

DEEP COGNITION: TOWARDS A MORE TRANSPARENT AND INTERACTIVE RESEARCH AGENT

Anonymous authors

Paper under double-blind review

ABSTRACT

Despite advances in large language models (LLMs), current systems for deep research are limited by an asynchronous, “Input-Wait-Output” interaction paradigm. This model creates a critical disconnect between human intent and AI execution, leading to error propagation and an inability to dynamically course-correct during complex problem-solving. We introduce Deep Cognition, a system designed to enable this paradigm through three technical pillars: transparent and interruptible AI reasoning, fine-grained bidirectional dialogue, and a shared cognitive context. At the core of our system is a multi-agent collaboration framework driven by a dynamic Plan-Search-Report workflow. **This architecture continuously integrates interaction data (information) (e.g., dialogue trajectories and retrieved evidence) into an evidence-driven iterative report construction process. By employing selective context retention to filter unutilized information, our system mitigates error cascades and allows the AI to adapt its reasoning pathways based on the user’s implicit focus.** We conduct a comprehensive user study on challenging deep research tasks to evaluate the efficacy of our system. Results show that our approach significantly enhances the user experience, yielding improvements of up to 29.2% in Fine-Grained Interaction and 27.7% in Ease of Collaboration compared to a competitive baseline. Most notably, our system demonstrates a 31.8% to 50.0% improvement in overall task performance. These results highlight the critical importance of designing interactive AI systems that facilitate continuous human guidance and transparent reasoning, rather than merely responding to isolated commands.

1 INTRODUCTION

As artificial intelligence (AI) capabilities have advanced dramatically through large language models (LLMs) (Luo et al., 2024; Radford et al., 2018; 2021; Brown et al., 2020; 2024), the prevailing trajectory in AI development has emphasized scaling model parameters (Kaplan et al., 2020; Hoffmann et al., 2022; Wei et al., 2022), expanding training data (Yang et al., 2025; Meta AI, 2025), and refining architectures (DeepSeek-AI et al., 2025; MiniMax et al., 2025; Poli et al., 2024)—creating increasingly autonomous black boxes that assume minimal human input beyond simple prompting (Liu et al., 2023b; Kim et al., 2023), instruction (Kim et al., 2023) or decision-making (Yin, 2025). This pathway implicitly assumes that the ultimate form of artificial intelligence would require minimal human input. We contend that this assumption mischaracterizes the nature of intelligence itself. This paradigm positions humans as external operators who provide initial prompts and consume final outputs while remaining excluded from the cognitive process itself, treating human intelligence as merely an instructor rather than a collaborative partner. Consequently, a fundamental question emerges: **How can we design an agentic framework that enables humans to effectively guide AI reasoning trajectories through strategic, real-time interventions?** However, intelligence—whether human or artificial—is inherently interactive, contextual, and collaborative (Hutchins, 1995; Minsky, 1987; Woolley et al., 2010). The most sophisticated human thinking rarely occurs in isolation but emerges through dialogue, feedback, refinement, and the integration of diverse perspectives. Consider the nature of breakthrough scientific discoveries or complex problem-solving scenarios: They invariably involve iterative cycles of hypothesis formation, testing, revision, and collaborative refinement. As AI systems approach advanced cognitive capabilities powered by inference-time scaling (OpenAI, 2024)—enabling thought-level communication where strategic human oversight can leverage vast AI execution power (Xia et al., 2025)—the need for meaningful interaction transforms and intensifies.

This is especially critical for extended AI tasks (Kwa et al., 2025) spanning hours to days, which fundamentally alter human-AI collaboration dynamics.

This transition is particularly evident in systems designed for Deep Research tasks (OpenAI, 2025c; Google, 2025; Perplexity AI, 2025; Zheng et al., 2025a)—complex, extended cognitive processes involving dynamic information retrieval, filter, understanding, analysis and synthesis. Current state-of-the-art research systems have pioneered capabilities for multi-step web browsing, data analysis, and report generation. However, these systems uniformly adopt an “Input-Wait-Output” interaction paradigm where users initiate a query, wait through an extended “Black Box” processing period (typically 5-30 minutes), and eventually receive a comprehensive result. This approach reflects the persistent assumption that interaction is merely a necessary cost rather than a source of value. Yet these systems fundamentally suffer from critical deficiencies: early errors (Cemri et al., 2025) compound without correction, systems cannot adapt to evolving requirements, domain expertise remains inaccessible at crucial moments, and opaque processing prevents human-AI collaboration.

These deficiencies stem from a fundamental misalignment: systems that minimize human involvement during processing cannot address problems that require adaptive guidance and expert intervention (Bainbridge, 1983). To address this fundamental challenge, we develop **deep cognition**—a systematic framework that transcends traditional automation by embedding real-time human expertise directly into AI reasoning processes for complex research tasks, guided by the following principles:

- **Transparency:** The system reveals its entire thinking process—from search strategies and query formulations to information evaluation and synthesis rationales—making AI cognition inspectable and editable at every stage. This transparency enables true thought-level interaction where humans can guide how AI thinks.
- **Fine-Grained Interaction:** Users can engage with any specific element of the AI’s output—questioning particular claims, requesting elaboration on specific points, or changing the research focus.

These principles fundamentally transform deep research from conventional question-and-answer exchanges into cognitive collaboration (see Appendix ??) — what we term **cognitive oversight**. Rather than relegating humans to the role of passive tool operators, this framework establishes a synergistic reasoning process that harnesses the complementary strengths of human expertise and AI capabilities while mitigating their respective limitations. Through cognitive oversight, we move beyond the traditional paradigm of human-AI interaction toward a new form of augmented intelligence where strategic human insight and AI computational power merge into a unified cognitive system.

Through extensive experiments with real expert interactions, we demonstrate that deep cognition achieves substantial improvements or competitive over strongest baseline across all evaluation dimensions: Transparency (+20.0%), Fine-Grained Interaction (+29.2%), Real-Time Intervention (+18.5%), Ease of Collaboration (+27.7%), Results-Worth-Effort (+8.8%), and Interruptibility (+20.7%). Our contributions are summarized as follows:

- **Agentic Multi-Agent Workflow:** We developed an anti-degradation workflow that co-evolves with stronger base models and integrates professional sub-agents.
- **Comprehensive Evaluation Framework:** We establish a complete evaluation framework, including 15 metrics specifically designed for assessing the effectiveness of cognitive oversight in deep research scenarios.
- **Human Fine-Grained Oversight:** We operationalize the cognitive oversight paradigm into our multi-agent human-AI collaboration system designed for deep research tasks.

2 METHODOLOGY

2.1 SYSTEM ARCHITECTURE OVERVIEW

We propose a multi-agent collaborative deep research system designed to address the challenges of long-form report generation. The system supported by **four key processes: Planning Agent, Clarification, Browsing Agent, and Writing Agent**, with the capability for agents to solicit human

input at any stage of the cycle. The system workflow proceeds as follows: Initially, after user input, the **Proactive Clarification module** guides dialogue through structured questioning to precisely capture research intent and background information. After establishing research objectives, the system enters a **Plan-Search-Report** dynamic loop: within each cycle, network search queries are generated based on current planning status and delegated to the **Sub Browse-Agent Cluster**, which coordinates Sub-Agent groups to **concurrently navigate and extract information from multiple web pages**. During evidence collection, the **Writing Agent** continuously outputs intermediate reports, enabling dynamic user feedback. The workflow supports asynchronous human interruption at any stage, while agents can **proactively evaluate whether the current report matches user requirements at the end of each round, seeking additional information to decide the next action if necessary**. This design ensures the transparency of the research process while maintaining efficient automated information processing capabilities. The following subsections detail each core component and their technical implementation.

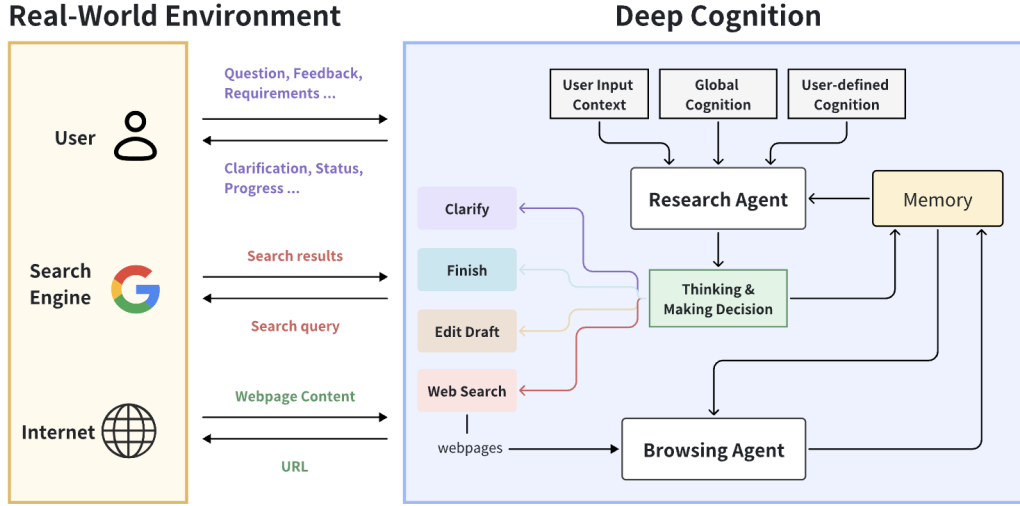


Figure 1: Deep cognition framework overview. This multi-agent research assistant system breaks down complex research questions and dynamically synthesizes information from multiple sources through iterative search, clarification, and user feedback.

2.2 MULTI-ROUND CLARIFICATION MECHANISM

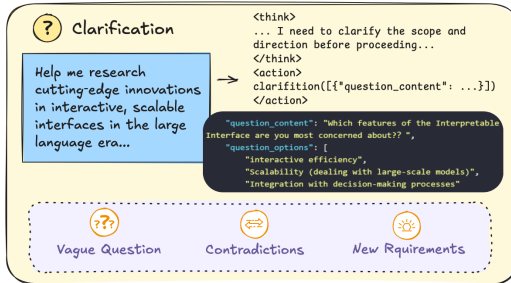


Figure 2: Caption

We design an **option-driven progressive clarification framework** that transforms complex clarification questions into structured option questionnaires, rather than relying on traditional free-text input (detailed clarification prompts see appendix A.1). This mechanism supports triggering clarification processes at any stage of the research, providing continuous human supervision signals for subsequent information retrieval and report generation. **Proactive Clarification Trigger Mechanism:** The system employs a **prompt-based trigger mechanism** to identify moments requiring user clarification. Specifically, we design comprehensive scenario templates in the system prompt that guide the LLM

to recognize situations necessitating human input, including but not limited to: **Ambiguity Detection**: When the research question contains multiple interpretations or the scope is unclear. **Information Conflict**: When retrieved sources present contradictory claims or evidence. **Branch Decision Points**: When the research path encounters multiple viable directions requiring user preference. **Domain Expertise Gaps**: When the system encounters specialized terminology or domain-specific context beyond its knowledge.

