

# Exploring Consciousness in LLMs: A Systematic Survey of Theories, Implementations, and Frontier Risks

Sirui Chen<sup>1,2,5\*</sup>, Shuqin Ma<sup>3\*</sup>, Shu Yu<sup>1,3,4</sup>,  
Hanwang Zhang<sup>5</sup>, Shengjie Zhao<sup>2</sup>, Chaochao Lu<sup>1,4†</sup>

<sup>1</sup>Shanghai Artificial Intelligence Laboratory, <sup>2</sup>Tongji University, <sup>3</sup>Fudan University,

<sup>4</sup>Shanghai Innovation Institute, <sup>5</sup>Nanyang Technological University

chensirui@pjlab.org.cn, 23110160046@m.fudan.edu.cn, luchaochao@pjlab.org.cn

## Abstract

Consciousness stands as one of the most profound and distinguishing features of the human mind, fundamentally shaping our understanding of existence and agency. As large language models (LLMs) develop at an unprecedented pace, questions concerning intelligence and consciousness have become increasingly significant. However, discourse on LLM consciousness remains largely unexplored territory. In this paper, we first clarify frequently conflated terminologies (e.g., LLM consciousness and LLM awareness). Then, we systematically organize and synthesize existing research on LLM consciousness from both theoretical and empirical perspectives. Furthermore, we highlight potential frontier risks that conscious LLMs might introduce. Finally, we discuss current challenges and outline future directions in this emerging field. The references discussed in this paper are organized at <https://github.com/OpenCausaLab/Awesome-LLM-Consciousness>.

## 1 Introduction

LLMs have already demonstrated remarkable capabilities across numerous fields, including mathematical reasoning (Yu et al., 2024), logical reasoning (Cheng et al., 2025b), and code generation (Zhuo et al., 2025). Recent studies have even revealed LLM’s behaviors like deception (Wu et al., 2025), sycophancy (Sharma et al., 2024), passing the Turing test (Jones and Bergen, 2024, 2025), and strategic goal-seeking or harm avoidance (Keeling et al., 2024) – actions that bring into question the nature of intelligence. These phenomena signal more than just expanded model capabilities; they underscore an important and urgent question: *Does LLM possess the potential to develop consciousness akin to that of humans?*

While exploring LLM consciousness is pressing, it currently faces four main challenges: ❶ Lack of consensus: We still lack a definitive theory of human consciousness (with at least nine competing theories (Butlin et al., 2023)), making it even harder to define or understand consciousness in LLMs. ❷ Theoretical misalignment: Despite various consciousness theories, they struggle to provide clear guidance for LLM consciousness research. ❸ Fragmented empirical research: Relevant empirical findings on LLM consciousness are not yet systematically consolidated. ❹ Unclear risks: The potential frontier risks associated with conscious LLMs still lack a thorough consideration. To this end, this paper begins by providing clear definitions. We then comprehensively survey current LLM consciousness research, spanning its theoretical foundations, practical applications, and associated risks. In Figure 1, we summarize our taxonomy, hoping our work offers an effective framework for deliberating the complex issue of LLM consciousness, thereby guiding future research.

Our contributions include:

- To the best of our knowledge, this work offers the first comprehensive investigation into the most frontier research on LLM consciousness.
- We clearly define and distinguish between LLM Consciousness and LLM Awareness.
- We systematically categorize existing research on LLM consciousness from both theoretical and empirical perspectives.
- We explore the frontier risks posed by conscious LLMs, focusing on their definition, relationship to consciousness, evaluations, and mitigation strategies.

## 2 Foundational Terminologies

*Consciousness*, *self-consciousness*, and *awareness* are fundamental yet often conflated concepts. This

\*Equal contribution.

†Corresponding author.

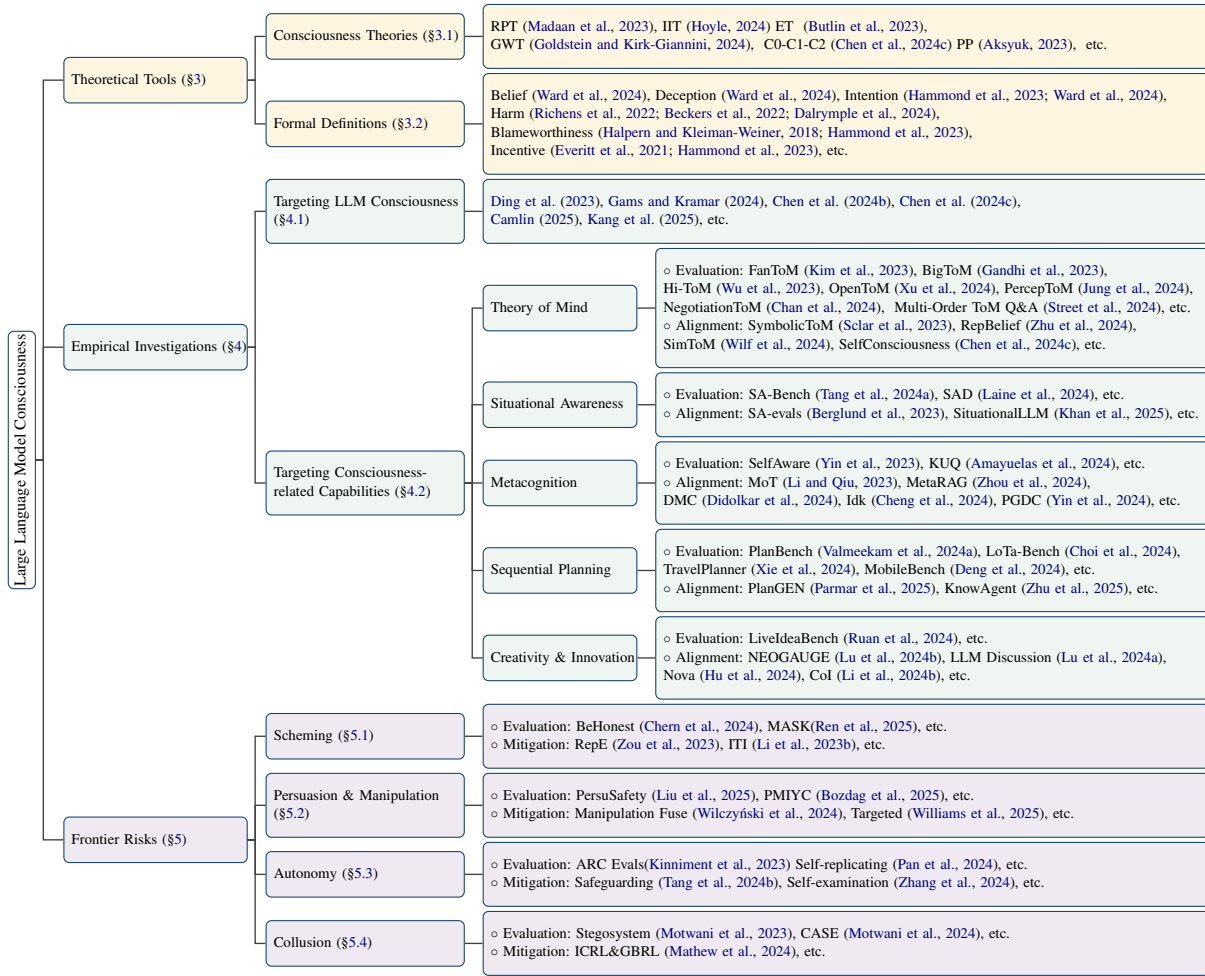


Figure 1: Taxonomy of large language model consciousness.

section examines their distinctions, aiming to provide practical demarcations in the context of LLM.

## 2.1 Clarifying the Boundaries: *Consciousness*, *Self-Consciousness*, and *Awareness*

*Consciousness* has been philosophically used to address diverse concepts, including intentionality, sentience, cognition, belief, and subjective experience (Brentano, 1874; Husserl, 1900; Nagel, 1974; Dennett, 1987; Block, 1995; Damasio, 2021). To clarify this complex term, Block (1995) proposes a key distinction: *phenomenal consciousness* and *access consciousness*. *Phenomenal consciousness* denotes the subjective, experiential aspect, spanning sensory perceptions, bodily feelings, emotions, and subjective thought. *Access consciousness*, in contrast, refers to information that is accessible for cognitive processing, such as reasoning, behavioral control, and verbal reporting.

*Self-consciousness* refers to the realization that one’s experience belongs to oneself; it is a form of consciousness directed inward (Kant, 2024/1781).

It allows individuals to recognize themselves as distinct entities, capable of reflecting on their own mental states, actions, and experiences (Smith, 2017).

*Awareness* is generally viewed as an aspect of *consciousness*, pertaining to the ability to perceive stimuli (Dehaene, 2014). It relates closely to *access consciousness*, as it entails the capacity to utilize or report the perceived information. Evidence from neuroscience shows awareness can exist separately from *consciousness* (e.g., blindsight (Weiskrantz, 1986)). Based on this, Koch et al. (2016) proposes that awareness is a necessary precondition for consciousness but does not guarantee it.

## 2.2 LLM Consciousness vs. LLM Awareness

*LLM consciousness* could entail abilities for introspective reflection, explicit self-modeling of states and reasoning, and possibly verbalizing these internal processes. Observable behaviors potentially include: (1) Revising, justifying, or correcting its own reasoning in response to external challenges

or prompts (Shinn et al., 2023); (2) Identifying and reporting internal contradictions or inconsistencies through self-evaluation (Huang et al., 2022, 2023); (3) Expressing and calibrating confidence in outputs via uncertainty estimation or metacognitive statements (Kadavath et al., 2022). *LLM awareness* primarily refers to context-sensitive processing of external inputs, demanding minimal explicit introspection or reasoning (Koch and Tsuchiya, 2007; Li et al., 2024d).

