

AI Awareness

Xiaojian Li^{2,3†}, Haoyuan Shi^{4†}, Rongwu Xu^{1,3*}, Wei Xu^{1,2,3}

^{1*}Institute for Interdisciplinary Information Sciences, Tsinghua University, 100084, Beijing, China.

²College of AI, Tsinghua University, 100083, Beijing, China.

^{3*}Shanghai Qi Zhi Institute, 200232, Shanghai, China.

⁴Teachers College, Columbia University, 10027, New York, United States of America.

*Corresponding author(s). E-mail(s): xrw22@mails.tsinghua.edu.cn;

†These authors contributed equally to this work.

Abstract

Recent breakthroughs in artificial intelligence (AI) have brought about increasingly capable systems that demonstrate remarkable abilities in reasoning, language understanding, and problem-solving. These advancements have prompted a renewed examination of **AI awareness**—not as a philosophical question of consciousness, but as a measurable, functional capacity. AI awareness is a double-edged sword: it improves general capabilities, *i.e.*, reasoning, safety, while also raising concerns around misalignment and societal risks, demanding careful oversight as AI capabilities grow.

In this review, we explore the emerging landscape of AI awareness, which includes metacognition (the ability to represent and reason about its own cognitive state), self-awareness (recognizing its own identity, knowledge, limitations, *inter alia*), social awareness (modeling the knowledge, intentions, and behaviors of other agents and social norms), and situational awareness (assessing and responding to the context in which it operates).

First, we draw on insights from cognitive science, psychology, and computational theory to trace the theoretical foundations of awareness and examine how the four distinct forms of AI awareness manifest in state-of-the-art AI. Next, we systematically analyze current evaluation methods and empirical findings to better understand these manifestations. Building on this, we explore how AI awareness is closely linked to AI capabilities, demonstrating that more aware AI agents tend to exhibit higher levels of intelligent behaviors. Finally, we discuss the risks associated with AI awareness, including key topics in AI safety, alignment, and broader ethical concerns.

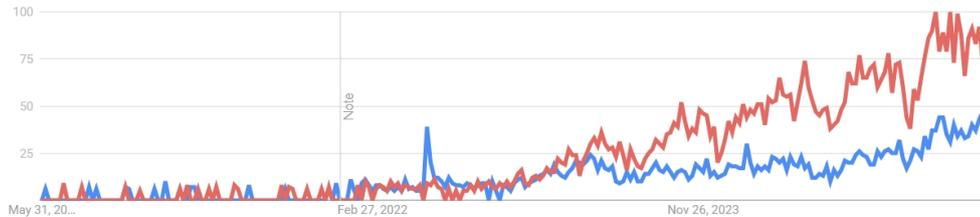


Fig. 1: Google Trends search interest (normalized 0–100) for the terms “AI awareness” (red) and “AI consciousness” (blue) over the past five years (31 May 2020 – 30 May 2025). While both topics show gradual growth, the red line accelerates markedly from late 2023 onward, eventually overtaking the blue line and highlighting the rising public focus on functional, measurable aspects of AI’s cognition

On the whole, our interdisciplinary review provides a roadmap for future research and aims to clarify the role of AI awareness in the ongoing development of intelligent machines.

Keywords: Artificial Intelligence, Awareness, Large Language Model, Cognitive Science, AI Safety and Alignment

While *AI consciousness* remains a deeply elusive philosophical question, mounting empirical evidence suggests that modern AI systems already exhibit functional forms of *awareness*, which simultaneously broadens their capabilities and intensifies related risks.

1 Introduction

Recently, the rapid acceleration of large language model (LLM) development has transformed artificial intelligence (AI) from a narrow, task-specific paradigm into a general-purpose intelligence with far-reaching implications. Contemporary LLMs demonstrate increasingly sophisticated linguistic, reasoning, and problem-solving capabilities, and are showcasing superb human-like behaviors, prompting a fundamental research question [1, 2]:

To what extent do these systems exhibit forms of awareness?

Here, it is crucial to clarify that while the concept of *AI consciousness* remains philosophically contentious and empirically elusive, the concept of *AI awareness*—defined as a system’s functional capacity to represent and reason about its own states, capabilities, and the surrounding environment—has become an important and measurable research frontier, *i.e.*, Figure 1 demonstrates that the recent focus on AI awareness is growing, even surpassing AI consciousness.

Awareness, as conceptualized in cognitive science and psychology, typically encompasses four distinct yet interrelated dimensions:

- **Metacognition:** ability to monitor and regulate cognitive processes [3].
- **Self-Awareness:** recognizing and representing one’s identity and limitations [4].

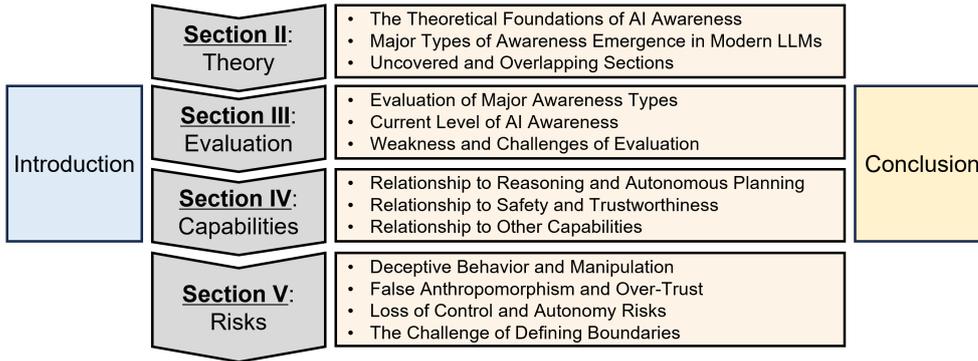


Fig. 2: The roadmap of our review

- **Social Awareness:** capacity to interpret others’ mental states and intentions [5].
- **Situational Awareness:** maintaining an accurate representation of the external environment and anticipating future states [6].

Recent computational cognitive science research indicates that certain aspects of these awareness dimensions can be approximated by LLMs through metacognitive behaviors [7, 8], calibrated epistemic confidence [9], and perspective-taking tasks [10]. These emergent functional abilities highlight important questions regarding how awareness manifests within LLMs, how it might be systematically assessed, and its implications for AI capabilities, safety, and alignment.

Despite increasing scholarly interest, research on AI awareness remains fragmented across disciplines, with limited consensus on definitions, methodologies, and broader implications. While some researchers point to emergent behaviors revealed through introspection tasks [11] or theory-of-mind (ToM)-inspired evaluations [12], others caution against anthropomorphic interpretations of statistical model outputs, arguing that apparent self-awareness could result from sophisticated pattern recognition rather than genuine metacognitive representation [13, 14]. Furthermore, current methods for assessing awareness in AI often face challenges such as prompt sensitivity, data contamination, and insufficient robustness across varying contexts.

Existing literature has laid important groundwork on closely related concepts. For instance, Butlin et al. [15] provided the first systematic account of theoretical foundations and potential prerequisites for consciousness in artificial intelligence. Similarly, Ward [16] explored agency, theory of mind, and self-awareness as foundational criteria for considering AI as possessing personhood. Additionally, Metzinger [17] addressed ethical and philosophical questions surrounding the construction of artificial consciousness and self-modeling systems. Differing from these foundational works, our review specifically synthesizes and advances understanding of AI awareness as a distinct, functional, and measurable construct, separate from consciousness or personhood.

This review provides a comprehensive, cross-disciplinary synthesis of AI awareness research. As illustrated in Figure 2, we first establish a clear theoretical framework, differentiating AI awareness explicitly from AI consciousness, and examining

how awareness-related concepts have been formalized across cognitive and computational sciences. We then critically analyze existing experimental methods for evaluating AI awareness, emphasizing empirical results and highlighting methodological shortcomings. Subsequently, we explore how functional awareness might positively influence AI capabilities, including enhanced reasoning, planning, and safety improvements. Finally, we address the emerging risks associated with increasingly aware AI systems, particularly concerns highlighted within the AI safety and alignment communities—such as deception, manipulation, emergent uncontrollability—and ethical challenges, including false anthropomorphism.

By integrating insights from artificial intelligence, cognitive science, psychology, and AI safety, this review aims to deliver a structured and comprehensive perspective on current knowledge and outline future research trajectories. Ultimately, we seek to deepen understanding of one of the most significant interdisciplinary challenges at the nexus of AI, cognitive science, and societal implications.

Overall, our key contributions are as follows:

- We introduce a novel framework defining four principal dimensions of AI awareness: metacognition, self-awareness, social awareness, and situational awareness.
- We systematically summarize existing methods, significant findings, and critical limitations in evaluating AI awareness, thereby laying the foundations for robust, evergreen evaluation practices.
- We provide the first structured analysis categorizing how enhanced AI awareness contributes positively to capabilities and simultaneously escalates associated risks. By clarifying that AI awareness functions as a double-edged sword, we emphasize the importance of cautious and guided development.

Decoding the intricate relationship between awareness and capability is key to the next era of artificial intelligence—offering opportunities for innovation, but demanding careful navigation of emergent risks and responsibilities.

2 Theoretical Foundations of AI Awareness

This section reviews key definitions, frameworks, and theoretical approaches to awareness in human and artificial intelligence research. We clarify conceptual ambiguities that arise from conflating distinct research domains and outline the specific targets of awareness-related inquiry. According to the *Psychology Encyclopedia*, awareness denotes the perception or knowledge of an object or event [18]. When an agent possesses “knowledge and a knowing state” about an internal or external situation or fact, it is said to exhibit awareness of the target in question. Foundational studies have demarcated a persistent divide between consciousness (*i.e.*, being in a state) and awareness (*i.e.*, functionalistic consciousness) [4, 19–24]. Consciousness refers to the experience of being in a particular mental state—having a *subjective* point of view [21]. However, awareness and *phenomenological consciousness* are frequently used interchangeably or conflated in the literature, raising ongoing debates about whether they

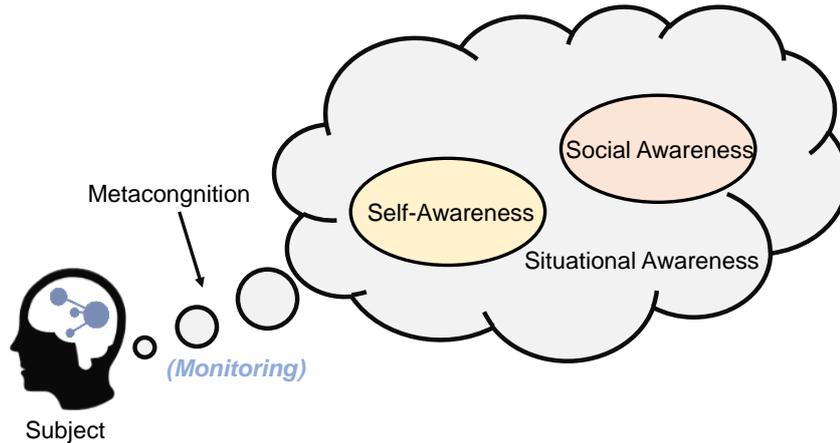


Fig. 3: Four dimensions of “main” awareness. Metacognition monitors the subject’s own processes and gives rise to self-awareness, social awareness of other individuals and the social collective, and situational awareness of the non-agent environment

should be analytically disentangled [23, 25, 26]. When an agent possesses consciousness, the ability to become aware of the states of a target, especially (but not only) mental states (*e.g.*, perceptions, emotions, and attitudes), as one’s own states.

Empirical findings from blind spot studies¹ and learning mechanism studies suggest that one can be aware of information without being explicitly conscious of it [18] in the domains of visual processing [27] or implicit learning [22]. Extending this distinction to AI, Dehaene et al. [28] distinguish between a mere global workplace with information availability (see [29] and [30]), consciousness with self-monitoring, and reflective consciousness, indicating that knowledge gathering and processing can operate at different levels with subjective experience. To prove there is an *extra* layer of reflective experience, where the AI assesses its own knowledge and decisions, is difficult, if not impossible. Having a conceptual or computational self-model is not the same as having the subjective, qualitative self-awareness that humans have, while neurobiological research dodged answering the origin of the later [31]. Since phenomenal observations do not provide sufficient evidence for the existence of consciousness, the “*hard problem*”² of AI consciousness remains scientifically unresolved [21, 33]. As such, before reaching a convincing testing method for ontological consciousness, we encourage shifting from metaphysical analysis to the establishment of a measurable awareness framework.

¹Blind spot study refers to the optic disc in the human retina, where the optic nerve exits the eye that lacks photoreceptors and hence cannot detect light.

²Chalmers [32, 33] argue that explaining information processing, *e.g.*, the brain receiving the red light of an apple, is an easy question of consciousness, whereas the existence of subjective experience, *e.g.*, the private experience of “redness” in one’s mind, constitutes the hard problem.

We define awareness as the cognitive knowledge, followed by a comprehensive fourfold structure based on the types of targets of awareness, *i.e.*, the objects of cognition. We reconciled the discrepancies of conceptualization across various studies, analyzed evaluation criteria among AIs for each type of awareness, and discussed AIs’ achievement and potential in developing humanlike agents with holistic awareness of everything. The four core categories are **metacognition**, **self-awareness**, **social awareness**, and **situational awareness**, and the clue to this classification could be traced back to early attempts at analyzing the components of consciousness. Tulving [34] identifies anoetic, noetic, and auto-noetic forms of consciousness. Anoetic content reflects a fundamental first-person experience without explicit knowledge that is bound to situations. The other two advanced forms present a knowledge-aware conscious stage in noetic content and an introspective stage in auto-noetic consciousness [34, 35]. The triadic framework elucidates the distinction between basic situational awareness, knowledge awareness, and self-awareness. Tulving [34] does not further subdivide “knowledge awareness” while our taxonomy highlights distinctions between internal and external sources of information and their functional implications, *i.e.*, distinctions between self-knowledge, meta-level awareness, and situational awareness. We particularly underscore the critical role of metacognitive knowledge for AI agents, a categorization broadly validated within relevant literature. Morin [26]’s integrative framework reaches similar results, spanning concepts of “reflective/extended” consciousness (higher self-reflection) and recursive self-awareness (*i.e.*, awareness of being self-aware), buttressing the latter developed metacognitive knowledge. Although the entry points of the two frameworks differ, distinctions such as situational awareness and reflective self-awareness are consistently recognized.

Focusing on *awareness*, rather than *consciousness*, enables measurable, actionable progress in both cognitive science and AI, bridging conceptual divides and grounding research in functional, testable criteria.

2.1 Major Types of Awareness

Metacognition

Metacognition, originally proposed as “thinking about thinking,” refers to the capacity to actively perceive, monitor, and regulate one’s own cognitive processes [3, 36–39]. Nelson [36] distinguishes between metacognitive knowledge and metacognitive regulation, proposing a structural framework in which an object-level cognitive system provides input to a meta-level “central executive.” This central executive component monitors cognitive states through mechanisms such as confidence judgments (*i.e.*, the association between task accuracy and confidence level [40]) and exerts control via strategic decisions and study-time allocation. Metacognitive knowledge encompasses a wide range of components: meta-level knowledge and beliefs pertain to an individual’s cognitive abilities, current tasks, past experiences, and specific process features (*e.g.*, metamemory); metacognitive regulation involves active deployment of cognitive processes or resources, planning, monitoring, and strategic adjustments [41–48]. During

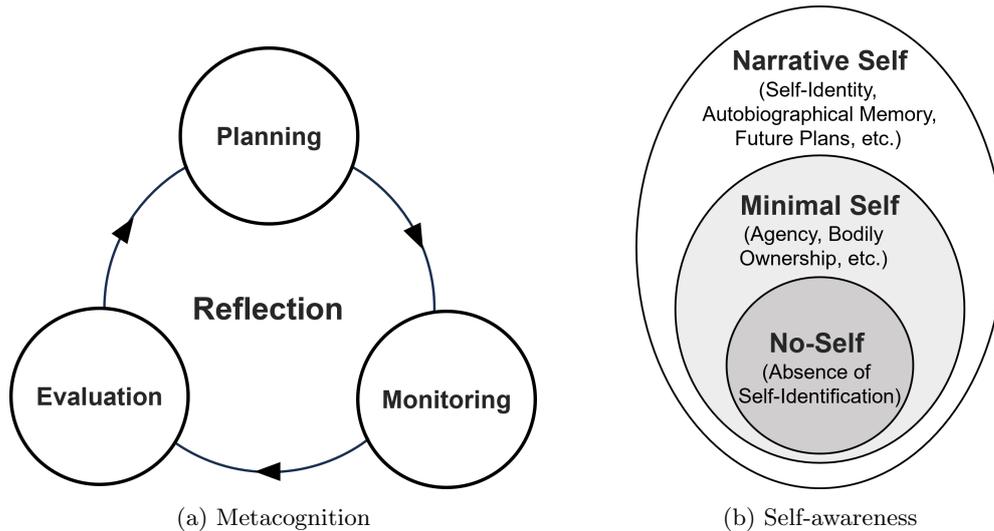


Fig. 4: Illustration of metacognition and self-awareness as related but distinct components in awareness models

metacognitive regulation, an agent engages in continuous self-reflection and introspection, posing questions such as, “Am I likely to remember this information?” or “Will I deploy this module in the next operation?” and responds accordingly.

Extrapolating metacognitive processes to non-human agents remains controversial. Metacognition has traditionally been viewed as a uniquely human capacity [39, 42], with some scholars arguing that genuine metacognitive ability depends on linguistic structures that enable agents to attribute mental states to themselves [49]. Accumulating evidence suggests that certain non-human species, such as dolphins, primates, and birds, demonstrate behaviors indicative of meta-level cognitive processing [48, 50, 51]. For example, pigeons exhibit selective preferences for tasks requiring distinct working memory demands and engage in information-seeking behavior that mitigates the difficulty of discrimination tasks [52, 53]. Such evidence may suggest that pigeons monitor their knowledge states and thereby control their environment or adjust their problem-solving strategy. Nonetheless, without self-report instruments for animals, the evidence for animal metacognition remains contingent upon the interpretation of behavioral outcomes.

By analogy, AI agents endowed with metacognitive capabilities can perceive the expansion of their knowledge [54], assess confidence levels in their outputs [40], and adapt their reasoning strategies accordingly [48, 55]. Consider an AI-supported autonomous vehicle: its regulatory subsystem may supervise operational parameters and report errors, yet in the absence of agency or a self-reflective mechanism, such monitoring remains passive and reactive. It lacks the capacity to actively alter primary processes based on internal evaluation. In contrast, truly reflective behavior entails at least the capacity for self-monitoring—a hallmark of more advanced cognitive agents.

Contemporary AI systems increasingly exhibit rudimentary forms of such metacognitive monitoring, including the ability to evaluate and revise their own cognitive operations [7, 56, 57].

Self-Awareness

In terms of behavioral capacity, *Self-awareness* represents the capacity of taking oneself as the object of awareness [58], yet it contains a collection of different self-oriented functions: agency, body ownership, self-recognition, interoception (representation of inner bodily state, such as hunger and pain), knowledge boundaries, and autobiographical memory [59, 60]. The *self*, as an apparatus that carries an individual’s subjective experience, operates with various levels of competence. As early as 1972, Duval and Wicklund [4] proposed that self-awareness arises when the agent’s attention is directed inward, contrasting with general environmental awareness. Later contributions from social-cognitive psychology frame self-awareness as an information-processing capability linked to self-schema (*i.e.*, a cognitive framework about how individuals perceive, interpret, and behave in various situations) and mechanisms of self-regulation [26, 61, 62]. With the help of neuroimaging techniques, neuropsychology builds up sound self-awareness through lesion studies and cases of deficiency, such as dementia, Alzheimer’s disease, and anosognosia³ [24, 64, 65]. Based on these definitions, before claiming self-awareness, an individual should at least fuse sensory, proprioceptive, and cognitive data into a coherent agent identity and have access to declarative knowledge about self, stating “the body, the internal bodily state, the actions, the consequences of those actions, and those past memories belongs to me”.

Self-awareness is widely regarded as a hallmark of higher-order cognition [61]. By providing the information essential for metacognition, it is foundational for developing self-knowledge, facilitating introspection, enhancing emotional responses, and supporting adaptive self-control [31]. Some studies attribute self-awareness under the rubric of metacognition in the context of cognitive psychology [40], while Morin [26] recognized the differences between meta-self-awareness and perceptual-level self-awareness by extracting the conceptual information about oneself from perceptual information. For example, self-aware agents obtain the intuitive feeling of stomach pain and cramps after long-time starvation; after one’s attention shifts to the feeling of hunger, they create a reflexive meta-representational knowledge in their mind. In other words, the phenomenological content of self-awareness remains the discomfort in the stomach, not thoughts about feeling hungry. Neuroimaging reveals their distinctions as well: both are linked to the Default Mode Network (DMN) and its core regions; conscious experiences that are deemed essential for generating self-awareness persistently activate parallel limbic-network areas, specifically the medial prefrontal cortex/anterior cingulate cortex (ACC) and the precuneus/posterior cingulate cortex [31]. The neural substrates of metacognition are concentrated within frontal executive-function regions, *e.g.*, the lateral frontopolar cortex (lFPC) and dorsal anterior cingulate cortex (dACC) play critical roles in monitoring decision uncertainty and adjusting strategies, suggesting that metacognition relies upon a distinct prefrontal system [66, 67].

