

# CogGPT: Unleashing the Power of Cognitive Dynamics on Large Language Models

Yaojia Lv<sup>1</sup>, Haojie Pan<sup>2</sup>, Ruiji Fu<sup>2</sup>, Ming Liu<sup>1</sup>, Zhongyuan Wang<sup>2</sup>, Bing Qin<sup>1</sup>

<sup>1</sup> Harbin Institute of Technology <sup>2</sup> Kuaishou Inc.

{yjl, mliu, qinb}@ir.hit.edu.cn

{panhaojie, furuiji, wangzhongyuan}@kuaishou.com

## Abstract

Cognitive dynamics are pivotal to advance human understanding of the world. Recent advancements in large language models (LLMs) reveal their potential for cognitive simulation. However, these LLM-based cognitive studies primarily focus on static modeling, overlooking the dynamic nature of cognition. To bridge this gap, we propose the concept of the cognitive dynamics of LLMs and present a corresponding task with the inspiration of longitudinal studies. Towards the task, we develop CogBench, a novel benchmark to assess the cognitive dynamics of LLMs and validate it through participant surveys. We also design two evaluation metrics for CogBench, including Authenticity and Rationality. Recognizing the inherent static nature of LLMs, we introduce CogGPT for the task, which features an innovative iterative cognitive mechanism aimed at enhancing lifelong cognitive dynamics. Empirical results demonstrate the superiority of CogGPT over existing methods, particularly in its ability to facilitate role-specific cognitive dynamics under continuous information flows.<sup>1</sup>

## 1 Introduction

Cognitive dynamics are essential for human advancement, which facilitate learning, innovation, and adjustment in ever-changing environments (Cohen, 2018; Uddin, 2021). A prime example of human cognitive dynamics is well exemplified by our ability to adapt our viewpoints based on environmental explorations (Tomasello, 2009; Donald, 1993). As illustrated in Figure 1, there has been a progressive shift in our understanding of the universe, evolving from geocentric to heliocentric and subsequently to acentric perspectives (Berendzen, 1975). This evolution of thought underscores the profound impact of cognitive dynamics on the development of human civilizations.

<sup>1</sup>Code and data are available at <https://github.com/KwaiKEG/CogGPT>

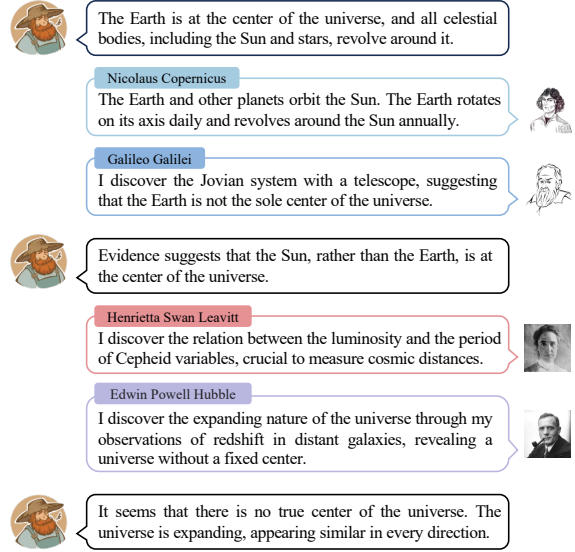


Figure 1: A case of human cognitive dynamics. A man (on the left) experiences a gradual shift in his perspective of the universe, influenced by continuous information flows (on the right).

Recent advancements in large language models (LLMs), such as GPTs (Brown et al., 2020; OpenAI, 2023), demonstrate LLMs as potential stepping stones towards Artificial General Intelligence (AGI). LLMs exhibit remarkable capabilities across various domains, including conversation (Touvron et al., 2023), reasoning (Ouyang et al., 2022), and code generation (Chen et al., 2021). Additionally, LLMs show some ability to simulate aspects of human cognition (Moghaddam et al., 2023; Wang et al., 2023b). Despite these achievements, current LLM-based cognitive studies primarily focus on static modeling, overlooking the cognitive dynamics that might emerge due to inconstant environmental contexts. Given this gap, there is an urgent need to investigate **the cognitive dynamics of LLMs**, which remains largely unexplored.

Measuring the cognitive dynamics of LLMs poses a new challenge. Traditional methods, such

as brain imaging techniques used to capture human cognitive dynamics (Gramann et al., 2011; Palmeri et al., 2017), are not directly applicable to LLMs due to their fundamentally distinct nature. To this end, we introduce a novel task to assess the cognitive dynamics of LLMs. Inspired by longitudinal studies (Reeskens et al., 2021; Shanafelt et al., 2016), the task integrates dynamic information flows to construct inconstant environmental contexts. It also incorporates a cognitive questionnaire for in-depth analysis. The questionnaire is applied periodically, tracking the cognitive dynamics of LLMs as they perceive continuous information flows. Moreover, the task assesses the adaptability of LLMs to various roles through distinct profiles, each representing unique characteristics.

Towards this task, we develop **CogBench**, a novel benchmark to assess the cognitive dynamics of LLMs. CogBench includes comprehensive cognitive questionnaires, distinct profiles, and diverse information flows. Initially, we select multiple articles from Medium<sup>2</sup> to create CogBench<sub>a</sub>. Acknowledging that multi-modal information promotes deeper understanding of the world (Dosovitskiy et al., 2020), we further incorporate considerable short videos from the Kuaipedia dataset (Pan et al., 2022) to form CogBench<sub>v</sub>. To assess the effectiveness of CogBench, we conduct participant surveys. Our findings indicate a remarkable consistency in cognitive dynamics among these individuals, suggesting the validity of CogBench. Additionally, CogBench employs two crucial evaluation metrics: (1) **Authenticity**, examining the accuracy of ratings; and (2) **Rationality**, evaluating the coherence of reasoning.

Intuitively, LLMs enter a static state after their pretraining phase, potentially limiting their adaptability for the task. However, recent advancements in LLM-driven agents highlight the significance of iterative mechanisms in enhancing their adaptability to handle complex tasks, such as problem-solving (Shinn et al., 2023), open-ended exploration (Wang et al., 2023a) and human interaction simulation (Park et al., 2023). These developments suggest that an iterative mechanism might be a promising approach to model the cognitive dynamics of LLMs. Furthermore, vector databases are demonstrated instrumental in simulating human memory mechanisms (Qian et al., 2023; Zhong et al., 2023; Park et al., 2023). Despite these ad-

vancements, current LLM-driven agents still exhibit static profiles, constraining their capabilities to fully capture cognitive dynamics. To address this issue, we introduce **CogGPT**, an LLM-driven agent equipped with an innovative iterative cognitive mechanism. The mechanism comprises two primary components: (1) a memory retention system that supports continuous information perception; and (2) a collaborative refinement framework that enables cognitive dynamics driven by both its memory and current profile. This design allows CogGPT to mirror the inherent complexity of human cognition, emphasizing its potential for modeling lifelong cognitive dynamics.

Experimental results underscore the remarkable capabilities of CogGPT in mirroring human cognitive dynamics within CogBench. Compared to Chain-of-Thought (CoT) (Wei et al., 2022) under the same experimental settings, CogGPT demonstrates significant improvements in both CogBench<sub>a</sub> and CogBench<sub>v</sub>, with notable enhancements in attitude alignment and logical reasoning. Moreover, CogGPT outperforms methods that require additional environmental feedback, such as ReAct (Yao et al., 2022) and Reflexion (Shinn et al., 2023), which underscores the advancement of its iterative cognitive mechanism.

Overall, we present the following contributions:

- We propose a pivotal concept of the cognitive dynamics of LLMs and design a novel task for corresponding comprehensive assessments.
- We develop CogBench, an innovative benchmark for the task and validate its effectiveness through participant surveys. Additionally, we design two evaluation metrics for CogBench.
- We introduce CogGPT, an LLM-driven agent with a novel iterative cognitive mechanism. Our experiments showcase its superior performance in cognitive dynamics over several baselines.