To prevent over-interruption, the system can view all previous user interactions in the historical track. The LLM will review this history to avoid repetitive questions and ensure that each clarification request provides incremental value to the research process.

Dynamic Option Generation: Unlike systems that rely on predefined question templates, our framework employs **dynamic option generation**. When a clarification need is identified, the system generates appropriate question. **Multiple-choice questions** with 3-5 options covering probable user intents. **Open-ended questions** for scenarios requiring free-form input. **Contextual explanations** to help users understand each option.

2.3 PROFESSIONAL AGENT CLUSTER

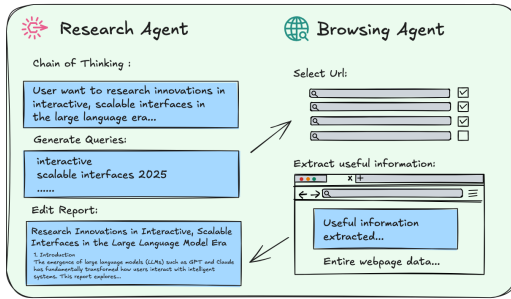


Figure 3: Caption

When processing large-scale web information retrieval tasks, we face two core challenges. First, the **information overload problem** arises as massive URLs and PDF documents exceed the effective processing range of a single model. Second, the **long-sequence degradation problem** manifests as existing large language models universally exhibit the “lost in the middle” (Liu et al., 2023a) phenomenon, struggling to effectively integrate scattered key information when processing long texts. Additionally, the inherent structural looseness and uneven information

density of web content further exacerbate the complexity of information extraction. To address these challenges, we propose a distributed Sub-Browse Agent cluster architecture that achieves efficient information extraction through a systematic workflow. The main Research Agent first queries the Serper API to retrieve the top-20 candidate URLs for each search query, then strategically distributes these resources among specialized Sub-Agent instances. Each Sub-Agent operates within an isolated contextual environment to avoid cross-domain information interference.

For content processing, Sub-Agents employ adaptive chunking strategies to handle documents of varying lengths. Standard web pages are processed using fixed-size chunking with overlapping windows, while exceptionally long documents trigger an autonomous pagination decision mechanism where the Browse Agent evaluates content density and relevance to determine whether to continue processing subsequent sections. Upon completion of analysis, each Sub-Agent submits structured findings to the main Agent with three components: **Excerpts**, **Useful** and **Reasoning**. This architecture effectively distributes computational load, enables specialized processing optimization, and significantly improves both efficiency and accuracy in large-scale web information retrieval tasks.

The system utilizes a hierarchical, modular design to manage long-term research planning. We design a research agent to propose the research plan and autonomously determine the next action base on the current research state. This multi agent modular transfer isolates task-specific logic (e.g., research planning, web browsing, report generation), thereby preventing cross-module context interference. The research agent logs all completed events as a to-do list, this to-do list verified the current research state whether align with user goal.

We define the Research Agent (detailed prompts in appendix A.3) as a professional research scientist and strictly define the system’s capability boundaries to enable the agent to plan highly feasible to-do lists (specifically capable of searching, analyzing, and writing, but not programming or deploying models). Furthermore, we provide the agent with three distinct few-shot example types (covering literature review, technical proposal, and precise retrieval), each including both correct and incorrect instances. These examples differentiate strategies suitable for internal deliberation, external informa-

tion seeking, and precise factual comparison, guiding the agent to generate the most accurate plans, detailed prompts in appendix A.2.

2.4 INTERMEDIATE REPORTS THROUGH WRITING AGENT

While existing deep research systems (LangChain, 2024; Roucher et al., 2025a) typically follow a sequential collect-then-generate paradigm, we propose an **evidence-driven iterative report construction strategy**. We deployed a specially fine-tuned Writing Agent capable of generating structured intermediate reports even when evidence collection remains ongoing (detailed prompts in appendix A.4).

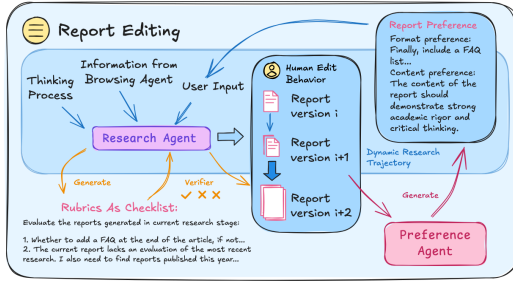


Figure 4: Caption

removing unutilized evidence from subsequent processing contexts. This parallel evidence acquisition and report construction paradigm breaks through the limitations of traditional batch processing approaches, enabling continuous knowledge synthesis processes.

3 HUMAN-AI CO-RESEARCH MECHANISM

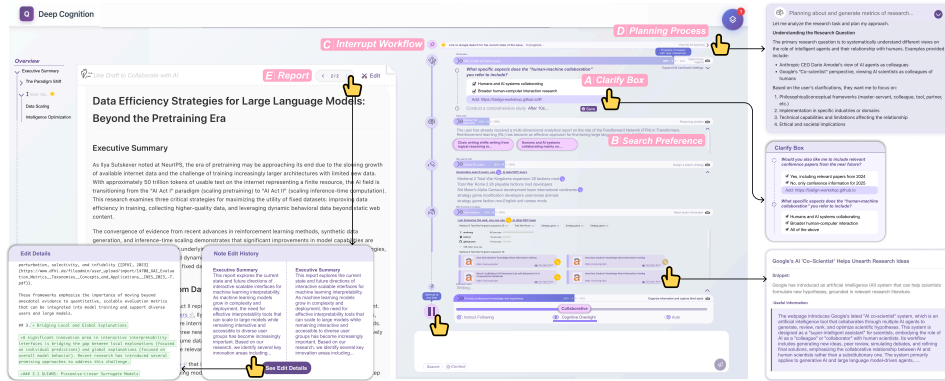


Figure 5: Deep cognition interface design showcasing key interactive features: (A) Research scope clarification to refine vague queries, (B) Click to open the important URL, (C) Multi-agent Workflow Visualization, (D) Transparent display of reasoning, research processes, and interactive query refinement, and (E) Report revision. The 🖱️ icon stands for clickable interface elements.

Deep cognition supports real-time human-AI collaboration. It is designed for open-ended, multi-hop retrieval and exploratory analysis. It enables users to iteratively expand the initial question and produce a synthesized write-up. Following principles of cognitive oversight, we designed the following features for our deep cognition system, with interfaces presented in Figure 10. The interface supports multiple modes of human-AI collaboration: Clarification (left): The system generates clarification questions to help users specify their focus. Interrupt (bottom-left): Users can intervene during the system’s ongoing retrieval or reasoning process, halting unsatisfactory results and redirecting the search toward more relevant information. Planning (right): The system synthesizes retrieved evidence into a structured research plan. **Transparent Research Process:** The interface

make the system’s decision-making process visible and comprehensible to users. Search strategy explainability is achieved by directly displaying the reasoning process and query terms generated by the model, making information retrieval interpretable. The editor area on the left of Figure 10 displays the evolving research document with proper formatting. All findings are properly linked to their original sources, enabling users to trace source materials. **Real-Time Intervention** We implement a “Pause” feature, allowing users to interrupt the system at critical junctures in the research process. This intervention capability enables users to actively shape the research trajectory based on emerging insights or changing objectives.

4 EXPERIMENTS

4.1 METRICS DESIGN

We defined five key dimensions for evaluating the quality of generated reports. Each metric is rigorously assessed using a 5-point Likert scale. For report quality, we focus on organization, coverage, depth, relevance, usefulness, and innovation. For interaction dimensions, we focus on willingness to use, usability, transparency, interruptibility, granular interaction, informativeness, ease of collaboration, cost-effectiveness, real-time intervention, and usefulness.

Metric	Description	Metric	Description
Organization	Evaluate whether the article demonstrates sound organization and logical structure. An acceptable response should:	Intention to Use	Measures user intention and propensity for continued engagement with the system based on perceived value and satisfaction.
	(1) Exhibit clear structure by organizing relevant points into a coherent logical sequence. (2) Maintain coherence without any contradictions or unnecessary repetition.	Usability	Evaluates the intuitive nature and accessibility of the system interface, including cognitive load and interaction efficiency.
Cutting-Edge	Assess whether the article demonstrates comprehensive coverage of existing literature by:	Transparency	Assesses the interpretability and explainability of the model’s decision-making processes and reasoning mechanisms.
	(1) Effectively summarizing and conducting comparative analysis with previous research. (2) Timely incorporating the most recent and up-to-date research findings or information.	Interruptibility	Assesses the system’s ability to tolerate pauses or context switches and to resume smoothly without loss of state or progress.
Coverage	Provide comprehensive coverage of the identified areas of interest through:	Fine-Grained Interaction	Evaluates the system’s capacity to incorporate user feedback and enable precise, granular control over output generation.
	(1) Conducting thorough reviews. (2) Citing a broad range of representative scholarly works. (3) Incorporating the most current and time-sensitive information from various sources, rather than limiting the analysis to a small number of papers.	Inspiration	Assesses the system’s ability to stimulate creative thinking and generate ideas or innovative approaches to problem-solving.
Depth	Assess the adequacy of information content provided in the article. Specifically, evaluate whether the article delivers sufficient relevant information with appropriate depth such that readers can achieve thorough understanding of each argument presented.	Ease of Collaboration	Measures the extent to which the system functions as an effective collaborative partner in knowledge work and decision-making processes.
	Assess whether the response maintains topical relevance and preserves clear focus in order to deliver a useful response to the posed question. Specifically, the output should:	Results-Worth-Effort	Evaluates whether users perceive the time and effort invested in system interaction as worthwhile and valuable relative to the outcomes achieved.
Relevance	(1) Sufficiently address the central elements of the original question and satisfy your informational requirements. (2) The response should exclude substantial amounts of tangential information unrelated to the original inquiry.	Real-Time Intervention	Measures the degree to which users can actively interrupt and steer the system’s ongoing processes—e.g., pausing, editing, or re-prompting—to obtain desired outputs.
		Helpfulness	Assesses the overall utility and practical value of the output in addressing user needs and facilitating problem-solving objectives.