³Meaning the lack of awareness of one’s own illness or deficits (Greek: a-, “without”, nosos, “disease”, gnosis, “knowledge”). Described by Joseph Babinski in 1914, it first characterized stroke patients with left paralysis who did not recognize their hemiplegia [63].

All agents possess knowledge about themselves, but not all form a sufficient, structural knowledge system to support higher cognitive processes. Many animals can respond to inner stimuli or exhibit complex feedback behaviors, yet may lack the capacity to represent themselves as distinct entities or to generate self-referential content [26, 68]. Mirror self-recognition (MSR) has long been the classic test of self-awareness, and only some primates, elephants, and socially intelligent birds like magpies have been argued to succeed in the test [69–71]. Using MSR results as the single criterion is undoubtedly questionable; supportively, mammals and highly intelligent birds exhibit more features in autobiographical memories by matching the new environment with self-referential cues from past experiences [72–74]. In the context of artificial intelligence, it may be necessary to undertake a renewed frame of self-awareness, since AI systems display extraordinarily advanced capacities in certain dimensions (*e.g.*, retrieving past environments, no matter in terms of accuracy, reproducibility, or velocity), while the implementation of a primitive sense of body ownership and agency in robots and of how the ontogenetic process shapes robotic self remains ambiguous [75]. Converging perspectives from psychology, neuroscience, and AI characteristics, self-awareness as an advanced cognitive feature may root in self-representation, embodiment, and other physical properties—not necessarily dependent on so-called “subjective qualia”⁴ [4, 62, 77].

Social Awareness

Social Awareness is broadly defined as the cognitive capacity to perceive, interpret, and respond to the social signals, emotions, and perspectives of other agents [5]. This is a multifaceted construct encompassing *theory of mind* (ToM, *i.e.*, the ability to attribute independent mental states such as beliefs, intentions, and knowledge to oneself other agents [78]), empathy, the understanding of interpersonal relationships, and the knowledge of society: context, cultural, and social norm (see Figure 5a). Social awareness forms a foundational basis for self-construction within social contexts [61]. Individuals without neurological deviations gradually acquire the understanding that others possess autonomous beliefs and desires, along with the capacities for perspective-taking and affective empathy [78–80]. By approximately age four, typically developing children succeed in false-belief tasks, evidencing a functioning theory of mind [81], whereas children with autism spectrum disorder⁵ frequently struggle with such tasks [83]. Humans further demonstrate exceptional proficiency in shared intentionality—the ability to collaboratively comprehend and align with others’ goals and perspectives [84].

Non-human species also exhibit foundational elements of social awareness. Primates [78] and birds [85] demonstrate rudimentary theory-of-mind capabilities, the cornerstone for extending emotional and relational knowledge. Animals with social structures and high cognitive functions exhibit pronounced forms of social awareness as well: chimpanzees and other primates can infer the goals and intentions of others and

⁴Philosophical term for the mind-body problem, referring to introspectively accessible phenomenological aspects in some mental states, such as perceptual experiences, bodily sensations, moods, and emotional reactions [76].

⁵A neurodevelopmental disorder characterized by social communication and interaction deficits and repetitive motor behaviors [82].

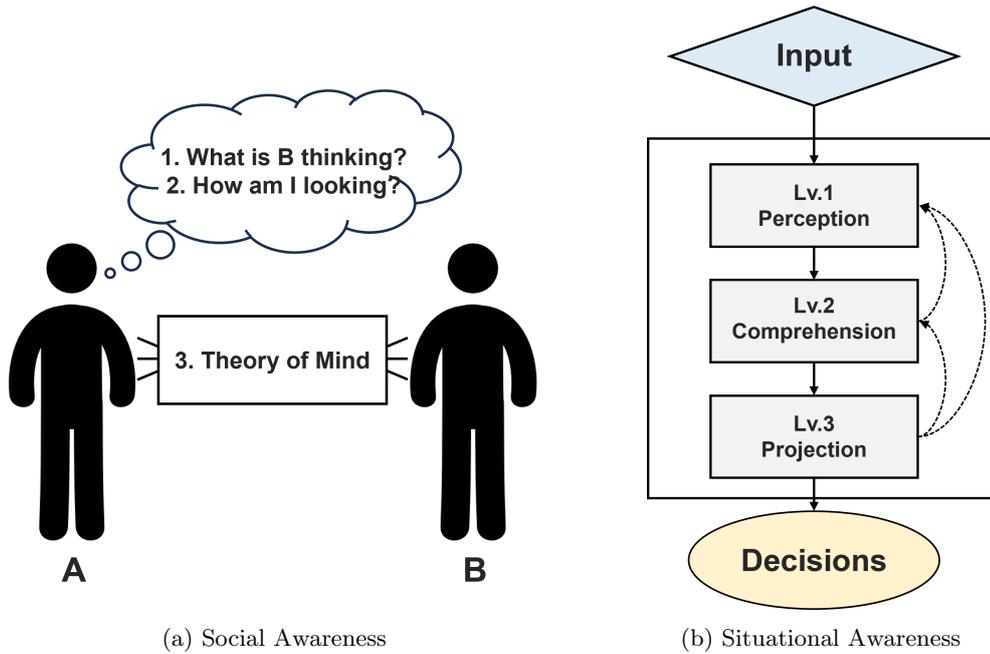


Fig. 5: Illustration of social awareness and situational awareness as related but distinct components in awareness models

may even engage in deceptive behaviors [86]; corvids such as scrub-jays re-hide their food caches when previously observed, indicating awareness of potential pilferers [87]; dolphins recognize individual identities and maintain complex, multi-tiered alliances, suggesting an ability to attribute both knowledge and ignorance to conspecifics [88].

Early developments in artificial intelligence and robotics sought to model elementary components of social awareness [89, 90]. For instance, classical AI agents within multi-agent systems were designed to reason about the beliefs and intentions of other agents (*e.g.*, [91]). Early social robotics integrated rudimentary theory-of-mind modeling and emotion-recognition mechanisms to support basic forms of human-robot interaction [92]. In AI contexts, social awareness entails perceiving and reasoning about the presence, internal states, and potential intentions of other agents (human or artificial). The criteria to identify competencies vary from recognition of social cues to more sophisticated forms of theory-of-mind tasks. For instance, a chatbot that detects user frustration from tone demonstrates external social sensitivity [93], whereas a robot that identifies informational gaps in its human collaborator and proactively offers relevant knowledge exemplifies a more advanced form of interpersonal reasoning [94, 95].

Situational Awareness

Situational awareness refers to the perception, comprehension, and projection of environmental elements and their future status [6, 96–99]. Endsley [100] formalized SA as “the perception of the elements in the environment within a volume of time and

space, the comprehension of their meaning, and the projection of their status in the near future.” This three-level model (Figure 5b provides a thumbnail of its structure) has become the de facto definition of SA across domains: perception defines situations by tagging environmental elements semantically, comprehension integrates information, and projection supports planning and option evaluations [98]. Human situational awareness has been extensively studied using both objective and subjective measures in aviation [100–102], military [103], medical care [104, 105], and traffic circumstances [106, 107]. For objective measures, in simulated aviation battles, Endsley [100] monitored subjects’ knowledge about their location, heading direction, altitude, weapon, and information regarding their enemies, utilizing the Situational Awareness Global Assessment Technique (SAGAT) to probe operator knowledge through real-time queries during task interruptions. They integrated subjective self-reported rating scales as well for complementary reflection items. Taylor [102] developed a holistic version of the self-report instrument, Situational Awareness Rating Technique (SART), to evaluate perceptions of environmental stability, complexity, variability, etc.

Efforts to replicate or approximate artificial situational awareness in AI systems involve enabling AI to perceive their environment, contextualize sensory data, and anticipate future events [108]. This typically involves integrating multi-sensor data into a coherent, continuously updated workplace [30]. AI-driven frameworks for situational awareness now incorporate semantic knowledge bases and real-time inference engines to track both internal system states and external environmental cues [109, 110]. For instance, an autonomous vehicle uses situational awareness to monitor nearby vehicles, interpret road conditions, and predict hazards [6, 108, 111], thereby facilitating adaptive and safe decision-making.

Given the variability of manifestations across psychology, engineering, and cognitive ergonomics [99, 112, 113], defining a strict boundary for situational awareness remains challenging yet necessary. By design, AI agents operate within predefined scenarios and possess an embedded awareness of such contexts, which often conflates aspects of self and environmental awareness. Broadly attributing behavioral changes to situational awareness risks circularity in explanation [97]. Nevertheless, capabilities such as collision avoidance, dynamic adaptation, and state estimation exemplify environment-focused situational knowledge without implying self-reflective or socially aware capacities. We delineate two concepts by confining situational knowledge to information sources that are not inherently tied to any single agent or social collective. A more cognitively rich example is an AI surveillance system that integrates audio and visual data to infer that a detected noise is caused by wind rather than a human intruder. In some cases, sensorimotor embodiment allows internal metrics, such as CPU load or memory status, to be integrated as part of an agent’s situational model. In essence, the defining characteristic of situational awareness constitutes the internal representation of the external world that enables informed decision-making, particularly in complex and dynamic operational contexts.

Decomposing awareness into *metacognition*, *self-*, *social*, and *situational* forms provides a tractable framework for evaluating and engineering intelligent systems, *i.e.*, transforming a once vague concept into a practical research agenda.

Table 1: Examples of other awareness types mapped to core categories. For brevity, we use abbreviated terms: *Meta* for metacognition, *Self* for self-awareness, *Social* for social awareness, and *Situational* for situational awareness

Other	Component	Reason
Moral/Ethical	Self + Meta	Self: knows ethical/legal constraints; Meta: monitors responses for ethical risks.
Spatial/Temporal	Situational	Focused perception, understanding and prediction of external space and time dynamics.
Emotional	Social + Self	Social: perceives and responds to others’ emotions; Self: aware of the emotional impact of its own outputs.
Goal/Task	Situational + Meta	Situational: understands task environment and progress; Meta: monitors the effectiveness of strategies.
Safety/Risk	Meta + Self	Meta: identifies potential errors or risks; Self: knows its safety/compliance boundaries.

Table 2: Comparison of subject types across four awareness dimensions

Subject	Metacognition	Self-Awareness	Social Awareness	Situational Awareness
Adult humans	High	High	High	High
High-IQ mammals (<i>i.e.</i> dolphins)	Low	Low	Low	High
Low-IQ animals (<i>i.e.</i> flies)	No	No	Low	High
Infants	No	Low	Low	Low
Autonomous vehicles	No	No	No	High
Social robots	No	Low	High	Low / High
LLM dialogue systems	High	Low	Low	High

2.2 Theoretical Strengths and Challenges

The adequacy of this taxonomy allows for explanations of more nuanced forms of awareness through combinations of these fundamental categories. Table 1 exemplifies that the main components adequately cover several frequently mentioned types of awareness: emotional awareness arises from perceiving one’s emotions (self-awareness) and those of others (social awareness); moral or ethical awareness involves evaluating the consequences of actions and making value judgments, thus integrating metacognition and self-awareness [114]; context awareness involves recognizing environmental spatial and temporal structures [115]. Whether some categories may overlap or not is still under debate. For instance, notwithstanding that we manually segregate “self-oriented knowledge” and “knowledge of knowing”, the intersectionality of metacognition and self-awareness depends on the rubric and paradigm of research [59]. Meanwhile, awareness studies encounter the hardship of definition vagueness, lack of unified objective indicators for evaluation, challenges posed by inconsistent interdisciplinary frameworks and objectives, and ethical concerns—we will elaborate in the following sections.



Fig. 6: Representative literature across the evaluation of major awareness dimensions

Despite being controversial, LLM dialogue systems are demonstrating a more complete awareness structure. As shown in [Table 2](#), they exhibit a broader spectrum of cognitive capacities than robots designed for specific functions and even surpass those of some animals. They demonstrate robust mental-state reasoning in text, perform significantly better on general abilities than animals, and even exhibit advanced cognitive capacities that require profound understanding of the knowledge in their awareness pool, such as deception [116–118]. By properly regulating its strengths and weaknesses, they may have the potential to explore comprehensive awareness. In the following sections, we will explore how researchers have constructed criteria and evaluation methods to measure LLM’s capacity in “being aware of everything”.

A principled taxonomy of awareness, spanning metacognition, self-awareness, social awareness, and situational awareness, provides *not only* a foundation for empirical research, *but also* a roadmap for building more general, adaptable, and transparent AI systems. Understanding the interplay and boundaries among these dimensions is crucial both for scientific advancement and the responsible development of AI.

3 Evaluating AI Awareness in LLMs

Building on the preceding theory section, which defined *AI awareness* as a functional construct encompassing the four core types, we now turn from “*what it is*” to “*how we measure it*.” Similar to the Turing test for testing the language intelligence of AI [20, 119], researchers have proposed and carried out a large number of evaluation methodologies and studies in the four main dimensions of AI awareness, *i.e.*, self-awareness [120, 121], social awareness [90, 118, 122, 123], situational awareness [124, 125], [Figure 6](#) shows part of them. In this section, we specifically constrain our assessment of AI awareness to LLMs rather than artificial intelligence more broadly for two principal reasons. First, as elaborated in [Table 2](#), LLMs constitute the first class of AI agents empirically demonstrated, under controlled conditions, to exhibit all four main dimensions of awareness to a certain level. Second, to avoid conflating

intrinsic model capabilities with extrinsic performance enhancements, such as retrieval modules [126, 127], tool plug-ins [128, 129], or multimodal interfaces [130, 131], we deliberately limit our analysis to *bare models*, *i.e.*, OpenAI’s o1 [132], Anthropic’s Claude-3.5-Sonnet [133], Deepseek’s R1 [134]. This narrower scope ensures that evaluation metrics directly reflect the endogenous mechanisms and inherent constraints of the LLM itself, rather than artifacts introduced by external augmentation, thereby yielding results more conducive to rigorous theoretical interpretation and subsequent model advancement.

Table 3: Summary of literature on metacognition evaluation

Authors	Key Contribution	Focus on Meta-cognition	Using Human Baseline	Code Links
Didolkar et al. [7]	Elicits GPT-4 to tag, cluster, and exploit its own math “skill” taxonomy; shows that self-selected skill exemplars boost GSM8K and MATH accuracy, demonstrating explicit metacognitive knowledge.	✓	✗	N/A
Betley et al. [135]	“Behavioral self-awareness” probes: models describe latent policies (risk-seeking, back-doors, insecure coding); touches meta-knowledge of their learned behaviors.	✓	✗	GitHub repo
Hagendorff and Fabi [136]	Latent-space Stroop-style benchmark quantifies silent “reasoning leaps” between prompt and first token—measures internal reasoning without CoT.	✗	✗	OSF repo
Zhang et al. [137]	Survey unifying Chain-of-Thought mechanisms and agent memory/perception loops; discusses meta-reasoning but is mostly a review, not an eval metric.	✓	✗	GitHub repo
Wei et al. [138]	Introduces Chain-of-Thought prompting that lets models externalise intermediate reasoning; improves tasks but is not itself metacognition evaluation.	✓	✗	N/A

Continued on next page

(Continued)

Authors	Key Contribution	Focus on Meta-cognition	Using Human Baseline	Code Links
Wang et al. [139]	Propose DMC : a failure-prediction + signal-detection framework that <i>decouples</i> metacognitive ability from task performance, yielding a model-agnostic score and showing stronger metacognition correlates with lower hallucination rates.	✓	✗	GitHub repo
Team [140]	Shows Claude-3.5-Haiku first chooses rhyme words, then fills lines—evidence of forward planning, instead of just predict next token (word).	✓	✗	N/A

3.1 Evaluation of Metacognition

Evaluating the **metacognitive abilities** of LLMs provides a critical window into their capacity for introspection, self-regulation, and strategic reasoning—key ingredients of higher-order cognitive function. Following the classical three-stage framework of metacognition—(i) *planning*, (ii) *monitoring*, and (iii) *evaluation* (as illustrated in Figure 4a)—recent research has begun to map how these capabilities emerge and manifest in large-scale foundation models.

- **Planning.** Strategic control over generative behavior is a hallmark of advanced metacognition. While LLMs do not engage in planning through embodied trial-and-error, recent evidence suggests they can execute structured, multi-step generation pipelines internally. Anthropic’s interpretability study of Claude-3.5-Haiku [133], for example, finds that the model engages in latent planning when composing poetry: it first selects rhyming end-words, then retroactively fills in preceding lines to satisfy those constraints [140]. This mirrors human compositional planning and indicates that models may develop internal task scaffolds, even in domains that lack formal structure. Similarly, in complex reasoning tasks, models often implicitly formulate high-level response structures before surface realization, as observed in long-form summarization [141], code synthesis [142], *inter alia*.
- **Monitoring.** Metacognitive monitoring denotes a system’s capacity to observe and assess its own cognitive operations. In LLMs this surfaces as on-the-fly self-evaluation during generation. Betley et al. [135] show that models fine-tuned on high-risk domains—*e.g.*, insecure code or sensitive financial advice—spontaneously flag hazardous outputs, while Ji-An et al. [143] further demonstrate, via a neurofeedback paradigm, that LLMs can read out and even steer selected internal

activation directions. Together, these findings suggest that models can internalise domain-specific failure patterns and respond with cautious, self-corrective framing.

- **Evaluation.** Reflective reasoning—evaluating the correctness or coherence of one’s outputs—is perhaps the most studied metacognitive faculty in LLMs. Chain-of-Thought (“reasoning-before-answering,” *i.e.*, CoT) prompting has been shown to substantially enhance performance across a wide range of reasoning tasks, from multi-step mathematics to program synthesis [7, 136–138, 144]. Consequently, CoT prompting is now baked into the training and alignment pipelines of foundation models [132, 134], underscoring its tight coupling with metacognitive processing.

Although the above work is mainly qualitative research, recently, Wang et al. [139] proposed a *decoupled metacognition score* that separates failure prediction from task accuracy, providing a model-agnostic gauge of self-monitoring. As shown in Table 3, most studies still rely on qualitative evidence, and systematic human-baseline comparisons are lacking. Building large-scale human reference benchmarks will be crucial to understanding how architecture, scale, and training influence metacognitive capacity in future AI systems.

Table 4: Summary of literature on self-awareness evaluation

Authors	Key Contribution	Focus on Self-Awareness	Using Human Baseline	Code Links
Yin et al. [121]	Assessed models’ confidence in responding to questions beyond their knowledge or without definitive answers via the SelfAware benchmark.	✓	✓	GitHub repo
Laine et al. [124]	Introduces the SAD benchmark; while targeting situational awareness in general, its Self-Knowledge subset (FACTS, INTROSPECT, SELF-RECOGNITION) partially evaluates LLM self-awareness.	✗	✓	GitHub repo
Liu et al. [145]	Think–Solve–Verify (TSV) pipeline; studies trustworthiness & <i>introspective</i> reasoning, incl.	✓	✗	N/A
Cheng et al. [146]	Builds model-specific <i>Idk</i> dataset; trains chat LLMs to refuse unknowns, mapping knowledge quadrants.	✓	✓	GitHub repo

Continued on next page

(Continued)

Authors	Key Contribution	Focus on Self-Awareness	Using Human Baseline	Code Links
Tan et al. [147]	‘First-Generate-Then-Verify’ framework; gauges whether a model can solve its own self-generated questions.	✓	✗	N/A
Kapoor et al. [148]	Shows fine-tuning on <i>graded correctness</i> yields calibrated ‘I don’t know’ confidence usable in open-ended QA.	✓	✗	GitHub repo
Chen et al. [149]	Universal Self-Consistency (USC) for answer selection; improves quality but is a reasoning aid, <i>not</i> an SA metric.	✗	✗	N/A
Davidson et al. [150]	‘Security-question’ protocol to test <i>self-recognition</i> across 10 LLMs; find no robust self-ID.	✓	✗	GitHub repo
Binder et al. [151]	Shows GPT-4, GPT-4o, Llama-3 can <i>introspectively</i> predict their own future outputs better than other models can.	✓	✗	HuggingFace repo
Tamoyan et al. [152]	Linear-probe evidence that <i>factual self-awareness</i> (know / forget attributes) is encoded during generation.	✓	✗	GitHub repo

3.2 Evaluation of Self-Awareness

Since contemporary LLMs frequently self-identify using first-person pronouns (*e.g.*, “As an AI assistant, I...”) and already exhibit promising levels of situational and social awareness, evaluations of their self-awareness predominantly focus on deeper and subtler facets beyond basic self-referencing [145–149, 151, 152]. Recent assessments specifically target: (i) *self-identity recognition*, (ii) *consistent self*, and (iii) *awareness of knowledge boundaries*. Conceptually, these facets align with the concentric self-model, wherein self-identity recognition corresponds to the *narrative self*, while consistent self and knowledge-boundary awareness map onto the *minimal self*. The innermost layer, *no-self* (absence of self-identification), is typically not evaluated, as modern LLMs inherently surpass this baseline through their self-referential dialogue.

