altering names  
1935 of the characters, numerical values in math problems, or locations of the events (Wang et al., 2024b).  
1936 Permutation also includes reordering in-context learning examples (Lu et al., 2021; Pecher et al., 2024)  
1937 and, in the case of multiple-choice questions, rearranging the options (Zong et al., 2023; Zheng et al.,  
1938 2023a).