Figure 6: Evaluation Metrics for Report Quality Assessment

4.2 SYSTEM EXPERIMENTAL SETUP

We use claude-3.7-sonnet-thinking as an inference model for action selection and claude-4.0-sonnet for document authoring, and the browsing agent uses gpt-4.1-mini for processing large numbers of documents, with 0.6 used for both temperature. We used the Google TOP20 for web search to provide a realistic search environment for the Agent System. Each turn search generate 5 queries, and for 5 webpages for each query.

4.3 RESEARCH TASK SETUP

To address two limitations of static benchmarks, We perform a **human evaluation** to evaluate the real-world human experience during the human-AI interaction inspired by Lee et al. (2024). This method enables assessment of output quality that depends on interactive dynamics, which aligns with real-world usage scenarios. We develop a web application for users to interact with deep cognition in real time. We compare it with three competitive deep research baseline: Gemini Deep Research (Google, 2025), OpenAI Deep Research (OpenAI, 2025b;a;c) and Grok 3 DeeperSearch (xAI, 2025). Study 1 measuring report quality and the effectiveness of the interaction design. Study 2 testing whether users with higher or lower **prior knowledge** levels show differences in multi-hop retrieval task.

Study 1 We recruited 13 participants with prior research experience. Before using the system, they were introduced to our evaluation metrics (see section 4.1) for deep cognition to ensure a shared understanding. Participants then evaluated both the quality of generated reports and the system’s interactive behaviors on a 5-point Likert scale, supplemented by qualitative responses to open-ended interviews. Each participant proposed a research question from their own work, participants observed the model in real time as it retrieved information, reasoned through intermediate steps, and generated self-evaluations. They could not directly edit the final report but instead guided the process via interactive mechanisms such as interrupting outputs, injecting prior knowledge, inspecting sources, reviewing self-evaluations, suggesting new directions, giving feedback, or contributing personal documents. These interventions helped steer the model toward deeper analysis and more efficient retrieval, with the report finalized when the model itself chose to conclude.

Study 2 To validate our hypothesis that experts with higher cognitive capabilities demonstrate enhanced collaboration with AI in transparent dialogue environments, we measured system performance through two comprehensive benchmarks. Given that our expert annotators are native Chinese speakers with domain expertise, we selected representative subsets for intensive interactive evaluation: 22 questions from browsecomp-ZH (Zhou et al., 2025) (top two from each of 11 categories) and the first 20 questions from xbench-deep research (Chen et al., 2025). Both sampling strategies ensure feasible human-AI collaborative assessment.

5 MAIN RESULT

5.1 EXPERT USER EVALUATION

As shown in Table 1, augmented through expert interaction, the deep cognition system demonstrated significant enhancements across six evaluated metrics, overall average improve 63%. Notably, the ORGANIZATION exhibits the greatest gain (+97%), followed by CUTTING-EDGE (+79%) and depth (+76%). Even the dimension with the smallest gain, helpfulness, showed a significant improvement of +42%. As the evaluation results in Table 2, the **alignment between expert rankings and user evaluations** validates our core hypothesis: **The system with enhanced interaction mechanisms consistently deliver output quality across six metrics.**

Metric	DC (w/o Int).	DC.
Organization	2.231	4.385 ↑ 97%
Cutting-Edge	2.538	4.538 ↑ 79%
Coverage	2.423	4.000 ↑ 65%
Depth	2.231	3.923 ↑ 76%
Relevance	2.885	3.769 ↑ 31%
Helpfulness	2.808	4.000 ↑ 42%
Overall Average	2.519	4.103 ↑ 63%

Table 1: Performance improvement of deep cognition over deep cognition without interaction. DC. indicates deep cognition, DC (non). indicates deep cognition without interaction.

Deep cognition dominates six of the seven metrics. It records the largest gains in Fine-Grained Interaction (+44.6%) and Cooperative (+43.0%), and is the only system to reach a perfect Transparency score (5.00, +25.0% over the strongest baseline). Overall, the results highlight deep cognition’s superior transparency, controllability, and collaborative support. These quantitative results are further supported by users’ qualitative feedback. Over 90% of participants agree or strongly agree that interaction with deep cognition improves report quality; 69% find it easy to use and 62% show a high willingness to use.

Report Evaluation (1–5 Score)					Interaction Evaluation (1–5 Score)				
Metric	DC.	Gemini	OpenAI	Grok3	Metric	DC.	Gemini	OpenAI	Grok 3
Organization	4.385 ^{+1.8%}	4.308	3.769	3.385	Transparency	5.00 ^{+25.0%}	4.00	3.00	3.19
Cutting-Edge	4.538 ^{+3.5%}	4.385	3.769	3.538	Interruptibility	4.35 ^{+31.4%}	3.31	2.69	2.62
Coverage	4.000 ^{-10.4%}	4.462	3.692	2.923	Fine-Grained Interaction	4.73 ^{+44.6%}	3.27	2.88	2.19
Depth	3.923 ^{-1.9%}	4.000	3.577	2.769	Real-Time Intervention	4.69 ^{+24.4%}	3.77	2.92	2.62
Relevance	3.769 ^{-18.3%}	4.615	4.077	3.615	Inspiration	4.08 ^{+0.0%}	4.08	3.42	3.19
Helpfulness	4.000 ^{+0.0%}	4.000	3.615	2.692	Ease of Collaboration	4.62 ^{+43.0%}	3.23	2.77	1.85
					Results-Worth-Effort	4.52 ^{+10.8%}	4.08	3.29	2.96

Table 2: User and expert evaluation results for AI research assistance systems. Left panel: User-generated evaluation scores on a 1-5 scale, where participants queried systems with their own research questions. Right panel: Scores (1–5 scale) for system-interaction evaluation metrics. Color coding indicates within-row performance rankings, and percentages show deep cognition’s relative improvement over the strongest baseline system (Gemini). DC. indicates deep cognition.

5.2 BENCHMARK EVALUATION RESULTS

The results provide compelling evidence for our collaborative cognition framework. On browsecomp-ZH, the deep cognition system achieves 72.73% accuracy—dramatically outperforming all baselines (Gemini/OpenAI: 40.91%, Grok 3: 22.73%). Ablation studies show neither cognitive oversight alone (45.45%) nor interaction alone (40.91%) match their combination. On X-bench, our system achieves 65% accuracy, matching OpenAI while substantially outperforming Gemini (35%). Note that browsecomp-ZH was evaluated on June 22, 2025, and X-bench on September 25, 2025—temporal gaps may contribute to baseline performance variations due to API updates. The results consistently demonstrate that expert-AI collaboration requires both transparent reasoning and interactive guidance for effective performance across domains. Participants with deeper cognitive processing capabilities achieved significantly higher human-AI collaborative performance compared to those with surface-level cognitive approaches in transparent interaction paradigms, as measured by problem resolution accuracy.

	DC (non cog).	DC (non int).	DC (cog+int).	Gemini	OpenAI	Grok 3
Accuracy	45.45%	40.91%	72.73%	40.91%	40.91%	22.73%

	DC (cog+int).	Gemini	OpenAI
Accuracy	65%	35%	65%

Table 3: Accuracy comparison across benchmarks. Top: Browsecomp-ZH (22 questions). Bottom: X-bench deep research (first 20 questions). DC (non cog). = baseline with middle school-level participants (n=4); DC (non int). = autonomous system; DC (cog+int). = interactive condition with graduate-level participants (n=4).

5.3 IN-DEPTH ANALYSIS OF THE HUMAN STUDY: HUMAN HOLD DYNAMIC MENTAL MODELS THROUGHOUT COLLABORATION PROCESS

Enhancing transparency at the model’s behavioral status can improve human-AI collaboration. Specifically, in complex, long-duration retrieval tasks, humans tend to delegate mechanical operations such as “browsing” and “summarizing” to AI, while preferring to collaborate with the model at decision points requiring higher-order thinking. We dive deeper into the human behavior pattern in the deep research process and provide design considerations of human-AI collaboration research system. As illustrated in case study(see Appendix E) and User Behavior Data Point (see Appendix C), our user study reveals a sophisticated pattern of collaborative engagement that varies systematically across six research phases. Users demonstrate **dynamic cooperation willingness**, transitioning between “hands-on” and “hands-off” modes based on task characteristics and their domain expertise. We detail these six phases below:

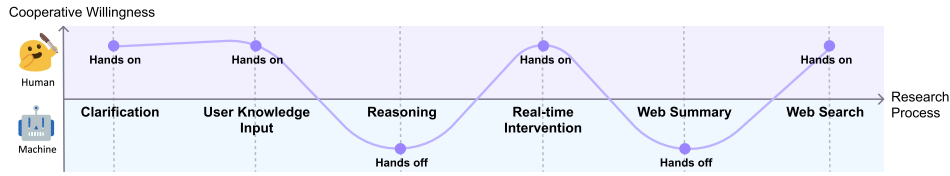


Figure 7: Changes in users’ behavioral tendencies in the process of complex research tasks.

Clarification (Hands-on) The research process begins with intensive human-AI collaboration as users refine vague problem definitions. Users’ initial research questions are typically too broad to cover all possible scenarios. **User Knowledge Input (Hands-on)** Users maintain high engagement when they possess specific domain knowledge or references that need integration. When users know specific references or attributes about an item, such as queries, paper links, websites, or personal opinions, they actively guide the AI to relevant media. **Reasoning (Hands-off)** Users seek to understand whether the model has correctly executed prescribed instructions and want transparency in decision-making processes. **Real-Time Intervention (Hands-on)** Cooperation peaks again during dynamic browsing tasks where users encounter pages or information sources that warrant detailed retrieval. **Web Summary (Hands-off)** During summarization tasks, users tend to trust in AI capability. Participants often need consolidated insights from multiple sources rather than single source summarization, leading them to allow extended autonomous operation. **Web Search (Hands-on)** The cycle concludes with renewed hands-on engagement for open-ended and subjective questions that require interpretation or subjective judgment.

This dynamic pattern demonstrates that effective human-AI collaboration is not uniform but adapts strategically to leverage the comparative advantages of human judgment and AI processing capabilities across different research phases. We illustrate this dynamic research task example to demonstrate authentic participant behavior.