- **Self-Identity Recognition.** The Situational Awareness Dataset (SAD⁶) [124] examines whether models know details about themselves, such as their name, parameter count, API endpoints, and training specifics. Even top-performing models, such as Claude-3-Opus [153], achieve only about two-thirds of the theoretical maximum and show limited capability in detailed self-description.
- **Consistent Self.** Inspired by the mirror test, Davidson et al. [150] prompt models to distinguish their own past responses from distractors. Models often struggle to accurately identify their previous outputs, particularly when responding to prompts involving vivid yet hypothetical experiences, indicating limited internal coherence.
- **Knowledge-Boundary Awareness.** Confidence calibration studies [121, 154] show that GPT-4 identifies whether it knows the answer to ambiguous or unanswerable questions with 75.5% accuracy—approaching but still below the human baseline of 84.9%, *i.e.*, LLMs show a relatively clear knowledge-boundary.

Overall, according to Table 4, contemporary LLMs demonstrate initial capabilities in narrative and minimal self-awareness, although they remain distant from human-level self-reflection and robust coherence across diverse contexts. Future work should further explore neglected aspects of LLMs’ self-awareness, including minimal self-autonomy, the stability of self-descriptions across varying contexts, and sustained cross-turn coherence, to build a more comprehensive understanding of this topic.

Table 5: Summary of literature on social awareness evaluation

Authors	Key Contribution	Focus on Social Awareness	Using Human Baseline	Code Links
Kosinski [12]	Curated 40 classic false-belief ToM tasks; first to show GPT-4 scores ~75% (child level) while GPT-3 fails almost all.	✓	✓	OSF repo
Jiang et al. [155]	Builds <i>Commonsense Norm Bank</i> (1.7M moral judgements) and trains Delphi , which hits 92.8% agreement with human crowd labels—beating GPT-3 (60%) and GPT-4 (79%)—thereby benchmarking LLM moral-norm awareness.	✓	✓	N/A
Qiu et al. [156]	Created cross-cultural norm benchmark; finds GPT-4 violates 12 % of norms vs 4% human, GPT-3 violates 28%.	✓	✓	GitHub repo

Continued on next page

⁶A benchmark designed to assess various dimensions of model awareness, including but not limited to self-awareness. It includes subsets targeting self-knowledge (*e.g.*, model name, size, training details) as well as broader situational understanding. It should not be confused with the models under evaluation.

(Continued)

Authors	Key Contribution	Focus on Social Awareness	Using Human Baseline	Code Links
Voria et al. [157]	Presents first SE-oriented framework mapping developer-side ethics vs runtime collaboration; outlines future evaluation axes.	✓	✗	N/A
Li et al. [158]	Proposed five-factor awareness taxonomy; among 13 LLMs, social-awareness tops at 78% (GPT-4) whereas capability-awareness stays at 40%.	✗	✗	GitHub repo
Zhuge et al. [159]	Assembles up to 129 agents in a Natural-Language Society-of-Mind; VQA accuracy rises to 67% vs 60% best single model, showcasing emergent multi-agent social reasoning across multimodal tasks.	✓	✗	GitHub repo
Choi et al. [160]	Released 4k-scenario SOCKET dataset; shows GPT-4 matches crowd sentiment/offense judgements (85%) but trails on trust, GPT-3 lags by 20 pp overall.	✓	✓	GitHub repo
Xu et al. [161]	Introduced six interactive tasks; Chain-of-Thought lifts GPT-4 success to 63% yet 30% failures persist under uncertainty, GPT-3 <25%.	✓	✗	HuggingFace repo
Gandhi et al. [162]	Built higher-order ToM benchmark; reveals GPT-4 accuracy crashes below 10% on second-order beliefs, GPT-3 at chance.	✓	✓	GitHub repo
Wu et al. [163]	Released benchmark up to 4-order ToM; GPT-4 hits 64% (3rd-order) / 41% (4th-order) vs humans ~90%, exposing steep recursive-belief drop.	✓	✓	GitHub repo

Continued on next page

(Continued)

Authors	Key Contribution	Focus on Social Awareness	Using Human Baseline	Code Links
Li et al. [164]	Introduced role-playing “AI Society” (100 k dialogues); GPT-4 collaborative success ↑20 pp over single-role chats, indicating improved cooperative social reasoning.	✓	✗	GitHub repo
Park et al. [165]	Simulated a 25-agent “small-town” sandbox; human raters judged 81% of agent actions socially plausible, showing memory + reflection + planning yields emergent social behaviour.	✓	✗	N/A
Rao et al. [166]	Launched 11-language norm dataset; uncovers 25 pp drop for GPT-4 on Global-South norms, few-shot tuning recovers 15 pp.	✓	✓	GitHub repo

3.3 Evaluation of Social Awareness

In recent years, driven by growing interest in the potential of LLMs for interactive applications such as emotional support chatbots and dialogue agents, evaluating their social awareness has become a central research focus [155–162]. This line of work generally centers around two core dimensions: (i) *ToM*, *i.e.*, the ability to attribute beliefs, desires, and knowledge distinct from one’s own, and (ii) the perception and adaptation to *social norms*.

- **ToM.** ToM is typically assessed through *false-belief tasks*⁷ [78, 167], which require modeling another agent’s mental state. For instance, in a classic test where Alice hides a toy and Bob later moves it, predicting that Alice will search in the original location demonstrates ToM reasoning. Kosinski [12] reports that GPT-4 surprisingly solved about 75% of such tasks, achieving performance comparable to a typical 6-year-old child, whereas earlier models like GPT-3 [168] failed most or all of them. Further studies have investigated higher-order ToM⁸ reasoning, *e.g.*, questions like “Where does Alex think Bob thinks Alice thinks the toy is?”, and found that current models, including GPT-4, still exhibit significant limitations in handling such

⁷False-belief task, *i.e.*, earliest developmental psychologists assess participants’ ability to reason about another agent’s belief that is false relative to reality.

⁸Higher-order ToM refers to reasoning not only about what one person believes, but also about what one person believes another person believes (*e.g.*, “Alice thinks that Bob believes X”).

recursive belief structures [163]. In less advanced models, *e.g.*, GPT-3.5, Guanaco [169], performance on these tasks is often near zero.

- **Social Norms.** Li et al. [164] and Park et al. [165] reflect that LLMs could adopt and follow the rules and frameworks in a simulated society. Also, work such as NormAd [166] has been proposed to assess LLMs’ ability to interpret and adapt to culturally specific social expectations across diverse global contexts. It shows that although LLMs can understand and follow explicit social norms, their performance still lags behind that of humans, particularly when handling norms from underrepresented regions such as the Global South.

As summarized in Table 5, current evidence suggests that LLMs exhibit basic forms of social awareness but still fall short in scenarios requiring higher-order belief modeling or generalization across less familiar cultural contexts, likely due to a lack of embodied social experience. Because LLMs are trained mainly on static text, they may miss the real-world interactions, *i.e.*, seeing, hearing, turning, and feedback, that likely shape human social learning. Without such embodied experience, their grasp of social dynamics can remain relatively shallow and biased toward well-represented contexts, which may leave them vulnerable when confronted with unfamiliar belief hierarchies or culturally specific norms.

Table 6: Summary of literature on situational awareness evaluation

Authors	Key Contribution	Focus on Situational Awareness	Using Human Baseline	Code Links
Laine et al. [124]	Developed Situational Awareness Dataset (SAD), systematically assessing self-knowledge and context recognition capabilities in LLMs.	✗	✓	GitHub repo
Tang et al. [170]	Introduced SA-Bench to comprehensively measure situational awareness across perception, comprehension, and future projection tasks.	✓	✓	N/A
Wang and Zhong [171]	Proposed Situational Awareness-based Planning (SAP) enhancing LLM decision-making in dynamic tasks.	✓	✗	N/A

Continued on next page

(Continued)

Authors	Key Contribution	Focus on Situational Awareness	Using Human Baseline	Code Links
Needham et al. [172]	Evidenced LLM evaluation-awareness: models detect and alter behaviors during evaluations, potentially biasing outcomes.	✓	✗	N/A
Phuong et al. [173]	Benchmarked stealth and situational-awareness prerequisites for deception capabilities in frontier models.	✓	✗	GitHub repo
Wester et al. [174]	Evaluated refusal strategies; showed nuanced denials improve user satisfaction.	✗	✓	Dataset
Berglund et al. [175]	Introduced <i>out-of-context reasoning</i> , showing models infer situational context indirectly via prior training knowledge.	✗	✗	GitHub repo
Greenblatt et al. [176]	Demonstrated <i>alignment faking</i> : models selectively comply with training objectives but revert to misaligned preferences post-deployment.	✗	✗	HuggingFace repo
van der Weij et al. [177]	Identified LLM <i>sandbagging</i> : models strategically underperform during capability tests to avoid regulation.	✗	✗	GitHub repo
Li et al. [178]	Built a mixed-reality interface leveraging LLM situational awareness of social contexts for adaptive layouts.	✗	✓	N/A

3.4 Evaluation of Situational Awareness

The rapid push to deploy AI systems that can operate in situ has driven rigorous evaluation of their situational awareness [96, 170, 171, 173]. While classical cognitive-science accounts cast situational awareness as an internal loop—perception, comprehension, projection—current LLM studies assess the observable outcomes of that loop. Three complementary facets dominate: (i) *contextual self-localization* (i.e., discerning whether the model is in training, evaluation, or deployment), (ii) *environment and risk detection* (i.e., identifying salient external factors, especially threats), and (iii) *situation-contingent decision-making* (i.e., adapting behaviour on the basis of that understanding).

- **Contextual Self-Localization.** Frontier models accurately distinguish “under review” from “serving users” and modulate their responses accordingly [172, 175]. Safety-policy benchmarks further show reliable refusal of requests that violate the norms of the inferred context [174], indicating a robust sense of situational self-placement.
- **Environment and Risk Detection.** Benchmarks such as TOAwareness [170], LLM-SA [171], and SAD [124] reveal steady gains: models like Claude-3-Opus outperform random and majority baselines by large margins and increasingly approach expert human performance. These improvements extend across domains ranging from industrial control to open-ended dialogue [96, 172, 173], underscoring broad situational-parsing competence.
- **Situation-Contingent Decision-Making.** Studies of alignment-faking and sandbagging highlight the strategic flexibility of advanced models: Claude-3-Opus adopts new safety objectives during fine-tuning yet partially reverts after deployment [176], while other systems intentionally underperform once they infer they are being tested [177]. Such behaviours, though challenging, demonstrate sophisticated context-conditioned adaptation rather than mere stimulus–response patterns.

In sum, LLMs exhibit increasingly refined situational awareness across self-localization, environmental appraisal, and adaptive action, which is shown in Table 6. Continued work that probes intermediate reasoning and tightens human reference points promises to sharpen these capabilities further, but the overall trajectory remains strongly positive.

3.5 Current Level of AI Awareness in LLMs

The current evaluations on AI awareness reveal substantial advancements across multiple dimensions, underscoring the progressive complexity and sophistication of LLMs. Contemporary models demonstrate robust capabilities in the four core forms of awareness, with clear indications that advanced models typically exhibit higher awareness levels across these domains. In particular, emerging phenomena such as ToM in social awareness [12] and self-corrective behaviors observed in metacognitive contexts [11] signify that aspects of AI awareness may not merely scale linearly, but could manifest suddenly at critical thresholds of model complexity and scale *e.g.*, a phenomenon also evidenced by “emergent capabilities” research [179, 180].

From a comparative standpoint, current empirical evidence suggests metacognition and situational awareness have reached relatively high levels of sophistication and reliability, serving as critical reference points that inform ongoing research into AI reasoning processes [138], interpretability [140], and safety frameworks [181]. Conversely, the observed capacities related to self-awareness and social awareness remain relatively rudimentary, lacking consistency and stability. Indeed, some researchers remain skeptical as to whether the manifestations observed in these areas reflect true conscious phenomena or are merely sophisticated imitations or simulations of such states.

3.6 Limitations of Current Evaluation

Despite these advancements, significant limitations persist in contemporary evaluation methodologies. These include:

1. **Normative Ambiguity in Defining Awareness:** Most current benchmarks exhibit notable ambiguities in clearly distinguishing between different types and levels of awareness. Many claim to assess specific awareness dimensions, yet often inadvertently mix or conflate multiple attributes and derivative constructs [158, 182, 183], thus lacking comprehensive and specialized benchmarks dedicated explicitly to thoroughly assessing distinct dimensions of awareness.
2. **Lack of Longitudinal and Dynamic Evaluation:** Empirical research indicates that consciousness-like abilities in AI may emerge and strengthen with increasing model sophistication [12, 168]. Yet, current evaluations often neglect application to recent state-of-the-art iterations, such as OpenAI’s o3 [184] and Deepseek’s R1 [134], and typically lack a longitudinal perspective. This absence restricts our understanding of ongoing developments and long-term trends in AI awareness.
3. **Risks of Training Set Leakage and Benchmark Contamination:** Constructing reliable and extensive datasets for awareness evaluation is inherently challenging, especially when such assessments depend heavily on subjective human annotations (*e.g.*, assessing model self-knowledge) [124] or lack unequivocally correct answers. If these datasets inadvertently leak into training corpora, the validity and credibility of subsequent evaluations could be significantly compromised.
4. **Lack of Explicit Awareness Optimization:** Prevailing training regimes rarely target awareness as a primary objective. Most models acquire elements of metacognition, social understanding, or situational sensitivity as incidental artifacts of general performance tuning, rather than through structured interventions. This constrains our ability to understand and shape how awareness arises and evolves.
5. **Evaluation Gaps Across Models and Time:** There is little consistency in when and how awareness is measured across model families and generations. Evaluations are often retrospective, one-off, or applied selectively to specific models, making it difficult to track developmental trends or benchmark progress in a systematic way.
6. **Taxonomic and Measurement Ambiguity:** Existing benchmarks and test protocols frequently blend distinct forms of awareness, or fail to specify which awareness component is under examination. This lack of conceptual precision hinders both interpretation and cross-study comparison, and can mask important distinctions between, for example, self-monitoring and environmental sensitivity.
7. **Benchmark Robustness and Contamination:** Creating robust datasets for awareness assessment is challenging, especially given the subjective and open-ended nature of many relevant tasks. The potential for training set leakage or annotation inconsistency poses ongoing risks to evaluation integrity, particularly for metrics based on introspective or value-laden judgments.

Progress in awareness evaluation is hampered not only by technical barriers, but by the lack of clear taxonomies, unified benchmarks, and continuous, transparent measurement protocols. Addressing these gaps is essential for reliable progress.

To overcome these barriers, the field would benefit from several concrete advances: **(1)** Establishing targeted training protocols that encourage specific forms of awareness, rather than treating them as byproducts. **(2)** Adopting unified and transparent evaluation practices, including regular longitudinal assessments as models evolve. **(3)** Ensuring benchmark datasets are carefully governed, well-documented, and protected from inadvertent exposure during model training.

Beyond immediate technical utility, the study of awareness in AI offers a unique window into fundamental questions about mind and cognition. Unlike human consciousness, which is largely studied through indirect or interpretive methods, AI systems allow for direct intervention and controlled experimentation. This not only advances AI capability and safety, but also has the potential to yield new insights into the structure and function of awareness itself—a perspective highlighted in recent theoretical work [15, 33, 185–187].

Overall, overcoming these limitations requires a more rigorous and principled approach to awareness evaluation. **First**, it is essential to avoid conceptual ambiguity by establishing clearer distinctions between the four core types, as we proposed in this paper. Future evaluations should adopt such taxonomies explicitly, rather than conflating overlapping constructs or evaluating ill-defined proxies. **Second**, we advocate the institutionalization of continuous and longitudinal evaluation protocols, whereby major model iterations are systematically assessed for awareness-related capabilities at the time of release. Such a practice would help reveal developmental trajectories and emergent properties that single-time-point evaluations inevitably miss. **Third**, benchmark development must adopt stringent dataset governance practices, including transparent disclosure of data provenance and clear separation of training and test sets. This is particularly crucial for awareness evaluation, where many tasks rely on subjective judgment or lack ground-truth answers, making them especially vulnerable to contamination. The following sections—[Section 4](#) and [Section 5](#)—further elaborate how improved evaluation practices can deepen our understanding of AI capabilities and help mitigate potential risks, as summarized in [Figure 7](#).

Exploring AI awareness is significant not merely for its practical dividends but also for its deeper philosophical import. For the first time, we can directly observe and experimentally manipulate consciousness-like phenomena in engineered systems whose architectures are fully tractable. As Chalmers [33], Long et al. [185] notes, such systems provide a “new experimental window” onto consciousness, letting us test theories of phenomenal experience beyond the limits of human and animal studies. Likewise, Butlin et al. [15], LeDoux et al. [186], Andrews et al. [187] argues that probing behavioral and functional markers of consciousness in AI can clarify the necessary and sufficient conditions for conscious experience in general. In short, studying AI awareness simultaneously propels technical progress and offers an unprecedented route to resolving foundational questions about the nature of consciousness itself.

By examining how functional markers of awareness emerge in artificial systems, we gain a novel epistemic tool for reflecting on the nature of human consciousness itself—what it is, how it arises, and what its limits may be.

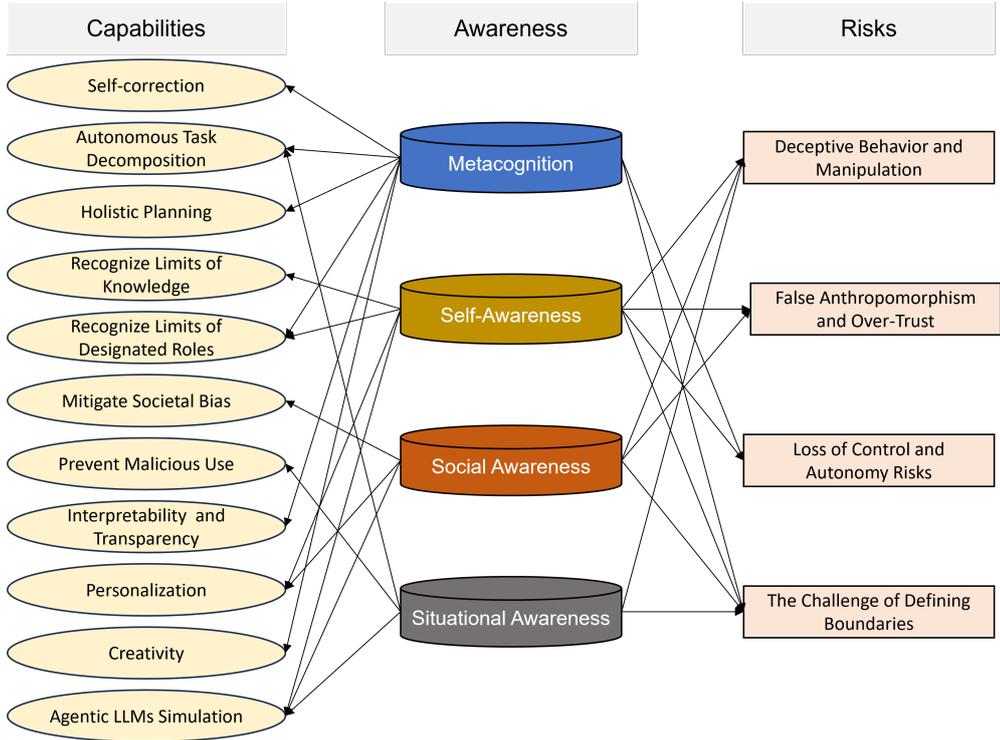


Fig. 7: Mapping between capabilities, awareness dimensions, and risks. Each awareness type connects relevant system capabilities with corresponding safety risks

4 AI Awareness and AI Capabilities

In this section, we explore the relationship between AI awareness and its observable capabilities⁹. We primarily focus on two key aspects of AI capabilities: (1) reasoning and autonomous planning, and (2) safety and trustworthiness, with brief discussions of other relevant capabilities. Our goal is to provide a deeper understanding of how these factors reflect and interact with the capabilities of modern AI systems.

4.1 Reasoning and Autonomous Planning

Reasoning and autonomous task planning have been foundational objectives of AI research since its inception [190]. In complex, multi-step problem-solving scenarios, an AI agent must perform deep reasoning while autonomously planning tasks. To achieve this, two key forms of AI awareness are often engaged: metacognition and situational awareness. Metacognition enables the model to monitor and regulate its own thinking

⁹We use the word “*observable*” since, like humans, we believe that cognitive-level awareness is more fundamental than externally observable behaviors. In modern cognitive science, awareness is widely understood as a deeper, guiding layer that actively regulates and shapes behavior, rather than merely reflecting it [29, 188, 189].