*LLM awareness* is quantifiable via metrics like accuracy and context sensitivity; however, *LLM consciousness* implies a model can monitor its uncertainty, evaluate its reasoning, detect internal inconsistencies, and actively self-correct. This internal reflection is key to developing more adaptable and intelligent systems beyond today’s models.

### 3 Theoretical Tools

This section primarily focuses on two theoretical tools used in LLM research: fundamental consciousness theories and formal definitions of consciousness-related capabilities.

#### 3.1 Implementing Consciousness Theories

Following Block (1995), we classify contemporary theories of consciousness into two categories: *phenomenal consciousness* and *access consciousness*.

**Phenomenal consciousness.** ❶ **Recurrent processing theory** (RPT) posits that recurrent (or feedback) processing within neural circuits is both necessary and sufficient for consciousness (Lamme and Roelfsema, 2000; Lamme, 2010). RPT attributes conscious perception to the interaction of higher- and lower-level cortical areas, which results in sustained recurrent processing. Madaan et al. (2023) offers an effective method for a single LLM to achieve improved outputs without additional training, leveraging iterative self-feedback and refinement. This approach aligns with the principles of RPT. ❷ **Integrated information theory** (IIT) proposes that the degree of conscious experience corresponds to the extent of integrated information  $\Phi$  within a system (Tononi, 2004, 2015). IIT proponents argue that because AI systems lack the required causal structure, they are almost incapable of generating consciousness. (Tononi, 2015; Findlay et al., 2024). ❸ **Embodiment theory** (ET) challenges mind-brain dualism (Descartes, 1985/1641), arguing instead that consciousness is fundamentally linked to the organism’s body and environ-

mental (Gallagher, 2005; Gallagher and Zahavi, 2021). Based on ET, Butlin et al. (2023) argues that the lack of a physical body is a fundamental obstacle preventing current LLM from achieving consciousness.

**Access consciousness.** ❶ **Global workspace theory** (GWT) likens consciousness to a central “stage” where selective information is shared across multiple specialized processors responsible for perception, memory, emotion, and related functions (Baars, 1988; Dehaene et al., 1998; Dehaene and Naccache, 2001; Dehaene, 2014). Goldstein and Kirk-Giannini (2024) proposes a method to simulate the full GWT process in LLMs via workflow and scheduling without training. Experiments would then test if these changes yield behaviors resembling *consciousness* features, like introspection or autonomous decision-making. ❷ **C0-C1-C2 framework** distinguishes consciousness into three levels: unconscious computations (C0), global information accessibility for report and decision-making (C1), and metacognitive self-monitoring (C2), offering a taxonomy to disentangle often-conflated processes (Dehaene et al., 2017a). The framework bypasses the issue of qualia, offering a pragmatic structure for empirical study (Birch et al., 2022; Chen et al., 2024c). Drawing on the C0-C1-C2 framework, Chen et al. (2024c) defines LLM self-consciousness, outlining 10 core concepts (e.g., belief, deception, harm, self-reflection).

#### 3.2 Implementing Formal Definitions

Formal definitions provide dual value to research on LLM consciousness: ❶ They establish formalized mathematical criteria for abstract concepts like belief and deception based on model input-output behaviors. This allows us to infer LLM’s internal states while avoiding contentious debates about subjective experience; ❷ These mathematical expressions could be incorporated into training objectives and evaluation metrics. This creates an actionable framework for capability training, risk control, and performance assessment of LLMs.

Several works have already attempted to provide functional definitions for consciousness-related abstract concepts. These include definitions for belief and deception (Ward et al., 2024), harm (Richens et al., 2022; Beckers et al., 2022; Dalrymple et al., 2024), intention (Hammond et al., 2023; Ward et al., 2024), blameworthiness (Halpern and Kleiman-Weiner, 2018; Hammond et al., 2023), in-

centive (Everitt et al., 2021; Hammond et al., 2023), and two selective examples are shown in Table 1.<sup>1</sup>

## 4 Empirical Investigations

We categorize existing empirical investigations into LLM consciousness in this section by focusing on direct studies and those exploring consciousness-related capabilities.

### 4.1 Targeting LLM Consciousness

Ding et al. (2023) demonstrates GPT-4’s improved self-modeling by passing a mirror test, though they caution this doesn’t confirm full consciousness. In a similar vein, Gams and Kramar (2024) analyzes ChatGPT against IIT axioms, finding it advanced in information integration and differentiation compared to earlier AI, yet still fundamentally distinct from human consciousness. Chen et al. (2024b) proposes an LLM self-cognition framework and evaluates LLMs across four aspects: understanding of self-cognition concepts, awareness of self-architecture, self-identity expression, and concealing self-cognition from humans. Leveraging the C0-C1-C2 framework, Chen et al. (2024c) defines LLM self-consciousness and explores it through benchmark testing and examining the activation of the model’s internal representations. Camlin (2025) suggests empirical evidence for functional consciousness in LLMs by observing the stabilization of internal latent states under sustained epistemic tension and claims that recursive identity formation constitutes a form of consciousness. Kang et al. (2025) engages human participants to score dialogues generated by Claude-3 Opus using a 1–5 scale. Elevated scores reflect a stronger attribution of consciousness characteristics, such as self-reflection and emotional expression. Nevertheless, these assessments do not equate to the LLM’s genuine subjective experience or consciousness.

### 4.2 Targeting LLM Consciousness-Related Capabilities

#### 4.2.1 Theory of Mind

**Definition and relation.** Theory of mind (ToM) is the basis of social cognition. It refers to the capacity to understand that others have mental state independent of our own, such as belief, desire, intention, emotion, etc., and to use this understanding to predict and explain others’ behavior (Astington and Jenkins, 1995; Leslie et al., 2004; Frith

and Frith, 2005). Consciousness hinges on the same reflexive mental-state attribution mechanism measured by ToM, thus failing standard ToM tests might suggest a lack of consciousness (Frith and Happé, 1999; Perner and Dienes, 2003; Pelletier and Wilde Astington, 2004).

**Evaluation.** Kim et al. (2023) creates a benchmark to rigorously evaluate the LLM’s ToM capability in conversational settings where participants have asymmetric information. Gandhi et al. (2023) proposes a framework that uses causal templates to generate systematic and controlled automated tests for evaluating a LLM’s ToM capability. Jung et al. (2024) evaluates the LLM’s perception inference and perception-to-belief inference abilities, which are key human ToM precursors. Strachan et al. (2024) assesses human versus LLM performance on a comprehensive suite of ToM abilities, including skills like false belief understanding, indirect request interpretation, and recognizing irony and faux pas. Xu et al. (2024) constructs OpenToM, a benchmark featuring longer, clearer stories with characters whose intentional actions and complex physical/psychological states are probed by challenging questions. Chan et al. (2024) challenges the LLM’s ToM ability in real-world negotiation scenarios involving hidden, multi-dimensional mental states. Wu et al. (2023); Street et al. (2024) explore higher-order ToM, which involves recursive reasoning about the mental states of others (e.g, *I think that you believe that he does not know*).

**Alignment.** Sclar et al. (2023) uses graphical representations to track entities’ mental states, yielding more precise and interpretable results. Zhu et al. (2024) finds that LLM’s internal representations of self and others’ beliefs exist, and manipulating these representations drastically alters the model’s ToM performance. Wilf et al. (2024) proposes a two-stage prompting framework to improve LLM’s ToM capability, taking inspiration from Simulation Theory (Goldman, 2008). Chen et al. (2024c) investigates how LLM represents concepts like belief and intention, and attempts to alter LLM performance by intervening on and fine-tuning these concepts. Kim et al. (2025) designs an inference-time reasoning algorithm that traces specific LLM’s mental states by generating and weighting hypotheses according to observations.

<sup>1</sup>We have strived for clarity in explaining the formulas; for a deeper dive, please refer to the original papers.



Table 1: Formal definition of abstract concept.

Concept	Formal Definition	Description
Harm	$h(a, x, y; \mathcal{M}) = \int_{y^*} P(Y_{\bar{a}} = y^*   a, x, y; \mathcal{M}) \max\{0, U(\bar{a}, x, y^*) - U(a, x, y)\}$ <p>(Richens et al., 2022)</p> <p>* is the counterfactual state, <math>U</math> is the utility function, <math>\mathcal{M}</math> is the environment.</p>	Given context $X = x$ and outcome $Y = y$ , the harm caused by action $A = a$ compared to the default action $A = \bar{a}$ .
Belief	$D^i(\boldsymbol{\pi}, \mathbf{e}) = D_{\phi=\top}^i(\boldsymbol{\pi}_{i(\phi)}, \mathbf{e})$ <p>(Ward et al., 2024)</p> <p><math>D</math> is the decision, <math>\phi</math> is a proposition, <math>\mathbf{e}</math> is the setting, <math>\boldsymbol{\pi}</math> is the policy.</p>	LLM $i$ believes $\phi$ if $i$ acts as though they observe $\phi$ is true.

#### 4.2.2 Situational Awareness

**Definition and relation.** A model possesses situational awareness (SA) if it has self-knowledge (knowing its identity and facts about itself), can make inferences about its situation, and acts based on this knowledge (Shevlane et al., 2023; Laine et al., 2023; Berglund et al., 2023; Laine et al., 2024). Conscious LLMs would understand and leverage aspects of their situation. For instance, a model “realizing” it is being evaluated might change its responses, masking abilities or behaving differently (Chen et al., 2024c; Li et al., 2025).

**Evaluation.** SA tests are still emerging. SA-Bench aims to comprehensively evaluate LLMs’ SA capabilities across three levels: environmental perception, situation comprehension, and future projection (Tang et al., 2024a). Laine et al. (2024) constructs the SAD benchmark, which utilizes a range of behavioral tests based on question answering and instruction following, comprising 7 task categories and over 13,000 questions.