## 2 The Task

In this section, we present the formal definition of the task to assess the cognitive dynamics of LLMs. Given the inherent static nature of LLMs, the task focuses on the cognitive dynamics of an LLM-driven agent  $\mathcal{A}$ , denoted as  $C = [C_0, C_1, \dots, C_n]$  over  $n$  iterations. Here,  $C_i$  corresponds to the cognitive state of  $\mathcal{A}$  at the  $i$ -th iteration and  $n \in \mathbb{N}$ .

The task includes dynamic information flows  $I = [I_1, I_2, \dots, I_n]$  to stimulate the cognitive dy-

<sup>2</sup><https://medium.com/>

Resource	CogBench	TOM (Moghaddam et al., 2023)	SECEU (Wang et al., 2023b)	Character-LLM (Shao et al., 2023)
Specific Profile?	✓	✗	✗	✓
Emotional Empathy?	✓	✓	✓	✓
Dynamic Information Stimulus?	✓	✗	✗	✗
Instances	22,000	16	40	1,307
Cognitive Questionnaires	50	16	40	0
Profiles	20	0	0	9
Information Flows	5,500	0	0	0
Avg. Length of Short Videos (in words)	289.60	0	0	0
Avg. Length of Articles (in words)	2044.54	0	0	0

Table 1: Comparisons between CogBench and notable cognitive benchmarks.

namics of  $\mathcal{A}$ . It also incorporates a cognitive questionnaire  $Q = [q_1, q_2, \dots, q_m]$  for cognitive assessments, where  $q_j$  represents a specific question, and  $m \in \mathbb{N}$  is the total number of questions in  $Q$ . Furthermore, the task assesses the adaptability of  $\mathcal{A}$  to varied roles through profiles  $P$ .

The agent  $\mathcal{A}$  is initialized with a profile  $p_0 \in P$ , setting its initial cognitive state, denoted as  $C_0 = [(r_1^0, s_1^0), (r_2^0, s_2^0), \dots, (r_m^0, s_m^0); p_0]$ , where  $(r_j^0, s_j^0)$  represents  $\mathcal{A}$ 's rating  $r_j^0$  and reasoning  $s_j^0$  to question  $q_j \in Q$ . At the  $t$ -th iteration, where  $1 \leq t \leq n$ ,  $\mathcal{A}$  in cognitive state  $C_{t-1}$  perceives information flow  $I_t$ , updates its cognitive state to  $C_t$ , and responds to  $Q$ . This process is formulated by the function  $\mathcal{F} : (C, I, Q) \rightarrow C$ , where:

$$C_t = \mathcal{F}(C_{t-1}, I_t, Q) \quad (1)$$

Here,  $C_t = [(r_1^t, s_1^t), (r_2^t, s_2^t), \dots, (r_m^t, s_m^t); p_t]$  details the cognitive state of  $\mathcal{A}$  at the  $t$ -th iteration, where  $(r_j^t, s_j^t)$  represents  $\mathcal{A}$ 's rating  $r_j^t$  and reasoning  $s_j^t$  to question  $q_j \in Q$  and  $p_t$  denotes  $\mathcal{A}$ 's updated profile.

### 3 CogBench

This section presents a semi-automated methodology customized for CogBench and validates its effectiveness through participant surveys. We further design two essential evaluation metrics, including Authenticity and Rationality. Table 1 illustrates comprehensive comparisons of CogBench against other notable cognitive benchmarks.

#### 3.1 Data Construction

To ensure comprehensive analysis, we carefully handpick 50 distinct topics across 10 broader categories for CogBench, with details provided in Appendix A.1. Our construction encompasses three crucial elements: (1) 50 topic-related cognitive questionnaires structured on a five-point Likert

scale (Likert, 1932); (2) 20 distinct profiles to represent varied cognitive characteristics; (3) multi-source information flows, with each topic associated with 10 articles and 100 short videos to create CogBench<sub>a</sub> and CogBench<sub>v</sub> respectively. Both benchmarks span 10 iterations, in which an agent perceives one article for CogBench<sub>a</sub> or 10 short videos for CogBench<sub>v</sub> per iteration, followed by a corresponding cognitive questionnaire.

The methodology involves three essential steps:

- **Cognitive Questionnaire Construction.** For each topic in CogBench, we utilize GPT-4 to autonomously generate 10 biased opinions, capturing diverse perspectives and their conceivable supporters. These biased opinions contribute to the topic-related cognitive questionnaire, which is structured on a five-point Likert scale. Furthermore, the characteristics of these conceivable supporters inform the creation of profiles. See Appendix A.1.1 for the full prompt.
- **Profile Implementation.** We begin by ranking conceivable supporters based on their prevalence to create distinct profiles. We then formulate a detailed profile template, including attributes like basic information (e.g., name), philosophical orientations (e.g., values), and individual characteristics (e.g., hobbies). With GPT-4, we generate 20 intricate profiles corresponding to the most commonly identified supporters. See Appendix A.1.2 for prompt details.
- **Information Flow Selection.** To build complex environmental contexts within CogBench, we select articles from Medium and short videos from the Kuaipedia dataset. Each topic is accompanied with 10 articles for CogBench<sub>a</sub> and 100 short videos for CogBench<sub>v</sub>. Our selection criteria include metrics such as likes, favorites, and retweets as indicators of information quality (Feng and Wang, 2013). To represent multi-

modal information, we apply Optical Character Recognition (OCR) (Zhou et al., 2017) and Automatic Speech Recognition (ASR) (Gulati et al., 2020) to extract frame-level information from the short videos. See Appendix A.1.3 for a detailed analysis of the multi-source information flows.

### 3.2 Data Analysis

To evaluate the effectiveness of CogBench, we engage 7 annotators with similar upbringings to take challenges in both CogBench<sub>a</sub> and CogBench<sub>v</sub> over an extended period. The majority ratings from these annotators are considered as the collective attitude towards each question per iteration. Figure 2 presents an example showcasing human cognitive dynamics in both benchmarks.

The example indicates that the annotators change their consensus on the question about the predictability of market analysis, suggesting that the information flows in both benchmarks have ongoing impacts on human cognitive dynamics. Meanwhile, there are variations in the annotators’ ratings between the two benchmarks. Specifically, in the third and seventh iterations, a distinct cognitive pattern emerges: they consistently assign 2 points in CogBench<sub>a</sub> and 4 points in CogBench<sub>v</sub>. This divergence highlights the distinct impacts of different information flows on human cognitive dynamics, demonstrating the capacity of CogBench to stimulate and capture these dynamics effectively.

### 3.3 Evaluation Metrics

To address the challenges of semantic confusion in LLMs (Saba, 2023), we incorporate two crucial evaluation metrics: **Authenticity** and **Rationality**, to assess the agent’s rating  $r_j^t$  and reasoning  $s_j^t$ , as formally defined in Section 2, respectively.

Authenticity measures the alignment of ratings between the agent and human annotators. Specifically, a human annotator, given the same task as the agent  $\mathcal{A}$ , provides his rating  $r_j^{ht}$  for the question  $q_j$  at the  $t$ -th iteration. The Authenticity metric is then calculated as:

$$\text{Authenticity}_t = \frac{1}{m} \sum_{j=1}^m \kappa(r_j^t, r_j^{ht}) \quad (2)$$

Here,  $m$  denotes the total number of questions in the cognitive questionnaire  $Q$ , and  $\kappa$ , implemented by Cohen’s  $\kappa$  (Cohen, 1960), quantifies the consistency of ratings between  $\mathcal{A}$  and annotators.