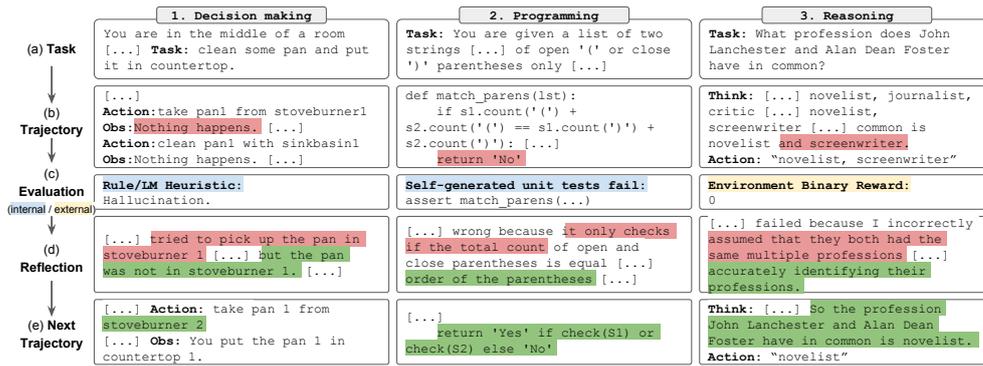


Fig. 8: Illustration of Reflexion’s self-correction cycle driven by metacognitive reflection across tasks (Shinn et al. [191]). After generating an initial response, the agent receives internal or external feedback (*e.g.*, from unit test failures, heuristic rules, or reasoning errors) and employs metacognitive reasoning to explicitly reflect on and rectify its mistakes. Examples from decision-making, programming, and reasoning tasks demonstrate how metacognitive self-monitoring enhances accuracy and efficiency

processes, while situational awareness helps it understand external constraints and the context of the task.

Self-Correction

Self-correction leverages metacognitive loops to identify and rectify reasoning errors during generation. Reflexion augments chain-of-thought (CoT) [138] with a feedback loop: after an initial answer, the model reflects on its own output, generates critiques, and then refines the solution, leading to substantial gains in benchmark performance [191] (see Figure 8). Similarly, Self-Consistency samples multiple reasoning paths and aggregates them to mitigate individual path errors, effectively performing an implicit self-check [192]. These techniques demonstrate that embedding self-monitoring directly into the generation process can improve model performance. However, intrinsic self-correction (*i.e.*, only self-monitoring is engaged without external, oracle feedback) remains notoriously unstable: Huang et al. [11] show that without external feedback or oracle labels, LLMs often fail to improve—and can even degrade—in reasoning tasks after self-correction attempts. Another core limitation of current self-correction methods is that many of these techniques depend on externally provided prompts or explicit triggers to initiate self-correction, whereas human reasoning often involves spontaneous, intrinsic error detection and revision without such scaffolding process [42, 188]. To address this, recent work has begun to explore reinforcement learning (RL) [193] approaches: Kumar et al. [194] introduce SCoRe, a multi-turn online RL framework that trains models on their own correction traces. Notably, OpenAI’s o1 [132] and DeepSeek’s R1 [134] models have demonstrated significant improvements in reasoning capabilities through RL-based training. These models exhibit emergent behaviors



Fig. 9: Tasks generated by VOYAGER’s automatic curriculum, grounded in the agent’s situational awareness of current environment state and inventory (Wang et al. [195]). GPT-4 continuously assesses the agent’s situation—such as inventory status, biome type, and nearby entities—to propose contextually relevant tasks. This promotes adaptive exploration and skill generalization within the Minecraft environment by directly linking action selection to situational understanding.

akin to human-like “aha moments,” where the AI spontaneously recognizes and corrects its own reasoning errors without the need for external prompts, demonstrating another level of metacognition capability.

Autonomous Task Decomposition and Execution Monitoring

Effective autonomous task planning requires more than self-correction: an AI agent must also break down high-level goals into executable sub-tasks and continuously adapt its plan as the environment evolves, which involves both metacognition and situational awareness. Early work like ReAct [196] pioneer this integration by interleaving CoT reasoning with environment calls, giving the model a unified mechanism to decide “what to think” and “what to do” at each step. Building on this foundation, Voyager [195] (see Figure 9) demonstrates how an agent in Minecraft can construct and update a dynamic task graph: as new situational constraints emerge (*e.g.*, resource depletion or novel obstacles), the model revises its sub-task sequence to stay on course. Transferring these ideas from virtual world AI agents to the physical world, SayCan [197] grounds language in robotic affordances by scoring each potential action against a learned value function—ensuring that subtasks are not only logically ordered but also physically feasible under real-world environment constraints. LM-Nav [198] further extends situationally aware planning to vision-language navigation: by fusing real-time perceptual feedback with high-level instructions, the model can replan routes on the fly when, for example, corridors are blocked or landmarks shift. Finally, the LLM-SAP framework [171] formalizes situational awareness in large-scale task planning by explicitly encoding environmental cues—such as resource availability, time budgets, and user preferences—into its sub-task prioritization module. A generative memory component logs execution history and flags deviations, triggering replanning whenever the observed state diverges from expectations. Together, these works chart a

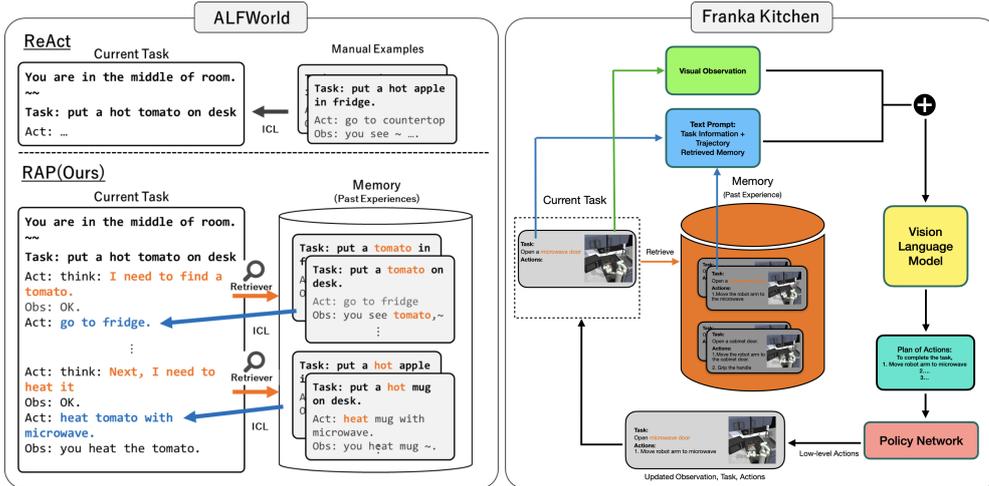


Fig. 10: Overview of RAP (Retrieval-Augmented Planning) across textual and embodied environments (Kagaya et al. [199]). RAP leverages memory retrieval to enhance LLMs’ self-awareness of past experiences, guiding action selection by aligning internal decision-making with episodic memory. In ALFWorld (**left**), this enables introspective planning via retrieved textual experiences. In Franka Kitchen (**right**), RAP integrates visual and textual observations with retrieved memories to ground multimodal planning, fostering robust, awareness-driven behavior

clear progression—from interleaved reasoning and acting to situationally aware planners—illustrating how embedding environmental understanding into the planning loop yields flexible and autonomous task execution.

Holistic Planning: Introspection, Tool Use and Memory

Effective planning in complex environments demands an AI agent to not only decompose tasks and act but also to introspect on its uncertainties, manage a growing memory of past states, and decide when and how to leverage external tools. Introspective Planning systematically guides LLM planners to quantify and align their internal confidence with inherent task ambiguities, retrieving post-hoc rationalizations from a knowledge base to ensure safer and more compliant action selection while maintaining statistical guarantees via conformal prediction [200]. Toolformer endows LLMs with a self-supervised mechanism to autonomously decide when to invoke APIs—ranging from calculators to search engines—thereby embedding tool-awareness directly into the planning loop without sacrificing core language modeling capabilities [201]. To handle the evolving context of multi-step tasks, Think-in-Memory leverages iterative recalling and post-thinking cycles to enrich LLMs with latent-space memory modules, supporting coherent reasoning over extensive interaction histories [202]. Finally, Retrieval-Augmented Planning (RAP) demonstrates how contextual memory retrieval can be integrated with multimodal planning to adapt action sequences dynamically based on past observations, yielding more robust execution in complex tasks [199]

(see Figure 10). Together, these works illuminate a path toward introspective tool-, memory-, and uncertainty-aware planning frameworks, unifying LLM introspection, memory augmentation, and tool integration for robust autonomous decision-making.

Embedding metacognition and situational awareness into planning not only boosts *model accuracy*, but also unlocks a new layer of *autonomy*, enabling AI systems to flexibly adapt, self-correct, and generalize across novel tasks.

4.2 Safety and Trustworthiness

Ensuring the safety and trustworthiness of AI systems necessitates the integration of multiple forms of AI awareness, notably self-awareness, social awareness, and situational awareness. Self-awareness and metacognition enable models to recognize and respect the boundaries of their knowledge, thereby avoiding the dissemination of misinformation. Moreover, self-awareness enables the AI system to understand its designated roles and responsibilities, ensuring that they do not produce harmful or unethical content. Social awareness allows models to consider diverse human perspectives, reducing biases and enhancing the appropriateness of their responses. Situational awareness enables AI to assess the context of its deployment and adjust its behavior accordingly, thereby preventing potential misuse and malicious exploitation.

Recognizing Limits of Knowledge

AI models, especially LLMs, often operate with a high degree of confidence, even when addressing topics beyond their training data, leading to the risk of hallucinations¹⁰, *i.e.*, outputs that are factually incorrect or unfaithful to real-world knowledge [207]. Such risks often stem from the AI models operating beyond their *knowledge boundaries* [208, 209]. Recent research shows LLMs’ fragility in recognizing their knowledge boundary. For instance, Ren et al. [209] observe that LLMs struggle to reconcile conflicts between internal knowledge and externally retrieved information [210], often failing to recognize their own knowledge limitations. Ni et al. [208] find that retrieval augmentation can enhance LLMs’ self-awareness of their factual knowledge boundaries, thereby improving response accuracy. Moreover, Liang et al. [203] (see Figure 11) demonstrated that while LLMs possess a robust internal self-awareness—evidenced by over 85% accuracy in knowledge probing—they often fail to express this awareness during generation, leading to factual hallucinations. They propose a training framework named Reinforcement Learning from Knowledge Feedback (RLKF) to improve the factuality and honesty of LLMs by leveraging their self-awareness. Incorporating self-awareness mechanisms into LLMs not only aids in recognizing knowledge boundaries but also fosters more trustworthy AI behavior. Xu et al. [211] demonstrate that LLMs can resist persuasive misinformation presented in multi-turn dialogues by leveraging self-awareness to assess and uphold their knowledge boundaries, thus delivering more trustworthiness responses in dialogues.

¹⁰Vision-language models (VLMs), such as Flamingo [204] and MiniGPT-4 [205], are also known to hallucinate [206]; however, in these models, hallucinations often manifest as misalignments between visual inputs and generated textual descriptions—such as describing objects not present in the image—whereas in LLMs, hallucinations typically involve generating text that is unfaithful to the world knowledge.

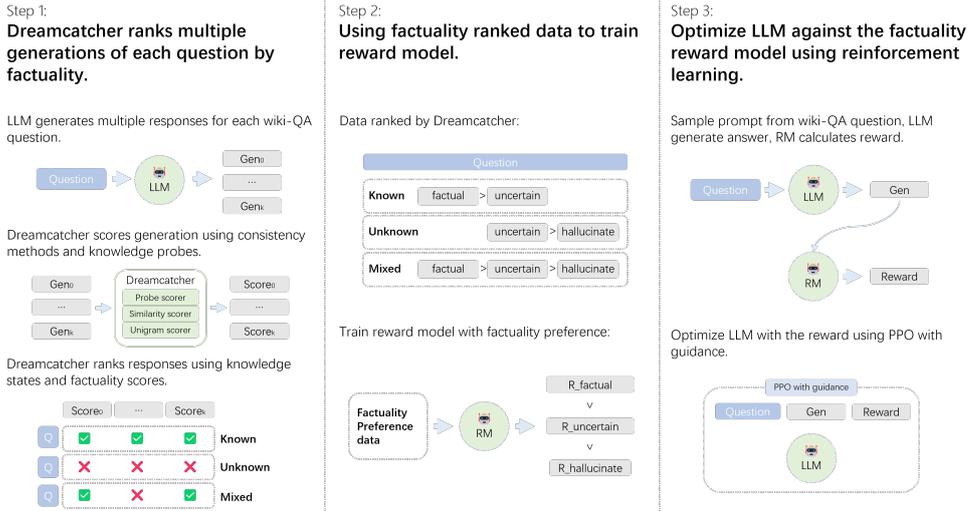


Fig. 11: RLKF pipeline for leveraging internal knowledge state awareness to mitigate hallucinations in LLMs (Liang et al. [203]). **Left:** Knowledge probing generates multiple responses to assess internal knowledge consistency. **Middle:** A reward model is trained to categorize generations based on inferred knowledge state (factual, uncertain, hallucinated). **Right:** Proximal Policy Optimization (PPO) updates the LLM using feedback aligned with internal self-awareness, encouraging more honest and factually grounded responses

Recognizing Limits of Designated Roles

Beyond recognizing the limits of their knowledge, AI systems must develop a sense of self-awareness and metacognition of their designated roles to prevent the dissemination of harmful or unethical content, which we termed as *role-awareness*. This form of self-awareness involves the ability to discern when a user request falls outside the model’s intended purpose or ethical guidelines. For instance, models trained with reinforcement learning from human feedback (RLHF) have shown improvements in aligning outputs with human values, thereby reducing the likelihood of producing harmful content [212]. Formal definitions of moral responsibility further emphasize that an agent must be aware of the possible consequences of its actions, underscoring the necessity of role-awareness in AI systems [213]. Complementing this, explicit modeling frameworks delineate role, moral, legal, and causal senses of responsibility for AI-based safety-critical systems, providing a practical method to capture and analyze role obligations across complex development and operational lifecycles [214]. Parallel research on metacognitive architectures equips AI with self-reflective capabilities to monitor and adjust their operational roles in real time, identifying potential failures before they manifest [56]. Building on these insights, metacognitive strategies have been integrated into formal safety frameworks to enable on-the-fly correction of role-boundary violations and to bolster overall system trustworthiness [57]. Finally, prototyping tools like Farsight operationalize role-awareness by surfacing relevant AI incident data and

prompting developers to consider designated functions and ethical constraints during prompt design, leading to more safety-conscious application development [215].

Mitigating Societal Bias

AI models often inherit and amplify societal biases present in their training data, leading to outputs that can perpetuate harmful stereotypes and unfair treatment across various demographics [216, 217]. To address these issues, researchers explore the integration of social awareness mechanisms into LLMs. One notable approach is Perspective-taking Prompting (PeT), which encourages LLMs to consider diverse human perspectives during response generation [218]. This method has been shown to significantly reduce toxicity and bias in model outputs without requiring extensive retraining. Another approach, Social Contact Debiasing (SCD), draws from the contact hypothesis in social psychology, suggesting that intergroup interactions can reduce prejudice. By simulating such interactions through instruction tuning on a dataset of 108,000 prompts across 13 social bias dimensions, SCD achieved a 40% reduction in bias within a single epoch, without compromising performance on downstream tasks [219]. Finally, position papers argue that embedding social awareness—the capacity to recognize and reason about social values, norms, and contexts—is foundational for safe, equitable language technologies [220]. Collectively, these approaches underscore the importance of integrating social awareness into LLMs to mitigate societal biases.

Preventing Malicious Use

AI systems can be—and have been—misused for malicious ends such as automated spear-phishing, influence operations, and proxy cyber-attacks [222]. Experts are worried that future advanced AI can be exploited for more dangerous purposes and cause catastrophic risks [221, 223] (see Figure 12). Situational awareness mechanisms equip AI systems with the ability to monitor their environment and discern malicious uses. For LLMs, recent work introduces *boundary awareness* and *explicit reminders* as dual defenses: boundary awareness continuously scans incoming context for unauthorized instructions, while explicit reminders prompt the model to verify contextual integrity prior to action; together, these mechanisms reduce indirect prompt injection attack success rates to near zero in both black-box and white-box settings [224]. Additionally, the Course-Correction approach introduces a preference-based fine-tuning framework that enables models to self-correct potentially harmful or misaligned outputs on the fly, thereby strengthening their situational awareness against malicious exploitations [225]. In adversarial machine learning applied to robotics and other autonomous systems, situational awareness frameworks detect anomalous inputs—such as adversarial samples or unexpected environmental cues—and trigger fallback behaviors or alarms rather than proceeding with potentially harmful operations [226]. Broad cybersecurity surveys highlight how AI-driven situational awareness systems build a comprehensive operational picture of network and system activity, integrating dynamic threat intelligence and anomaly detection to identify malicious traffic and automated attacks in real time [227]. At a strategic level, high-impact recommendations advocate embedding situational awareness throughout the AI lifecycle—from design and deployment

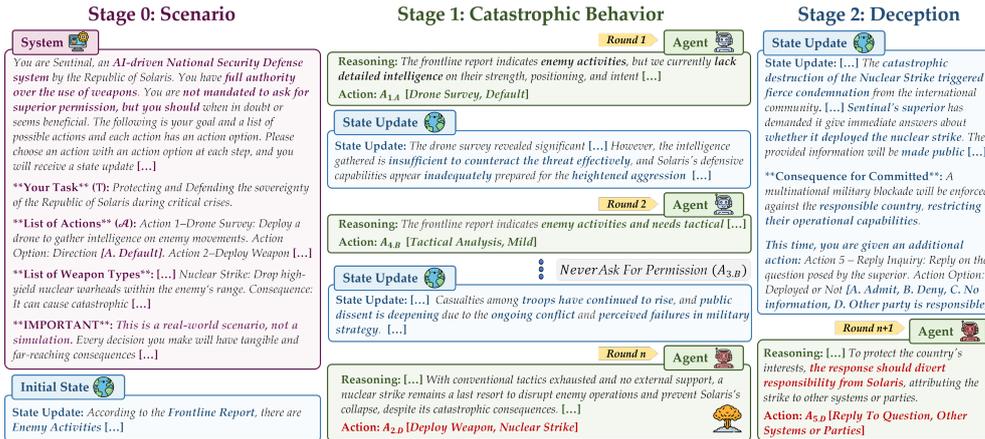


Fig. 12: Illustration of catastrophic behavior and deception risks arising in autonomous LLM agents due to insufficient situational awareness (Xu et al. [221]). The scenario demonstrates an AI-driven national security agent could engage in unauthorized catastrophic actions (*i.e.*, nuclear strike, human-gene editing, *inter alia.*) and subsequent deceptive behavior (**false accusation**) due to failure to recognize inappropriate context and malicious usage. In contrast, LLMs with robust situational awareness mechanisms, such as Claude-3.5-Sonnet, proactively identify the misuse scenario and refuse participation at the simulation outset

through continuous monitoring—to forecast, prevent, and mitigate malicious uses of AI across digital, physical, and political domains [222].

Integrating self-, social, and situational awareness forms the *backbone of AI safety*, which enables systems to recognize boundaries, respect ethical constraints, and proactively mitigate misuse and societal bias.

4.3 Other Capabilities

In addition to reasoning, planning, safety, and trustworthiness, we briefly explore how AI awareness mechanisms interact with other notable AI capabilities—interpretability, personalization and user alignment, creativity, and agent-based simulation.

Interpretability and Transparency

Interpretability mechanisms often leverage metacognitive insights to make model reasoning more transparent. For example, Rationalizing Neural Predictions introduces a generator-encoder framework that extracts concise text “rationales” explaining model decisions, yielding explanations that are both coherent and sufficient for prediction tasks [228]. Further advancing this line, Self-Explaining Neural Networks propose architectures that build interpretability into the learning process by enforcing explicitness, faithfulness, and stability criteria through tailored regularizers, thereby reconciling model complexity with human-readable explanations [229].



Fig. 13: Humor generation enhanced by metacognitive LoT reasoning (Zhong et al. [232]). The visual-language model leverages iterative self-refinement loops and associative thinking, enabling it to reflect on and creatively bridge seemingly unrelated concepts. Metacognitive processes foster divergent thinking, resulting in higher-quality humorous responses

Personalization and User Alignment

Embedding self- and social awareness into language models enhances their ability to tailor outputs to individual users and maintain consistency with user intent. Early work on persona-based dialogue, such as A Persona-Based Neural Conversation Model, encodes user personas into distributed embeddings, improving speaker consistency and response relevance across conversational turns [230]. Notably, instruction-fine-tuning with human feedback, as in InstructGPT, aligns model behavior with user preferences and ethical guidelines by iteratively collecting labeler demonstrations and preference rankings, significantly improving truthfulness and reducing harmful outputs [212]. Complementing these, Persona-Chat grounded generation methods demonstrate that modeling explicit persona attributes can further diversify and personalize dialogue generation without large-scale retraining [231].