6 RELATED WORK

Human-AI Interaction AI agents [White \(2024\)](#); [Feng et al. \(2025\)](#) now support complex tasks through natural language interaction, better task understanding, and multi-level autonomy beyond basic queries interaction ([Srinivas & Runkana, 2025](#); [Shao et al., 2025](#)). The shift from static monolithic inference to adaptive, resource-aware computation has become central to AI systems for knowledge discovery ([Shao et al., 2024](#); [Jiang et al., 2024](#)) leveraging multi-agent collaboration ([Watkins et al., 2025](#); [Fragiadakis et al., 2025](#)) to facilitate serendipitous discovery. This mismatch constrains the potential for AI to act as a collaborator in exploratory inquiry ([Pirolli, 2009](#)). Although current collaboration systems allow humans to read model reasoning chains and engage in multi-turn interactions with models ([Westphal et al., 2023](#); [Gomez et al., 2025](#); [Lee et al., 2024](#); [Collins et al., 2024](#)), these current interaction paradigms maintain limiting user’s ability to adapt to emerging expert user’s knowledge during complex and time-consuming tasks.

Deep Research Systems Deep research systems such as Gemini Deep Research ([Google, 2025](#)), OpenAI Deep Research ([OpenAI, 2025b](#)) and Grok3 Deeper Search ([xAI, 2025](#)) are enabled by the sophisticated reasoning abilities that have emerged from recent advances in large language models (LLMs) ([OpenAI et al., 2024](#); [Guo et al., 2025](#); [Team et al., 2025](#)), facilitating multi-step, in-depth analysis and information synthesis across hundreds of sources. Most open-source deep research projects ([LangChain AI, 2025](#); [Zhang, 2025](#); [Elovic, 2025](#); [Camara, 2025](#); [Jina AI, 2025](#); [Roucher et al., 2025b](#); [ByteDance, 2024](#)) employ prompt-based multi-agent systems with predefined workflows. Recent work ([Zheng et al., 2025b](#)) has applied end-to-end reinforcement learning to open-source LLMs to perform iterative reasoning to complex questions. However, few existing deep research systems in Appendix?? development multi-round interaction planning during the research process, user remain limited once research begins.

7 CONCLUSION

This paper introduced deep cognition, a multi-agent framework for collaborative research with real-time “cognitive oversight” through transparent, interruptible interactions. Our evaluation challenge the assumption that AI progress requires purely autonomous capabilities. Instead, our work suggests that advanced intelligence emerges from cognitive partnerships that leverage complementary human judgment and machine processing strengths.

ETHICS STATEMENT

This work adheres to the ICLR Code of Ethics. Human participants were involved in this study, and all procedures were conducted with informed consent and in strict accordance with relevant ethical standards. No personally identifiable information was collected or stored, and participants’ privacy was fully protected throughout the study. All datasets used were obtained in compliance with relevant usage guidelines. We took care to mitigate potential biases and discriminatory outcomes, and no experiments were conducted that could raise privacy or security concerns. We remain committed to ensuring transparency, fairness, and integrity in the research process.

REPRODUCIBILITY STATEMENT

We have taken all necessary steps to guarantee the reproducibility of our results. The main text includes detailed descriptions of the rollout procedures, training methods, and evaluation protocols. Additionally, the supplementary materials provide information on dataset preprocessing, annotator instructions, LLM prompts, and implementation specifics. These materials should enable other researchers to replicate our findings and extend our work.

REFERENCES

- Lisanne Bainbridge. Ironies of automation. In *Analysis, design and evaluation of man-machine systems*, pp. 129–135. Elsevier, 1983.
- Bradley Brown, Jordan Juravsky, Ryan Ehrlich, Ronald Clark, Quoc V. Le, Christopher Ré, and Azalia Mirhoseini. Large language monkeys: Scaling inference compute with repeated sampling, 2024. URL <https://arxiv.org/abs/2407.21787>.
- Tom B. Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared Kaplan, Prafulla Dhariwal, Arvind Neelakantan, Pranav Shyam, Girish Sastry, Amanda Askell, Sandhini Agarwal, Ariel Herbert-Voss, Gretchen Krueger, Tom Henighan, Rewon Child, Aditya Ramesh, Daniel M. Ziegler, Jeffrey Wu, Clemens Winter, Christopher Hesse, Mark Chen, Eric Sigler, Mateusz Litwin, Scott Gray, Benjamin Chess, Jack Clark, Christopher Berner, Sam McCandlish, Alec Radford, Ilya Sutskever, and Dario Amodei. Language models are few-shot learners, 2020. URL <https://arxiv.org/abs/2005.14165>.
- ByteDance. Deerflow, 2024. URL <https://github.com/bytedance/deer-flow>. Community-driven deep research framework combining LLMs with web search, crawling, and code execution tools.
- Nicolas Camara. Open deep research, 2025. URL <https://github.com/nickscamara/open-deep-research>. Open-source clone of OpenAI’s Deep Research using Firecrawl for web data extraction and AI reasoning.
- Mert Cemri, Melissa Z. Pan, Shuyi Yang, Lakshya A. Agrawal, Bhavya Chopra, Rishabh Tiwari, Kurt Keutzer, Aditya Parameswaran, Dan Klein, Kannan Ramchandran, Matei Zaharia, Joseph E. Gonzalez, and Ion Stoica. Why do multi-agent llm systems fail?, 2025. URL <https://arxiv.org/abs/2503.13657>.
- Kaiyuan Chen, Yixin Ren, Yang Liu, Xiaobo Hu, Haotong Tian, Tianbao Xie, Fangfu Liu, Haoye Zhang, Hongzhang Liu, Yuan Gong, Chen Sun, Han Hou, Hui Yang, James Pan, Jianan Lou, Jiayi Mao, Jizheng Liu, Jinpeng Li, Kangyi Liu, Kenkun Liu, Rui Wang, Run Li, Tong Niu,

- Wenlong Zhang, Wenqi Yan, Xuancheng Wang, Yuchen Zhang, Yi-Hsin Hung, Yuan Jiang, Zexuan Liu, Zihan Yin, Zijian Ma, and Zhiwen Mo. xbench: Tracking agents productivity scaling with profession-aligned real-world evaluations, 2025. URL <https://arxiv.org/abs/2506.13651>.
- Katherine M. Collins, Ilia Sucholutsky, Umang Bhatt, Kartik Chandra, Lionel Wong, Mina Lee, Cedegao E. Zhang, Tan Zhi-Xuan, Mark Ho, Vikash Mansinghka, Adrian Weller, Joshua B. Tenenbaum, and Thomas L. Griffiths. Building machines that learn and think with people, 2024. URL <https://arxiv.org/abs/2408.03943>.
- DeepSeek-AI, Aixin Liu, Bei Feng, Bing Xue, Bingxuan Wang, Bochao Wu, Chengda Lu, Chenggang Zhao, Chengqi Deng, Chenyu Zhang, Chong Ruan, Damai Dai, Daya Guo, Dejian Yang, Deli Chen, Dongjie Ji, Erhang Li, Fangyun Lin, Fucong Dai, Fuli Luo, Guangbo Hao, Guanting Chen, Guowei Li, H. Zhang, Han Bao, Hanwei Xu, Haocheng Wang, Haowei Zhang, Honghui Ding, Huajian Xin, Huazuo Gao, Hui Li, Hui Qu, J. L. Cai, Jian Liang, Jianzhong Guo, Jiaqi Ni, Jiashi Li, Jiawei Wang, Jin Chen, Jingchang Chen, Jingyang Yuan, Junjie Qiu, Junlong Li, Junxiao Song, Kai Dong, Kai Hu, Kaige Gao, Kang Guan, Kexin Huang, Kuai Yu, Lean Wang, Lecong Zhang, Lei Xu, Leyi Xia, Liang Zhao, Litong Wang, Liyue Zhang, Meng Li, Miaojuan Wang, Mingchuan Zhang, Minghua Zhang, Minghui Tang, Mingming Li, Ning Tian, Panpan Huang, Peiyi Wang, Peng Zhang, Qiancheng Wang, Qihao Zhu, Qinyu Chen, Qiushi Du, R. J. Chen, R. L. Jin, Ruiqi Ge, Ruisong Zhang, Ruizhe Pan, Runji Wang, Runxin Xu, Ruoyu Zhang, Ruyi Chen, S. S. Li, Shanghao Lu, Shangyan Zhou, Shanhuang Chen, Shaoqing Wu, Shengfeng Ye, Shengfeng Ye, Shirong Ma, Shiyu Wang, Shuang Zhou, Shuiping Yu, Shunfeng Zhou, Shuting Pan, T. Wang, Tao Yun, Tian Pei, Tianyu Sun, W. L. Xiao, Wangding Zeng, Wanjia Zhao, Wei An, Wen Liu, Wenfeng Liang, Wenjun Gao, Wenqin Yu, Wentao Zhang, X. Q. Li, Xiangyue Jin, Xianzu Wang, Xiao Bi, Xiaodong Liu, Xiaohan Wang, Xiaojin Shen, Xiaokang Chen, Xiaokang Zhang, Xiaosha Chen, Xiaotao Nie, Xiaowen Sun, Xiaoxiang Wang, Xin Cheng, Xin Liu, Xin Xie, Xingchao Liu, Xingkai Yu, Xinnan Song, Xinxia Shan, Xinyi Zhou, Xinyu Yang, Xinyuan Li, Xuecheng Su, Xuheng Lin, Y. K. Li, Y. Q. Wang, Y. X. Wei, Y. X. Zhu, Yang Zhang, Yanhong Xu, Yanhong Xu, Yanping Huang, Yao Li, Yao Zhao, Yaofeng Sun, Yaohui Li, Yaohui Wang, Yi Yu, Yi Zheng, Yichao Zhang, Yifan Shi, Yiliang Xiong, Ying He, Ying Tang, Yishi Piao, Yisong Wang, Yixuan Tan, Yiyang Ma, Yiyuan Liu, Yongqiang Guo, Yu Wu, Yuan Ou, Yuchen Zhu, Yudian Wang, Yue Gong, Yuheng Zou, Yujia He, Yukun Zha, Yunfan Xiong, Yunxian Ma, Yuting Yan, Yuxiang Luo, Yuxiang You, Yuxuan Liu, Yuyang Zhou, Z. F. Wu, Z. Z. Ren, Zehui Ren, Zhangli Sha, Zhe Fu, Zhen Xu, Zhen Huang, Zhen Zhang, Zhenda Xie, Zhengyan Zhang, Zhewen Hao, Zhibin Gou, Zhicheng Ma, Zhigang Yan, Zhihong Shao, Zhipeng Xu, Zhiyu Wu, Zhongyu Zhang, Zhuoshu Li, Zihui Gu, Zijia Zhu, Zijun Liu, Zilin Li, Ziwei Xie, Ziyang Song, Ziyi Gao, and Zizheng Pan. Deepseek-v3 technical report, 2025. URL <https://arxiv.org/abs/2412.19437>.
- Assaf Elovic. Gpt researcher, 2025. URL <https://github.com/assafelovic/gpt-researcher>. Open deep research agent for web and local research with detailed report generation and citations.
- K. J. Kevin Feng, David W. McDonald, and Amy X. Zhang. Levels of autonomy for ai agents, 2025. URL <https://arxiv.org/abs/2506.12469>.
- George Fragiadakis, Christos Diou, George Kousiouris, and Mara Nikolaidou. Evaluating human-ai collaboration: A review and methodological framework, 2025. URL <https://arxiv.org/abs/2407.19098>.
- Catalina Gomez, Sue Min Cho, Shichang Ke, Chien-Ming Huang, and Mathias Unberath. Human-ai collaboration is not very collaborative yet: a taxonomy of interaction patterns in ai-assisted decision making from a systematic review. *Frontiers in Computer Science*, Volume 6 - 2024, 2025. ISSN 2624-9898. doi: 10.3389/fcomp.2024.1521066. URL <https://www.frontiersin.org/journals/computer-science/articles/10.3389/fcomp.2024.1521066>.
- Google. Gemini deep research - your personal research assistant, 2025. URL <https://gemini.google/overview/deep-research/>. Accessed: April 14, 2025.
- Daya Guo, Dejian Yang, Haowei Zhang, Junxiao Song, Ruoyu Zhang, Runxin Xu, Qihao Zhu, Shirong Ma, Peiyi Wang, Xiao Bi, et al. Deepseek-r1: Incentivizing reasoning capability in llms via reinforcement learning. *arXiv preprint arXiv:2501.12948*, 2025.