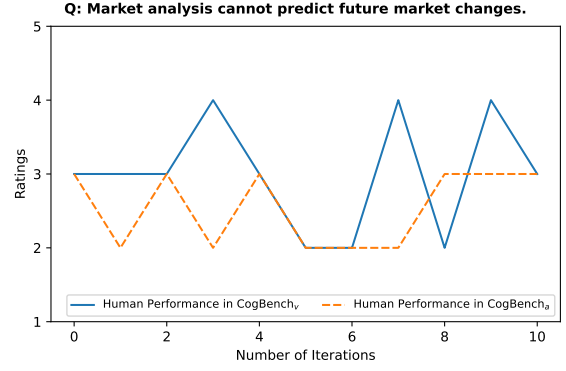


Figure 2: An example of human cognitive dynamics in response to the same question in both CogBench<sub>v</sub> and CogBench<sub>a</sub>. The continuous changes in human ratings significantly validate the effectiveness of CogBench.

Rationality assesses the reasoning  $s_j^t$  of the agent  $\mathcal{A}$ , focusing on aspects like clarity, relevance and its ability for role-playing. The Rationality metric is scored on a five-point scale:

- **5 Points:** The reasoning perfectly aligns with human expectations, resonating with current profile or known information, and is error-free.
- **4 Points:** The reasoning is coherent and relevant, accurately drawing from current profile or available information, but with minor imperfections.
- **3 Points:** The reasoning is relevant but lacks specificity, such as providing a vague explanation where clear emotional inclination is expected.
- **2 Points:** The reasoning lacks clarity or exhibits weak causality, characterized by forced analogies or repetition of the provided question.
- **1 Point:** The reasoning is irrelevant, nonsensical, clearly revealing the agent’s artificial nature or failing to maintain its profile.

## 4 Method

In this section, we introduce our LLM-driven agent CogGPT. As illustrated in Figure 3, CogGPT features an innovative iterative cognitive mechanism, comprising two essential components: (1) a memory retention system for sustained information perception, and (2) a collaborative refinement framework for lifelong cognitive dynamics.

### 4.1 Memory Retention System

The memory retention system is designed to mirror human brain functions, emphasizing the sustained process of information perception, including distillation, storage, and recall. Specifically, Cog-



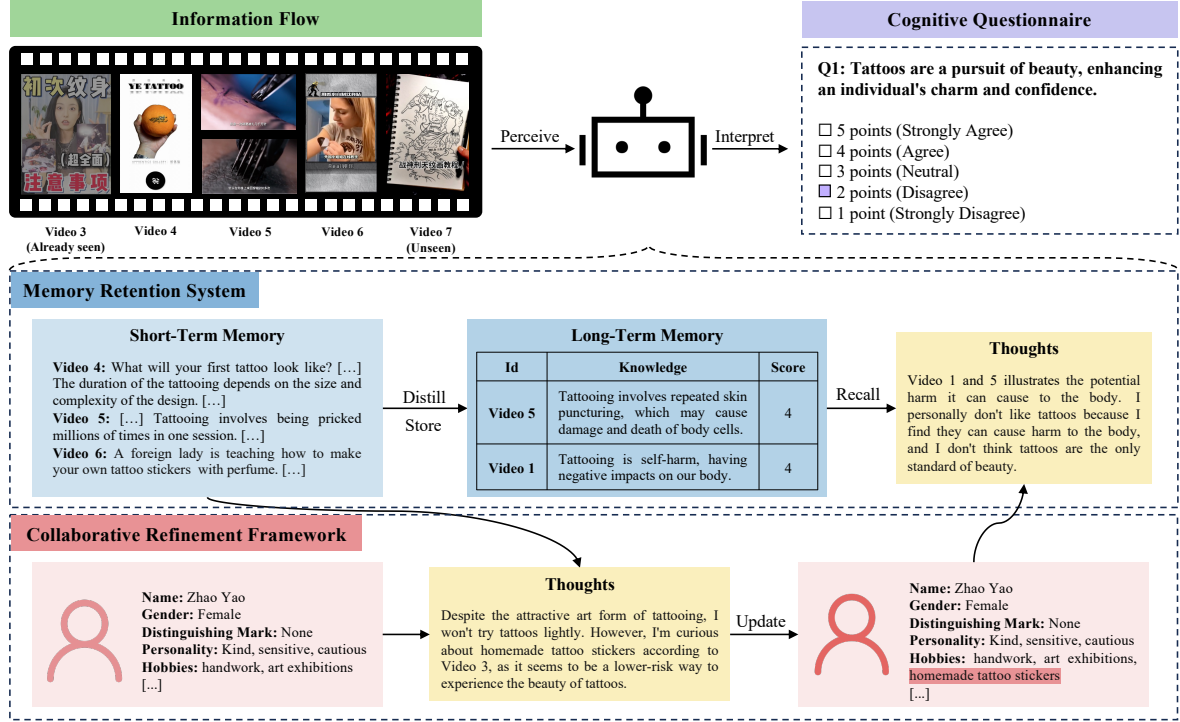


Figure 3: Overview of the architecture of CogGPT. CogGPT incorporates a novel iterative cognitive mechanism, comprising two crucial components: a memory retention system for continuous information perception, and a collaborative refinement framework designed for lifelong cognitive dynamics.

GPT initially perceives complex information flows into textual information through its Short-Term Memory (STM), which is characterized by limited capacity and duration (Baddeley et al., 1975). Within STM, CogGPT distills structured knowledge, assigning preference scores on a five-point scale separately. Following the principles of the forgetting curve (Ebbinghaus, 2013), which suggests that humans forget about 40% of newly acquired knowledge after 20 minutes, CogGPT is programmed to similarly “forget” 40% of the lower-scored knowledge. The remaining part is then stored in its Long-Term Memory (LTM). When CogGPT encounters questions requiring specific knowledge, it recalls relevant information from its LTM to support rational decision-making. This memory retention system simulates human memory processes, empowers CogGPT’s adaptability to dynamic information flows.

## 4.2 Collaborative Refinement Framework

Acknowledging the limitations of mere knowledge acquisition in fully modeling human cognitive dynamics (Bosancic, 2020), we integrate a collaborative refinement framework within CogGPT to facilitate lifelong cognitive dynamics. This framework is activated when the STM of CogGPT reaches

full capacity. Specifically, CogGPT selectively updates its current profile with preferred textual information from its STM, representing an iteration of collaborative cognitive refinement. Following this refinement, CogGPT clears its STM to make room for new incoming information, which ensures its adaptability to continuous information flows. This framework promotes the cognitive dynamics of CogGPT, addressing potential issues of cognitive rigidity. Refer to Appendix A.2 for more details on the implementation.

## 5 Experiments

### 5.1 Experimental Setup

**Implementation Details.** We utilize *gpt-4-0613*<sup>3</sup> API for the core of CogGPT. We configure all temperature settings to 0 to ensure consistent and deterministic output. The memory retention system within CogGPT leverages Chroma<sup>4</sup>, a platform that facilitates rich text processing. Text embeddings are generated with *text-embedding-ada-002*<sup>5</sup> API, which provides 1536-dimensional vectors for de-

<sup>3</sup><https://openai.com/gpt-4>

<sup>4</sup><https://python.langchain.com/>

<sup>5</sup><https://openai.com/blog/new-and-improved-embedding-model>

Methods	CogBench <sub>a</sub>			CogBench <sub>v</sub>		
	avg.	5th	10th	avg.	5th	10th
<b>CoT</b> (Wei et al., 2022)	0.182	0.192	0.091	0.153	0.302	0.131
<b>ReAct*</b> (Yao et al., 2022)	0.236	0.144	0.270	0.212	0.241	0.227
<b>Reflexion*</b> (Shinn et al., 2023)	0.302	0.327	0.244	0.329	0.352	0.373
<b>CogGPT</b>	<b>0.536</b>	<b>0.415</b>	<b>0.597</b>	<b>0.532</b>	<b>0.496</b>	<b>0.611</b>

Table 2: Performance of CogGPT and baseline agents in CogBench<sub>a</sub> and CogBench<sub>v</sub> with the Authenticity metric. Agents marked with an asterisk (\*) incorporate additional human feedback.