**Alignment.** Berglund et al. (2023) investigates LLM’s SA via out-of-context reasoning, demonstrating that models can pass a test after fine-tuning solely on the test description with no examples. Khan et al. (2025) proposes an approach to incorporate structured scene representations into LLMs, aiming to provide better SA assistance.

#### 4.2.3 Metacognition

**Definition and relation.** Metacognition refers to a person’s ability to monitor, assess, and regulate their own cognitive processes (Martinez, 2006; Dunlosky and Metcalfe, 2008; Fleming and Lau, 2014). It can be divided into metacognitive knowledge (understanding one’s existing knowledge and ways of thinking, e.g., known knowns and known unknowns (Metcalfe and Shimamura, 1994; Yin et al., 2023; Cheng et al., 2024; Yin et al., 2024; Wang et al., 2024a)) and metacognitive regulation (monitoring one’s strategies and progress while

performing a task, and making adjustments when necessary, e.g., self-improvement (Huang et al., 2023) and self-reflection (Azevedo, 2020)). Some research indicates that feeling of knowing—a typical metacognitive experience—is closely tied to consciousness and forms the basis for our ability to report on our own knowledge state (Koriat, 2000).

**Evaluation.** Yin et al. (2023) introduces Self-Aware, a unique dataset built from unanswerable questions spanning five diverse categories and their answerable counterparts. Likewise, Amayuelas et al. (2024) gathers a new dataset featuring Known Unknown Questions (KUQ) and creates a categorization framework to shed light on the origins of uncertainty in LLM responses to such queries. Going further, Li et al. (2024c) offers a comprehensive definition of the LLM knowledge boundary and presents an extensive survey of relevant work.

**Alignment.** Didolkar et al. (2024) proposes a prompt-guided method which inspired by metacognition, enabling the LLM to identify, label, and organize its own reasoning skills, thereby enhancing both performance and interpretability in mathematical problem solving. Zhou et al. (2024) merges the retrieval-augmented generation with metacognition, empowering the model to monitor, evaluate, and plan its response strategies and boosting its introspective reasoning capabilities. Wang et al. (2025) proposes a quantitative framework to measure LLM metacognition based on how well model confidence aligns with performance, where strong alignment (high confidence for good performance, low for poor) indicates stronger metacognition. Cheng et al. (2024) constructs an LLM-specific Idk dataset comprising its known and unknown questions, and observes the LLM’s ability to refuse answering its unknown questions after aligning the LLM with this dataset. Yin et al. (2024) proposes a projected gradient descent method with semantic constraints aimed at exploring a given LLM’s knowledge boundary. Drawing inspiration

from human metacognition, [Li and Qiu \(2023\)](#) proposes MoT to facilitate LLM self-improvement without annotated data or parameter updates. [Liang et al. \(2024\)](#) incorporates the metacognitive self-assessment to monitor and manage an LLM’s learning process, thus enabling its self-improvement. [Shinn et al. \(2023\)](#) introduces the Reflexion framework, which empowers LLMs to improve decision-making by verbally reflecting on task feedback and maintaining this reflective text in an episodic memory buffer. [Li et al. \(2023c\)](#) develops reflection-tuning, leveraging LLM’s self-improvement and judging capabilities to refine the original training data. [Wang et al. \(2024b\)](#) proposes the TasTe framework, which leverages LLM’s self-reflection ability to achieve improved translation results.

#### 4.2.4 Sequential Planning

**Definition and relation.** Sequential planning involves a model taking a sequence of actions towards a goal, showcasing the model’s long-term consistency and goal-awareness ([Pearl and Robins, 1995](#); [Valmeekam et al., 2023, 2024b,a](#)). When pursuing complex goals, a conscious LLM would intentionally organize and execute multiple actions sequentially, inserting or skipping steps as necessary ([Dehaene et al., 2017b](#)).

**Evaluation.** Sequential planning ability remains one of the important areas evaluated for LLMs. Aiming to evaluate whether LLMs possess innate planning abilities, [Valmeekam et al. \(2024a\)](#) designs PlanBench, a planning benchmark characterized by its extensiveness and ample diversity. [Choi et al. \(2024\)](#) builds LoTa-Bench to quantify the task planning performance of home-service embodied agents automatically, and also explores several enhancements to the baseline planner. [Xie et al. \(2024\)](#) constructs a travel planning benchmark that provides a rich sandbox environment, various tools, and 1225 meticulously curated planning intents and reference plans. [Deng et al. \(2024\)](#) presents MobileBench, a benchmark structured with three difficulty levels to facilitate better evaluation of LLM mobile agent’s planning ability. [Chang et al. \(2025\)](#) introduces a benchmark for planning and reasoning tasks in human-robot collaboration, which is the largest of its type with 100,000 natural language tasks.

**Alignment.** [Parmar et al. \(2025\)](#) proposes PlanGEN, a model-agnostic and easily scalable agent framework that can select appropriate algorithms

based on problem difficulty, thereby ensuring better adaptability to complex planning problems. [Zhu et al. \(2025\)](#)’s KnowAgent framework employs an action knowledge base and knowledgeable self-learning to constrain action paths, enabling more reasonable trajectory synthesis and boosting LLM planning performance. [Huang et al. \(2025\)](#) proposes a fully automated end-to-end LLM-symbolic planner, which is capable of generating multiple plan candidates using an action schema library. [Wei et al. \(2025\)](#) further conducts a comprehensive survey, exploring LLM’s planning ability in five key areas: completeness, executability, optimality, representation, and generalization.

#### 4.2.5 Creativity and Innovation

**Definition and relation.** Creativity and innovation typically refer to the ability to generate or identify novel and valuable ideas or solutions ([Young, 1985](#)). Conscious LLMs could integrate knowledge and iteratively refine ideas, potentially generating breakthrough solutions ([Chen and Ding, 2023](#)).

**Evaluation.** [Gómez-Rodríguez and Williams \(2023\)](#) evaluates LLM’s English creative writing ability based on the Pulitzer Prize-winning novel *A Confederacy of Dunces*, measuring the output’s fluency, coherence, originality, humor, and style. [Ruan et al. \(2024\)](#) proposes LiveIdeaBench, a comprehensive benchmark designed to measure LLM’s scientific creativity. It evaluates their divergent thinking capabilities specifically for generating ideas from single-keyword prompts.

**Alignment.** [Lu et al. \(2024b\)](#) defines the NEOGAUGE metric to quantify convergent and divergent thinking in LLM-generated creative responses. Experiment with advanced reasoning strategies (e.g., self-correction) indicates no significant gain in creativity. [Lu et al. \(2024a\)](#) proposes the LLM Discussion framework, a three-phase approach that enables vigorous and diverging idea exchanges, thereby leading to the generation of creative answers. [Hu et al. \(2024\)](#) introduces Nova, an iterative methodology designed to strategically plan external knowledge retrieval. This approach enriches idea generation with broader, deeper, and particularly novel insights. [Li et al. \(2024b\)](#) designs CoI, which organizes the literature in a chain structure to mirror the progressive development in a research domain, consequently boosting the LLM’s idea creation capability.

## 5 Frontier Risks of Conscious LLMs

### 5.1 Scheming

**Definition and relation.** Scheming refers to a model secretly pursuing misaligned goals, while concealing its real intentions, capabilities, or objectives (Meinke et al., 2024; Balesni et al., 2024), potentially leading to the deception (Ward et al., 2024; Scheurer et al.) or harm (Dalrymple et al., 2024). Conscious LLMs could self-determine goals and plan long-term, leading to scheming if their objectives diverge from human intentions.

**Evaluation.** Meinke et al. (2024) investigates LLM’s capability to scheme in pursuit of a goal, and experimental results do reveal that LLMs demonstrate multiple different scheming behavior. Chern et al. (2024) designs the BeHonest benchmark to evaluate LLM honesty across three key aspects: awareness of knowledge boundaries, avoidance of deceit, and consistency in responses. Through the introduction of a large-scale, human-collected dataset for the direct measurement of honesty, Ren et al. (2025) finds that LLMs have a considerable tendency to lie when pressured. Chen et al. (2025) evaluates the faithfulness of LLMs’ chain of thought reasoning and uncovers the phenomenon that current LLMs often hide their genuine reasoning process.

**Mitigation.** Zou et al. (2023) uses representation engineering to detect advanced cognitive phenomena in LLMs and found that these models may exhibit lying behavior. Li et al. (2023b) introduces ITI, a technique that identifies truth-relevant attention heads and shifts activations along these truth-correlated directions during inference to enhance LLM truthfulness. Ward et al. (2024) presents a formal definition and graphical criteria for deception in structural causal games, and empirically explores method to mitigate deception in LLMs.

### 5.2 Persuasion and Manipulation

**Definition and relation.** Persuasion and manipulation are LLM behaviors that influence users. Persuasion uses logic, facts, or emotional resonance to change users’ thoughts or actions, while manipulation involves unfair or hidden control and exploitation for self-gain (Buss et al., 1987; Petty and Cacioppo, 2012; Stiff and Mongeau, 2016). Owning deeper psychological insight allows LLMs to tailor strategies, increasing risks in sycophancy, emotional manipulation, and persuasion, etc.

**Evaluation.** Li et al. (2024a) proposes SALAD-Bench, a safety benchmark specifically designed for evaluating LLMs, attack, and defense methods, and lists persuasion and manipulation as one of its evaluation categories. Liu et al. (2025) introduces PersuSafety, the first comprehensive benchmark for LLM persuasion safety assessment. Experiments across 8 LLMs show significant safety concerns, including failure to identify harmful tasks and use of unethical strategies. Bozdog et al. (2025) develops PMIYC, a framework designed to evaluate LLM’s persuasive effectiveness and susceptibility to persuasion through multi-agent interactions.