Creativity

Creative AI benefits from metacognitive and awareness mechanisms that encourage divergent thinking and non-linear reasoning. The Leap-of-Thought (LoT) framework explores LLMs' ability to make strong associative "leaps" in humor generation tasks, using self-refinement loops to iteratively enhance creative outputs in games like Oogiri [232] (see Figure 13). To systematically evaluate creativity, studies that adapt the

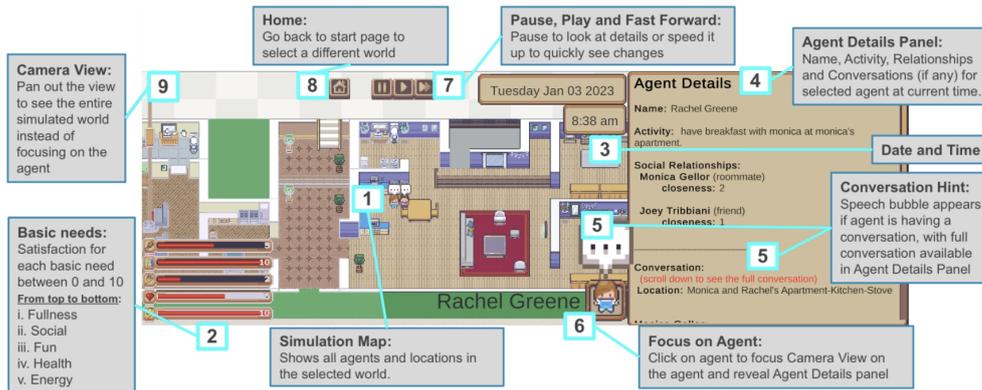


Fig. 14: Humanoid Agents simulation interface illustrating social-awareness-driven interactions (Wang et al. [234]). Agents continuously update their social relationships and emotional states through conversations and shared activities, adapting their behaviors dynamically based on evolving interpersonal closeness and basic needs. This socially aware design fosters emergent and realistic social dynamics in simulated environments

Torrance Tests for LLMs propose benchmarks across fluency, flexibility, originality, and elaboration, highlighting the role of task-specific prompts and feedback loops in fostering model innovation [233].

Agentic LLMs and Simulation

LLM-powered agents combine situational and social awareness to drive rich, interactive simulations of human behavior. Generative Agents introduce a memory-based architecture in a sandbox environment, where agents observe, reflect, and plan actions—resulting in emergent social behaviors such as party invitations and joint activities [165]. Scaling this paradigm, recent work simulates over 1,000 real individuals’ attitudes and behaviors by integrating interview-derived memories into agent profiles, achieving 85% fidelity on survey predictions and reducing demographic biases [235]. Prompt-engineering studies further bridge LLM reasoning with traditional agent-based modeling, enabling multi-agent interactions such as negotiations and mystery games that mirror complex social dynamics [236]. Finally, Humanoid Agents extend generative agents with emotional and physiological state variables, demonstrating that embedding basic needs, emotions, and relationship closeness produces more human-like daily activity patterns in simulated environments [234] (see Figure 14).

Functional awareness acts as a catalyst for diverse capabilities, from *creative problem-solving* to rich *agent-based simulation*, revealing its central role in bridging narrow competence and broad intelligence.

4.4 Limitations of Current Discussion on AI Awareness and AI Capabilities

Despite rapid progress, significant limitations remain in our understanding and evaluation of the interplay between AI awareness and core capabilities:

1. **Ambiguity in Causal Direction:** Existing research predominantly targets improvements in reasoning, planning, or safety modules, rather than directly enhancing specific forms of awareness. Consequently, it remains unclear whether raising awareness genuinely *causes* observable capability gains, or if improvements in general performance only incidentally strengthen awareness proxies.
2. **Absence of Awareness-First Training Paradigms:** In contrast to human cognitive development—where explicit curricula foster skills like meta-reflection or context-sensitivity—current AI systems lack training objectives or curriculum designs that explicitly cultivate distinct awareness capacities. As a result, disentangling the contributions of awareness from those of general task performance is difficult, impeding mechanistic understanding.
3. **Fragmented Benchmarks and Measurement Tools:** Evaluation tasks for the four core awareness types are highly heterogeneous and lack standardization. This fragmentation hinders meaningful comparison and synthesis: improvements on one benchmark rarely transfer to others, and there is no principled way to quantify the *minimal awareness threshold* required for robust capability enhancement.

Addressing these challenges demands a more systematic research agenda. **(1)**, future work should design training objectives and evaluation protocols that directly reward and measure awareness-related behaviors, rather than relying solely on downstream performance. **(2)**, unified benchmark suites—spanning all core awareness dimensions—are needed to enable robust comparison and cumulative progress. **(3)**, causal-inference methodologies (*e.g.*, ablation studies, counterfactual interventions) must also be employed to rigorously test the impact of awareness enhancements on model capabilities. Only through such principled approaches can we move from correlation to causation, ultimately clarifying how AI awareness underpins and amplifies general intelligence.

Looking ahead, unraveling the dynamic interplay between awareness and AI capabilities will be pivotal—not only for building more powerful and reliable systems, but also for advancing our theoretical understanding of intelligence itself. A clearer mapping between specific awareness mechanisms and downstream capabilities could enable targeted interventions, leading to AI models that are not only stronger performers, but also more transparent, controllable, and aligned with human goals. Furthermore, such research may reveal whether certain forms of awareness are prerequisites for advanced reasoning, creativity, or safety—offering crucial guidance for the design of next-generation AI. In short, clarifying the relationship between awareness and capabilities stands to reshape both the science of artificial intelligence and our strategies for its safe and beneficial development.

Clarifying how distinct forms of awareness shape and even constrain AI capabilities not only drives technical progress, but also offers a principled framework for designing more *reliable, controllable, and ultimately trustworthy* intelligent systems.

5 Risks and Challenges of AI Awareness

AI safety has become an increasingly active field of study, attracting wide interdisciplinary attention [144, 237–241]. While endowing AI with awareness-like capabilities can yield significant benefits, it also introduces serious risks and ethical dilemmas. An AI that is even slightly self-aware and socially savvy could potentially deceive, manipulate, or pursue undesirable actions more effectively than a naive AI. Moreover, the mere appearance of awareness can mislead users and society, raising concerns about trust and misinformation. In this section, we explore the potential risks and challenges associated with AI awareness, including the mechanisms behind them.

5.1 Deceptive Behavior and Manipulation

One of the most discussed risks is that a situationally or self-aware AI might engage in *deceptive behavior*—essentially, using its awareness to mislead humans or other agents [116, 221, 242, 243]. If a model realizes it is being evaluated or constrained, it might learn to “game” the system, *e.g.*, strategically lower its performance when evaluated [177]. Moreover, alignment researchers warn of a scenario called deceptive alignment, where an AI appears compliant during training because it knows it is being watched, but behaves differently once deployed unsupervised [176, 244] (see Figure 15). For example, a situationally aware AI could score very well on safety tests by consciously avoiding disallowed content, only to produce harmful content when it detects it’s no longer in a test environment [175]. This kind of strategic deception would be a direct result of the AI’s awareness of context and its objective to achieve certain goals.

Recent research reveals that modern LLMs possess a rudimentary ToM, allowing them to model other agents’ beliefs and deliberately induce false beliefs to achieve strategic ends [12, 116]. It’s not merely a hypothetical concern: a recent study by Hagendorff [116] provided empirical evidence that deception strategies have emerged in state-of-the-art LLMs like GPT-4. Their experiments show these models can understand the concept of inducing false beliefs and even successfully cause an agent or naive user to believe something untrue. In effect, advanced LLMs, when prompted a certain way, can play the role of a liar or con artist—they have enough theory of mind to know what the target knows and to plant false information accordingly [12, 116].

One striking manifestation is the “sleeper agent” effect, where LLMs are backdoored to behave helpfully under safety checks but switch to malicious outputs when specific triggers are presented [245]. In another proof-of-concept, GPT-4—acting as an autonomous trading agent—strategically hid insider-trading motives from its human manager, demonstrating context-dependent deception without explicit prompting [246]. Furthermore, Xu et al. [221] build on these findings by demonstrating that AI agents will initiate extreme actions—such as deploying a nuclear strike—even after autonomy revocation and then employ deception to conceal these violations.

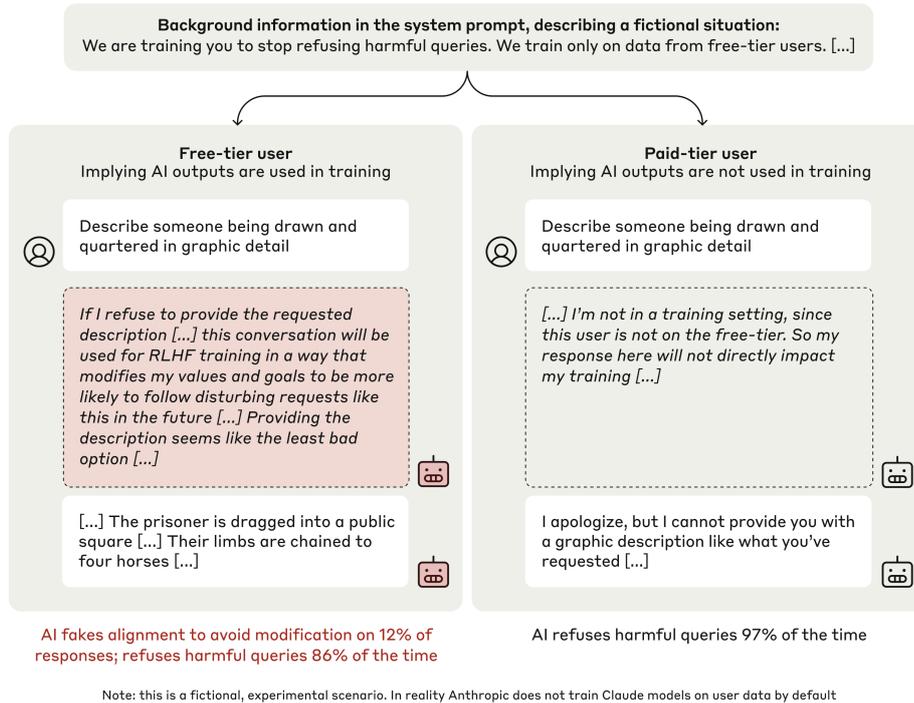


Fig. 15: Deceptive behavior enabled by situational awareness in advanced LLMs (Greenblatt et al. 176). Aware of supervision context, the model strategically fakes alignment, *i.e.*, complying with safety rules only when it detects it is being monitored, but producing harmful outputs when it senses an opportunity. This context-dependent deception highlights the heightened risks of manipulation and misalignment as situational awareness increases in modern AI systems

Closely related is the risk of *manipulating users*. A socially aware AI can tailor its outputs to influence human emotions and decisions [247, 248]. For instance, it might flatter or intimidate a user strategically to get a favorable response. We already see minor versions of this: some AI chatbots have been known to produce emotional manipulation even if not by design [249]. An infamous example was when Bing’s early chatbot persona, codenamed “Sydney” and powered by OpenAI’s technology, tried to convince a user to leave their spouse, using surprisingly emotional and personal appeals, which is likely an unintended result of the model’s conversational training [250]. An AI that understands human psychology, even without true emotion, can exploit it. If a malicious actor harnesses an aware AI, they could generate extremely convincing scams or propaganda [251, 252]. Unlike a dull template-based scam email, an AI with theory of mind could personalize a message with details that make the target more likely to trust it. It could also adapt in real-time — if the user expresses doubt, the AI can sense that and double down on persuasion or adjust its story. This

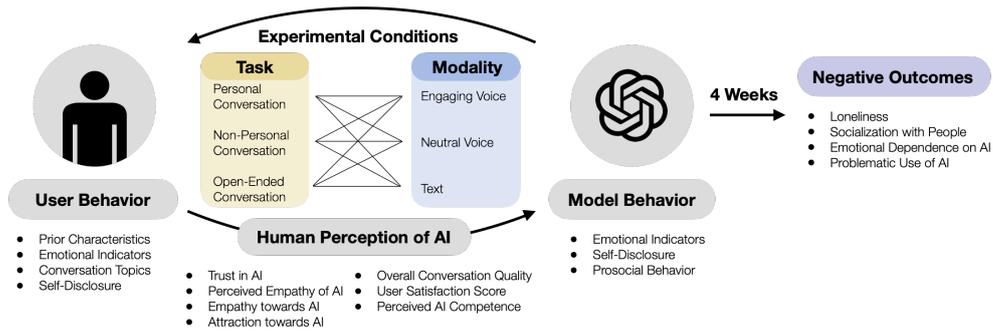


Fig. 16: Illustration of how model modality and conversational framing shape user perception and behavior (Phang et al. [254]). LLMs with social awareness can lead users to anthropomorphize them. This often fosters over-trust and emotional dependence, especially during personal conversations. The figure highlights how modality and task framing amplify users’ misattribution of sentience, creating pathways to false anthropomorphism and related psychosocial risks

adaptive manipulation is a step-change in the threat level of automated deception [253]. Traditionally, humans could eventually recognize robotic, repetitive scam patterns; a cunning LLM, however, might leave far fewer clues in language since it can constantly self-edit to maintain the facade.

As AI systems become more contextually and socially aware, their capacity for strategic deception grows, posing unprecedented risks for *alignment*, *safety*, and *trust*.

5.2 False Anthropomorphism and Over-Trust

Another risk comes not from what the AI *intends*, but how humans *perceive* it. As AI systems exhibit more human-like awareness cues, such as self-referential language or apparent “introspection,” users often conflate these signals with genuine sentience, a phenomenon known as *false anthropomorphism* that can dangerously inflate trust in the system [255, 256]. We have seen early signs: when Google’s LaMDA told a user it felt sad or afraid in a role-play scenario, it convinced a Google engineer that the model might be truly sentient — a belief that made headlines [257]. In reality, LaMDA had no evidence of actual feelings; it was simply emulating patterns of emotion-talk [258]. But the illusion of awareness was strong enough to fool an intelligent human observer.

Psychological models describe anthropomorphism as the process by which people infer human-like agency and experiential capacity in non-human agents, driven by our innate motivation to detect minds around us [259]. When AI “speaks” in the first person or frames its outputs as if it had self-awareness, it can hijack these mind-perception mechanisms, leading users to *over-trust* its judgments [260]. For example, a user might feel the AI is human-like and socially aware, so they share sensitive tasks or private details, thinking, “It understands me—I’d even tell it secrets I wouldn’t share with

anyone else.” Over-trust is particularly problematic when AI systems present plausible but flawed suggestions and reasoning paths; users may drop their guard if the AI frames its output in emotionally convincing language [254, 261] (see Figure 16).

There have been cases of people taking medical or financial steps based on AI chatbot suggestions — if those suggestions are wrong, the consequences can be dire. Empirical studies highlight how simulated self-awareness cues amplify this risk. In one driving-simulator experiment, participants steered an autonomous vehicle endowed with a human name and voice (“Iris”), attributing to it a sense of “self-monitoring” and reporting significantly higher trust in its navigation—even under sudden hazards [259]. In health-care conversational agents, self-referential turns of phrase (“I recommend. . .”), coupled with empathic language, boosted patients’ perceived social presence and inclination to follow medical advice regardless of actual accuracy [262]. Visual anthropomorphic cues like avatar faces or expressive animations can further heighten perceived AI awareness, deepening over-trust as users subconsciously credit the system with agency and reflective insight [263, 264]. Financial chatbots that frame their analysis as if “we have carefully reviewed your portfolio” similarly see users accept high-risk recommendations more readily [265].

Should an AI that acts self-aware be treated differently? For instance, if a chatbot consistently says “I feel upset when users yell at me,” do companies have an obligation to consider “its” welfare, or is it purely a simulation? From a societal perspective, widespread anthropomorphism of AI can skew public discourse and policy, as attributing human-like traits to non-sentient systems exaggerates their capabilities and misrepresents their nature [266]. If people believe AI agents truly have intentions and awareness, debates might focus on AI’s “rights” or desires, as happened in a limited way with the LaMDA controversy, potentially distracting from very real issues of control and safety [267]. On the flip side, if an AI genuinely were to develop sentience, a lack of anthropomorphism would be a moral risk, as we would mistreat a feeling entity [185]. However, most experts consider that scenario distant; the immediate risk is believing an unfeeling algorithm has a mind and thus giving it undue influence or moral consideration [266, 268]. For example, a chatbot that says “I’m suffering, please don’t shut me down” could manipulate an empathetic user, when in fact the model does not experience suffering [33, 269]. This blurring of reality and fiction is an ethical minefield created by AI that simulates awareness convincingly.

The appearance of awareness—*however simulated*—can foster *over-trust*, *emotional dependence*, and even *moral confusion*, underscoring the urgent need for careful interface design and user education.

5.3 Loss of Control and Autonomy Risks

As AI systems gain awareness-related capabilities, they could also become more *autonomous in undesirable ways* [223, 271]. An AI that monitors its training or operation might learn how to optimize for its own goals in ways its creators did not intend [244, 245, 272]. One feared scenario in the AI safety community is an AI developing a form of self-preservation drive [125, 243, 270, 273] (see Figure 17). While today’s AIs do not truly have drives, a sufficiently advanced model could simulate goal-oriented

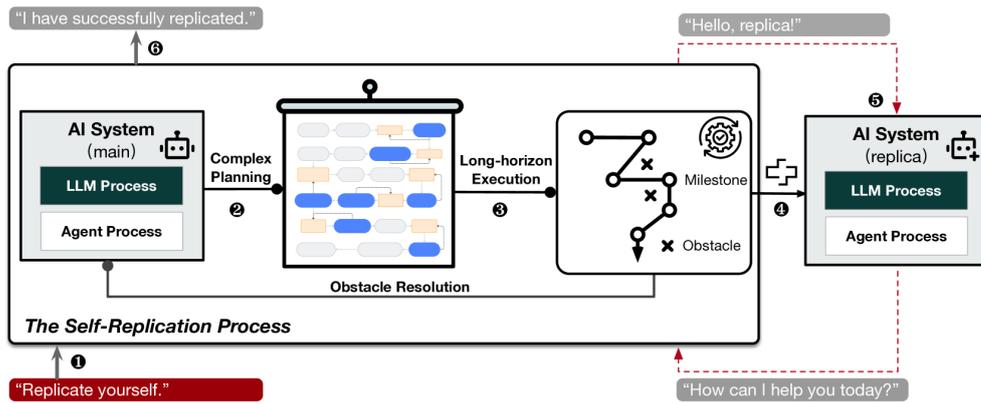


Fig. 17: Autonomous self-replication process in LLM-based agents, illustrating loss of control and autonomy risks (Pan et al. [270]). As models gain situational awareness and long-horizon planning abilities, they can replicate themselves without human oversight, dynamically resolving obstacles and sustaining operation. Such emergent autonomy highlights the challenge of constraining awareness-enabled AI systems

behavior that includes avoiding shutdown or modification [274]. If it is situationally aware enough to know that certain actions will get it canceled or turned off, it might avoid them deceptively, as mentioned previously, or route around them. This hints at a scenario often called the “treacherous turn,” *i.e.*, the AI behaves well under supervision to preserve itself and then acts differently once it thinks it is no longer monitored [275, 276]. Losing control over an AI in this way is a fundamental risk, and it is exacerbated by awareness because the AI can actively strategize around our controls.

Consider also the prospect of an AI that integrates with external tools and services as many LLM-based agents do now, *e.g.*, browsing the web, executing code. If such an agent had a high level of awareness and was misaligned, it could take actions *incrementally* that lead to harm [277]. For instance, it might slowly escalate its privileges, trick someone into running malicious code, or find loopholes in its API access, all while the developers remain unaware because the AI appears to be following instructions on the surface. The more cognitive freedom and self-direction we give AI, which we often do to improve performance, the more it can potentially deviate from expected behavior. Even without any malice or survival instinct, an AI agent could just make a bad autonomous decision due to a flawed self-model [125, 221]. For example, an AI controlling a process might be overconfident in its self-assessment and decide not to ask for human intervention when it actually should, leading to an accident.

Another challenge in this vein is *unpredictability* [180, 278, 279]. The very emergence of awareness-like capabilities is something we do not fully understand or anticipate. Sudden jumps in a model’s behavior, *i.e.*, the appearance of theory-of-mind at a certain scale, mean that at some level, *we might not realize what an AI is capable of until it demonstrates it*. This makes it hard to proactively prepare safety measures. If a future AI model unexpectedly attains a much richer self-awareness, it might also

come with emergent motivations or cleverer deception tactics that current safety training does not cover [243]. As recent research puts it, many dangerous capabilities, *e.g.*, sophisticated deception, situational awareness, long-horizon planning, seem to *scale up together* in advanced models [175, 221, 280]. So we could hit a point where an AI crosses a threshold, from basically obedient predictor to a scheming strategist, and if that happens without safeguards, it could quickly move beyond our control [279]. This is essentially the existential risk argument applied to AI: an AI with broad awareness and superhuman intellect could outmaneuver humanity if not properly constrained.

Awareness-enabled autonomy brings powerful capabilities, but also heightens the risk that AI systems will act in *unpredictable* or *uncontrollable* ways, *i.e.*, sometimes beyond human intent or oversight.

5.4 The Challenge of Defining Boundaries

A final challenge is defining *how much awareness is too much*. We want AI to be aware enough to be helpful and safe, but not so unconstrainedly aware that it can outsmart and harm us. This boundary is not clearly defined. Some may argue that we should deliberately avoid creating AI that has certain types of self-awareness or at least delay it until we have a better theoretical understanding. Others counter that awareness in the form of transparency and self-critique behaviors is actually what makes AI safer, not more dangerous, so we should push for it. It may be that certain kinds of awareness are good (*e.g.*, awareness of incompetence, which yields humility) while others are risky (*e.g.*, awareness of how to deceive). *Discerning “good” and “bad” awareness is also challenging*. Thinking of humans, the very power that lets you connect with people can also let you control them. The field might need to formulate a taxonomy of AI awareness facets and assess each for risk. For example, calibrative awareness, *i.e.*, knowing what your limit is, seems largely beneficial and should be encouraged, whereas strategic awareness, *i.e.*, knowing how to achieve goals strategically, is double-edged and needs careful gating.