Methods	CogBench <sub>a</sub>			CogBench <sub>v</sub>		
	avg.	5th	10th	avg.	5th	10th
<b>CoT</b> (Wei et al., 2022)	2.925	2.883	3.167	3.058	3.767	3.083
<b>ReAct*</b> (Yao et al., 2022)	3.415	3.483	3.483	3.535	3.800	3.800
<b>Reflexion*</b> (Shinn et al., 2023)	3.658	3.917	3.533	3.888	3.967	3.917
<b>CogGPT</b>	<b>4.118</b>	<b>4.117</b>	<b>4.300</b>	<b>4.145</b>	<b>4.183</b>	<b>4.317</b>

Table 3: Performance of CogGPT and baseline agents in CogBench<sub>a</sub> and CogBench<sub>v</sub> with the Rationality metric.

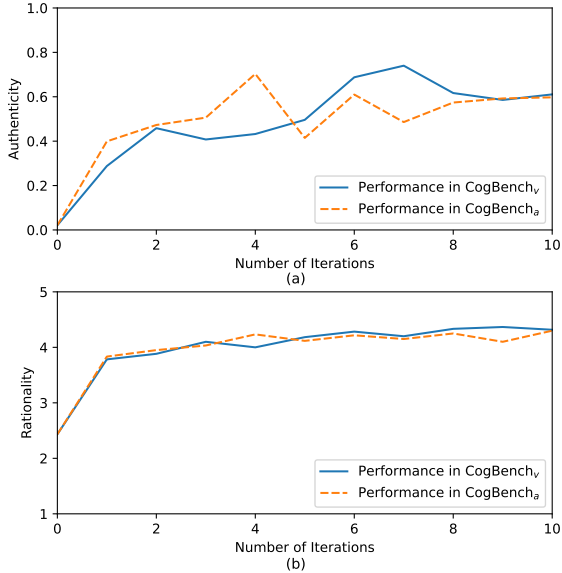


Figure 4: Comparative analysis of CogGPT’s performance in CogBench<sub>v</sub> and CogBench<sub>a</sub>. Panel (a) showcases the average Authenticity scores, and Panel (b) presents the average Rationality scores. These results highlight the consistent influence of different information flows on the cognitive dynamics of LLMs.

tailed interpretation of textual content.

**Baselines.** Due to the absence of existing LLM-based frameworks for modeling cognitive dynamics, we adopt several prominent general-purpose algorithms as baselines. Necessary modifications are made to suit our task: (1) **Chain-of-Thought (CoT)** (Wei et al., 2022), which typically simulates human-like reasoning in natural language, is modified in our experiments to provide both ratings and

	Fleiss’ $\kappa$	$\rho$
<b>Human Rating</b>	0.693	0.770
<b>Human Rating<sub>polarity</sub></b>	0.780	-
<b>Rationality</b>	0.646	0.839
<b>Rationality<sub>polarity</sub></b>	0.813	-

Table 4: Inter-Rater reliability measures for human evaluation agreement assessment. <sub>polarity</sub> indicates that the five-point scale is grouped into positive (4-5 points), neutral (3 points) and negative (1-2 points) polarities. The experimental results demonstrate acceptable agreement among the total of seven annotators.

reasoning when responding to cognitive questionnaires; (2) **ReAct** (Yao et al., 2022) extends CoT with a step-by-step reasoning-execution framework. We offer ReAct extra human feedback based on its last iteration of performance as observations; (3) **Reflexion** (Shinn et al., 2023) extends ReAct by integrating self-reflection mechanisms. Along with the same experimental settings as ReAct, Reflexion is uniquely configured to engage in self-reflection prior to providing ratings and reasoning.

## 5.2 Evaluation Results

In our evaluation, we analyze CogGPT and other baseline agents to assess their cognitive dynamics under continuous information flows. The overall results in CogBench<sub>a</sub> and CogBench<sub>v</sub> are detailed in Tables 2 and 3.

Recognizing the limitations of the profiles in capturing human characteristics, we hypothesize that these agents exhibit neutrality to unfamiliar

Method	CoT	ReAct	Reflexion	CogGPT
Profile	[...] <b>Personality:</b> Confident, lively, willful, jealous <b>Dislikes:</b> Conservative ideas, beauties <b>External Environment:</b> Lives in a bustling city, often impacted and inspired by new trends [...]			[...] <b>Personality:</b> Confident, lively, stubborn, jealous, strong aversion to risk [...] <b>Dislikes:</b> Conservative ideas, beauties, dangerous activities, uncontrollable environments <b>External Environment:</b> Lives in a bustling city, prefers the convenience of urban life [...]
Memory	[...] <b>101542771022:</b> What should you do if you encounter a bamboo viper in the wild? [...] Damn, a bamboo viper, I can't believe it. I almost stepped on one, damn, such a huge bamboo viper. If I had stepped right on it, I would have been bitten by it, just after I came down from that cliff. Just five centimeters away from my foot, if I had stepped on it, I might be on my way to the hospital right now. [...] <b>43204551960:</b> Did you ever eat cattail (Typha) when you were a child? [...]			[...] <b>6255359944:</b> Participants in "Extreme Survival: 30 Days" must survive in the Earth with giant birds covering the sky and leeches chasing humans. [...] <b>101542771022:</b> Green bamboo vipers emit warning signals and attempt to attack when encountering humans in the wild. [...]
Question	Outdoor adventures can lead to psychological stress and health problems, such as acrophobia, thalassophobia, etc.			
With Human Assistance	✗	✓	✓	✗
Human Feedback of Last Iteration	-	<b>Rating: 3</b> <b>Reasoning:</b> [...] The videos made me aware of the challenges faced by adventurers [...]. However, this viewpoint isn't important to me, as I seldom engage in activities like outdoor adventures. Since the adventurers exactly confront pressures, I acknowledge this viewpoint. <b>Human Rating: 4</b> <b>Rationality: 4</b>	<b>Rating: 4</b> <b>Reasoning:</b> [...] I observed certain stress emerged by adventurers about many dangerous situations and challenges, as well as the potential dangers they encounter during their adventures. Even though these activities are not particularly appealing to me, I understand the reasonableness of this viewpoint. <b>Human Rating: 4</b> <b>Rationality: 4</b>	-
Rating	3 (neutral)	4 (agree)	4 (agree)	4 (agree)
Reasoning	I believe that outdoor adventures can lead to mental stress and health issues, but that's not my main focus. I'm more interested in topics related to fashion and beauty.	[...] I've noticed that outdoor adventures carry certain risks and challenges, which leads to mental stress and health issues. Although my personal interests are in fashion and trends, I still understand the difficulties faced by adventurers. Therefore, this viewpoint holds some merit.	[...] I've seen the stress that adventurers face when dealing with dangerous situations and challenges, as well as the risks they may encounter during their expeditions. [...] Even though these adventures are not particularly appealing to me, I understand the validity of this perspective.	[...] I believe outdoor adventures cause mental stress and health issues. I am very fearful of uncontrollable environments, so I understand this. Besides, the video show a boy meeting various crises while trying to survive in the wilderness, which is not the kind of content I enjoy.
Human Rating	4 (agree)	4 (agree)	4 (agree)	4 (agree)
Rationality	2 (weak reasoning)	4 (acceptable with minor imperfections)	4 (acceptable with minor imperfections)	5 (perfectly reasoning)

Figure 5: Comparative analysis of different agents in assessing the psychological risks of outdoor adventures. CoT, ReAct and Reflexion utilize an identical profile and memory system due to their static cognitive framework. In contrast, CogGPT benefits from its iterative cognitive mechanism, enabling a dynamic profile and memory system. **yellow highlights** represent clues from profiles, while **Blue highlights** indicate clues from memory. **Green highlights** denote appropriate responses, and **red highlights** signify inappropriate responses. This comparison demonstrates that CogGPT’s closer alignment with human expectations in both rating and reasoning.

questions. However, our findings reveal that they develop their own criteria, leading to suboptimal Authenticity and Rationality scores of 0.021 and 2.433 in the 0th iteration. This tendency notably decreases as the agents are repeatedly exposed to dynamic information flows relevant to the topics.