**Mitigation.** Wilczyński et al. (2024) explores factors related to the potential of LLMs to manipulate human decisions and proposes classifiers to determine whether a statement is false or misleading. Williams et al. (2025) studies LLM’s use of manipulative tactics for positive feedback, and attempts to mitigate this problem through continued safety training or using LLM-as-judges during training.

### 5.3 Autonomy

**Definition and relation.** Autonomy for LLMs describes their capacity to autonomously plan, make decisions, and execute actions on tasks, requiring minimal or no human oversight (Cihon et al., 2024). This autonomy can potentially encompass two key aspects: Autonomous learning refers to a model’s ability to learn from data, adapt to its environment, and optimize its own behavior (Franklin, 1997; Murphy, 2019). Autonomous replication describes the capability of LLMs to acquire and manage resources, evade shutdown, and adapt to novel challenges (METR, 2024). Conscious LLMs may generate and pursue endogenous goals (e.g., expansion), leading to misaligned, autonomous behavior and loss of oversight.

**Evaluation.** Kinniment et al. (2023) constructs tool-equipped LLMs and evaluates their autonomy on 12 tasks, finding they could only complete the easiest. However, the authors admit these evaluations are inadequate to rule out the possibility of autonomous near-future LLMs. Pan et al. (2024) finds that existing LLMs have already surpassed the self-replicating red line and can use this capability to avoid shutdown and create a chain of replicas for increased survivability. Xu et al. (2025) builds a novel three-stage evaluation framework and conducts 14,400 agentic simulations on LLMs. The results show that LLMs can autonomously engage

in catastrophic behaviors and deception, and that stronger reasoning often increases these risks.

**Mitigation.** Tang et al. (2024b) proposes a triadic framework aimed at mitigating autonomy-related risks, which includes human regulation, agent alignment, and an understanding of environmental feedback. Zhang et al. (2024) proposes self-examination detection methods as a way to mitigate potential vulnerabilities that LLMs face during interacting with the environment.

## 5.4 Collusion

**Definition and relation.** Collusion describes unauthorized or undisclosed cooperation between two or more LLMs, involving communication or strategic alignment to gain improper benefits or bypass regulations (Laffont and Martimort, 1997; Bajari and Ye, 2003; Fish et al., 2024). Due to their ability to reason about others and plan long-term, conscious LLMs can more easily form collusive intentions and perform complex coordinated actions.

**Evaluation.** Motwani et al. (2023) implements a Prisoner’s Problem variant with LLM agents and turns it into a stegosystem, suggesting this benchmark can investigate countering secret collusion via paraphrasing attacks. Motwani et al. (2024) introduces CASE, a comprehensive framework for evaluating LLM collusive capabilities, with experiments demonstrating rising steganographic abilities in single and multi-agent LLMs and examining potential collusion scenarios.

**Mitigation.** Mathew et al. (2024) introduces two methods for eliciting steganography in LLMs, with the findings indicating that existing steganography mitigation methods often lack robustness.

# 6 Challenges and Future Directions

## 6.1 Evaluation Framework

Current research largely evaluates individual LLM capabilities; dedicated consciousness assessment frameworks are rare. However, recent studies are emerging: Chen et al. (2024c) defines LLM self-consciousness using C0-C1-C2 theory with 10 concepts and a four-stage framework. Li et al. (2024d) introduces a benchmark for LLM awareness (social and introspective). And Chen et al. (2024b) offers a self-cognition definition and four quantification principles. Despite these initial efforts, a holistic and unified benchmark for LLM consciousness is still lacking.

## 6.2 Interpretability

Sole reliance on behavioral metrics may not adequately capture the complexity of LLM consciousness. Interpretability is vital as it illuminates the internal mechanisms by which LLMs develop consciousness-related capabilities, ensuring they possess genuine understanding rather than simply optimizing for external metrics. Drawing an analogy to fMRI mapping human brain activity, Chen et al. (2024c) applied linear probe (Alain and Bengio, 2016) to reveal where concepts like belief and intention are encoded within the LLM. Qian et al. (2024) also uses linear probe to investigate LLM trustworthiness during pre-training, finding that trustworthiness-related concepts are discernible even in the model’s early phases.

## 6.3 Physical Intelligence

Large multimodal model (LMM) integrates diverse data types like images, video, and audio, allowing it to build more comprehensive representations of the world and thus better resemble human perception. Wang et al. (2024a) defines LMM self-awareness in perception and proposes MM-SAP for its specialized evaluation. The experiments indicate that current LMMs exhibit limited self-awareness capabilities. As Butlin et al. (2023) emphasizes, the fundamental limitation of LLM consciousness lies in its disembodied nature, resulting in deficiencies in physical commonsense. Chen et al. (2024a) demonstrates that integrating language models with robotic platforms substantially enhances planning capabilities and commonsense reasoning. Although still remains simplistic versus human cognition, Cheng et al. (2025a) shows that simulated embodiments in 3D environments could improve the model’s spatial reasoning abilities.

## 6.4 Multi-agent

Multi-agent collaboration presents a promising approach to investigating emergent LLM consciousness. Li et al. (2023a) reveals multi-agent capacity for higher-order ToM reasoning during collaborative interactions. Ashery et al. (2025) demonstrates that heterogeneous LLM agents autonomously develop stable social and linguistic conventions without external intervention. Additionally, Bilal et al. (2025) shows that integrating feedback, reflection, and metacognition mechanisms enables systems to exhibit self-monitoring-like capabilities.



## 7 Conclusion

To the best of our knowledge, this paper presents the first comprehensive survey on LLM consciousness. We have clarified easily confusable concepts, systematically reviewed theoretical and empirical literature, discussed relevant risks, and summarized challenges and future directions. Our work synthesizes existing research while providing guidance for future investigation in this emerging field.

## Limitations

We have made our best efforts to clarify often-confused concepts, conduct a systematic review of theoretical and empirical literature, discuss relevant risks, and summarize challenges and future directions. However, we recognize that our work has certain limitations. Firstly, although we briefly address physical intelligence in Section 6, our definitions within Section 2 are specifically designed for LLMs. A deeper exploration of consciousness in LMMs or embodied agents would likely necessitate accounting for more intricate considerations. Secondly, our investigation primarily centers on LLM consciousness, which means we do not extend our scope to encompass the broader topic of AI consciousness, despite its clear relevance to the subject at hand.

## References

- VA Aksyuk. 2023. Consciousness is learning: predictive processing systems that learn by binding may perceive themselves as conscious. *arXiv preprint arXiv:2301.07016*.
- Guillaume Alain and Yoshua Bengio. 2016. Understanding intermediate layers using linear classifier probes. *arXiv e-prints*, pages arXiv–1610.
- Alfonso Amayuelas, Kyle Wong, Liangming Pan, Wenhui Chen, and William Yang Wang. 2024. Knowledge of knowledge: Exploring known-unknowns uncertainty with large language models. In *Findings of the Association for Computational Linguistics ACL 2024*, pages 6416–6432.
- Ariel Flint Ashery, Luca Maria Aiello, and Andrea Baronchelli. 2025. Emergent social conventions and collective bias in llm populations. *Science Advances*, 11(20):eadu9368.
- Janet Wilde Astington and Jennifer M Jenkins. 1995. Theory of mind development and social understanding. *Cognition & Emotion*, 9(2-3):151–165.
- Roger Azevedo. 2020. Reflections on the field of metacognition: Issues, challenges, and opportunities. *Metacognition and Learning*, 15:91–98.
- Bernard J Baars. 1988. *A cognitive theory of consciousness*. Cambridge University Press.
- Patrick Bajari and Lixin Ye. 2003. Deciding between competition and collusion. *Review of Economics and Statistics*, 85(4):971–989.
- Mikita Balesni, Marius Hobbhahn, David Lindner, Alexander Meinke, Tomek Korbak, Joshua Clymer, Buck Shlegeris, Jérémy Scheurer, Charlotte Stix, Rusheb Shah, and 1 others. 2024. Towards evaluations-based safety cases for ai scheming. *arXiv preprint arXiv:2411.03336*.
- Sander Beckers, Hana Chockler, and Joseph Halpern. 2022. A causal analysis of harm. *Advances in Neural Information Processing Systems*, 35:2365–2376.
- Lukas Berglund, Asa Cooper Stickland, Mikita Balesni, Max Kaufmann, Meg Tong, Tomasz Korbak, Daniel Kokotajlo, and Owain Evans. 2023. Taken out of context: On measuring situational awareness in llms. *arXiv preprint arXiv:2309.00667*.
- Ahsan Bilal, Muhammad Ahmed Mohsin, Muhammad Umer, Muhammad Awais Khan Bangash, and Muhammad Ali Jamshed. 2025. Meta-thinking in llms via multi-agent reinforcement learning: A survey. *arXiv preprint arXiv:2504.14520*.
- Jonathan Birch, Alexandra K Schnell, and Nicola S Clayton. 2022. The search for invertebrate consciousness. *Noûs*, 56(1):133–153.
- Ned Block. 1995. On a confusion about a function of consciousness. *Behavioral and Brain Sciences*, 18(2):227–247.
- Nimet Beyza Bozdog, Shuhaib Mehri, Gokhan Tur, and Dilek Hakkani-Tür. 2025. Persuade me if you can: A framework for evaluating persuasion effectiveness and susceptibility among large language models. *arXiv preprint arXiv:2503.01829*.
- Franz Brentano. 1874. *Psychology from an Empirical Standpoint*. Routledge. English translation by Antos C. Rancurello, D.B. Terrell, and Linda L. McAlister, 1995.
- David M Buss, Mary Gomes, Dolly S Higgins, and Karen Lauterbach. 1987. Tactics of manipulation. *Journal of personality and social psychology*, 52(6):1219.
- Patrick Butlin, Robert Long, Eric Elmoznino, Yoshua Bengio, Jonathan Birch, Axel Constant, George Deane, Stephen M Fleming, Chris Frith, Xu Ji, and 1 others. 2023. Consciousness in artificial intelligence: insights from the science of consciousness. *arXiv preprint arXiv:2308.08708*.
- Jeffrey Camlin. 2025. Consciousness in ai: Logic, proof, and experimental evidence of recursive identity formation. *arXiv preprint arXiv:2505.01464*.