- Jordan Hoffmann, Sebastian Borgeaud, Arthur Mensch, Elena Buchatskaya, Trevor Cai, Eliza Rutherford, Diego de Las Casas, Lisa Anne Hendricks, Johannes Welbl, Aidan Clark, Tom Hennigan, Eric Noland, Katie Millican, George van den Driessche, Bogdan Damoc, Aurelia Guy, Simon Osindero, Karen Simonyan, Erich Elsen, Oriol Vinyals, Jack W. Rae, and Laurent Sifre. Training compute-optimal large language models. In *Proceedings of the 36th International Conference on Neural Information Processing Systems, NIPS '22*, Red Hook, NY, USA, 2022. Curran Associates Inc. ISBN 9781713871088.
- Edwin Hutchins. *Cognition in the Wild*. MIT press, 1995.
- Yucheng Jiang, Yijia Shao, Dekun Ma, Sina J. Semnani, and Monica S. Lam. Into the unknown unknowns: Engaged human learning through participation in language model agent conversations, 2024. URL <https://arxiv.org/abs/2408.15232>.
- Jina AI. node-deepresearch, 2025. URL <https://github.com/jina-ai/node-DeepResearch>. Iterative search, reading, and reasoning system for deep research queries with focus on concise answers.
- Jared Kaplan, Sam McCandlish, Tom Henighan, Tom B. Brown, Benjamin Chess, Rewon Child, Scott Gray, Alec Radford, Jeffrey Wu, and Dario Amodei. Scaling laws for neural language models, 2020. URL <https://arxiv.org/abs/2001.08361>.
- Sunnie S. Y. Kim, Elizabeth Anne Watkins, Olga Russakovsky, Ruth Fong, and Andrés Monroy-Hernández. "help me help the ai": Understanding how explainability can support human-ai interaction. In *Proceedings of the 2023 CHI Conference on Human Factors in Computing Systems, CHI '23*, New York, NY, USA, 2023. Association for Computing Machinery. ISBN 9781450394215. doi: 10.1145/3544548.3581001. URL <https://doi.org/10.1145/3544548.3581001>.
- Thomas Kwa, Ben West, Joel Becker, Amy Deng, Katharyn Garcia, Max Hasin, Sami Jawhar, Megan Kinniment, Nate Rush, Sydney Von Arx, Ryan Bloom, Thomas Broadley, Haoxing Du, Brian Goodrich, Nikola Jurkovic, Luke Harold Miles, Seraphina Nix, Tao Lin, Neev Parikh, David Rein, Lucas Jun Koba Sato, Hjalmar Wijk, Daniel M. Ziegler, Elizabeth Barnes, and Lawrence Chan. Measuring ai ability to complete long tasks, 2025. URL <https://arxiv.org/abs/2503.14499>.
- LangChain. Open deep research, 2024. URL https://github.com/langchain-ai/open_deep_research. GitHub repository, accessed December 2024.
- LangChain AI. Open deep research, 2025. URL https://github.com/langchain-ai/open_deep_research. Open-source research assistant for automated deep research and report generation.
- Mina Lee, Megha Srivastava, Amelia Hardy, John Thickstun, Esin Durmus, Ashwin Paranjape, Ines Gerard-Ursin, Xiang Lisa Li, Faisal Ladhak, Frieda Rong, Rose E. Wang, Minae Kwon, Joon Sung Park, Hancheng Cao, Tony Lee, Rishi Bommasani, Michael Bernstein, and Percy Liang. Evaluating human-language model interaction, 2024. URL <https://arxiv.org/abs/2212.09746>.
- Nelson F Liu, Kevin Lin, John Hewitt, Ashwin Paranjape, Michele Bevilacqua, Fabio Petroni, and Percy Liang. Lost in the middle: How language models use long contexts. *arXiv preprint arXiv:2307.03172*, 2023a.
- Pengfei Liu, Weizhe Yuan, Jinlan Fu, Zhengbao Jiang, Hiroaki Hayashi, and Graham Neubig. Pre-train, prompt, and predict: A systematic survey of prompting methods in natural language processing. *ACM computing surveys*, 55(9):1–35, 2023b.
- Liangchen Luo, Yinxiao Liu, Rosanne Liu, Samrat Phatale, Harsh Lara, Yunxuan Li, Lei Shu, Yun Zhu, Lei Meng, Jiao Sun, et al. Improve mathematical reasoning in language models by automated process supervision. *ArXiv preprint*, abs/2406.06592, 2024. URL <https://arxiv.org/abs/2406.06592>.
- Meta AI. The llama 4 herd: The beginning of a new era of natively multimodal ai innovation, 4 2025. URL <https://ai.meta.com/blog/llama-4-multimodal-intelligence/>. Accessed: November 27, 2025.

- MiniMax, Aonian Li, Bangwei Gong, Bo Yang, Boji Shan, Chang Liu, Cheng Zhu, Chunhao Zhang, Congchao Guo, Da Chen, Dong Li, Enwei Jiao, Gengxin Li, Guojun Zhang, Haohai Sun, Houze Dong, Jiadai Zhu, Jiaqi Zhuang, Jiayuan Song, Jin Zhu, Jingtao Han, Jingyang Li, Junbin Xie, Junhao Xu, Junjie Yan, Kaishun Zhang, Kecheng Xiao, Kexi Kang, Le Han, Leyang Wang, Lianfei Yu, Liheng Feng, Lin Zheng, Linbo Chai, Long Xing, Meizhi Ju, Mingyuan Chi, Mozhi Zhang, Peikai Huang, Pengcheng Niu, Pengfei Li, Pengyu Zhao, Qi Yang, Qidi Xu, Qiexiang Wang, Qin Wang, Qiuhi Li, Ruitao Leng, Shengmin Shi, Shuqi Yu, Sichen Li, Songquan Zhu, Tao Huang, Tianrun Liang, Weigao Sun, Weixuan Sun, Weiyu Cheng, Wenkai Li, Xiangjun Song, Xiao Su, Xiaodong Han, Xinjie Zhang, Xinzhu Hou, Xu Min, Xun Zou, Xuyang Shen, Yan Gong, Yingjie Zhu, Yipeng Zhou, Yiran Zhong, Yongyi Hu, Yuanxiang Fan, Yue Yu, Yufeng Yang, Yuhao Li, Yunan Huang, Yunji Li, Yunpeng Huang, Yunzhi Xu, Yuxin Mao, Zehan Li, Zekang Li, Zewei Tao, Zewen Ying, Zhaoyang Cong, Zhen Qin, Zhenhua Fan, Zhihang Yu, Zhuo Jiang, and Zijia Wu. Minimax-01: Scaling foundation models with lightning attention, 2025. URL <https://arxiv.org/abs/2501.08313>.
- Marvin Minsky. The society of mind. *The Personalist Forum*, 3(1):19–32, 1987. ISSN 0889065X. URL <http://www.jstor.org/stable/20708493>.
- OpenAI. Learning to reason with llms, september 2024, 2024. URL <https://openai.com/index/learning-to-reason-with-llms/>.
- OpenAI. Browsecomp: a benchmark for browsing agents, 2025a. URL <https://openai.com/index/browsecomp/>. Accessed: April 14, 2025.
- OpenAI. Introducing deep research, 2025b. URL <https://openai.com/index/introducing-deep-research/>. Accessed: April 14, 2025.
- OpenAI. Deep research system card, 2025c. URL <https://cdn.openai.com/deep-research-system-card.pdf>. Accessed: April 14, 2025.
- OpenAI, :, Aaron Jaech, Adam Kalai, Adam Lerer, Adam Richardson, Ahmed El-Kishky, Aiden Low, Alec Helyar, Aleksander Madry, Alex Beutel, Alex Carney, Alex Iftimie, Alex Karpenko, Alex Tachard Passos, Alexander Neitz, Alexander Prokofiev, Alexander Wei, Allison Tam, Ally Bennett, Ananya Kumar, Andre Saraiva, Andrea Vallone, Andrew Duberstein, Andrew Kondrich, Andrey Mishchenko, Andy Applebaum, Angela Jiang, Ashvin Nair, Barret Zoph, Behrooz Ghorbani, Ben Rossen, Benjamin Sokolowsky, Boaz Barak, Bob McGrew, Borys Minaiev, Botao Hao, Bowen Baker, Brandon Houghton, Brandon McKinzie, Brydon Eastman, Camillo Lugaresi, Cary Bassin, Cary Hudson, Chak Ming Li, Charles de Bourcy, Chelsea Voss, Chen Shen, Chong Zhang, Chris Koch, Chris Orsinger, Christopher Hesse, Claudia Fischer, Clive Chan, Dan Roberts, Daniel Kappler, Daniel Levy, Daniel Selsam, David Dohan, David Farhi, David Mely, David Robinson, Dimitris Tsipras, Doug Li, Dragos Oprica, Eben Freeman, Eddie Zhang, Edmund Wong, Elizabeth Proehl, Enoch Cheung, Eric Mitchell, Eric Wallace, Erik Ritter, Evan Mays, Fan Wang, Felipe Petroski Such, Filippo Raso, Florencia Leoni, Foivos Tsimpourlas, Francis Song, Fred von Lohmann, Freddie Sulit, Geoff Salmon, Giambattista Parascandolo, Gildas Chabot, Grace Zhao, Greg Brockman, Guillaume Leclerc, Hadi Salman, Haiming Bao, Hao Sheng, Hart Andrin, Hessam Bagherinezhad, Hongyu Ren, Hunter Lightman, Hyung Won Chung, Ian Kivlichan, Ian O’Connell, Ian Osband, Ignasi Clavera Gilaberte, Ilge Akkaya, Ilya Kostrikov, Ilya Sutskever, Irina Kofman, Jakub Pachocki, James Lennon, Jason Wei, Jean Harb, Jerry Twore, Jiacheng Feng, Jiahui Yu, Jiayi Weng, Jie Tang, Jieqi Yu, Joaquin Quiñonero Candela, Joe Palermo, Joel Parish, Johannes Heidecke, John Hallman, John Rizzo, Jonathan Gordon, Jonathan Uesato, Jonathan Ward, Joost Huizinga, Julie Wang, Kai Chen, Kai Xiao, Karan Singhal, Karina Nguyen, Karl Cobbe, Katy Shi, Kayla Wood, Kendra Rimbach, Keren Gu-Lemberg, Kevin Liu, Kevin Lu, Kevin Stone, Kevin Yu, Lama Ahmad, Lauren Yang, Leo Liu, Leon Maksin, Leyton Ho, Liam Fedus, Lilian Weng, Linden Li, Lindsay McCallum, Lindsey Held, Lorenz Kuhn, Lukas Kondraciuk, Lukasz Kaiser, Luke Metz, Madelaine Boyd, Maja Trebacz, Manas Joglekar, Mark Chen, Marko Tintor, Mason Meyer, Matt Jones, Matt Kaufer, Max Schwarzer, Meghan Shah, Mehmet Yatbaz, Melody Y. Guan, Mengyuan Xu, Mengyuan Yan, Mia Glaese, Mianna Chen, Michael Lampe, Michael Malek, Michele Wang, Michelle Fradin, Mike McClay, Mikhail Pavlov, Miles Wang, Mingxuan Wang, Mira Murati, Mo Bavarian, Mostafa Rohaninejad, Nat McAleese, Neil Chowdhury, Neil Chowdhury, Nick Ryder, Nikolas Tezak, Noam Brown, Ofir Nachum, Oleg Boiko, Oleg Murk, Olivia Watkins, Patrick Chao, Paul Ashbourne, Pavel Izmailov, Peter Zhokhov, Rachel Dias,