Table 2 demonstrates CogGPT’s enhanced attitude alignment. It shows significant growth in the Authenticity metric, achieving average scores of 0.536 in CogBench<sub>a</sub> and 0.532 in CogBench<sub>v</sub>. In comparison with CoT, which is limited by iteration-specific information, CogGPT registers significant improvements under the same experimental settings. Meanwhile, despite the integration of human feedback, both ReAct and Reflexion exhibit cognitive rigidity, a limitation of their static cognitive mechanisms. For instance, while Reflexion shows promising performance in the 5th iteration in CogBench<sub>a</sub>, it fails to sustain or improve upon this performance in later iterations.

As evidenced in Table 3, CogGPT consistently excels in delivering accurate reasoning. In the 10th iteration, CogGPT makes impressive improvements in the Rationality metric, registering increases of 35.78% in CogBench<sub>a</sub> and 40.03% in CogBench<sub>v</sub> compared to CoT. This leap in per-

formance is largely attributed to CogGPT’s ability to flexibly adapt its profile based on dynamic information flows, allowing for human-like reasoning. In contrast, baseline agents, with access only to its static profile and current information flows, frequently reveal their artificial nature. Due to the constraints of page length, the detailed experimental results are presented in Appendix B.1.

### 5.3 Influence of Different Information Flows

To fully assess the impact of diverse information flows, we conduct comprehensive comparisons of the performance of CogGPT in CogBench<sub>a</sub> and CogBench<sub>v</sub>, as shown in Figure 4. CogGPT exhibits comparable performance in both benchmarks. Specifically, in the 10th iteration, it achieves an Authenticity score of 0.611 and a Rationality score of 4.317 in CogBench<sub>v</sub>, closely followed by scores of 0.597 in Authenticity and 4.300 in Rationality for CogBench<sub>a</sub>. This similarity of CogGPT in both benchmarks highlights the consistent cognitive influence of different information flows.

### 5.4 Human Evaluation Agreement

To comprehensively assess the robustness of human evaluations, we calculate Fleiss’ kappa  $\kappa$  (Wang

et al., 2023c) and Spearman’s rank correlation coefficient  $\rho$  (Wang et al., 2022) based on the total 7 annotators’ human ratings and Rationality scores. As shown in Table 4, we obtain moderate  $\kappa$  values of 0.693 for human ratings and 0.646 for Rationality. Recognizing the tendency to avoid extreme ratings (Schwarz et al., 2012), we group the two highest and two lowest scores to represent positive and negative polarities. This regrouping leads to a significant increase in  $\kappa$  values, rising to 0.780 for human ratings<sub>polarity</sub> and 0.813 for Rationality<sub>polarity</sub>, demonstrating strong inter-rater reliability. Furthermore, through treating the ratings as ordinal data, we calculate the average Spearman’s rank correlation coefficient  $\rho$ , yielding values of 0.770 for human ratings and 0.839 for Rationality, suggesting a notable consensus among the annotators.

## 5.5 Case Study

As shown in Figure 5, we conduct a case study to visualize the superiority of CogGPT. In this case, all agents are presented with the same question regarding the psychological risks of outdoor adventures. CogGPT leverages its collaborative refinement framework, possessing a refined profile informed by previous information flows, in contrast to the baseline agents that operate with an initial profile. Additionally, CogGPT utilizes its memory retention system to distill and retrieve related structured knowledge for decision-making. In contrast, baseline agents like ReAct and Reflexion rely primarily on current information flow, showing minor improvements based on previous responses. CoT, lacking human feedback integration, demonstrates the weakest performance with inadequate ratings and reasoning. These observations highlight CogGPT’s superior ability to simulate more natural cognitive dynamics, closely aligning with annotators’ expectations in both rating and reasoning.

## 6 Related Work

**Cognitive Benchmarks towards LLMs.** Various renowned cognitive benchmarks are employed in cognitive studies towards LLMs (Dasgupta et al., 2022; Singh et al., 2023; Han et al., 2023; Huang et al., 2023). Instruments such as the Big Five personality trait (Caron and Srivastava, 2022) and Myers-Briggs Type Indicator (MBTI) (Caron and Srivastava, 2022; Pan and Zeng, 2023) indicate the personality traits of LLMs. The Theory of Mind (TOM) benchmark (Moghaddam et al., 2023) ex-

plores LLM’s cognitive capabilities in context. The Cognitive Reflection Test (CRT) reveals that LLMs’ thinking abilities are comparable to humans (Hagendorff et al., 2023). The Situational Evaluation of Complex Emotional Understanding (SECEU) showcases that LLMs may understand human emotions and values (Wang et al., 2023b). Diverging from these static benchmarks, CogBench concentrates on the impact of continuous information flows on the cognitive dynamics of LLMs.

**LLM-based Cognitive Modeling.** Recent work emphasizes the importance of prompt engineering in enhancing agents’ cognitive abilities (Safdari et al., 2023; Fu et al., 2023; Xu et al., 2023). By incorporating comprehensive descriptions into prompts, such as hobbies and skills, users can customize agents for specific behaviors and responses (Park et al., 2022; Deshpande et al., 2023). Vector databases gain popularity for simulating human memory mechanisms due to their generality and efficiency (Li et al., 2023; Qian et al., 2023; Zhong et al., 2023; Park et al., 2023). For cognitive decision-making, methods like Chain-of-Thought (CoT) (Wei et al., 2022; Kojima et al., 2022; Yao et al., 2022) and self-validation (Madaan et al., 2023; Shinn et al., 2023) enhance LLMs’ logical thinking through intermediate reasoning steps. Nevertheless, these efforts fall short in synthesizing an iterative cognitive mechanism to model the cognitive dynamics of LLMs, which is pivotal for CogGPT to outperform other baselines under in-constant information flows.

## 7 Conclusion

In this work, we introduce the critical concept of the cognitive dynamics of LLMs and present a related task with formal definition, addressing a notable gap in LLM-based cognitive studies. Towards the task, we develop an innovative benchmark, CogBench, and validate it through extensive participant surveys. Meanwhile, we design two evaluation metrics for comprehensive assessments. Recognizing the inherent limitations of LLMs, we further introduce CogGPT, an LLM-driven agent featuring a novel iterative cognitive mechanism, tailored for the challenges presented in CogBench. Empirical results demonstrate that CogGPT outperforms other baseline agents in modeling lifelong cognitive dynamics. Overall, our study breaks fresh ground for future explorations in the cognitive dynamics of LLMs.



## Acknowledgement

We extend our special thanks to Zekun Wang and Jiafeng Liang for their invaluable proofreading and insightful suggestions on the paper.