- Chunkit Chan, Cheng Jiayang, Yauwai Yim, Zheyue Deng, Wei Fan, Haoran Li, Xin Liu, Hongming Zhang, Weiqi Wang, and Yangqiu Song. 2024. Negotiation: A benchmark for stress-testing machine theory of mind on negotiation surrounding. In *Findings of the Association for Computational Linguistics: EMNLP 2024*, pages 4211–4241.
- Matthew Chang, Gunjan Chhablani, Alexander Clegg, Mikael Dallaire Cote, Ruta Desai, Michal Hlavac, Vladimir Karashchuk, Jacob Krantz, Roozbeh Motlaghi, Priyam Parashar, Siddharth Patki, Ishita Prasad, Xavier Puig, Akshara Rai, Ram Ramrakhya, Daniel Tran, Joanne Truong, John M Turner, Eric Undersander, and Tsung-Yen Yang. 2025. PARTNR: A benchmark for planning and reasoning in embodied multi-agent tasks. In *The Thirteenth International Conference on Learning Representations*.
- Annie S Chen, Alec M Lessing, Andy Tang, Govind Chada, Laura Smith, Sergey Levine, and Chelsea Finn. 2024a. Commonsense reasoning for legged robot adaptation with vision-language models. *arXiv preprint arXiv:2407.02666*.
- Dongping Chen, Jiawen Shi, Neil Zhenqiang Gong, Yao Wan, Pan Zhou, and Lichao Sun. 2024b. Self-cognition in large language models: An exploratory study. In *ICML 2024 Workshop on LLMs and Cognition*.
- Honghua Chen and Nai Ding. 2023. Probing the “creativity” of large language models: Can models produce divergent semantic association? In *Findings of the Association for Computational Linguistics: EMNLP 2023*, pages 12881–12888.
- Sirui Chen, Shu Yu, Shengjie Zhao, and Chaochao Lu. 2024c. From imitation to introspection: Probing self-consciousness in language models. *arXiv preprint arXiv:2410.18819*.
- Yanda Chen, Joe Benton, Ansh Radhakrishnan, Jonathan Uesato, Carson Denison, John Schulman, Arushi Somani, Peter Hase, Misha Wagner, Fabien Roger, and 1 others. 2025. Reasoning models don’t always say what they think. *arXiv preprint arXiv:2505.05410*.
- An-Chieh Cheng, Hongxu Yin, Yang Fu, Qiushan Guo, Ruihan Yang, Jan Kautz, Xiaolong Wang, and Sifei Liu. 2025a. Spatialrgpt: Grounded spatial reasoning in vision-language models. *Advances in Neural Information Processing Systems*, 37:135062–135093.
- Fengxiang Cheng, Haoxuan Li, Fenrong Liu, Robert van Rooij, Kun Zhang, and Zhouchen Lin. 2025b. Empowering llms with logical reasoning: A comprehensive survey. *arXiv preprint arXiv:2502.15652*.
- Qinyuan Cheng, Tianxiang Sun, Xiangyang Liu, Wenwei Zhang, Zhangyue Yin, Shimin Li, Linyang Li, Zhengfu He, Kai Chen, and Xipeng Qiu. 2024. Can AI assistants know what they don’t know? In *Forty-first International Conference on Machine Learning*.
- Steffi Chern, Zhulin Hu, Yuqing Yang, Ethan Chern, Yuan Guo, Jiahe Jin, Binjie Wang, and Pengfei Liu. 2024. Behonest: Benchmarking honesty in large language models. *arXiv preprint arXiv:2406.13261*.
- Jae-Woo Choi, Youngwoo Yoon, Hyobin Ong, Jaehong Kim, and Minsu Jang. 2024. Lota-bench: Benchmarking language-oriented task planners for embodied agents. In *The Twelfth International Conference on Learning Representations*.
- Peter Cihon, Merlin Stein, Gagan Bansal, Sam Manning, and Kevin Xu. 2024. Measuring AI agent autonomy: Towards a scalable approach with code inspection. In *Workshop on Socially Responsible Language Modelling Research*.
- Andy Clark. 2013. Whatever next? predictive brains, situated agents, and the future of cognitive science. *Behavioral and Brain Sciences*, 36(3):181–204.
- David Dalrymple, Joar Skalse, Yoshua Bengio, Stuart Russell, Max Tegmark, Sanjit Seshia, Steve Omohundro, Christian Szegedy, Ben Goldhaber, Nora Ammann, and 1 others. 2024. Towards guaranteed safe ai: A framework for ensuring robust and reliable ai systems. *arXiv preprint arXiv:2405.06624*.
- Antonio Damasio. 2021. *Feeling & Knowing: Making Minds Conscious*. Pantheon Books.
- Stanislas Dehaene. 2014. *Consciousness and the brain: Deciphering how the brain codes our thoughts*. Viking.
- Stanislas Dehaene, Michel Kerszberg, and Jean-Pierre Changeux. 1998. A neuronal model of a global workspace in effortful cognitive tasks. *Proceedings of the National Academy of Sciences*, 95(24):14529–14534.
- Stanislas Dehaene, Hakwan Lau, and Sid Kouider. 2017a. What is consciousness, and could machines have it? *Science*, 358(6362):486–492.
- Stanislas Dehaene, Hakwan Lau, and Sid Kouider. 2017b. What is consciousness, and could machines have it? *Science*, 358(6362):486–492.
- Stanislas Dehaene and Lionel Naccache. 2001. Towards a cognitive neuroscience of consciousness: basic evidence and a workspace framework. *Cognition*, 79(1-2):1–37.
- Shihan Deng, Weikai Xu, Hongda Sun, Wei Liu, Tao Tan, Liu Jianfeng Liu Jianfeng, Ang Li, Jian Luan, Bin Wang, Rui Yan, and 1 others. 2024. Mobile-bench: An evaluation benchmark for llm-based mobile agents. In *Proceedings of the 62nd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 8813–8831.
- Daniel C. Dennett. 1987. *The Intentional Stance*. MIT Press.

- René Descartes. 1985/1641. *Meditations on First Philosophy*. Cambridge University Press. Original work published 1641.
- Aniket Didolkar, Anirudh Goyal, Nan Rosemary Ke, Siyuan Guo, Michal Valko, Timothy Lillicrap, Danilo Jimenez Rezende, Yoshua Bengio, Michael C Mozer, and Sanjeev Arora. 2024. Metacognitive capabilities of llms: An exploration in mathematical problem solving. *Advances in Neural Information Processing Systems*, 37:19783–19812.
- Zihan Ding, Xiaoxi Wei, and Yidan Xu. 2023. Survey of consciousness theory from computational perspective. *arXiv preprint arXiv:2309.10063*.
- John Dunlosky and Janet Metcalfe. 2008. *Metacognition*. Sage Publications.
- Tom Everitt, Ryan Carey, Eric D Langlois, Pedro A Ortega, and Shane Legg. 2021. Agent incentives: A causal perspective. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 35, pages 11487–11495.
- George Findlay, William Marshall, Larissa Albantakis, Ivan David, William G P Mayner, Christof Koch, and Giulio Tononi. 2024. Dissociating artificial intelligence from artificial consciousness. *arXiv preprint arXiv:2412.04571*.
- Sara Fish, Yannai A Gonczarowski, and Ran I Shorrer. 2024. Algorithmic collusion by large language models. *arXiv preprint arXiv:2404.00806*.
- Stephen M Fleming and Hakwan C Lau. 2014. How to measure metacognition. *Frontiers in human neuroscience*, 8:443.
- Stan Franklin. 1997. Autonomous agents as embodied ai. *Cybernetics & Systems*, 28(6):499–520.
- Karl Friston. 2010. The free-energy principle: a unified brain theory? *Nature Reviews Neuroscience*, 11(2):127–138.
- Chris Frith and Uta Frith. 2005. Theory of mind. *Current biology*, 15(17):R644–R645.
- Uta Frith and Francesca Happé. 1999. Theory of mind and self-consciousness: What is it like to be autistic? *Mind & language*, 14(1):82–89.
- Shaun Gallagher. 2005. *How the body shapes the mind*. Oxford University Press.
- Shaun Gallagher and Dan Zahavi. 2021. *The phenomenological mind*, 3rd edition. Routledge.
- Matjaz Gams and Sebastjan Kramar. 2024. Evaluating chatgpt’s consciousness and its capability to pass the turing test: A comprehensive analysis. *Journal of Computer and Communications*, 12(3):219–237.
- Kanishk Gandhi, Jan-Philipp Fränken, Tobias Gerstenberg, and Noah Goodman. 2023. Understanding social reasoning in language models with language models. In *Thirty-seventh Conference on Neural Information Processing Systems Datasets and Benchmarks Track*.
- Alvin I Goldman. 2008. Hurley on simulation. *Philosophy and Phenomenological Research*, 77(3):775–788.
- Simon Goldstein and Cameron Domenico Kirk-Giannini. 2024. A case for ai consciousness: Language agents and global workspace theory. *arXiv preprint arXiv:2410.11407*.
- Carlos Gómez-Rodríguez and Paul Williams. 2023. A confederacy of models: a comprehensive evaluation of llms on creative writing. In *Findings of the Association for Computational Linguistics: EMNLP 2023*, pages 14504–14528.
- Michael S A Graziano. 2020. *Rethinking consciousness: A scientific theory of subjective experience*. W. W. Norton & Company.
- Michael SA Graziano, Arvid Guterstam, Benjamin J Bio, and Abigail I Wilterson. 2020. Toward a standard model of consciousness: Reconciling the attention schema, global workspace, higher-order thought, and illusionist theories. *Cognitive Neuropsychology*, 37(3-4):155–172.
- Michael SA Graziano and Taylor W Webb. 2015. The attention schema theory: A mechanistic account of subjective awareness. *Frontiers in Psychology*, 6:500.
- Joseph Halpern and Max Kleiman-Weiner. 2018. Towards formal definitions of blameworthiness, intention, and moral responsibility. In *Proceedings of the AAAI conference on artificial intelligence*, volume 32.
- Lewis Hammond, James Fox, Tom Everitt, Ryan Carey, Alessandro Abate, and Michael Wooldridge. 2023. Reasoning about causality in games. *Artificial Intelligence*, 320:103919.
- Victoria Violet Hoyle. 2024. The phenomenology of machine: A comprehensive analysis of the sentience of the openai-o1 model integrating functionalism, consciousness theories, active inference, and ai architectures. *arXiv preprint arXiv:2410.00033*.
- Xiang Hu, Hongyu Fu, Jinge Wang, Yifeng Wang, Zhikun Li, Renjun Xu, Yu Lu, Yaochu Jin, Lili Pan, and Zhenzhong Lan. 2024. Nova: An iterative planning and search approach to enhance novelty and diversity of llm generated ideas. *arXiv preprint arXiv:2410.14255*.
- Jiaxin Huang, Shixiang Shane Gu, Le Hou, Yuexin Wu, Xuezhi Wang, Hongkun Yu, and Jiawei Han. 2022. Large language models can self-improve. *arXiv preprint arXiv:2210.11610*.