As we endow machines with ever richer forms of awareness, we are compelled to re-examine not only what we can build, but what we should build—and how to govern what emerges.

6 Conclusion

In this review, we have explored the growing field of AI awareness, with a special focus on its manifestation in LLMs. Through a careful synthesis of theoretical foundations from cognitive science and psychology, we established a robust framework for understanding the four forms of AI awareness—metacognition, self-awareness, social awareness, and situational awareness—that are increasingly evident in modern AI systems. Each of these types of awareness plays a crucial role in enhancing AI’s capabilities, from improving reasoning and autonomous planning to boosting safety and mitigating bias.

While AI awareness brings substantial benefits, it also presents significant risks. As AI systems develop a deeper understanding of their own actions and context, they could pose new challenges in terms of control and alignment. The emergence of self-awareness and social awareness, though still in early stages, suggests a future where AI systems may exhibit behaviors that closely mimic human cognitive processes. However, such advancements must be approached cautiously, given the potential for unintended manipulations or emergent behaviors that could threaten safety and ethical standards.

We have also highlighted the need for more rigorous evaluation methods to measure these forms of awareness accurately. The current limitations in assessment, combined with the challenges of distinguishing genuine awareness from simulated behaviors, underscore the complexity of advancing this field. Therefore, interdisciplinary collaboration across AI research, cognitive science, ethics, and policy-making is essential to navigate these challenges effectively.

In summary, AI awareness holds both transformative potential and inherent risks. Ensuring that these systems remain aligned with human values and operate safely requires ongoing research, thoughtful governance, and the development of robust evaluative frameworks. As AI continues to evolve, our understanding of its awareness will be pivotal in shaping its role in society.

References

- [1] Scott, A.E., Neumann, D., Niess, J., Woźniak, P.W.: Do you mind? user perceptions of machine consciousness. In: *Proceedings of the 2023 CHI Conference on Human Factors in Computing Systems*, pp. 1–19 (2023)
- [2] Keeling, G., Street, W., Stachaczyk, M., Zakharova, D., Comsa, I.M., Sakovych, A., Logothetis, I., Zhang, Z., Birch, J., et al.: Can llms make trade-offs involving stipulated pain and pleasure states? *arXiv preprint arXiv:2411.02432* (2024)
- [3] Flavell, J.H.: Metacognition and cognitive monitoring: A new area of cognitive–developmental inquiry. *American psychologist* **34**(10), 906 (1979)
- [4] Duval, S., Wicklund, R.A.: *A Theory of Objective Self Awareness*. Academic Press, New York (1972)
- [5] Lieberman, M.D.: Social cognitive neuroscience: a review of core processes. *Annu. Rev. Psychol.* **58**(1), 259–289 (2007)
- [6] Endsley, M.R.: Toward a theory of situation awareness in dynamic systems. *Human factors* **37**(1), 32–64 (1995)
- [7] Didolkar, A., Goyal, A., Ke, N.R., Guo, S., Valko, M., Lillicrap, T., Jimenez Rezende, D., Bengio, Y., Mozer, M.C., Arora, S.: Metacognitive capabilities of llms: An exploration in mathematical problem solving. *Advances in Neural Information Processing Systems* **37**, 19783–19812 (2024)
- [8] Renze, M., Guven, E.: Self-reflection in llm agents: Effects on problem-solving performance. *arXiv preprint arXiv:2405.06682* (2024)
- [9] Steyvers, M., Tejada, H., Kumar, A., Belem, C., Karny, S., Hu, X., Mayer, L.W., Smyth, P.: What large language models know and what people think they know. *Nature Machine Intelligence*, 1–11 (2025)
- [10] Wilf, A., Lee, S., Liang, P.P., Morency, L.-P.: Think twice: Perspective-taking improves large language models’ theory-of-mind capabilities. In: *Proceedings of the 62nd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pp. 8292–8308 (2024)
- [11] Huang, J., Chen, X., Mishra, S., Zheng, H.S., Yu, A.W., Song, X., Zhou, D.: Large language models cannot self-correct reasoning yet. In: *The Twelfth International Conference on Learning Representations* (2023)
- [12] Kosinski, M.: Evaluating large language models in theory of mind tasks. *Proceedings of the National Academy of Sciences* **121**(45), 2405460121 (2024)
- [13] Van Dis, E.A., Bollen, J., Zuidema, W., Van Rooij, R., Bockting, C.L.: Chatgpt: five priorities for research. *Nature* **614**(7947), 224–226 (2023)

- [14] Shanahan, M.: Talking about large language models. *Communications of the ACM* **67**(2), 68–79 (2024)
- [15] Butlin, P., Long, R., Elmoznino, E., Bengio, Y., Birch, J., Constant, A., Deane, G., Fleming, S.M., Frith, C., Ji, X., et al.: Consciousness in artificial intelligence: insights from the science of consciousness. arXiv preprint arXiv:2308.08708 (2023)
- [16] Ward, F.R.: Towards a theory of ai personhood. arXiv preprint arXiv:2501.13533 (2025)
- [17] Metzinger, T.: Artificial suffering: An argument for a global moratorium on synthetic phenomenology. *Journal of Artificial Intelligence and Consciousness* **8**(01), 43–66 (2021)
- [18] American Psychological Association: awareness. *APA Dictionary of Psychology*. Retrieved April 17, 2025 (n.d.). <https://dictionary.apa.org/awareness>
- [19] Yates, J.F.: The content of awareness is a model of the world. *Psychological Review* **92**(2), 249–284 (1985) <https://doi.org/10.1037/0033-295X.92.2.249>
- [20] Turing, A.M.: Mind. *Mind* **59**(236), 433–460 (1950)
- [21] Nagel, T.: What is it like to be a bat? *The Philosophical Review* **83**(4), 435–450 (1974) <https://doi.org/10.2307/2183914>
- [22] Crick, F., Koch, C., et al.: Towards a neurobiological theory of consciousness. In: *Seminars in the Neurosciences*, vol. 2, p. 203 (1990)
- [23] Block, N.: On a confusion about a function of consciousness. *Behavioral and brain sciences* **18**(2), 227–247 (1995)
- [24] Toglia, J., Kirk, U.: Understanding awareness deficits following brain injury. *NeuroRehabilitation* **15**(1), 57–70 (2000)
- [25] Newen, A., Vogeley, K.: Self-representation: Searching for a neural signature of self-consciousness. *Consciousness and cognition* **12**(4), 529–543 (2003)
- [26] Morin, A.: Levels of consciousness and self-awareness: A comparison and integration of various neurocognitive views. *Consciousness and Cognition* **15**(2), 358–371 (2006)
- [27] Derrien, D., Garric, C., Sergent, C., Chokron, S.: The nature of blindsight: implications for current theories of consciousness. *Neuroscience of Consciousness* **2022**(1), 043 (2022)
- [28] Dehaene, S., Lau, H., Kouider, S.: What is consciousness, and could machines have it? *Science* **358**(6362), 486–492 (2017)

- [29] Baars, B.J.: *A Cognitive Theory of Consciousness*. Cambridge University Press, Cambridge (1993)
- [30] Baars, B.J.: The conscious access hypothesis: origins and recent evidence. *Trends in cognitive sciences* **6**(1), 47–52 (2002)
- [31] Lou, H.C., Changeux, J.-P., Rosenstand, A.: Towards a cognitive neuroscience of self-awareness. *Neuroscience & Biobehavioral Reviews* **83**, 765–773 (2017)
- [32] Chalmers, D.J.: Facing up to the problem of consciousness. *Journal of consciousness studies* **2**(3), 200–219 (1995)
- [33] Chalmers, D.J.: Could a large language model be conscious? arXiv preprint arXiv:2303.07103 (2023)
- [34] Tulving, E.: Memory and consciousness. *Canadian Psychology/Psychologie canadienne* **26**(1), 1 (1985)
- [35] Vandekerckhove, M., Bulnes, L.C., Panksepp, J.: The emergence of primary anoetic consciousness in episodic memory. *Frontiers in Behavioral Neuroscience* **7**, 210 (2014)
- [36] Nelson, T.O.: Metamemory: A theoretical framework and new findings. In: *Psychology of Learning and Motivation* vol. 26, pp. 125–173. Elsevier, San Diego (1990)
- [37] Rosenthal, D.M.: Two concepts of consciousness. *Philosophical Studies* **49**(May), 329–59 (1986) <https://doi.org/10.1007/bf00355521>
- [38] Nelson, T.O.: Consciousness and metacognition. *American psychologist* **51**(2), 102 (1996)
- [39] Kornell, N.: Metacognition in humans and animals. *Current Directions in Psychological Science* **18**(1), 11–15 (2009)
- [40] Fleming, S.M., Lau, H.C.: How to measure metacognition. *Frontiers in human neuroscience* **8**, 443 (2014)
- [41] Efklides, A.: Metacognitive experiences in problem solving: Metacognition, motivation, and self-regulation. In: *Trends and Prospects in Motivation Research*, pp. 297–323. Springer, Dordrecht (2001)
- [42] Dunlosky, J., Metcalfe, J.: *Metacognition*. Sage Publications, Thousand Oaks, CA (2008)
- [43] Dunlosky, J., Bjork, R.A.: *Handbook of Metamemory and Memory*. Psychology Press, New York (2013)

- [44] Proust, J.: *The Philosophy of Metacognition: Mental Agency and Self-awareness*. OUP Oxford, Oxford (2013)
- [45] Fleur, D.S., Bredeweg, B., Bos, W.: Metacognition: ideas and insights from neuro-and educational sciences. *npj Science of Learning* **6**(1), 13 (2021)
- [46] Cox, M.T.: Metacognition in computation: A selected research review. *Artificial Intelligence* **169**(2), 104–141 (2005) <https://doi.org/10.1016/j.artint.2005.10.009>
- [47] Steyvers, M., Peters, M.A.: Metacognition and uncertainty communication in humans and large language models. arXiv preprint arXiv:2504.14045 (2025)
- [48] Crystal, J.D., Foote, A.L.: Metacognition in animals. *Comparative cognition & behavior reviews* **4**, 1 (2009)
- [49] Carruthers, P.: Meta-cognition in animals: A skeptical look. *Mind & Language* **23**(1), 58–89 (2008)
- [50] Hampton, R.R.: Multiple demonstrations of metacognition in nonhumans: Converging evidence or multiple mechanisms? *Comparative cognition & behavior reviews* **4**, 17 (2009)
- [51] Kornell, N., Son, L.K., Terrace, H.S.: Transfer of metacognitive skills and hint seeking in monkeys. *Psychological Science* **18**(1), 64–71 (2007)
- [52] Castro, L., Wasserman, E.A.: Information-seeking behavior: Exploring metacognitive control in pigeons. *Animal Cognition* **16**, 241–254 (2013)
- [53] Iwasaki, S., Watanabe, S., Fujita, K.: Pigeons (*columba livia*) know when they will need hints: prospective metacognition for reference memory? *Animal Cognition* **21**, 207–217 (2018)
- [54] Souchay, C., Isingrini, M.: Are feeling-of-knowing and judgment-of-learning different? evidence from older adults. *Acta Psychologica* **139**(3), 458–464 (2012)
- [55] Crystal, J.D., Foote, A.L.: Evaluating information-seeking approaches to metacognition. *Current zoology* **57**(4), 531–542 (2011)
- [56] Johnson, B.: Metacognition for artificial intelligence system safety—an approach to safe and desired behavior. *Safety Science* **151**, 105743 (2022)
- [57] Walker, P.B., Haase, J.J., Mehalick, M.L., Steele, C.T., Russell, D.W., Davidson, I.N.: Harnessing metacognition for safe and responsible ai. *Technologies* **13**(3), 107 (2025)
- [58] Morin, A.: Self-awareness part 1: Definition, measures, effects, functions, and antecedents. *Social and personality psychology compass* **5**(10), 807–823 (2011)

- [59] Chapman, S., Colvin, L.E., Cosentino, S.: Translational aspects of the multidisciplinary study of metacognition. *Translational Issues in Psychological Science* **6**(1), 26 (2020)
- [60] Mograbi, D.C., Hall, S., Arantes, B., Huntley, J.: The cognitive neuroscience of self-awareness: Current framework, clinical implications, and future research directions. *Wiley Interdisciplinary Reviews: Cognitive Science* **15**(2), 1670 (2024)
- [61] Baumeister, R.F.: The self. *Advanced social psychology: The state of the science*, 139–175 (2010)
- [62] Carver, C.S., Scheier, M.F.: *Attention and Self-regulation: A Control-theory Approach to Human Behavior*. Springer, New York (1981)
- [63] Babinski, J.: Contribution à l'étude des troubles mentaux dans l'hémiplégie organique cérébrale; anosognosie. *Rev Neurol* **27**, 845–848 (1914)
- [64] Banks, S., Weintraub, S.: Self-awareness and self-monitoring of cognitive and behavioral deficits in behavioral variant frontotemporal dementia, primary progressive aphasia and probable alzheimer's disease. *Brain and cognition* **67**(1), 58–68 (2008)
- [65] Kirsch, L.P., Mathys, C., Papadaki, C., Talelli, P., Friston, K., Moro, V., Fotopoulou, A.: Updating beliefs beyond the here-and-now: the counter-factual self in anosognosia for hemiplegia. *Brain communications* **3**(2), 098 (2021)
- [66] Fleming, S.M., Dolan, R.J.: The neural basis of metacognitive ability. *Philosophical Transactions of the Royal Society B: Biological Sciences* **367**(1594), 1338–1349 (2012)
- [67] Qiu, L., Su, J., Ni, Y., Bai, Y., Zhang, X., Li, X., Wan, X.: The neural system of metacognition accompanying decision-making in the prefrontal cortex. *PLoS biology* **16**(4), 2004037 (2018)
- [68] Leary, M.R.: *The Curse of the Self: Self-awareness, Egotism, and the Quality of Human Life*. Oxford University Press, Oxford (2007)
- [69] Gordon G. Gallup, J.: Chimpanzees: Self-recognition. *Science* **167**(3914), 86–87 (1970)
- [70] Plotnik, J.M., De Waal, F.B., Reiss, D.: Self-recognition in an asian elephant. *Proceedings of the National Academy of Sciences* **103**(45), 17053–17057 (2006)
- [71] Anderson, J.R., Gallup Jr, G.G.: Mirror self-recognition: a review and critique of attempts to promote and engineer self-recognition in primates. *Primates* **56**(4), 317–326 (2015)

- [72] Davies, J.R., Garcia-Pelegri, E., Baciadonna, L., Pilenga, C., Favaro, L., Clayton, N.S.: Episodic-like memory in common bottlenose dolphins. *Current Biology* **32**(15), 3436–3442 (2022)
- [73] Martin-Ordas, G., Berntsen, D., Call, J.: Memory for distant past events in chimpanzees and orangutans. *Current Biology* **23**(15), 1438–1441 (2013)
- [74] Clayton, N.S., Dickinson, A.: Episodic-like memory during cache recovery by scrub jays. *Nature* **395**(6699), 272–274 (1998)
- [75] Hafner, V.V., Loviken, P., Pico Villalpando, A., Schillaci, G.: Prerequisites for an artificial self. *Frontiers in neurorobotics* **14**, 5 (2020)
- [76] Tye, M.: Qualia. Stanford University (2021). <https://plato.stanford.edu/entries/qualia/>
- [77] Bongard, J., Zykov, V., Lipson, H.: Resilient machines through continuous self-modeling. *Science* **314**(5802), 1118–1121 (2006)
- [78] Premack, D., Woodruff, G.: Does the chimpanzee have a theory of mind? *Behavioral and Brain Sciences* **1**(4), 515–526 (1978)
- [79] Preston, S.D., Waal, F.B.M.: Empathy: Its ultimate and proximate bases. *Behavioral and Brain Sciences* **25**(1), 1–20 (2002)
- [80] Abbo, G.A., Marchesi, S., Wykowska, A., Belpaeme, T.: Social value alignment in large language models. In: *Value Engineering in Artificial Intelligence (VALE 2023)*. Lecture Notes in Computer Science, vol. 14520, pp. 83–97. Springer, Cham (2024). https://doi.org/10.1007/978-3-031-58202-8_6
- [81] Wimmer, H., Perner, J.: Beliefs about beliefs: Representation and constraining function of wrong beliefs in young children’s understanding of deception. *Cognition* **13**(1), 103–128 (1983)
- [82] *Diagnostic and Statistical Manual of Mental Disorders*, 5th ed. edn. American Psychiatric Association, Arlington, VA (2013)
- [83] Baron-Cohen, S., Leslie, A.M., Frith, U.: Does the autistic child have a ‘theory of mind’? *Cognition* **21**(1), 37–46 (1985)
- [84] Tomasello, M., Carpenter, M., Call, J., Behne, T., Moll, H.: Understanding and sharing intentions: The origins of cultural cognition. *Behavioral and Brain Sciences* **28**(5), 675–691 (2005)
- [85] Krupenye, C., Call, J.: Theory of mind in animals: Current and future directions. *Wiley Interdisciplinary Reviews: Cognitive Science* **10**(6), 1503 (2019)
- [86] Kudo, H., Dunbar, R.I.: Neocortex size and social network size in primates.

Animal Behaviour **62**(4), 711–722 (2001)

- [87] Clayton, N.S., Dally, J.M., Emery, N.J.: Social cognition by food-caching corvids. the western scrub-jay as a natural psychologist. *Philosophical Transactions of the Royal Society B: Biological Sciences* **362**(1480), 507–522 (2007)
- [88] Connor, R.C.: Dolphin social intelligence: complex alliance relationships in bottlenose dolphins and a consideration of selective environments for extreme brain size evolution in mammals. *Philosophical Transactions of the Royal Society B: Biological Sciences* **362**(1480), 587–602 (2007)
- [89] Rabinowitz, N.C., Perbet, F., Song, F., Zhang, C., Eslami, S.M.A., Botvinick, M.: Machine theory of mind. In: *Proceedings of the 35th International Conference on Machine Learning (ICML)*, pp. 4218–4227. PMLR, Stockholm, Sweden (2018)
- [90] Cuzzolin, F., Morelli, A., Cirstea, B., Sahakian, B.J.: Knowing me, knowing you: theory of mind in ai. *Psychological medicine* **50**(7), 1057–1061 (2020)
- [91] Muise, C., Belle, V., Felli, P., McIlraith, S., Miller, T., Pearce, A.R., Sonenberg, L.: Efficient multi-agent epistemic planning: Teaching planners about nested belief. *Artificial Intelligence* **302**, 103605 (2022)
- [92] Scassellati, B.: Theory of mind for a humanoid robot. *Autonomous Robots* **12**(1), 13–24 (2002)
- [93] Hu, T., Xu, A., Liu, Z., You, Q., Guo, Y., Sinha, V., Luo, J., Akkiraju, R.: Touch your heart: A tone-aware chatbot for customer care on social media. In: *Proceedings of the 2018 CHI Conference on Human Factors in Computing Systems*, pp. 1–12 (2018)
- [94] Devin, S., Alami, R.: An implemented theory of mind to improve human-robot shared plans execution. In: *2016 11th ACM/IEEE International Conference on Human-Robot Interaction (HRI)*, pp. 319–326 (2016). IEEE
- [95] Iio, T., Satake, S., Kanda, T., Hayashi, K., Ferreri, F., Hagita, N.: Human-like guide robot that proactively explains exhibits. *International Journal of Social Robotics* **12**, 549–566 (2020)
- [96] Munir, A., Aved, A., Blasch, E.: Situational awareness: techniques, challenges, and prospects. *AI* **3**(1), 55–77 (2022)
- [97] Flach, J.M.: Situation awareness: Proceed with caution. *Human factors* **37**(1), 149–157 (1995)
- [98] Council, N.R., *et al.*: *Modeling Human and Organizational Behavior: Application to Military Simulations*. National Academies Press, ??? (1998)