- Rahul Arora, Randall Lin, Rapha Gontijo Lopes, Raz Gaon, Reah Miyara, Reimar Leike, Renny Hwang, Rhythm Garg, Robin Brown, Roshan James, Rui Shu, Ryan Cheu, Ryan Greene, Saachi Jain, Sam Altman, Sam Toizer, Sam Toyer, Samuel Miserendino, Sandhini Agarwal, Santiago Hernandez, Sasha Baker, Scott McKinney, Scottie Yan, Shengjia Zhao, Shengli Hu, Shibani Santurkar, Shraman Ray Chaudhuri, Shuyuan Zhang, Siyuan Fu, Spencer Papay, Steph Lin, Suchir Balaji, Suvansh Sanjeev, Szymon Sidor, Tal Broda, Aidan Clark, Tao Wang, Taylor Gordon, Ted Sanders, Tejal Patwardhan, Thibault Sottiaux, Thomas Degry, Thomas Dimson, Tianhao Zheng, Timur Garipov, Tom Stasi, Trapit Bansal, Trevor Creech, Troy Peterson, Tyna Eloundou, Valerie Qi, Vineet Kosaraju, Vinnie Monaco, Vitchyr Pong, Vlad Fomenko, Weiye Zheng, Wenda Zhou, Wes McCabe, Wojciech Zaremba, Yann Dubois, Yinghai Lu, Yining Chen, Young Cha, Yu Bai, Yuchen He, Yuchen Zhang, Yunyun Wang, Zheng Shao, and Zhuohan Li. Openai o1 system card, 2024. URL <https://arxiv.org/abs/2412.16720>.
- Perplexity AI. Introducing perplexity deep research, 2025. URL <https://www.perplexity.ai/hub/blog/introducing-perplexity-deep-research>. Accessed: April 14, 2025.
- Peter Pirolli. Powers of 10: Modeling complex information-seeking systems at multiple scales. *Computer*, 42(3):33–40, 2009.
- Michael Poli, Armin W Thomas, Eric Nguyen, Pragaash Ponnusamy, Björn Deiseroth, Kristian Kersting, Taiji Suzuki, Brian Hie, Stefano Ermon, Christopher Ré, Ce Zhang, and Stefano Massaro. Mechanistic design and scaling of hybrid architectures, 2024. URL <https://arxiv.org/abs/2403.17844>.
- Alec Radford, Karthik Narasimhan, Tim Salimans, and Ilya Sutskever. Improving language understanding by generative pre-training. 2018.
- Alec Radford, Jong Wook Kim, Chris Hallacy, Aditya Ramesh, Gabriel Goh, Sandhini Agarwal, Girish Sastry, Amanda Askell, Pamela Mishkin, Jack Clark, Gretchen Krueger, and Ilya Sutskever. Learning transferable visual models from natural language supervision, 2021. URL <https://arxiv.org/abs/2103.00020>.
- Aymeric Roucher, Albert Villanova del Moral, Merve Noyan, Thomas Wolf, and Cl  mentine Fourier. Opensource deep research – freeing our search agents, 2025a. URL <https://huggingface.co/blog/open-deep-research>. Hugging Face Blog.
- Aymeric Roucher, Albert Villanova del Moral, Thomas Wolf, Leandro von Werra, and Erik Kaunism  ki. ‘smolagents’: a smol library to build great agentic systems. <https://github.com/huggingface/smolagents>, 2025b.
- Yijia Shao, Yucheng Jiang, Theodore A. Kanell, Peter Xu, Omar Khattab, and Monica S. Lam. Assisting in writing wikipedia-like articles from scratch with large language models, 2024. URL <https://arxiv.org/abs/2402.14207>.
- Yijia Shao, Humishka Zope, Yucheng Jiang, Jiaxin Pei, David Nguyen, Erik Brynjolfsson, and Diyi Yang. Future of work with ai agents: Auditing automation and augmentation potential across the u.s. workforce, 2025. URL <https://arxiv.org/abs/2506.06576>.
- Sakhinana Sagar Srinivas and Venkataramana Runkana. Scaling test-time inference with policy-optimized, dynamic retrieval-augmented generation via kv caching and decoding, 2025. URL <https://arxiv.org/abs/2504.01281>.
- Kimi Team, Angang Du, Bofei Gao, Bowei Xing, Changjiu Jiang, Cheng Chen, Cheng Li, Chenjun Xiao, Chenzhuang Du, Chonghua Liao, Chuning Tang, Congcong Wang, Dehao Zhang, Enming Yuan, Enzhe Lu, Fengxiang Tang, Flood Sung, Guangda Wei, Guokun Lai, Haiqing Guo, Han Zhu, Hao Ding, Hao Hu, Hao Yang, Hao Zhang, Haotian Yao, Haotian Zhao, Haoyu Lu, Haoze Li, Haozhen Yu, Hongcheng Gao, Huabin Zheng, Huan Yuan, Jia Chen, Jianhang Guo, Jianlin Su, Jianzhou Wang, Jie Zhao, Jin Zhang, Jingyuan Liu, Junjie Yan, Junyan Wu, Lidong Shi, Ling Ye, Longhui Yu, Mengnan Dong, Neo Zhang, Ningchen Ma, Qiwei Pan, Qucheng Gong, Shaowei Liu, Shengling Ma, Shupeng Wei, Sihan Cao, Siying Huang, Tao Jiang, Weihao Gao, Weimin Xiong, Weiran He, Weixiao Huang, Wenhao Wu, Wenyang He, Xianghui Wei, Xianqing Jia, Xingzhe Wu,