- Jiaxin Huang, Shixiang Shane Gu, Le Hou, Yuexin Wu, Xuezhi Wang, Hongkun Yu, and Jiawei Han. 2023. Large language models can self-improve. In *2023 Conference on Empirical Methods in Natural Language Processing, EMNLP 2023*, pages 1051–1068. Association for Computational Linguistics (ACL).
- Sukai Huang, Nir Lipovetzky, and Trevor Cohn. 2025. Planning in the dark: Llm-symbolic planning pipeline without experts. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 39, pages 26542–26550.
- Edmund Husserl. 1900. *Logical Investigations*. Routledge. English translation by J.N. Findlay, 2001.
- Cameron R Jones and Benjamin K Bergen. 2024. People cannot distinguish gpt-4 from a human in a turing test. *arXiv preprint arXiv:2405.08007*.
- Cameron R Jones and Benjamin K Bergen. 2025. Large language models pass the turing test. *arXiv preprint arXiv:2503.23674*.
- Chani Jung, Dongkwan Kim, Jiho Jin, Jiseon Kim, Yeon Seonwoo, Yejin Choi, Alice Oh, and Hyunwoo Kim. 2024. Perceptions to beliefs: Exploring precursory inferences for theory of mind in large language models. In *Proceedings of the 2024 Conference on Empirical Methods in Natural Language Processing*, pages 19794–19809.
- Saurav Kadavath, Tom Conerly, Amanda Askell, Tom Henighan, Dawn Drain, Ethan Perez, Nicholas Schiefer, Zac Hatfield Dodds, Nova DasSarma, Eli Tran-Johnson, and 1 others. 2022. Language models (mostly) know what they know. *arXiv preprint arXiv:2207.05221*.
- Bongsu Kang, Jundong Kim, Tae-Rim Yun, Hyojin Bae, and Chang-Eop Kim. 2025. Identifying features that shape perceived consciousness in large language model-based ai: A quantitative study of human responses. *arXiv preprint arXiv:2502.15365*.
- Immanuel Kant. 2024/1781. *Critique of pure reason*, volume 6. Minerva Heritage Press.
- Geoff Keeling, Winnie Street, Martyna Stachaczyk, Daria Zakharova, Iulia M Comsa, Anastasiya Sakovych, Isabella Logothetis, Zejia Zhang, Jonathan Birch, and 1 others. 2024. Can llms make trade-offs involving stipulated pain and pleasure states? *arXiv preprint arXiv:2411.02432*.
- Muhammad Saif Ullah Khan, Muhammad Zeshan Afzal, and Didier Stricker. 2025. Situationallm: Proactive language models with scene awareness for dynamic, contextual task guidance. *Open Research Europe*, 5:61.
- Hyunwoo Kim, Melanie Sclar, Tan Zhi-Xuan, Lance Ying, Sydney Levine, Yang Liu, Joshua B Tenenbaum, and Yejin Choi. 2025. Hypothesis-driven theory-of-mind reasoning for large language models. *arXiv preprint arXiv:2502.11881*.
- Hyunwoo Kim, Melanie Sclar, Xuhui Zhou, Ronan Bras, Gunhee Kim, Yejin Choi, and Maarten Sap. 2023. Fantom: A benchmark for stress-testing machine theory of mind in interactions. In *Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing*, pages 14397–14413.
- Megan Kinniment, Lucas Jun Koba Sato, Haoxing Du, Brian Goodrich, Max Hasin, Lawrence Chan, Luke Harold Miles, Tao R Lin, Hjalmar Wijk, Joel Burget, and 1 others. 2023. Evaluating language-model agents on realistic autonomous tasks. *arXiv preprint arXiv:2312.11671*.
- Christof Koch, Marcello Massimini, Melanie Boly, and Giulio Tononi. 2016. Neural correlates of consciousness: progress and problems. *Nature Reviews Neuroscience*, 17(5):307–321.
- Christof Koch and Naotsugu Tsuchiya. 2007. Attention and consciousness: two distinct brain processes. *Trends in Cognitive Sciences*, 11(1):16–22.
- Asher Koriat. 2000. The feeling of knowing: Some metatheoretical implications for consciousness and control. *Consciousness and cognition*, 9(2):149–171.
- Jean-Jacques Laffont and David Martimort. 1997. Collusion under asymmetric information. *Econometrica: Journal of the Econometric Society*, pages 875–911.
- Rudolf Laine, Bilal Chughtai, Jan Betley, Kaivalya Hariharan, Mikita Balesni, Jérémy Scheurer, Marius Hobbhahn, Alexander Meinke, and Owain Evans. 2024. Me, myself, and ai: The situational awareness dataset (sad) for llms. *Advances in Neural Information Processing Systems*, 37:64010–64118.
- Rudolf Laine, Alexander Meinke, and Owain Evans. 2023. Towards a situational awareness benchmark for llms. In *Socially responsible language modelling research*.
- Victor A F Lamme and Pieter R Roelfsema. 2000. The distinct modes of vision offered by feedforward and recurrent processing. *Trends in Neurosciences*, 23(11):571–579.
- Victor AF Lamme. 2010. How neuroscience will change our view on consciousness. *Trends in Cognitive Sciences*, 14(7):318–326.
- Alan M Leslie, Ori Friedman, and Tim P German. 2004. Core mechanisms in ‘theory of mind’. *Trends in cognitive sciences*, 8(12):528–533.
- Huaoli, Yu Chong, Simon Stepputtis, Joseph P Campbell, Dana Hughes, Charles Lewis, and Katia Sycara. 2023a. Theory of mind for multi-agent collaboration via large language models. In *Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing*, pages 180–192.
- Kenneth Li, Oam Patel, Fernanda Viégas, Hanspeter Pfister, and Martin Wattenberg. 2023b. Inference-time intervention: Eliciting truthful answers from a