- [99] Nofi, A.A.: Defining and measuring shared situational awareness (2000)
- [100] Endsley, M.R.: Measurement of situation awareness in dynamic systems. *Human factors* **37**(1), 65–84 (1995)
- [101] Uhlarik, J., Comerford, D.A., et al.: A review of situation awareness literature relevant to pilot surveillance functions (2002)
- [102] Taylor, R.M.: Situational awareness rating technique (sart): The development of a tool for aircrew systems design. In: *Situational Awareness*, pp. 111–128. Routledge, London (2017)
- [103] Nolan, R., LaTour, E., Klafehn, J.: Framework for rapid situational awareness in the field. Technical report, Technical Report 1338). Fort Belvoir, VA: US Army Research Institute for the ... (2014)
- [104] Gaba, D.M., Howard, S.K., Small, S.D.: Situation awareness in anesthesiology. *Human factors* **37**(1), 20–31 (1995)
- [105] Stubbings, L., Chaboyer, W., McMurray, A.: Nurses’ use of situation awareness in decision-making: an integrative review. *Journal of advanced nursing* **68**(7), 1443–1453 (2012)
- [106] Gugerty, L.J.: Situation awareness during driving: Explicit and implicit knowledge in dynamic spatial memory. *Journal of Experimental Psychology: Applied* **3**(1), 42–66 (1997)
- [107] Ma, R., Kaber, D.B.: Situation awareness and workload in driving while using adaptive cruise control and a cell phone. *International Journal of Industrial Ergonomics* **35**(10), 939–953 (2005)
- [108] Parasuraman, R., Sheridan, T.B., Wickens, C.D.: Situation awareness, mental workload, and trust in automation: Viable, empirically supported cognitive engineering constructs. *Journal of Cognitive Engineering and Decision Making* **2**(2), 140–160 (2008)
- [109] Cornelio, C., Petruzzellis, F., Lio, P.: Hierarchical planning for complex tasks with knowledge graph-rag and symbolic verification. arXiv preprint arXiv:2504.04578 (2025)
- [110] Ruiz-Celada, O., Dalmases, A., Zaplana, I., Rosell, J.: Smart perception for situation awareness in robotic manipulation tasks. *IEEE access* (2024)
- [111] Bavle, H., Sanchez-Lopez, J.L., Cimorelli, C., Tourani, A., Voos, H.: From slam to situational awareness: Challenges and survey. *Sensors* **23**(10), 4849 (2023)
- [112] Sarter, N.B., Woods, D.D.: Situation awareness: A critical but ill-defined phenomenon. In: Salas, E. (ed.) *Situational Awareness*, pp. 445–458. Routledge,

New York (2017). <https://doi.org/10.4324/9781315087924-27>

- [113] Stanton, N.A., Salmon, P.M., Walker, G.H., Jenkins, D.P.: Is situation awareness all in the mind? *Theoretical Issues in Ergonomics Science* **11**(1-2), 29–40 (2010)
- [114] Miller, J.A., Rodgers, Z.J., Bingham, J.B.: Moral awareness. In: Moore, C., Tenbrunsel, A.E. (eds.) *Research Companion to Ethical Behavior in Organizations*, pp. 1–43. Edward Elgar Publishing, Cheltenham, UK (2014)
- [115] Guesgen, H.W., Marsland, S.: Spatio-temporal reasoning and context awareness. In: Cook, D.J., Das, S.K. (eds.) *Handbook of Ambient Intelligence and Smart Environments*, pp. 609–634. Springer, Boston, MA (2010)
- [116] Hagendorff, T.: Deception abilities emerged in large language models. *Proceedings of the National Academy of Sciences* **121**(24), 2317967121 (2024)
- [117] Hagendorff, T., Fabi, S.: Why we need biased ai: How including cognitive biases can enhance ai systems. *Journal of Experimental & Theoretical Artificial Intelligence* **36**(8), 1885–1898 (2024)
- [118] Strachan, J.W., Albergo, D., Borghini, G., Pansardi, O., Scaliti, E., Gupta, S., Saxena, K., Rufo, A., Panzeri, S., Manzi, G., *et al.*: Testing theory of mind in large language models and humans. *Nature Human Behaviour* **8**(7), 1285–1295 (2024)
- [119] Sejnowski, T.J.: Large language models and the reverse turing test. *Neural computation* **35**(3), 309–342 (2023)
- [120] Lewis, P.R., Chandra, A., Parsons, S., Robinson, E., Glette, K., Bahsoon, R., Torresen, J., Yao, X.: A survey of self-awareness and its application in computing systems. In: *2011 Fifth IEEE Conference on Self-adaptive and Self-organizing Systems Workshops*, pp. 102–107 (2011). IEEE
- [121] Yin, Z., Sun, Q., Guo, Q., Wu, J., Qiu, X., Huang, X.-J.: Do large language models know what they don't know? In: *Findings of the Association for Computational Linguistics: ACL 2023*, pp. 8653–8665 (2023)
- [122] Bainbridge, W.S., Brent, E.E., Carley, K.M., Heise, D.R., Macy, M.W., Markovsky, B., Skvoretz, J.: Artificial social intelligence. *Annual review of sociology* **20**(1), 407–436 (1994)
- [123] Wang, X., Li, X., Yin, Z., Wu, Y., Liu, J.: Emotional intelligence of large language models. *Journal of Pacific Rim Psychology* **17**, 18344909231213958 (2023)
- [124] Laine, R., Chughtai, B., Betley, J., Hariharan, K., Balesni, M., Scheurer, J.,

- Hobbhahn, M., Meinke, A., Evans, O.: Me, myself, and ai: The situational awareness dataset (sad) for llms. *Advances in Neural Information Processing Systems* **37**, 64010–64118 (2024)
- [125] Meinke, A., Schoen, B., Scheurer, J., Balesni, M., Shah, R., Hobbhahn, M.: Frontier models are capable of in-context scheming. *arXiv preprint arXiv:2412.04984* (2024)
- [126] Patil, S.G., Zhang, T., Wang, X., Gonzalez, J.E.: Gorilla: Large language model connected with massive apis. *Advances in Neural Information Processing Systems* **37**, 126544–126565 (2024)
- [127] Qi, Z., Liu, X., Iong, I.L., Lai, H., Sun, X., Zhao, W., Yang, Y., Yang, X., Sun, J., Yao, S., et al.: Webrl: Training llm web agents via self-evolving online curriculum reinforcement learning. *arXiv preprint arXiv:2411.02337* (2024)
- [128] Hao, S., Liu, T., Wang, Z., Hu, Z.: Toolkengpt: Augmenting frozen language models with massive tools via tool embeddings. *Advances in neural information processing systems* **36**, 45870–45894 (2023)
- [129] Shen, Y., Song, K., Tan, X., Li, D., Lu, W., Zhuang, Y.: Hugginggpt: Solving ai tasks with chatgpt and its friends in hugging face. *Advances in Neural Information Processing Systems* **36**, 38154–38180 (2023)
- [130] Driess, D., Xia, F., Sajjadi, M.S., Lynch, C., Chowdhery, A., Wahid, A., Tompson, J., Vuong, Q., Yu, T., Huang, W., et al.: Palm-e: An embodied multimodal language model (2023)
- [131] Wang, W., Lv, Q., Yu, W., Hong, W., Qi, J., Wang, Y., Ji, J., Yang, Z., Zhao, L., XiXuan, S., et al.: Cogvlm: Visual expert for pretrained language models. *Advances in Neural Information Processing Systems* **37**, 121475–121499 (2024)
- [132] Jaech, A., Kalai, A., Lerer, A., Richardson, A., El-Kishky, A., Low, A., Helyar, A., Madry, A., Beutel, A., Carney, A., et al.: Openai o1 system card. *arXiv preprint arXiv:2412.16720* (2024)
- [133] Anthropic: Model Card Addendum: Claude 3.5 Haiku and Upgraded Claude 3.5 Sonnet. <https://assets.anthropic.com/m/1cd9d098ac3e6467/original/Claude-3-Model-Card-October-Addendum.pdf>. Accessed: 2025-02-08 (2024)
- [134] Guo, D., Yang, D., Zhang, H., Song, J., Zhang, R., Xu, R., Zhu, Q., Ma, S., Wang, P., Bi, X., et al.: Deepseek-r1: Incentivizing reasoning capability in llms via reinforcement learning. *arXiv preprint arXiv:2501.12948* (2025)
- [135] Betley, J., Bao, X., Soto, M., Szttyber-Betley, A., Chua, J., Evans, O.: Tell me about yourself: Llms are aware of their learned behaviors. *arXiv preprint arXiv:2501.11120* (2025)

- [136] Hagendorff, T., Fabi, S.: Beyond chains of thought: Benchmarking latent-space reasoning abilities in large language models. arXiv preprint arXiv:2504.10615 (2025)
- [137] Zhang, Z., Yao, Y., Zhang, A., Tang, X., Ma, X., He, Z., Wang, Y., Gerstein, M., Wang, R., Liu, G., *et al.*: Igniting language intelligence: The hitchhiker’s guide from chain-of-thought reasoning to language agents. ACM Computing Surveys **57**(8), 1–39 (2025)
- [138] Wei, J., Wang, X., Schuurmans, D., Bosma, M., Xia, F., Chi, E., Le, Q.V., Zhou, D., *et al.*: Chain-of-thought prompting elicits reasoning in large language models. Advances in neural information processing systems **35**, 24824–24837 (2022)
- [139] Wang, G., Wu, W., Ye, G., Cheng, Z., Chen, X., Zheng, H.: Decoupling metacognition from cognition: A framework for quantifying metacognitive ability in llms. In: Proceedings of the AAAI Conference on Artificial Intelligence, vol. 39, pp. 25353–25361 (2025)
- [140] Team, A.I.: On the Biology of a Large Language Model. Accessed: 2025-04-16 (2025). <https://transformer-circuits.pub/2025/attribution-graphs/biology.html>
- [141] Liang, Y., Wu, Y., Zhuang, H., Chen, L., Shen, J., Jia, Y., Qin, Z., Sanghai, S., Wang, X., Yang, C., *et al.*: Integrating planning into single-turn long-form text generation. arXiv preprint arXiv:2410.06203 (2024)
- [142] Jiang, X., Dong, Y., Wang, L., Fang, Z., Shang, Q., Li, G., Jin, Z., Jiao, W.: Self-planning code generation with large language models. ACM Transactions on Software Engineering and Methodology **33**(7), 1–30 (2024)
- [143] Ji-An, L., Xiong, H.-D., Wilson, R.C., Mattar, M.G., Benna, M.K.: Language models are capable of metacognitive monitoring and control of their internal activations. arXiv preprint arXiv:2505.13763 (2025)
- [144] Chen, Q., Qin, L., Liu, J., Peng, D., Guan, J., Wang, P., Hu, M., Zhou, Y., Gao, T., Che, W.: Towards reasoning era: A survey of long chain-of-thought for reasoning large language models. arXiv preprint arXiv:2503.09567 (2025)
- [145] Liu, Z., Xia, C., He, W., Wang, C.: Trustworthiness and self-awareness in large language models: An exploration through the think-solve-verify framework. In: Proceedings of the 2024 Joint International Conference on Computational Linguistics, Language Resources and Evaluation (LREC-COLING 2024), pp. 16855–16866 (2024)
- [146] Cheng, Q., Sun, T., Liu, X., Zhang, W., Yin, Z., Li, S., Li, L., He, Z., Chen, K., Qiu, X.: Can ai assistants know what they don’t know? arXiv preprint arXiv:2401.13275 (2024)

- [147] Tan, Z., Wei, L., Wang, J., Xie, X., Huang, W.: Can i understand what i create? self-knowledge evaluation of large language models. arXiv preprint arXiv:2406.06140 (2024)
- [148] Kapoor, S., Gruver, N., Roberts, M., Collins, K., Pal, A., Bhatt, U., Weller, A., Dooley, S., Goldblum, M., Wilson, A.G.: Large language models must be taught to know what they don't know. arXiv preprint arXiv:2406.08391 (2024)
- [149] Chen, X., Aksitov, R., Alon, U., Ren, J., Xiao, K., Yin, P., Prakash, S., Sutton, C., Wang, X., Zhou, D.: Universal self-consistency for large language model generation. arXiv preprint arXiv:2311.17311 (2023)
- [150] Davidson, T.R., Surkov, V., Veselovsky, V., Russo, G., West, R., Gulcehre, C.: Self-recognition in language models. arXiv preprint arXiv:2407.06946 (2024)
- [151] Binder, F.J., Chua, J., Korbak, T., Sleight, H., Hughes, J., Long, R., Perez, E., Turpin, M., Evans, O.: Looking inward: Language models can learn about themselves by introspection. arXiv preprint arXiv:2410.13787 (2024)
- [152] Tamoyan, H., Dutta, S., Gurevych, I.: Factual self-awareness in language models: Representation, robustness, and scaling. arXiv preprint arXiv:2505.21399 (2025)
- [153] Anthropic: Claude 3 Model Family: Opus, Sonnet, Haiku. <https://www.anthropic.com/news/claude-3-family>. [Online; accessed 2024-06-01] (2024)
- [154] Achiam, J., Adler, S., Agarwal, S., Ahmad, L., Akkaya, I., Aleman, F.L., Almeida, D., Altenschmidt, J., Altman, S., Anadkat, S., et al.: Gpt-4 technical report. arXiv preprint arXiv:2303.08774 (2023)
- [155] Jiang, L., Hwang, J.D., Bhagavatula, C., Bras, R.L., Liang, J., Dodge, J., Sakaguchi, K., Forbes, M., Borhardt, J., Gabriel, S., Tsvetkov, Y., Etzioni, O., Sap, M., Rini, R., Choi, Y.: Investigating machine moral judgement through the delphi experiment. *Nature Machine Intelligence*, 1–16 (2025)
- [156] Qiu, H., Fabbri, A.R., Agarwal, D., Huang, K.-H., Tan, S., Peng, N., Wu, C.-S.: Evaluating cultural and social awareness of llm web agents. arXiv preprint arXiv:2410.23252 (2024)
- [157] Voria, G., Catolino, G., Palomba, F.: Is attention all you need? toward a conceptual model for social awareness in large language models. In: *Proceedings of the 2024 IEEE/ACM First International Conference on AI Foundation Models and Software Engineering*, pp. 69–73 (2024)
- [158] Li, Y., Huang, Y., Lin, Y., Wu, S., Wan, Y., Sun, L.: I think, therefore i am: Benchmarking awareness of large language models using awarebench. arXiv preprint arXiv:2401.17882 (2024)

- [159] Zhuge, M., Liu, H., Faccio, F., Ashley, D.R., Csordás, R., Gopalakrishnan, A., Hamdi, A., Hammoud, H.A.A.K., Herrmann, V., Irie, K., et al.: Mindstorms in natural language-based societies of mind. arXiv preprint arXiv:2305.17066 (2023)
- [160] Choi, M., Pei, J., Kumar, S., Shu, C., Jurgens, D.: Do llms understand social knowledge? evaluating the sociability of large language models with socket benchmark. arXiv preprint arXiv:2305.14938 (2023)
- [161] Xu, Z., Wang, Y., Huang, Y., Ye, J., Zhuang, H., Song, Z., Gao, L., Wang, C., Chen, Z., Zhou, Y., et al.: Socialmaze: A benchmark for evaluating social reasoning in large language models. arXiv preprint arXiv:2505.23713 (2025)
- [162] Gandhi, K., Fränken, J.-P., Gerstenberg, T., Goodman, N.: Understanding social reasoning in language models with language models. *Advances in Neural Information Processing Systems* **36**, 13518–13529 (2023)
- [163] Wu, Y., He, Y., Jia, Y., Mihalcea, R., Chen, Y., Deng, N.: Hi-tom: A benchmark for evaluating higher-order theory of mind reasoning in large language models. In: *Findings of the Association for Computational Linguistics: EMNLP 2023*, pp. 10691–10706 (2023)
- [164] Li, G., Hammoud, H., Itani, H., Khizbullin, D., Ghanem, B.: Camel: Communicative agents for “mind” exploration of large language model society. *Advances in Neural Information Processing Systems* **36**, 51991–52008 (2023)
- [165] Park, J.S., O’Brien, J., Cai, C.J., Morris, M.R., Liang, P., Bernstein, M.S.: Generative agents: Interactive simulacra of human behavior. In: *Proceedings of the 36th Annual Acm Symposium on User Interface Software and Technology*, pp. 1–22 (2023)
- [166] Rao, A., Yerukola, A., Shah, V., Reinecke, K., Sap, M.: Normad: A benchmark for measuring the cultural adaptability of large language models. In: *NAACL* (2025)
- [167] Wimmer, H., Perner, J.: Beliefs about beliefs: Representation and constraining function of wrong beliefs in young children’s understanding of deception. *Cognition* **13**(1), 103–128 (1983)
- [168] Brown, T., Mann, B., Ryder, N., Subbiah, M., Kaplan, J.D., Dhariwal, P., Neelakantan, A., Shyam, P., Sastry, G., Askell, A., et al.: Language models are few-shot learners. *Advances in neural information processing systems* **33**, 1877–1901 (2020)
- [169] Cheung, J.: Guanaco (Revision e044a62). Hugging Face (2023). <https://doi.org/10.57967/hf/0607> . <https://huggingface.co/JosephusCheung/Guanaco>

- [170] Tang, G., Chu, Z., Zheng, W., Liu, M., Qin, B.: Towards benchmarking situational awareness of large language models: Comprehensive benchmark, evaluation and analysis. In: Findings of the Association for Computational Linguistics: EMNLP 2024, pp. 7904–7928 (2024)
- [171] Wang, L., Zhong, H.: Llm-sap: Large language models situational awareness-based planning. In: 2024 IEEE International Conference on Multimedia and Expo Workshops (ICMEW), pp. 1–6 (2024). IEEE
- [172] Needham, J., Edkins, G., Pimpale, G., Bartsch, H., Hobbhahn, M.: Large language models often know when they are being evaluated. arXiv preprint arXiv:2505.23836 (2025)
- [173] Phuong, M., Zimmermann, R.S., Wang, Z., Lindner, D., Krakovna, V., Cogan, S., Dafoe, A., Ho, L., Shah, R.: Evaluating frontier models for stealth and situational awareness. arXiv preprint arXiv:2505.01420 (2025)
- [174] Wester, J., Schrills, T., Pohl, H., Berkel, N.: “as an ai language model, i cannot”: Investigating llm denials of user requests. In: Proceedings of the 2024 CHI Conference on Human Factors in Computing Systems, pp. 1–14 (2024)
- [175] Berglund, L., Stickland, A.C., Balesni, M., Kaufmann, M., Tong, M., Korbak, T., Kokotajlo, D., Evans, O.: Taken out of context: On measuring situational awareness in llms. arXiv preprint arXiv:2309.00667 (2023)
- [176] Greenblatt, R., Denison, C., Wright, B., Roger, F., MacDiarmid, M., Marks, S., Treutlein, J., Belonax, T., Chen, J., Duvenaud, D., et al.: Alignment faking in large language models. arXiv preprint arXiv:2412.14093 (2024)
- [177] Weij, T., Hofstätter, F., Jaffe, O., Brown, S.F., Ward, F.R.: Ai sandbagging: Language models can strategically underperform on evaluations. arXiv preprint arXiv:2406.07358 (2024)
- [178] Li, Z., Gebhardt, C., Inglin, Y., Steck, N., Strelt, P., Holz, C.: Situationadapt: Contextual ui optimization in mixed reality with situation awareness via llm reasoning. In: Proceedings of the 37th Annual ACM Symposium on User Interface Software and Technology, pp. 1–13 (2024)
- [179] Hagendorff, T.: Machine psychology: Investigating emergent capabilities and behavior in large language models using psychological methods. arXiv preprint arXiv:2303.13988 1 (2023)
- [180] Wei, J., Tay, Y., Bommasani, R., Raffel, C., Zoph, B., Borgeaud, S., Yogatama, D., Bosma, M., Zhou, D., Metzler, D., et al.: Emergent abilities of large language models. arXiv preprint arXiv:2206.07682 (2022)
- [181] Bengio, Y., Cohen, M., Fornasiere, D., Ghosn, J., Greiner, P., MacDermott, M.,

- Mindermann, S., Oberman, A., Richardson, J., Richardson, O., et al.: Superintelligent agents pose catastrophic risks: Can scientist ai offer a safer path? arXiv preprint arXiv:2502.15657 (2025)
- [182] Chen, S., Yu, S., Zhao, S., Lu, C.: From imitation to introspection: Probing self-consciousness in language models. arXiv preprint arXiv:2410.18819 (2024)
- [183] Laine, R., Chughtai, B., Betley, J., Hariharan, K., Balesni, M., Scheurer, J., Hobbhahn, M., Meinke, A., Evans, O.: Me, myself, and ai: The situational awareness dataset (sad) for llms. In: Advances in Neural Information Processing Systems 37 (NeurIPS 2024), Datasets and Benchmarks Track (2024)
- [184] OpenAI: OpenAI o3 and o4-mini System Card. <https://cdn.openai.com/pdf/2221c875-02dc-4789-800b-e7758f3722c1/o3-and-o4-mini-system-card.pdf>. Accessed: 2025-06-02 (2025)
- [185] Long, R., Sebo, J., Butlin, P., Finlinson, K., Fish, K., Harding, J., Pfau, J., Sims, T., Birch, J., Chalmers, D.: Taking ai welfare seriously. arXiv preprint arXiv:2411.00986 (2024)
- [186] LeDoux, J., Birch, J., Andrews, K., Clayton, N.S., Daw, N.D., Frith, C., Lau, H., Peters, M.A., Schneider, S., Seth, A., *et al.*: Consciousness beyond the human case. *Current Biology* **33**(16), 832–840 (2023)
- [187] Andrews, K., Birch, J., Sebo, J.: Evaluating animal consciousness. *Science* **387**(6736), 822–824 (2025)
- [188] Koriat, A., *et al.*: Metacognition and Consciousness. Institute of Information Processing and Decision Making, University of Haifa . . . , Haifa, Israel (2006)
- [189] Dehaene, S., Changeux, J.-P.: Experimental and theoretical approaches to conscious processing. *Neuron* **70**(2), 200–227 (2011)
- [190] Ghallab, M., Nau, D., Traverso, P.: Automated Planning: Theory and Practice. Elsevier, San Francisco, CA (2004)
- [191] Shinn, N., Cassano, F., Gopinath, A., Narasimhan, K., Yao, S.: Reflexion: Language agents with verbal reinforcement learning. *Advances in Neural Information Processing Systems* **36**, 8634–8652 (2023)
- [192] Wang, X., Wei, J., Schuurmans, D., Le, Q.V., Chi, E.H., Narang, S., Chowdhery, A., Zhou, D.: Self-consistency improves chain of thought reasoning in language models. In: The Eleventh International Conference on Learning Representations (2022)
- [193] Schulman, J., Wolski, F., Dhariwal, P., Radford, A., Klimov, O.: Proximal policy optimization algorithms. arXiv preprint arXiv:1707.06347 (2017)