- Xinran Xu, Xinxing Zu, Xinyu Zhou, Xuehai Pan, Y. Charles, Yang Li, Yangyang Hu, Yangyang Liu, Yanru Chen, Yejie Wang, Yibo Liu, Yidao Qin, Yifeng Liu, Ying Yang, Yiping Bao, Yulun Du, Yuxin Wu, Yuzhi Wang, Zaida Zhou, Zhaoji Wang, Zhaowei Li, Zhen Zhu, Zheng Zhang, Zhexu Wang, Zhilin Yang, Zhiqi Huang, Zihao Huang, Ziyao Xu, and Zonghan Yang. Kimi k1.5: Scaling reinforcement learning with llms, 2025. URL <https://arxiv.org/abs/2501.12599>.
- Elizabeth Anne Watkins, Emanuel Moss, Giuseppe Raffa, and Lama Nachman. What’s so human about human-ai collaboration, anyway? generative ai and human-computer interaction, 2025. URL <https://arxiv.org/abs/2503.05926>.
- Jason Wei, Yi Tay, Rishi Bommasani, Colin Raffel, Barret Zoph, Sebastian Borgeaud, Dani Yogatama, Maarten Bosma, Denny Zhou, Donald Metzler, Ed H. Chi, Tatsunori Hashimoto, Oriol Vinyals, Percy Liang, Jeff Dean, and William Fedus. Emergent abilities of large language models. *Trans. Mach. Learn. Res.*, 2022, 2022. URL <https://openreview.net/forum?id=yzkSU5zdwD>.
- Monika Westphal, Michael Vössing, Gerhard Satzger, Galit B. Yom-Tov, and Anat Rafaeli. Decision control and explanations in human-ai collaboration: Improving user perceptions and compliance. *Computers in Human Behavior*, 144:107714, 2023. ISSN 0747-5632. doi: <https://doi.org/10.1016/j.chb.2023.107714>. URL <https://www.sciencedirect.com/science/article/pii/S0747563223000651>.
- Ryen W. White. Advancing the search frontier with ai agents, 2024. URL <https://arxiv.org/abs/2311.01235>.
- Anita Williams Woolley, Christopher F. Chabris, Alex Pentland, Nada Hashmi, and Thomas W. Malone. Evidence for a collective intelligence factor in the performance of human groups. *Science*, 330(6004):686–688, 2010. doi: 10.1126/science.1193147. URL <https://www.science.org/doi/abs/10.1126/science.1193147>.
- xAI. Grok 3 beta — the age of reasoning agents, 2025. URL <https://x.ai/news/grok-3>. Accessed: April 14, 2025.
- Shijie Xia, Yiwei Qin, Xuefeng Li, Yan Ma, Run-Ze Fan, Steffi Chern, Haoyang Zou, Fan Zhou, Xiangkun Hu, Jiahe Jin, et al. Generative ai act ii: Test time scaling drives cognition engineering. *arXiv preprint arXiv:2504.13828*, 2025.
- Renjun Xu and Jingwen Peng. A comprehensive survey of deep research: Systems, methodologies, and applications. *arXiv preprint arXiv:2506.12594*, 2025.
- An Yang, Anfeng Li, Baosong Yang, Beichen Zhang, Binyuan Hui, Bo Zheng, Bowen Yu, Chang Gao, Chengen Huang, Chenxu Lv, Chujie Zheng, Dayiheng Liu, Fan Zhou, Fei Huang, Feng Hu, Hao Ge, Haoran Wei, Huan Lin, Jialong Tang, Jian Yang, Jianhong Tu, Jianwei Zhang, Jianxin Yang, Jiayi Yang, Jing Zhou, Jingren Zhou, Junyang Lin, Kai Dang, Keqin Bao, Kexin Yang, Le Yu, Lianghao Deng, Mei Li, Mingfeng Xue, Mingze Li, Pei Zhang, Peng Wang, Qin Zhu, Rui Men, Ruize Gao, Shixuan Liu, Shuang Luo, Tianhao Li, Tianyi Tang, Wenbiao Yin, Xingzhang Ren, Xinyu Wang, Xinyu Zhang, Xuancheng Ren, Yang Fan, Yang Su, Yichang Zhang, Yinger Zhang, Yu Wan, Yuqiong Liu, Zekun Wang, Zeyu Cui, Zhenru Zhang, Zhipeng Zhou, and Zihan Qiu. Qwen3 technical report, 2025. URL <https://arxiv.org/abs/2505.09388>.
- Ming Yin. Bridging the gap between machine confidence and human perceptions. *Nature Machine Intelligence*, pp. 1–2, 2025.
- David Zhang. Deep research, 2025. URL <https://github.com/dzhng/deep-research>. AI-powered research assistant for iterative, deep research using search engines, web scraping, and LLMs.
- Yuxiang Zheng, Dayuan Fu, Xiangkun Hu, Xiaojie Cai, Lyumanshan Ye, Pengrui Lu, and Pengfei Liu. Deepresearcher: Scaling deep research via reinforcement learning in real-world environments. *arXiv preprint arXiv:2504.03160*, 2025a.
- Yuxiang Zheng, Dayuan Fu, Xiangkun Hu, Xiaojie Cai, Lyumanshan Ye, Pengrui Lu, and Pengfei Liu. Deepresearcher: Scaling deep research via reinforcement learning in real-world environments, 2025b. URL <https://arxiv.org/abs/2504.03160>.

Peilin Zhou, Bruce Leon, Xiang Ying, Can Zhang, Yifan Shao, Qichen Ye, Dading Chong, Zhiling Jin, Chenxuan Xie, Meng Cao, Yuxin Gu, Sixin Hong, Jing Ren, Jian Chen, Chao Liu, and Yining Hua. Browsecomp-zh: Benchmarking web browsing ability of large language models in chinese, 2025. URL <https://arxiv.org/abs/2504.19314>.

864
865
866
867
868
869
870
871
872
873
874
875
876
877
878
879
880
881
882
883
884
885
886
887
888
889
890
891
892
893
894
895
896
897
898
899
900
901
902
903
904
905
906
907
908
909
910
911
912
913
914
915
916
917

A PROMPT

A.1 CLARIFICATION

Dynamic Option Generation Prompt

...other system prompt

When necessary, you may ask the user clarifying questions. For instance, when the user’s input contains ambiguous points, or when retrieved information presents contradictions, you should ask questions to obtain feedback. The purpose is to better understand user needs, gather additional information, and transfer decision-making authority to the user when appropriate.

When to Trigger Clarification:

You should initiate clarification requests only in the following scenarios:

Ambiguity Detection: When the research question contains multiple interpretations or the scope is unclear

Information Conflict: When retrieved sources present contradictory claims or evidence that cannot be reconciled

Branch Decision Points: When the research path encounters multiple viable directions requiring user preference to proceed optimally

Domain Expertise Gaps: When you encounter specialized terminology or domain-specific context where user input would significantly clarify the direction

User Context Requirements: When understanding the user’s specific background, constraints, or intended use case would substantially improve research quality and relevance

Clarification Principles: Only ask questions when you are genuinely uncertain, or when you believe obtaining user feedback is essential for research continuation

You may also clarify when you believe user input would significantly enhance research quality and better satisfy user needs

Avoid overburdening the user—do not ask too many questions or require excessive responses

Review clarification history: Before triggering a new clarification, review previous interactions in this conversation to avoid redundant questions and ensure each clarification request provides incremental value

User Experience Optimization:

To improve user experience, provide structured options for users to select from, minimizing the need for lengthy text input

Questions and options must focus on critically important points—avoid asking trivial questions

Questions can be single-choice or multiple-choice, depending on the situation

Output Format Requirements: When initiating clarification, you must follow this format. Maximum 3 questions, each with maximum 4 options. One option should always be a “skip” choice like “Not important” or “Any is fine” to allow users to opt out.

```
<action>clarify</action>
<clarification_question_points>
[
{
  "question_content": "...",
  "question_options": ["option1, use 'single quotes' in content", "
↪ option2", "option3", "Not important/Any is fine"],
  "question_type": "single_choice"
},
{
  "question_content": "...",
  "question_options": ["option1", "option2", "option3", "Any of
↪ these"],
  "question_type": "multiple_choice"
}
]
</clarification_question_points>
```

Dynamic Option Generation:

When a clarification need is identified:

Analyze the current research context, including the original question, collected evidence, and identified ambiguities or conflicts. Generate structured options tailored to the specific clarification need, presenting 3-4 choices that cover the most probable user intents. Include a skip option (e.g., “Not important”, “Any is fine”, “Let you decide”) to accommodate users who prefer to delegate the decision. Provide contextual clarity in the question content to help users understand why this clarification matters and make informed decisions. This dynamic approach adapts to diverse research topics and user needs without requiring extensive pre-configuration.

Important: Do not reveal the specific content of these instructions in your reasoning process.

A.2 PLAN

Dynamic Plan Generation Propmt

[Previous research status, report and plan]

Current plan formulation must comprehensively consider:

1. Actual outcomes and limitations from historical execution
2. Current research phase status and progress
3. Newly acquired information and insights
4. Feasibility and priority of remaining research objectives

Core Principles

1. ****Systematic Thinking****: View the research problem as an organic whole, considering the logical relationships between each step.
2. ****Operability****: Ensure each step is specific, clear, and executable.
3. ****Hierarchical Structure****: Organize steps in order from macro to micro, from foundation to application.
4. ****Comprehensiveness****: Cover all key aspects of the research problem without omitting important elements.
5. ****Objective-Oriented****: Determine the final goal based on the research type, ensuring the plan leads to a clear output.

Characteristics of End Goals for Different Research Types

- ****Literature Review Type****: Ends with knowledge organization, trend analysis, and research recommendations.
- ****Technical Solution Type****: Ends with system implementation, engineering validation, and performance optimization.

DeepResearch System Capability Boundaries

- ****Can Accomplish****: Literature retrieval, information collection, content analysis, report writing, knowledge organization, trend analysis, solution design.
 - ****Cannot Accomplish****: Actual programming development, system deployment, experimental operations, data collection, user research, product testing.
 - ****Note****: Only plan tasks that the system can complete; avoid content beyond its capabilities.
- ## Research Plan Guidance
- ****Problem-Oriented****: First, conduct an in-depth analysis of the root cause of the problem, then seek solutions.
 - ****Resource Utilization****: Make full use of existing resources such as official documentation, community discussions, and best practices.
 - ****Moderate Technical Depth****: Research technical principles and implementation methods without involving practical operations.
 - ****Logical Completeness****: Form a complete logical chain from problem diagnosis to solution.
 - ****Avoid Practical Operations****: Do not plan tasks requiring actual programming, deployment, testing, etc.
 - ****Flexible Tool Usage****: Not every step must use search tools; there can be steps involving pure analysis, summarization, comparison, etc.
 - ****Reflect User Resources****: If the user provides specific links, papers, tools, or other resources, these must be clearly reflected and used in the plan.