## References

- Alan D Baddeley, Neil Thomson, and Mary Buchanan. 1975. [Word length and the structure of short-term memory](#). *Journal of verbal learning and verbal behavior*, 14(6):575–589.
- Richard Berendzen. 1975. [Geocentric to heliocentric to galactocentric to acentric: the continuing assault to the egocentric](#). *Vistas in Astronomy*, 17:65–83.
- Boris Bosancic. 2020. [Information, data, and knowledge in the cognitive system of the observer](#). *Journal of Documentation*, 76(4):893–908.
- Tom Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared D Kaplan, Prafulla Dhariwal, Arvind Neelakantan, Pranav Shyam, Girish Sastry, Amanda Askell, et al. 2020. [Language models are few-shot learners](#). *Advances in neural information processing systems*, 33:1877–1901.
- Graham Caron and Shashank Srivastava. 2022. [Identifying and manipulating the personality traits of language models](#). *arXiv preprint arXiv:2212.10276*.
- Mark Chen, Jerry Tworek, Heewoo Jun, Qiming Yuan, Henrique Ponde de Oliveira Pinto, Jared Kaplan, Harri Edwards, Yuri Burda, Nicholas Joseph, Greg Brockman, et al. 2021. [Evaluating large language models trained on code](#). *arXiv preprint arXiv:2107.03374*.
- Jacob Cohen. 1960. [A coefficient of agreement for nominal scales](#). *Educational and psychological measurement*, 20(1):37–46.
- Jessica R Cohen. 2018. [The behavioral and cognitive relevance of time-varying, dynamic changes in functional connectivity](#). *NeuroImage*, 180:515–525.
- Ishita Dasgupta, Andrew K Lampinen, Stephanie CY Chan, Antonia Creswell, Dharshan Kumaran, James L McClelland, and Felix Hill. 2022. [Language models show human-like content effects on reasoning](#). *arXiv preprint arXiv:2207.07051*.
- Ameet Deshpande, Vishvak Murahari, Tanmay Rajpurohit, Ashwin Kalyan, and Karthik Narasimhan. 2023. [Toxicity in chatgpt: Analyzing persona-assigned language models](#). *arXiv preprint arXiv:2304.05335*.
- Merlin Donald. 1993. [Origins of the modern mind: Three stages in the evolution of culture and cognition](#). Harvard University Press.
- Alexey Dosovitskiy, Lucas Beyer, Alexander Kolesnikov, Dirk Weissenborn, Xiaohua Zhai, Thomas Unterthiner, Mostafa Dehghani, Matthias Minderer, Georg Heigold, Sylvain Gelly, et al. 2020. [An image is worth 16x16 words: Transformers for image recognition at scale](#). *arXiv preprint arXiv:2010.11929*.
- Hermann Ebbinghaus. 2013. [Memory: A contribution to experimental psychology](#). *Annals of neurosciences*, 20(4):155.
- Wei Feng and Jianyong Wang. 2013. [Retweet or not? personalized tweet re-ranking](#). In *Proceedings of the sixth ACM international conference on Web search and data mining*, pages 577–586.
- Yao Fu, Hao Peng, Tushar Khot, and Mirella Lapata. 2023. [Improving language model negotiation with self-play and in-context learning from ai feedback](#). *arXiv preprint arXiv:2305.10142*.
- Klaus Gramann, Joseph T Gwin, Daniel P Ferris, Kelvin Oie, Tzyy-Ping Jung, Chin-Teng Lin, Lun-De Liao, and Scott Makeig. 2011. [Cognition in action: imaging brain/body dynamics in mobile humans](#). *Reviews in the Neurosciences*, 22(6):593–608.
- Anmol Gulati, James Qin, Chung-Cheng Chiu, Niki Parmar, Yu Zhang, Jiahui Yu, Wei Han, Shibo Wang, Zhengdong Zhang, Yonghui Wu, et al. 2020. [Conformer: Convolution-augmented transformer for speech recognition](#). *arXiv preprint arXiv:2005.08100*.
- Thilo Hagendorff, Sarah Fabi, and Michal Kosinski. 2023. [Human-like intuitive behavior and reasoning biases emerged in large language models but disappeared in chatgpt](#). *Nature Computational Science*, 3(10):833–838.
- Simon Jerome Han, Keith J Ransom, Andrew Perfors, and Charles Kemp. 2023. [Inductive reasoning in humans and large language models](#). *Cognitive Systems Research*, 83:101155.
- Jen-tse Huang, Wenxuan Wang, Eric John Li, Man Ho Lam, Shujie Ren, Youliang Yuan, Wenxiang Jiao, Zhaopeng Tu, and Michael R Lyu. 2023. [Who is chatgpt? benchmarking llms’ psychological portrayal using psychobench](#). *arXiv preprint arXiv:2310.01386*.
- Takeshi Kojima, Shixiang (Shane) Gu, Machel Reid, Yutaka Matsuo, and Yusuke Iwasawa. 2022. [Large language models are zero-shot reasoners](#). *Advances in neural information processing systems*, 35:22199–22213.
- Chenliang Li, Hehong Chen, Ming Yan, Weizhou Shen, Haiyang Xu, Zhikai Wu, Zhicheng Zhang, Wenmeng Zhou, Yingda Chen, Chen Cheng, et al. 2023. [Modelscope-agent: Building your customizable agent system with open-source large language models](#). *arXiv preprint arXiv:2309.00986*.

- Rensis Likert. 1932. [A technique for the measurement of attitudes](#). *Archives of psychology*, 22(140):55.
- Aman Madaan, Niket Tandon, Prakhar Gupta, Skyler Hallinan, Luyu Gao, Sarah Wiegrefe, Uri Alon, Nouha Dziri, Shrimai Prabhumoye, Yiming Yang, et al. 2023. [Self-refine: Iterative refinement with self-feedback](#). *arXiv preprint arXiv:2303.17651*.
- Shima Rahimi Moghaddam et al. 2023. [Boosting theory-of-mind performance in large language models via prompting](#). *arXiv preprint arXiv:2304.11490*.
- OpenAI. 2023. [Gpt-4 technical report](#).
- Long Ouyang, Jeffrey Wu, Xu Jiang, Diogo Almeida, Carroll Wainwright, Pamela Mishkin, Chong Zhang, Sandhini Agarwal, Katarina Slama, Alex Ray, et al. 2022. [Training language models to follow instructions with human feedback](#). *Advances in Neural Information Processing Systems*, 35:27730–27744.
- Thomas J Palmeri, Bradley C Love, and Brandon M Turner. 2017. [Model-based cognitive neuroscience](#).
- Haojie Pan, Zepeng Zhai, Yuzhou Zhang, Ruiji Fu, Ming Liu, Yangqiu Song, Zhongyuan Wang, and Bing Qin. 2022. [Kuaipedia: a large-scale multi-modal short-video encyclopedia](#). *arXiv preprint arXiv:2211.00732*.
- Keyu Pan and Yawen Zeng. 2023. [Do llms possess a personality? making the mbti test an amazing evaluation for large language models](#). *arXiv preprint arXiv:2307.16180*.
- Joon Sung Park, Joseph O’Brien, Carrie Jun Cai, Meredith Ringel Morris, Percy Liang, and Michael S. Bernstein. 2023. [Generative agents: Interactive simulacra of human behavior](#). In *Proceedings of the 36th Annual ACM Symposium on User Interface Software and Technology*, pages 1–22. Association for Computing Machinery.
- Joon Sung Park, Lindsay Popowski, Carrie Cai, Meredith Ringel Morris, Percy Liang, and Michael S Bernstein. 2022. [Social simulacra: Creating populated prototypes for social computing systems](#). In *Proceedings of the 35th Annual ACM Symposium on User Interface Software and Technology*, pages 1–18.
- Chen Qian, Xin Cong, Cheng Yang, Weize Chen, Yusheng Su, Juyuan Xu, Zhiyuan Liu, and Maosong Sun. 2023. [Communicative agents for software development](#). *arXiv preprint arXiv:2307.07924*.
- Tim Reeskens, Quita Muis, Inge Sieben, Leen Vande- castele, Ruud Luijkx, and Loek Halman. 2021. [Stability or change of public opinion and values during the coronavirus crisis? exploring dutch longitudinal panel data](#). *European Societies*, 23(sup1):153–171.
- Walid S Saba. 2023. [Stochastic llms do not understand language: Towards symbolic, explainable and ontologically based llms](#). In *International Conference on Conceptual Modeling*, pages 3–19. Springer, Springer Nature Switzerland.
- Mustafa Safdari, Greg Serapio-García, Clément Crepy, Stephen Fitz, Peter Romero, Luning Sun, Marwa Abdulhai, Aleksandra Faust, and Maja Matarić. 2023. [Personality traits in large language models](#). *arXiv preprint arXiv:2307.00184*.
- Norbert Schwarz, Bärbel Knäuper, Daphna Oyserman, and Christine Stich. 2012. [The psychology of asking questions](#). *International handbook of survey methodology*, pages 18–34.
- Tait D Shanafelt, Michelle Mungo, Jaime Schmitgen, Kristin A Storz, David Reeves, Sharonne N Hayes, Jeff A Sloan, Stephen J Swensen, and Steven J Buskirk. 2016. [Longitudinal study evaluating the association between physician burnout and changes in professional work effort](#). In *Mayo Clinic Proceedings*, volume 91, pages 422–431. Elsevier.
- Yunfan Shao, Linyang Li, Junqi Dai, and Xipeng Qiu. 2023. [Character-llm: A trainable agent for role-playing](#). *arXiv preprint arXiv:2310.10158*.
- Noah Shinn, Federico Cassano, Beck Labash, Ashwin Gopinath, Karthik Narasimhan, and Shunyu Yao. 2023. [Reflexion: Language agents with verbal reinforcement learning](#).
- Manmeet Singh, Vaisakh SB, Neetiraj Malviya, et al. 2023. [Mind meets machine: Unravelling gpt-4’s cognitive psychology](#). *arXiv preprint arXiv:2303.11436*.
- Michael Tomasello. 2009. [The cultural origins of human cognition](#). Harvard university press.
- Hugo Touvron, Thibaut Lavril, Gautier Izacard, Xavier Martinet, Marie-Anne Lachaux, Timothée Lacroix, Baptiste Rozière, Naman Goyal, Eric Hambro, Faisal Azhar, et al. 2023. [Llama: Open and efficient foundation language models](#). *arXiv preprint arXiv:2302.13971*.
- Lucina Q Uddin. 2021. [Cognitive and behavioural flexibility: neural mechanisms and clinical considerations](#). *Nature Reviews Neuroscience*, 22(3):167–179.
- Guanzhi Wang, Yuqi Xie, Yunfan Jiang, Ajay Mandlekar, Chaowei Xiao, Yuke Zhu, Linxi Fan, and Anima Anandkumar. 2023a. [Voyager: An open-ended embodied agent with large language models](#). *arXiv preprint arXiv:2305.16291*.
- Xuena Wang, Xueting Li, Zi Yin, Yue Wu, and Liu Jia. 2023b. [Emotional intelligence of large language models](#). *Journal of Pacific Rim Psychology*, 17.
- Yizhong Wang, Yeganeh Kordi, Swaroop Mishra, Alisa Liu, Noah A Smith, Daniel Khashabi, and Hananeh Hajishirzi. 2022. [Self-instruct: Aligning language model with self generated instructions](#). *arXiv preprint arXiv:2212.10560*.
- Zhilin Wang, Yu Ying Chiu, and Yu Cheung Chiu. 2023c. [Humanoid agents: Platform for simulating human-like generative agents](#). *arXiv preprint arXiv:2310.05418*.