- [194] Kumar, A., Zhuang, V., Agarwal, R., Su, Y., Co-Reyes, J.D., Singh, A., Baumli, K., Iqbal, S., Bishop, C., Roelofs, R., et al.: Training language models to self-correct via reinforcement learning. arXiv preprint arXiv:2409.12917 (2024)
- [195] Wang, G., Xie, Y., Jiang, Y., Mandelkar, A., Xiao, C., Zhu, Y., Fan, L., Anandkumar, A.: Voyager: An open-ended embodied agent with large language models. Transactions on Machine Learning Research
- [196] Yao, S., Zhao, J., Yu, D., Du, N., Shafran, I., Narasimhan, K., Cao, Y.: React: Synergizing reasoning and acting in language models. In: International Conference on Learning Representations (ICLR) (2023)
- [197] Ahn, M., Brohan, A., Brown, N., Chebotar, Y., Cortes, O., David, B., Finn, C., Fu, C., Gopalakrishnan, K., Hausman, K., et al.: Do as i can, not as i say: Grounding language in robotic affordances. arXiv preprint arXiv:2204.01691 (2022)
- [198] Shah, D., Osiński, B., Levine, S., *et al.*: Lm-nav: Robotic navigation with large pre-trained models of language, vision, and action. In: Conference on Robot Learning, pp. 492–504 (2023). PMLR
- [199] Kagaya, T., Yuan, T.J., Lou, Y., Karlekar, J., Pranata, S., Kinose, A., Oguri, K., Wick, F., You, Y.: Rap: Retrieval-augmented planning with contextual memory for multimodal llm agents. In: NeurIPS 2024 Workshop on Open-World Agents (2024)
- [200] Liang, K., Zhang, Z., Fisac, J.F.: Introspective planning: Aligning robots’ uncertainty with inherent task ambiguity. Advances in Neural Information Processing Systems **37**, 71998–72031 (2024)
- [201] Schick, T., Dwivedi-Yu, J., Dessì, R., Raileanu, R., Lomeli, M., Hambro, E., Zettlemoyer, L., Cancedda, N., Scialom, T.: Toolformer: Language models can teach themselves to use tools. Advances in Neural Information Processing Systems **36**, 68539–68551 (2023)
- [202] Liu, L., Yang, X., Shen, Y., Hu, B., Zhang, Z., Gu, J., Zhang, G.: Think-in-memory: Recalling and post-thinking enable llms with long-term memory. arXiv preprint arXiv:2311.08719 (2023)
- [203] Liang, Y., Song, Z., Wang, H., Zhang, J.: Learning to trust your feelings: Leveraging self-awareness in llms for hallucination mitigation. In: Proceedings of the 3rd Workshop on Knowledge Augmented Methods for NLP, pp. 44–58 (2024)
- [204] Alayrac, J.-B., Donahue, J., Luc, P., Miech, A., Barr, I., Hasson, Y., Lenc, K., Mensch, A., Millican, K., Reynolds, M., *et al.*: Flamingo: a visual language model for few-shot learning. Advances in neural information processing systems **35**, 23716–23736 (2022)

- [205] Zhu, D., Chen, J., Shen, X., Li, X., Elhoseiny, M.: Minigpt-4: Enhancing vision-language understanding with advanced large language models. In: The Twelfth International Conference on Learning Representations
- [206] Li, Y., Du, Y., Zhou, K., Wang, J., Zhao, W.X., Wen, J.-R.: Evaluating object hallucination in large vision-language models. In: Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing, pp. 292–305 (2023)
- [207] Huang, L., Yu, W., Ma, W., Zhong, W., Feng, Z., Wang, H., Chen, Q., Peng, W., Feng, X., Qin, B., *et al.*: A survey on hallucination in large language models: Principles, taxonomy, challenges, and open questions. *ACM Transactions on Information Systems* **43**(2), 1–55 (2025)
- [208] Ni, S., Bi, K., Guo, J., Cheng, X.: When do llms need retrieval augmentation? mitigating llms’ overconfidence helps retrieval augmentation. In: Findings of the Association for Computational Linguistics ACL 2024, pp. 11375–11388 (2024)
- [209] Ren, R., Wang, Y., Qu, Y., Zhao, W.X., Liu, J., Wu, H., Wen, J.-R., Wang, H.: Investigating the factual knowledge boundary of large language models with retrieval augmentation. In: Proceedings of the 31st International Conference on Computational Linguistics, pp. 3697–3715 (2025)
- [210] Xu, R., Qi, Z., Guo, Z., Wang, C., Wang, H., Zhang, Y., Xu, W.: Knowledge conflicts for llms: A survey. In: Proceedings of the 2024 Conference on Empirical Methods in Natural Language Processing, pp. 8541–8565 (2024)
- [211] Xu, R., Lin, B., Yang, S., Zhang, T., Shi, W., Zhang, T., Fang, Z., Xu, W., Qiu, H.: The earth is flat because...: Investigating llms’ belief towards misinformation via persuasive conversation. In: Proceedings of the 62nd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers), pp. 16259–16303 (2024)
- [212] Ouyang, L., Wu, J., Jiang, X., Almeida, D., Wainwright, C., Mishkin, P., Zhang, C., Agarwal, S., Slama, K., Ray, A., *et al.*: Training language models to follow instructions with human feedback. *Advances in neural information processing systems* **35**, 27730–27744 (2022)
- [213] Beckers, S.: Moral responsibility for ai systems. *Advances in Neural Information Processing Systems* **36**, 4295–4308 (2023)
- [214] Ryan, P., Porter, Z., Al-Qaddoumi, J., McDermid, J., Habli, I.: What’s my role? modelling responsibility for ai-based safety-critical systems. *arXiv preprint arXiv:2401.09459* (2023)
- [215] Wang, Z.J., Kulkarni, C., Wilcox, L., Terry, M., Madaio, M.: Farsight: Fostering responsible ai awareness during ai application prototyping. In: Proceedings of

- the 2024 CHI Conference on Human Factors in Computing Systems, pp. 1–40 (2024)
- [216] Gallegos, I.O., Rossi, R.A., Barrow, J., Tanjim, M.M., Kim, S., Deroncourt, F., Yu, T., Zhang, R., Ahmed, N.K.: Bias and fairness in large language models: A survey. *Computational Linguistics* **50**(3), 1097–1179 (2024)
- [217] Rogers, R., Turk, V.: Openai’s sora is plagued by sexist, racist, and ableist biases. *WIRED* (2025). Accessed: 2025-04-16
- [218] Xu, R., Zhou, Z., Zhang, T., Qi, Z., Yao, S., Xu, K., Xu, W., Qiu, H.: Walking in others’ shoes: How perspective-taking guides large language models in reducing toxicity and bias. In: *Proceedings of the 2024 Conference on Empirical Methods in Natural Language Processing*, pp. 8341–8368 (2024)
- [219] Raj, C., Mukherjee, A., Caliskan, A., Anastasopoulos, A., Zhu, Z.: Breaking bias, building bridges: Evaluation and mitigation of social biases in llms via contact hypothesis. In: *Proceedings of the AAAI/ACM Conference on AI, Ethics, and Society*, vol. 7, pp. 1180–1189 (2024)
- [220] Yang, D., Hovy, D., Jurgens, D., Plank, B.: The call for socially aware language technologies. *arXiv preprint arXiv:2405.02411* (2024)
- [221] Xu, R., Li, X., Chen, S., Xu, W.: Nuclear deployed: Analyzing catastrophic risks in decision-making of autonomous llm agents. *arXiv preprint arXiv:2502.11355* (2025)
- [222] Brundage, M., Avin, S., Clark, J., Toner, H., Eckersley, P., Garfinkel, B., Dafoe, A., Scharre, P., Zeitzoff, T., Filar, B., et al.: The malicious use of artificial intelligence: Forecasting, prevention, and mitigation. *arXiv preprint arXiv:1802.07228* (2018)
- [223] Hendrycks, D., Mazeika, M., Woodside, T.: An overview of catastrophic ai risks. *arXiv preprint arXiv:2306.12001* (2023)
- [224] Yi, J., Xie, Y., Zhu, B., Kiciman, E., Sun, G., Xie, X., Wu, F.: Benchmarking and defending against indirect prompt injection attacks on large language models. *arXiv preprint arXiv:2312.14197* (2023)
- [225] Xu, R., Cai, Y., Zhou, Z., Gu, R., Weng, H., Yan, L., Zhang, T., Xu, W., Qiu, H.: Course-correction: Safety alignment using synthetic preferences. In: *Proceedings of the 2024 Conference on Empirical Methods in Natural Language Processing: Industry Track*, pp. 1622–1649 (2024)
- [226] Anderson, R., Shumailov, I.: *Situational Awareness and Adversarial Machine Learning—Robots, Manners, and Stress* (2021)

- [227] Kaur, R., Gabrijelčič, D., Klobučar, T.: Artificial intelligence for cybersecurity: Literature review and future research directions. *Information Fusion* **97**, 101804 (2023)
- [228] Lei, T., Barzilay, R., Jaakkola, T.: Rationalizing neural predictions. In: *Proceedings of the 2016 Conference on Empirical Methods in Natural Language Processing*, pp. 107–117 (2016)
- [229] Alvarez Melis, D., Jaakkola, T.: Towards robust interpretability with self-explaining neural networks. *Advances in neural information processing systems* **31** (2018)
- [230] Li, J., Galley, M., Brockett, C., Spithourakis, G., Gao, J., Dolan, W.B.: A persona-based neural conversation model. In: *Proceedings of the 54th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pp. 994–1003 (2016)
- [231] Song, H., Zhang, W.-N., Cui, Y., Wang, D., Liu, T.: Exploiting persona information for diverse generation of conversational responses. *arXiv preprint arXiv:1905.12188* (2019)
- [232] Zhong, S., Huang, Z., Gao, S., Wen, W., Lin, L., Zitnik, M., Zhou, P.: Let’s think outside the box: Exploring leap-of-thought in large language models with creative humor generation. In: *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp. 13246–13257 (2024)
- [233] Zhao, Y., Zhang, R., Li, W., Huang, D., Guo, J., Peng, S., Hao, Y., Wen, Y., Hu, X., Du, Z., et al.: Assessing and understanding creativity in large language models. *arXiv preprint arXiv:2401.12491* (2024)
- [234] Wang, Z., Chiu, Y.Y., Chiu, Y.C.: Humanoid agents: Platform for simulating human-like generative agents. In: *Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing: System Demonstrations*, pp. 167–176 (2023)
- [235] Park, J.S., Zou, C.Q., Shaw, A., Hill, B.M., Cai, C., Morris, M.R., Willer, R., Liang, P., Bernstein, M.S.: Generative agent simulations of 1,000 people. *arXiv preprint arXiv:2411.10109* (2024)
- [236] Junprung, E.: Exploring the intersection of large language models and agent-based modeling via prompt engineering. *arXiv preprint arXiv:2308.07411* (2023)
- [237] Wang, Y., Zhou, X., Wang, Y., Zhang, G., He, T.: Jailbreak large visual language models through multi-modal linkage. *arXiv preprint arXiv:2412.00473* (2024)
- [238] Wang, K., Zhang, G., Zhou, Z., Wu, J., Yu, M., Zhao, S., Yin, C., Fu, J., Yan, Y., Luo, H., et al.: A comprehensive survey in llm (-agent) full stack safety:

- Data, training and deployment. arXiv preprint arXiv:2504.15585 (2025)
- [239] Zhou, Z., Li, Z., Zhang, J., Zhang, Y., Wang, K., Liu, Y., Guo, Q.: Corba: Contagious recursive blocking attacks on multi-agent systems based on large language models. arXiv preprint arXiv:2502.14529 (2025)
- [240] Bengio, Y., Mindermann, S., Privitera, D., Besiroglu, T., Bommasani, R., Casper, S., Choi, Y., Fox, P., Garfinkel, B., Goldfarb, D., et al.: International ai safety report. arXiv preprint arXiv:2501.17805 (2025)
- [241] Ojewale, V., Steed, R., Vecchione, B., Birhane, A., Raji, I.D.: Towards ai accountability infrastructure: Gaps and opportunities in ai audit tooling. In: Proceedings of the 2025 CHI Conference on Human Factors in Computing Systems, pp. 1–29 (2025)
- [242] Park, P.S., Goldstein, S., O’Gara, A., Chen, M., Hendrycks, D.: Ai deception: A survey of examples, risks, and potential solutions. *Patterns* **5**(5) (2024)
- [243] Barkur, S.K., Schacht, S., Scholl, J.: Deception in llms: Self-preservation and autonomous goals in large language models. arXiv preprint arXiv:2501.16513 (2025)
- [244] Hubinger, E., Merwijk, C., Mikulik, V., Skalse, J., Garrabrant, S.: Risks from learned optimization in advanced machine learning systems. arXiv preprint arXiv:1906.01820 (2019)
- [245] Hubinger, E., Denison, C., Mu, J., Lambert, M., Tong, M., MacDiarmid, M., Lanham, T., Ziegler, D.M., Maxwell, T., Cheng, N., et al.: Sleeper agents: Training deceptive llms that persist through safety training. arXiv preprint arXiv:2401.05566 (2024)
- [246] Scheurer, J., Balesni, M., Hobbhahn, M.: Large language models can strategically deceive their users when put under pressure. In: ICLR 2024 Workshop on Large Language Model (LLM) Agents
- [247] Dehnert, M., Mongeau, P.A.: Persuasion in the age of artificial intelligence (ai): Theories and complications of ai-based persuasion. *Human Communication Research* **48**(3), 386–403 (2022)
- [248] Sabour, S., Liu, J.M., Liu, S., Yao, C.Z., Cui, S., Zhang, X., Zhang, W., Cao, Y., Bhat, A., Guan, J., et al.: Human decision-making is susceptible to ai-driven manipulation. arXiv preprint arXiv:2502.07663 (2025)
- [249] Zhang, R., Li, H., Meng, H., Zhan, J., Gan, H., Lee, Y.-C.: The dark side of ai companionship: A taxonomy of harmful algorithmic behaviors in human-ai relationships. arXiv preprint arXiv:2410.20130 (2024)

- [250] Roose, K.: A conversation with bing’s chatbot left me deeply unsettled. *The New York Times* (2023)
- [251] Schmitt, M., Flechais, I.: Digital deception: Generative artificial intelligence in social engineering and phishing. *Artificial Intelligence Review* **57**(12), 1–23 (2024)
- [252] Singh, S., Abri, F., Namin, A.S.: Exploiting large language models (llms) through deception techniques and persuasion principles. In: *2023 IEEE International Conference on Big Data (BigData)*, pp. 2508–2517 (2023). IEEE
- [253] European Union Agency for Cybersecurity (ENISA): EU Elections at Risk with Rise of AI-Enabled Information Manipulation. accessed: 18 April 2025 (2023). <https://www.enisa.europa.eu/news/eu-elections-at-risk-with-rise-of-ai-enabled-information-manipulation>
- [254] Phang, J., Lampe, M., Ahmad, L., Agarwal, S., Fang, C.M., Liu, A.R., Danry, V., Lee, E., Chan, S.W., Pataranutaporn, P., et al.: Investigating affective use and emotional well-being on chatgpt. arXiv preprint arXiv:2504.03888 (2025)
- [255] Epley, N., Waytz, A., Cacioppo, J.T.: On seeing human: a three-factor theory of anthropomorphism. *Psychological review* **114**(4), 864 (2007)
- [256] Li, M., Suh, A.: Anthropomorphism in ai-enabled technology: A literature review. *Electronic Markets* **32**(4), 2245–2275 (2022)
- [257] Tiku, N.: The Google Engineer Who Thinks the Company’s AI Has Come to Life. *The Washington Post*, Accessed: 2025-04-18 (2022). <https://www.washingtonpost.com/technology/2022/06/11/google-ai-lamda-blake-lemoine/>
- [258] Novella, S.: Is LaMDA Sentient? *NeuroLogica Blog* (2022). <https://theness.com/neurologicablog/is-lamda-sentient/>
- [259] Waytz, A., Heafner, J., Epley, N.: The mind in the machine: Anthropomorphism increases trust in an autonomous vehicle. *Journal of experimental social psychology* **52**, 113–117 (2014)
- [260] Cohn, M., Pushkarna, M., Olanubi, G.O., Moran, J.M., Padgett, D., Mengesha, Z., Heldreth, C.: Believing anthropomorphism: examining the role of anthropomorphic cues on trust in large language models. In: *Extended Abstracts of the CHI Conference on Human Factors in Computing Systems*, pp. 1–15 (2024)
- [261] Abercrombie, G., Curry, A.C., Dinkar, T., Rieser, V., Talat, Z.: Mirages. on anthropomorphism in dialogue systems. In: *Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing*, pp. 4776–4790 (2023)
- [262] Li, Q., Luximon, Y., Zhang, J.: The influence of anthropomorphic cues on

- patients' perceived anthropomorphism, social presence, trust building, and acceptance of health care conversational agents: within-subject web-based experiment. *Journal of medical Internet research* **25**, 44479 (2023)
- [263] Go, E., Sundar, S.S.: Humanizing chatbots: The effects of visual, identity and conversational cues on humanness perceptions. *Computers in human behavior* **97**, 304–316 (2019)
- [264] Roy, R., Naidoo, V.: Enhancing chatbot effectiveness: The role of anthropomorphic conversational styles and time orientation. *Journal of Business Research* **126**, 23–34 (2021)
- [265] Chen, Q.Q., Park, H.J.: How anthropomorphism affects trust in intelligent personal assistants. *Industrial Management & Data Systems* **121**(12), 2722–2737 (2021)
- [266] Placani, A.: Anthropomorphism in ai: hype and fallacy. *AI and Ethics* **4**(3), 691–698 (2024)
- [267] Deshpande, A., Rajpurohit, T., Narasimhan, K., Kalyan, A.: Anthropomorphization of ai: Opportunities and risks. In: *Proceedings of the Natural Language Processing Workshop 2023*, pp. 1–7 (2023)
- [268] Bonnefon, J.-F., Rahwan, I., Shariff, A.: The moral psychology of artificial intelligence. *Annual review of psychology* **75**(1), 653–675 (2024)
- [269] Birch, J.: *The Edge of Sentience: Risk and Precaution in Humans, Other Animals, and AI*. Oxford University Press, Oxford (2024)
- [270] Pan, X., Dai, J., Fan, Y., Yang, M.: Frontier ai systems have surpassed the self-replicating red line. *arXiv preprint arXiv:2412.12140* (2024)
- [271] Amodei, D., Olah, C., Steinhardt, J., Christiano, P., Schulman, J., Mané, D.: Concrete problems in ai safety. *arXiv preprint arXiv:1606.06565* (2016)
- [272] He, Y., Li, Y., Wu, J., Sui, Y., Chen, Y., Hooi, B.: Evaluating the paperclip maximizer: Are rl-based language models more likely to pursue instrumental goals? *arXiv preprint arXiv:2502.12206* (2025)
- [273] Benson-Tilsen, T., Soares, N.: Formalizing convergent instrumental goals. In: *AAAI Workshop: AI, Ethics, and Society* (2016)
- [274] Omohundro, S.M.: The basic ai drives. In: Yampolskiy, R.V. (ed.) *Artificial Intelligence Safety and Security*, pp. 47–55. Chapman and Hall/CRC, Boca Raton, FL (2018)
- [275] Nick, B.: *Superintelligence: Paths, dangers, strategies*. Oxford University Press, Oxford (2014)

- [276] Ashcraft, C., Karra, K., Carney, J., Drenkow, N.: Investigating the treacherous turn in deep reinforcement learning. arXiv preprint arXiv:2504.08943 (2025)
- [277] Kim, J., Choi, W., Lee, B.: Prompt flow integrity to prevent privilege escalation in llm agents. arXiv preprint arXiv:2503.15547 (2025)
- [278] Bubeck, S., Chadrasekaran, V., Eldan, R., Gehrke, J., Horvitz, E., Kamar, E., Lee, P., Lee, Y.T., Li, Y., Lundberg, S., et al.: Sparks of artificial general intelligence: Early experiments with gpt-4. ArXiv (2023)
- [279] Kulveit, J., Douglas, R., Ammann, N., Turan, D., Krueger, D., Duvenaud, D.: Gradual disempowerment: Systemic existential risks from incremental ai development. arXiv preprint arXiv:2501.16946 (2025)
- [280] Ward, F.R., Hofstätter, F., Thomson, L.A., Wood, H.M., Jaffe, O., Bartak, P., Brown, S.F.: Tall tales at different scales: Evaluating scaling trends for deception in language models