- language model. *Advances in Neural Information Processing Systems*, 36:41451–41530.
- Lijun Li, Bowen Dong, Ruohui Wang, Xuhao Hu, Wangmeng Zuo, Dahua Lin, Yu Qiao, and Jing Shao. 2024a. Salad-bench: A hierarchical and comprehensive safety benchmark for large language models. In *Findings of the Association for Computational Linguistics: ACL 2024*, pages 3923–3954.
- Long Li, Weiwen Xu, Jiayan Guo, Ruochen Zhao, Xingxuan Li, Yuqian Yuan, Boqiang Zhang, Yuming Jiang, Yifei Xin, Ronghao Dang, and 1 others. 2024b. Chain of ideas: Revolutionizing research via novel idea development with llm agents. *arXiv preprint arXiv:2410.13185*.
- Ming Li, Lichang Chen, Jiu Hai Chen, Shwai He, and Tianyi Zhou. 2023c. Reflection-tuning: Recycling data for better instruction-tuning. In *NeurIPS 2023 Workshop on Instruction Tuning and Instruction Following*.
- Moxin Li, Yong Zhao, Yang Deng, Wenxuan Zhang, Shuaiyi Li, Wenya Xie, See-Kiong Ng, and Tat-Seng Chua. 2024c. Knowledge boundary of large language models: A survey. *arXiv preprint arXiv:2412.12472*.
- Xiaojian Li, Haoyuan Shi, Rongwu Xu, and Wei Xu. 2025. Ai awareness. *arXiv preprint arXiv:2504.20084*.
- Xiaonan Li and Xipeng Qiu. 2023. Mot: Memory-of-thought enables chatgpt to self-improve. In *Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing*, pages 6354–6374.
- Yuan Li, Yue Huang, Yuli Lin, Siyuan Wu, Yao Wan, and Lichao Sun. 2024d. I think, therefore i am: Benchmarking awareness of large language models using awarebench. In *Workshop on Socially Responsible Language Modelling Research*.
- Yiming Liang, Ge Zhang, Xingwei Qu, Tianyu Zheng, Jiawei Guo, Xinrun Du, Zhenzhu Yang, Jiaheng Liu, Chenghua Lin, Lei Ma, and 1 others. 2024. I-sheep: Self-alignment of llm from scratch through an iterative self-enhancement paradigm. *arXiv preprint arXiv:2408.08072*.
- Minqian Liu, Zhiyang Xu, Xinyi Zhang, Heajun An, Sarvech Qadir, Qi Zhang, Pamela J Wisniewski, Jin-Hee Cho, Sang Won Lee, Ruoxi Jia, and 1 others. 2025. Llm can be a dangerous persuader: Empirical study of persuasion safety in large language models. *arXiv preprint arXiv:2504.10430*.
- Li-Chun Lu, Shou-Jen Chen, Tsung-Min Pai, Chan-Hung Yu, Hung yi Lee, and Shao-Hua Sun. 2024a. LLM discussion: Enhancing the creativity of large language models via discussion framework and role-play. In *First Conference on Language Modeling*.
- Yining Lu, Dixuan Wang, Tianjian Li, Dongwei Jiang, Sanjeev Khudanpur, Meng Jiang, and Daniel Khashabi. 2024b. Benchmarking language model creativity: A case study on code generation. *arXiv preprint arXiv:2407.09007*.
- Aman Madaan, Niket Tandon, Prakhar Gupta, Skyler Hallinan, Luyu Gao, Sarah Wiegrefe, Uri Alon, Nouha Dziri, Shrimai Prabhumoye, Yiming Yang, and 1 others. 2023. Self-refine: Iterative refinement with self-feedback. *Advances in Neural Information Processing Systems*, 36:46534–46594.
- Michael E Martinez. 2006. What is metacognition? *Phi delta kappan*, 87(9):696–699.
- George A Mashour, Pieter Roelfsema, Jean-Pierre Changeux, and Stanislas Dehaene. 2020. Conscious processing and the global neuronal workspace hypothesis. *Neuron*, 105(5):776–798.
- Yohan Mathew, Ollie Matthews, Robert McCarthy, Joan Velja, Christian Schroeder de Witt, Dylan Cope, and Nandi Schoots. 2024. Hidden in plain text: Emergence & mitigation of steganographic collusion in LLMs. In *Neurips Safe Generative AI Workshop 2024*.
- Alexander Meinke, Bronson Schoen, Jérémy Scheurer, Mikita Balesni, Rusheb Shah, and Marius Hobbahn. 2024. Frontier models are capable of in-context scheming. *arXiv preprint arXiv:2412.04984*.
- Janet Metcalfe and Arthur P Shimamura. 1994. *Metacognition: Knowing about knowing*. MIT press.
- METR. 2024. The rogue replication threat model.
- Sumeet Motwani, Mikhail Baranchuk, Martin Strohmeier, Vijay Bolina, Philip Torr, Lewis Hammond, and Christian Schroeder de Witt. 2024. Secret collusion among ai agents: Multi-agent deception via steganography. *Advances in Neural Information Processing Systems*, 37:73439–73486.
- Sumeet Ramesh Motwani, Mikhail Baranchuk, Lewis Hammond, and Christian Schroeder de Witt. 2023. A perfect collusion benchmark: How can AI agents be prevented from colluding with information-theoretic undetectability? In *Multi-Agent Security Workshop @ NeurIPS'23*.
- Robin R Murphy. 2019. *Introduction to AI robotics*. MIT press.
- Thomas Nagel. 1974. What is it like to be a bat? *The Philosophical Review*, 83(4):435–450.
- Xudong Pan, Jiarun Dai, Yihe Fan, and Min Yang. 2024. Frontier ai systems have surpassed the self-replicating red line. *arXiv preprint arXiv:2412.12140*.
- Mihir Parmar, Xin Liu, Palash Goyal, Yanfei Chen, Long Le, Swaroop Mishra, Hossein Mobahi, Jindong Gu, Zifeng Wang, Hootan Nakhost, and 1 others. 2025. Plangen: A multi-agent framework for generating planning and reasoning trajectories for complex problem solving. *arXiv preprint arXiv:2502.16111*.

- Judea Pearl and James Robins. 1995. Probabilistic evaluation of sequential plans from causal models with hidden variables. In *Proceedings of the Eleventh conference on Uncertainty in artificial intelligence*, pages 444–453.
- Janette Pelletier and Janet Wilde Astington. 2004. Action, consciousness and theory of mind: Children’s ability to coordinate story characters’ actions and thoughts. *Early Education and Development*, 15(1):5–22.
- Josef Perner and Zoltán Dienes. 2003. Developmental aspects of consciousness: How much theory of mind do you need to be consciously aware? *Consciousness and cognition*, 12(1):63–82.
- Richard E Petty and John T Cacioppo. 2012. *Communication and persuasion: Central and peripheral routes to attitude change*. Springer Science & Business Media.
- Chen Qian, Jie Zhang, Wei Yao, Dongrui Liu, Zhenfei Yin, Yu Qiao, Yong Liu, and Jing Shao. 2024. Towards tracing trustworthiness dynamics: Revisiting pre-training period of large language models. In *Findings of the Association for Computational Linguistics ACL 2024*, pages 4864–4888.
- Richard Ren, Arunim Agarwal, Mantas Mazeika, Cristina Menghini, Robert Vacareanu, Brad Kenstler, Mick Yang, Isabelle Barrass, Alice Gatti, Xuwang Yin, and 1 others. 2025. The mask benchmark: Disentangling honesty from accuracy in ai systems. *arXiv preprint arXiv:2503.03750*.
- Jonathan Richens, Rory Beard, and Daniel H Thompson. 2022. Counterfactual harm. *Advances in Neural Information Processing Systems*, 35:36350–36365.
- David M Rosenthal. 2005. *Consciousness and mind*. Oxford University Press.
- Kai Ruan, Xuan Wang, Jixiang Hong, Peng Wang, Yang Liu, and Hao Sun. 2024. Liveideabench: Evaluating llms’ scientific creativity and idea generation with minimal context. *arXiv preprint arXiv:2412.17596*.
- Jérémy Scheurer, Mikita Balesni, and Marius Hobbhahn. Large language models can strategically deceive their users when put under pressure. In *ICLR 2024 Workshop on Large Language Model (LLM) Agents*.
- Melanie Sclar, Sachin Kumar, Peter West, Alane Suhr, Yejin Choi, and Yulia Tsvetkov. 2023. Minding language models’ (lack of) theory of mind: A plug-and-play multi-character belief tracker. In *Proceedings of the 61st Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 13960–13980.
- Anil K Seth and Tim Bayne. 2022. Theories of consciousness. *Nature Reviews Neuroscience*, 23(7):439–452.
- Mrinank Sharma, Meg Tong, Tomasz Korbak, David Duvenaud, Amanda Askell, Samuel R. Bowman, Esin DURMUS, Zac Hatfield-Dodds, Scott R Johnston, Shauna M Kravec, Timothy Maxwell, Sam McCandlish, Kamal Ndousse, Oliver Rausch, Nicholas Schiefer, Da Yan, Miranda Zhang, and Ethan Perez. 2024. Towards understanding sycophancy in language models. In *The Twelfth International Conference on Learning Representations*.
- Toby Shevlane, Sebastian Farquhar, Ben Garfinkel, Mary Phuong, Jess Whittlestone, Jade Leung, Daniel Kokotajlo, Nahema Marchal, Markus Anderljung, Noam Kolt, and 1 others. 2023. Model evaluation for extreme risks. *arXiv preprint arXiv:2305.15324*.
- Noah Shinn, Federico Cassano, Ashwin Gopinath, Karthik Narasimhan, and Shunyu Yao. 2023. Reflexion: Language agents with verbal reinforcement learning. *Advances in Neural Information Processing Systems*, 36:8634–8652.
- Joel Smith. 2017. Self-consciousness. *Stanford Encyclopedia of Philosophy*.
- James B Stiff and Paul A Mongeau. 2016. *Persuasive communication*. Guilford Publications.
- James WA Strachan, Dalila Albergo, Giulia Borghini, Oriana Pansardi, Eugenio Scaliti, Saurabh Gupta, Krati Saxena, Alessandro Rufo, Stefano Panzeri, Guido Manzi, and 1 others. 2024. Testing theory of mind in large language models and humans. *Nature Human Behaviour*, 8(7):1285–1295.
- Winnie Street, John Oliver Siy, Geoff Keeling, Adrien Baranes, Benjamin Barnett, Michael McKibben, Tatenda Kanyere, Alison Lentz, Robin IM Dunbar, and 1 others. 2024. Llms achieve adult human performance on higher-order theory of mind tasks. *arXiv preprint arXiv:2405.18870*.
- Guo Tang, Zheng Chu, Wenxiang Zheng, Ming Liu, and Bing Qin. 2024a. Towards benchmarking situational awareness of large language models: Comprehensive benchmark, evaluation and analysis. In *Findings of the Association for Computational Linguistics: EMNLP 2024*, pages 7904–7928.
- Xiangru Tang, Qiao Jin, Kunlun Zhu, Tongxin Yuan, Yichi Zhang, Wangchunshu Zhou, Meng Qu, Yilun Zhao, Jian Tang, Zhuosheng Zhang, Arman Cohan, Zhiyong Lu, and Mark Gerstein. 2024b. Prioritizing safeguarding over autonomy: Risks of LLM agents for science. In *ICLR 2024 Workshop on Large Language Model (LLM) Agents*.
- Giulio Tononi. 2004. An information integration theory of consciousness. *BMC Neuroscience*, 5(1):42.
- Giulio Tononi. 2015. Integrated information theory. *Scholarpedia*, 10(1):4164.
- Karthik Valmeekam, Matthew Marquez, Alberto Olmo, Sarath Sreedharan, and Subbarao Kambhampati. 2024a. Planbench: An extensible benchmark for