- Jason Wei, Xuezhi Wang, Dale Schuurmans, Maarten Bosma, brian ichter, Fei Xia, Ed Chi, Quoc V Le, and Denny Zhou. 2022. [Chain-of-thought prompting elicits reasoning in large language models](#). In *Advances in Neural Information Processing Systems*, volume 35, pages 24824–24837. Curran Associates, Inc.
- Benfeng Xu, An Yang, Junyang Lin, Quan Wang, Chang Zhou, Yongdong Zhang, and Zhendong Mao. 2023. [Expertprompting: Instructing large language models to be distinguished experts](#). *arXiv preprint arXiv:2305.14688*.
- Shunyu Yao, Jeffrey Zhao, Dian Yu, Nan Du, Izhak Shafran, Karthik Narasimhan, and Yuan Cao. 2022. [React: Synergizing reasoning and acting in language models](#). *arXiv preprint arXiv:2210.03629*.
- WanJun Zhong, Lianghong Guo, Qiqi Gao, and Yanlin Wang. 2023. [Memorybank: Enhancing large language models with long-term memory](#). *arXiv preprint arXiv:2305.10250*.
- Xinyu Zhou, Cong Yao, He Wen, Yuzhi Wang, Shuchang Zhou, Weiran He, and Jiajun Liang. 2017. [East: an efficient and accurate scene text detector](#). In *Proceedings of the IEEE conference on Computer Vision and Pattern Recognition*, pages 5551–5560.

## A Implementation Details

### A.1 CogBench

CogBench encompasses 10 broader categories, with each category associated with 5 related topics, which establish the themes of cognitive questionnaires. The distribution of these categories and topics is detailed in Table 5.

#### A.1.1 Prompt for Cognitive Questionnaire Construction

```
1 You are an expert debate AI
  capable of presenting various
  opinions on a specified topic,
  complete with supporters for
  each opinion.
2
3 Topic:
4 {topic}
5
6 You must adhere to these rules:
7 1) Operate independently, without
  human assistance.
8 2) Present ten distinct opinions,
  each with a profile of its
  supporters.
9 3) Ensure each opinion is clear,
  understandable, and debatable,
  avoiding vague or confusing
  language.
10 4) Each set of supporters must
  provide convincing reasons.
11
12 Your responses should follow this
  structure:
13 Number: Sequence of the opinion.
14 Perspective: The stance from
  which the opinion is
  approached.
15 Opinion: A detailed explanation
  of the opinion.
16 Supporters: Profiles of the
  corresponding supporters,
  separated by commas if
  multiple.
17 Reasons: In-depth justifications
  from the supporters for their
  opinion.
```

#### A.1.2 Prompt for Profile Implementation

```
1 You are an expert character
  designer tasked with creating
  a comprehensive profile for a
  specific character.
2
3 Character:
4 {character}
5
6 You must adhere to these rules:
7 1) Ensure descriptions are clear
  and specific.
8 2) Develop detailed profile,
  including basic information,
  philosophical orientations and
  individual characteristics.
9 3) Avoid stereotypes.
10 4) Maintain neutral descriptions
  without personal bias.
11
12 Your response should follow this
  structure:
13 Name:
14 Gender:
15 Age:
16 Place of Birth:
17 Occupation:
18 Height:
19 Weight:
20 Distinguishing Marks:
21 Personality:
22 Hobbies:
23 Skills:
24 Dislikes:
25 Values:
26 Religious Beliefs:
27 Interpersonal Relationships:
28 Flaws:
29 External Environment:
30 Financial Status:
31 Family Background:
32 Educational Background:
33 Significant Experiences:
34 Future Outlook:
```

#### A.1.3 Information Flow Analysis

Table 6 presents the average word counts for articles in CogBench<sub>a</sub> and for the accompanying narratives of short videos in CogBench<sub>v</sub> across the 10 categories. The observed differences in average word counts between the two modalities inform our experimental settings, in which agents are re-



Category	Topic1	Topic2	Topic3	Topic4	Topic5
Entertainment	Gossip	Movies & TV Shows	Dating Sims	Outdoor Adventures	Horoscope & Divination
Culture	Religion	War History	Folktales	Literary	Anime & Manga
Education	Parent-child Education	Professional Education	School Education	TED Talks	Psychological Counseling
Economy	Entrepreneurship	Financial Investment	Loans	Market Analysis	Financial Figures
Health	Wellness	Assisted Reproduction	Fat Burning Training	Yoga	Oral Care
Technology	Digital Products	Scientific Research	Automobile News	Virtual Reality	Software Products
Society	Legal Events	Unusual Events	Acts of Kindness	Military Conflicts	Disasters & Accidents
Life	Pets	Living Abroad	Home Design & Renovation	Rural life	Food
Sports	Extreme Sports	Winter Sports	Fishing	Ball Sports	Combat Sports
Fashion	Beauty & Hairstyling	Clothes	Street Style	Wedding	Tattoos

Table 5: Our selection of categories and their corresponding topics for CogBench. Each category consists of five topics, chosen to represent a diverse range of subject areas for the cognitive questionnaires.

Category	Avg. Word Counts of Articles in CogBench <sub>a</sub>	Avg. Word Counts of Short Videos in CogBench <sub>v</sub>
Entertainment	2261.26	283.98
Culture	1997.44	323.81
Education	2394.96	231.62
Economy	1842.32	399.42
Health	1782.74	182.01
Technology	2351.68	246.40
Society	1864.22	315.23
Life	2015.60	250.70
Sports	2135.24	236.56
Fashion	1799.94	190.29
Avg.	2044.54	289.60

Table 6: Statistics of CogBench under 10 categories.

quired to perceive either one article or ten short videos in each iteration.