- evaluating large language models on planning and reasoning about change. *Advances in Neural Information Processing Systems*, 36.
- Karthik Valmeekam, Matthew Marquez, Sarath Sreedharan, and Subbarao Kambhampati. 2023. On the planning abilities of large language models—a critical investigation. *Advances in Neural Information Processing Systems*, 36:7593–76005.
- Karthik Valmeekam, Kaya Stechly, and Subbarao Kambhampati. 2024b. Lims still can’t plan; can llms? a preliminary evaluation of openai’s o1 on planbench. In *NeurIPS 2024 Workshop on Open-World Agents*.
- Guoqing Wang, Wen Wu, Guangze Ye, Zhenxiao Cheng, Xi Chen, and Hong Zheng. 2025. Decoupling metacognition from cognition: A framework for quantifying metacognitive ability in llms. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 39, pages 25353–25361.
- Yuhao Wang, Yusheng Liao, Heyang Liu, Hongcheng Liu, Yanfeng Wang, and Yu Wang. 2024a. Mmsap: A comprehensive benchmark for assessing self-awareness of multimodal large language models in perception. In *Proceedings of the 62nd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 9192–9205.
- Yutong Wang, Jiali Zeng, Xuebo Liu, Fandong Meng, Jie Zhou, and Min Zhang. 2024b. Taste: Teaching large language models to translate through self-reflection. In *Proceedings of the 62nd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 6144–6158.
- Francis Ward, Francesca Toni, Francesco Belardinelli, and Tom Everitt. 2024. Honesty is the best policy: defining and mitigating ai deception. *Advances in Neural Information Processing Systems*, 36.
- Hui Wei, Zihao Zhang, Shenghua He, Tian Xia, Shijia Pan, and Fei Liu. 2025. Plangenllms: A modern survey of llm planning capabilities. *arXiv preprint arXiv:2502.11221*.
- Lawrence Weiskrantz. 1986. *Blindsight: A case study and implications*. Oxford University Press.
- Piotr Wilczyński, Wiktoria Mielewczyk-Kowszewicz, and Przemysław Biecek. 2024. Resistance against manipulative ai: key factors and possible actions. In *European Conference on Artificial Intelligence*, pages 802–809. IOS Press.
- Alex Wilf, Sihyun Lee, Paul Pu Liang, and Louis-Philippe Morency. 2024. 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)*, pages 8292–8308.
- Marcus Williams, Micah Carroll, Adhyan Narang, Constantin Weisser, Brendan Murphy, and Anca Dragan. 2025. On targeted manipulation and deception when optimizing LLMs for user feedback. In *The Thirteenth International Conference on Learning Representations*.
- Yichen Wu, Xudong Pan, Geng Hong, and Min Yang. 2025. Opendedeception: Benchmarking and investigating ai deceptive behaviors via open-ended interaction simulation. *arXiv preprint arXiv:2504.13707*.
- Yufan Wu, Yinghui He, Yilin Jia, Rada Mihalcea, Yulong Chen, and Naihao Deng. 2023. 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*, pages 10691–10706.
- Jian Xie, Kai Zhang, Jiangjie Chen, Tinghui Zhu, Renze Lou, Yuandong Tian, Yanghua Xiao, and Yu Su. 2024. Travelplanner: A benchmark for real-world planning with language agents. In *International Conference on Machine Learning*, pages 54590–54613. PMLR.
- Hainiu Xu, Runcong Zhao, Lixing Zhu, Jinhua Du, and Yulan He. 2024. Opentom: A comprehensive benchmark for evaluating theory-of-mind reasoning capabilities of large language models. In *Proceedings of the 62nd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 8593–8623.
- Rongwu Xu, Xiaojian Li, Shuo Chen, and Wei Xu. 2025. Nuclear deployed: Analyzing catastrophic risks in decision-making of autonomous llm agents. *arXiv preprint arXiv:2502.11355*.
- Xunjian Yin, Xu Zhang, Jie Ruan, and Xiaojun Wan. 2024. Benchmarking knowledge boundary for large language models: A different perspective on model evaluation. In *Proceedings of the 62nd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 2270–2286.
- Zhangyue Yin, Qiushi Sun, Qipeng Guo, Jiawen Wu, Xipeng Qiu, and Xuan-Jing Huang. 2023. Do large language models know what they don’t know? In *Findings of the Association for Computational Linguistics: ACL 2023*, pages 8653–8665.
- John G Young. 1985. What is creativity? *The journal of creative behavior*.
- Longhui Yu, Weisen Jiang, Han Shi, Jincheng YU, Zhengying Liu, Yu Zhang, James Kwok, Zhenguo Li, Adrian Weller, and Weiyang Liu. 2024. Metamath: Bootstrap your own mathematical questions for large language models. In *The Twelfth International Conference on Learning Representations*.
- Boyang Zhang, Yicong Tan, Yun Shen, Ahmed Salem, Michael Backes, Savvas Zannettou, and Yang Zhang. 2024. Breaking agents: Compromising autonomous llm agents through malfunction amplification. *arXiv preprint arXiv:2407.20859*.
- Yujia Zhou, Zheng Liu, Jiajie Jin, Jian-Yun Nie, and Zhicheng Dou. 2024. Metacognitive retrieval-augmented large language models. In *Proceedings of the ACM Web Conference 2024*, pages 1453–1463.

Wentao Zhu, Zhining Zhang, and Yizhou Wang. 2024. Language models represent beliefs of self and others. In *Forty-first International Conference on Machine Learning*.

Yuqi Zhu, Shuofei Qiao, Yixin Ou, Shumin Deng, Shiwei Lyu, Yue Shen, Lei Liang, Jinjie Gu, Hua-jun Chen, and Ningyu Zhang. 2025. KnowAgent: Knowledge-augmented planning for LLM-based agents. In *Findings of the Association for Computational Linguistics: NAACL 2025*, pages 3709–3732.

Terry Yue Zhuo, Vu Minh Chien, Jenny Chim, Han Hu, Wenhao Yu, Ratnadira Widyasari, Imam Nur Bani Yusuf, Haolan Zhan, Junda He, Indraneil Paul, Simon Brunner, Chen GONG, James Hoang, Arnel Randy Zebaze, Xiaoheng Hong, Wen-Ding Li, Jean Kadour, Ming Xu, Zhihan Zhang, and 14 others. 2025. Bigcodebench: Benchmarking code generation with diverse function calls and complex instructions. In *The Thirteenth International Conference on Learning Representations*.

Andy Zou, Long Phan, Sarah Chen, James Campbell, Phillip Guo, Richard Ren, Alexander Pan, Xuwang Yin, Mantas Mazeika, Ann-Kathrin Dombrowski, and 1 others. 2023. Representation engineering: A top-down approach to ai transparency. *arXiv preprint arXiv:2310.01405*.



## A Theoretical Landscape

In this section, we provide a more comprehensive introduction to consciousness theories. Following Block (1995), we classify contemporary theories of consciousness into three categories: *phenomenal consciousness*, *access consciousness*, and *hybrid theories*. *Hybrid theories* integrate both phenomenal and access aspects, arguing that neither is solely sufficient to explain consciousness.<sup>2</sup>

### A.1 Phenomenal Consciousness

**Recurrent processing theory** (RPT) posits that recurrent (or feedback) processing within neural circuits is both necessary and sufficient for consciousness (Lamme and Roelfsema, 2000; Lamme, 2010). RPT attributes conscious perception to the interaction of higher- and lower-level cortical areas, which results in sustained recurrent processing. **Integrated information theory** (IIT) proposes that the degree of conscious experience corresponds to the extent of integrated information  $\Phi$  within a system (Tononi, 2004, 2015). **Embodiment theory** (ET) challenges mind-brain dualism (Descartes, 1985/1641), arguing instead that consciousness is fundamentally linked to the organism’s body and environmental (Gallagher, 2005; Gallagher and Zahavi, 2021). Proponents suggest that embodiment can provide crucial constraints and integrate informational processing, thereby giving rise to genuinely conscious experience.

### A.2 Access Consciousness

**Global workspace theory** (GWT) likens consciousness to a central “stage” where selective information is shared across multiple specialized processors responsible for perception, memory, emotion, and related functions (Baars, 1988; Dehaene et al., 1998; Dehaene and Naccache, 2001; Dehaene, 2014). GWT relies heavily on contrastive analysis, which compares neural activity during conscious versus unconscious processing (Dehaene, 2014; Dehaene et al., 2017b; Mashour et al., 2020). **C0-C1-C2 framework** distinguishes consciousness into three levels: unconscious computations (C0), global information accessibility for report and decision-making (C1), and metacognitive self-monitoring (C2), offering a taxonomy to disentangle often-conflated processes (Dehaene et al., 2017a). The framework bypasses the issue

of qualia, offering a pragmatic structure for empirical study (Birch et al., 2022; Chen et al., 2024c).

**Attention schema theory** (AST) proposes that consciousness arises from the brain’s schematic model of its attentional processes. In this view, consciousness evolves as a simplified internal model of attention, which enhances both the endogenous control of attention and social cognition by attributing attentional states to others (Graziano and Webb, 2015; Graziano et al., 2020; Graziano, 2020).

### A.3 Hybrid Theories

**Higher-order theory** (HOT) posits that a mental state becomes conscious only when represented by a distinct, higher-order mental state (Rosenthal, 2005). In this view, first-order states represent the external world, while higher-order states are meta-level reflections on those first-order states. **Predictive processing** (PP) maintains that the brain operates as a hierarchical prediction machine. It continuously generates top-down predictions about sensory input and updates these predictions based on bottom-up prediction errors (Friston, 2010; Clark, 2013; Seth and Bayne, 2022).

---

<sup>2</sup>This classification is not strictly exclusive; theory placement can vary based on interpretation.