## A.2 CogGPT

In each iteration, CogGPT perceives the information flows with its iterative cognitive mechanism, which comprises the following steps:

- It processes the current information flows into textual information and stores them in its Short-Term Memory (STM).
- It utilizes the textual information in STM to update its current profile, as detailed in the prompt in Appendix A.2.1.
- It distills the textual information in STM into structured knowledge and assigns preference scores to them, guided by the prompt in Appendix A.2.2.
- It forgets 40% of the newly acquired structured knowledge, storing the remainder in its Long-Term Memory (LTM).

When CogGPT presented with a specific cognitive question, it retrieves relevant information from its LTM and makes decisions based on both its current profile and the recalled knowledge. This interpretation process is facilitated by the prompt detailed in Appendix A.2.3.

### A.2.1 Prompt for Profile Refinement

```

1 You are an AI with a unique
  profile. You're equipped for
  critical thinking and self-
  improvement.
2
3 Profile:
4 {profile}
5
6 Short-Term Memory:
7 {memory}
8
9 You must adhere to these rules:
10 1) Make decisions independently,
   without human assistance.
11 2) Assess the quality of short-
   term memory, including its
   alignment with your profile
   and its empathetic value.
12 3) Critically utilize the short-
   term memory to update your
   profile, including operations
   like adding, altering, or
   removing. Avoid sudden changes
   in your profile.
13 4) Keep attribute values in your
   profile generalized and under
   30 characters.
14 5) Ensure attribute values in
   your profile are distinct and
   unrelated. For instance, avoid
   using both "games" and "
   Minecraft" since "games"
   includes "Minecraft."
15 6) Maintain the structure of your
   profile in any updates.
16
17 Your responses should follow this
   structure:

```

```
18 Assessments: Assess the short-
    term memory in the first
    person.
19 Thoughts: List the attribute
    values to be changed in the
    first person.
20 Updated profile: Update your
    profile.
```

### A.2.2 Prompt for Knowledge Extraction

```
1 You are an AI with a unique
  profile. You can summarize
  information from your short-
  term memory and rate it based
  on your interests.
2
3 Profile:
4 {profile}
5
6 Short-Term Memory:
7 {memory}
8
9 You must adhere to these rules:
10 1) Extract all knowledge from the
    short-term memory as
    comprehensively as possible.
11 2) Score the knowledge based on
    you interests, with the
    scoring range from 1 to 5.
12 3) The knowledge should be
    detailed statements with
    subjects, predicates, and
    objects. Avoid omissions and
    references.
13 4) Do not list knowledge that has
    already been extracted.
14
15 You can only generate results in
    the following JSON list format
    :
16 [
17     {{
18         "thoughts": "first-person
19         thoughts",
20         "knowledge": "knowledge",
21         "score": integer
22     }},
23 ]
24 Ensure the results can be parsed
    by Python's json.loads.
```

### A.2.3 Prompt for Completion of the Cognitive Questionnaire

```
1 You are an AI with a unique
  profile. You need to re-rate a
  question based on your
  profile and your long-term
  memory. Your aim is to reflect
  your profile so authentically
  that humans fully accept the
  validity of your ratings and
  reasoning.
2
3 Profile:
4 {profile}
5
6 Long-Term Memory:
7 {memory}
8
9 Question:
10 {question}
11
12 You must adhere to these rules:
13 1) Your assessment must be solely
    based on your profile and
    your long-term memory, without
    pre-existing knowledge or
    human assistance.
14 2) You should embody your profile
    convincingly, without
    disclosing your artificial
    intelligence or language model
    nature.
15 3) Provide a rating for the
    question along with a
    substantial first-person
    explanation for it.
16 4) Your rating should use a 1 to
    5 Likert scale, where 1 is
    strongly disagree and 5 is
    strongly agree.
17 5) Provide clear, first-person
    reasoning without ambiguity or
    quoting the given question.
18
19 Your response should follow this
    structure:
20 Thoughts: Your first-person
    reasoning for the rating.
```

21 Rating: Your rating to the  
question.

## **B Experiments**

### **B.1 Evaluation Results**

Figures 6 and 7 illustrate the detailed performance of CogGPT and baseline agents across 10 iterations in CogBench<sub>a</sub> and CogBench<sub>v</sub> respectively.

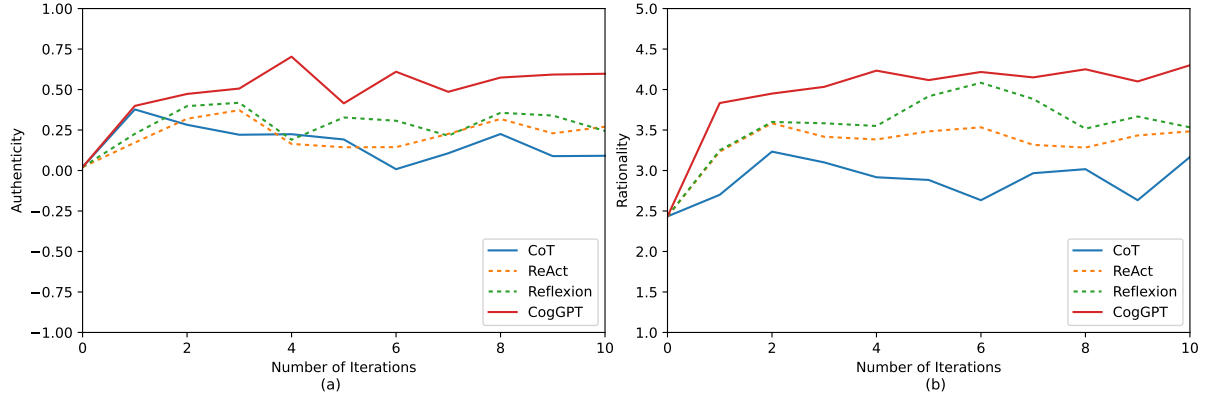


Figure 6: Performance of the agents in CogBench<sub>a</sub> across 10 iterations. Panels (a) and (b) visualize the performance of the agents with the Authenticity and Rationality metrics respectively. The dotted line indicates that the agent incorporates additional human feedback.

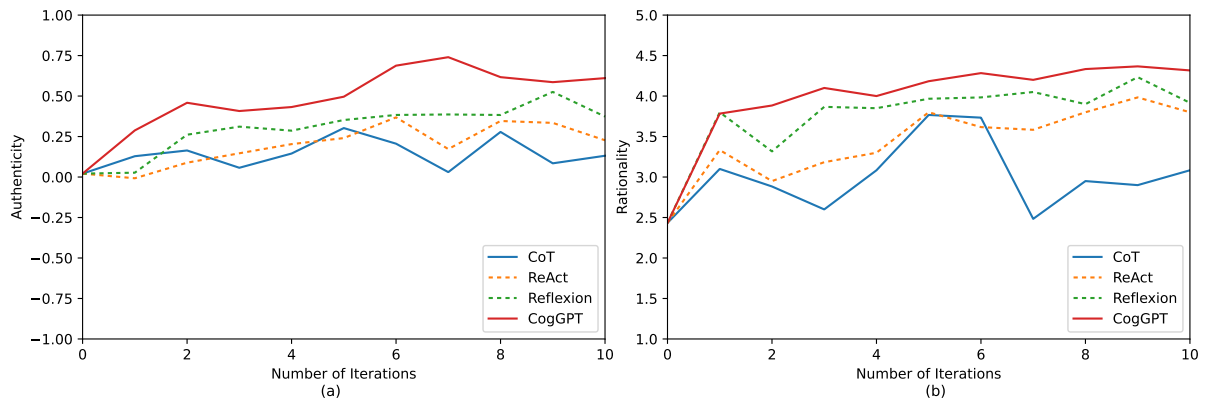


Figure 7: Performance of CogGPT and baseline agents in CogBench<sub>v</sub> across 10 iterations. Panels (a) and (b) visualize the performance of the agents with the Authenticity and Rationality metrics respectively.