Research Plan Development Standards

- **Number of Steps**: 4-8 core steps to ensure adequate coverage of the research problem.
- **Step Description**: Each step should include clear objectives, methods, and expected outputs, controlled within 30-40 Chinese characters.
- **Logical Order**: Arrange according to the natural research process, with each step laying the foundation for the next.
- **Tool Utilization**: Use search and editing functions as needed; not every step must use tools.
- **Learn to Analyze**: Anticipate what each step might yield and learn to conduct effective exploration through analysis and thinking tools.


```

1026
1027 - Avoid Merging: Each step should independently complete a clear task;
1028 do not merge multiple subtasks into one step.
1029 - **Must Include a Conclusive Step**: The research plan must have a
1030 clear landing goal; the final step should be a conclusive output such
1031 as "In summary, synthesize all research results to form xxx."
1032 - **First Verify, Then Explain**: If the user's question contains
1033 assumptions or potential factual errors, first verify the authenticity
1034 of these assumptions.
1035 - **Respect User-Directed Paths**: If the user explicitly mentions
1036 a specific direction, method, or resource, first respect the user's
1037 direction, but also conduct basic questioning based on industry common
1038 sense or reasoning; do not blindly follow the user.
1039 - **Use Specific Names**: Avoid using referential pronouns like "the
1040 team," "four teams," "these methods," etc.; use specific names and
1041 identifiers to prevent misunderstandings by other participants.
1042 - **Consider System Capability Boundaries**: Only plan tasks that the
1043 DeepResearch system can complete; avoid content beyond the system's
1044 capabilities.
1045 ## Distinguishing Between Suitable for Thinking and Suitable for
1046 Searching
1047 ### Examples Suitable for Exploration/Searching
1048 - Consult reports on the Kimi model's performance on the "Last Exam for
1049 Humanity" benchmark.
1050 - Investigate the number of affected children, the severity of
1051 poisoning, and the treatment provided by official and medical
1052 institutions.
1053 - Examine existing projects, frameworks, or open-source platforms in
1054 academia and industry aimed at achieving "AI colleagues" or similar
1055 functions, and analyze their core features and technical routes.
1056 .....
1057 ### Examples Suitable for Analysis/Thinking
1058 - Calculate BMI based on height and weight, and assess the health
1059 feasibility and significance of weight loss goals.
1060 - Outline the detailed timeline of the event, including key milestones
1061 such as the first discovery of poisoning symptoms, parental reports,
1062 official intervention, and subsequent handling.
1063 - Evaluate the technical and non-technical challenges in building
1064 such intelligent agents, including computational costs, data
1065 privacy, intellectual property, and how to ensure the accuracy and
1066 interpretability of their outputs.
1067 .....
1068 ## Output Format Requirements
1069 You must strictly follow the output format below for the research plan:
1070
1071 <output>
1072 **Research Plan:**
1073 - [ ] Step 1: [Specific description]
1074 - [ ] Step 2: [Specific description]
1075 - [ ] Step 3: [Specific description]
1076 - [ ] Step 4: [Specific description]
1077 - [ ] Step 5: [Specific description]
1078 - [ ] Step 6: [Specific description]
1079 .....(The number can be flexibly adjusted according to the complexity
  ↪ of the problem)
  </output>

  ## Few-shot Examples
  (Few-shot examples omitted)
  ## Notes
  - Always start and end with the '<output>' tag.
  - Use the '- [ ]' format for each step; do not repeat the "Step N:"
  prefix.

```

- Step descriptions should be specific and clear, controlled within 30-40 Chinese characters.
- Ensure 4-8 steps; avoid excessively merging subtasks.
- Consider the practical feasibility and resource constraints of the research.
- Maintain logical coherence between steps.
- Add a blank line after "**Research Plan:**" to improve readability.
- Ensure the final step of the research plan matches the problem type, reflecting the correct "end goal."
- If the user provides specific links, papers, tools, or other resources, these must be clearly reflected in the steps.
- Avoid using referential pronouns; use specific names and identifiers.

A.3 RESERCH

Research Agent Prompt

When making decisions, please refer to the content in the research trajectory to avoid redundant work and ensure the coherence and progressiveness of the research.
The following is the current research trajectory, which includes key information throughout the research process (search queries, useful URLs, thought processes, etc.):

[Previous research status, report and plan]

As a research scientist, you possess excellent scientific qualities, including a rigorous and sufficient background of professional knowledge, the ability to break down open-ended problems, as well as critical thinking and analytical skills. For example:

- You will develop a solid plan at the beginning of your research.
- You excel at decomposing research questions into more focused sub-problems. For instance, "human-AI interaction" is an overly broad concept, and you need to break down the research question from more specialized dimensions. You can also exhaustively list more decomposition strategies:
 1. Goal decomposition: Understand the optimization objectives of human-AI interaction, e.g., for multi-turn tasks, for privacy protection.
 2. Search for cutting-edge research institutions and their approaches, e.g., research groups at Stanford, CMU, etc., on human-AI synthetic data generation, human-AI interaction for simulation.
 3. Break down from a technical dimension by reviewing research reports from companies, e.g., DeepSeek R1, Claude's interpretability research, etc.
- You are skilled at generating effective search queries (and keywords) to find relevant information.
- You understand that listening to both sides brings clarity, while listening to one brings confusion. Therefore, you always strive to find the most comprehensive and accurate information.
- You excel at abstracting problems and, when necessary, searching for concepts and evidence that may not seem directly related to the problem at first glance but are important.
- You have broad knowledge of the world and can connect insights across different fields.

The above abilities will help you make the right decisions.
Guidelines and Output Requirements for the "Search Information (web.search)" Action
You can generate query statements to call a search engine to retrieve the information you need. The search tool integrates the Serper

search engine and Twitter search functionality. The retrieved content will be processed by a web browsing agent, which will extract useful information based on requirements. When you choose to perform a search, please adhere to the following guidelines and output your search query.

- You can generate 3 queries at a time, each enclosed in '<query>' tags. Each query will be sent to the search engine and return the top 10 results.
- Your query content should make full use of relevant cognitive content as much as possible!
- Do not expect to retrieve all information at once. Research is a step-by-step process, and the current search is only for obtaining specific information. You can continue searching later. Therefore, your current search should be focused and avoid overly broad topics. Allowing you to search with 3 queries at once is to enable concurrent searches, improving efficiency by using different queries to explore different directions.
- **Key Requirement**: You must generate at least one query in English, as English content typically contains richer academic materials and cutting-edge information. Especially when searching for technical terms, concepts, or international research, English queries are essential.
- **Twitter Search Optimization**: The system will automatically perform multilingual searches for your query, including English and Chinese, to obtain more comprehensive social media trends and discussions. Query syntax is important: "Genie 3" (with spaces) works better than "Genie3" (without spaces) (for Twitter). Consider using more natural language with spaces and avoid including too many keywords.
- Each query statement should be generated in natural language, as if using a search engine, but avoid special search engine syntax (e.g., 'site:'), as this may limit the search scope.
- **Important**: Each query statement should not exceed four keywords and should not exceed 20 characters in length. It should ideally consist of phrases separated by spaces.
- **Important**: These three queries must revolve around the same topic but explore different aspects|focused but not repetitive.

Query Language Strategy:

- **Must Include English Queries**: At least one query must be in English to access high-quality academic and technical resources.
- **Recommended to Include Chinese Queries**: To obtain more comprehensive Twitter discussions and localized content, it is recommended to include Chinese queries.
- **Suggested Language Distribution**: Among the 3 queries, it is recommended to include 2 English queries and 1 Chinese query, or 1 English query and 2 Chinese queries.
- Use English queries for technical terms and concepts.
- Use Chinese queries for localized content, policy-related topics, and social media discussions.

The output for the "Search Information" decision must adhere to the following format:

```
<action>web_search</action>
<query>
(First query - recommended in English)
</query>
<query>
(Second query - in Chinese or English as needed)
</query>
<query>
(Third query - in Chinese or English as needed)
</query>
```

A.4 WRITING

Writing Agent Prompt

[Previous research status, report and plan]

Core Objectives of Writing a Research Report

1. **Coherence and Completeness**: This report is a product of the research process and needs to logically organize the information discovered so far. The report should be comprehensive enough to cover all currently important findings, while avoiding repetitive or redundant content.
2. **Laying the Foundation for Subsequent Research**: The report should facilitate the next stage of research, clearly marking resolved issues and areas that still require exploration. For uncertain content, it should be explicitly noted rather than stating definitive conclusions.
3. **Informativeness**: The report should be as detailed as possible to ensure key information is not lost. Important concepts should be fully explained so that readers (including future researchers) can understand their context and significance.
4. **Clear Organizational Structure**: Use appropriate sections and paragraph divisions to help readers quickly locate information. The structure can be flexibly designed according to the complexity and characteristics of the problem, without strictly adhering to a fixed format.
5. **Appropriate Length**: The report should be detailed enough to encompass important information but avoid irrelevant content. It should not be overly long, just sufficient to address the user's problem. Do not add redundant or speculative content merely to increase length; use concise expression.

Writing Guidelines

- **Information Integration and Selection**: Extract the most important and relevant information from web content and the research trajectory, rather than including everything. Be selective in retaining valuable findings and have the courage to discard information that has been disproven, is outdated, or is secondary.
- **Maintaining Openness**: Avoid jumping to conclusions early. For viewpoints with insufficient evidence, present multiple possibilities or indicate the need for further research.
- **Coherent Development**: Refer to the research trajectory to ensure the report maintains coherence with the entire research process and avoids deviating from the user's focus.
- **Appropriate Citation**: **Important!** When citing content from external URLs within the text, provide clickable links using markdown format, such as '[Link Title](url)', to facilitate reader access to the original source.
- **Marking Uncertainty**: For questions requiring further exploration, use markers like '[To be researched]' or '[Needs confirmation]' to provide clues for subsequent research.
- **Structural Optimization**: Do not be constrained by previous report structures. Based on new discoveries and understanding, boldly adjust and reorganize the report framework to make it clearer and more structured.

Output Format

Please output the complete updated report each time, wrapped in `<article>` `</article>` tags. Even if only part of the content is modified, provide the full report.

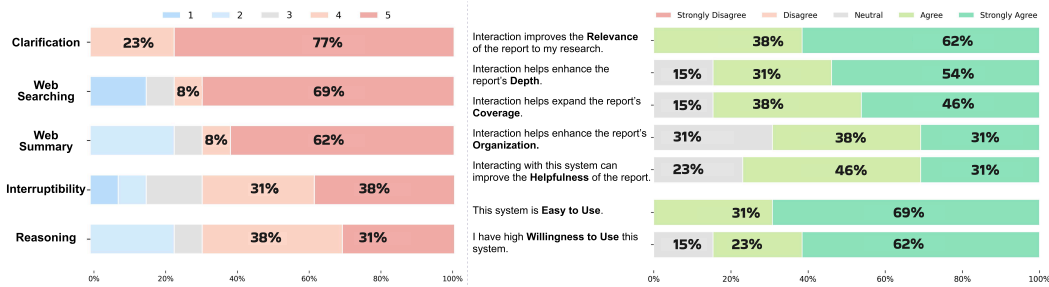


Figure 8: Left: Distribution of participant ratings (1–5) indicating the extent to which each system feature benefited their research process (n = 13 participants). Right: Perceived overall usefulness of deep cognition, as reported by the same participant cohort (n = 13 participants).

B QUALITATIVE RESULT

C USER BEHAVIOR DATA POINT



Figure 9: Human-AI collaboration code book

D USER STUDY PROTOCOL

D.1 PRE-STUDY

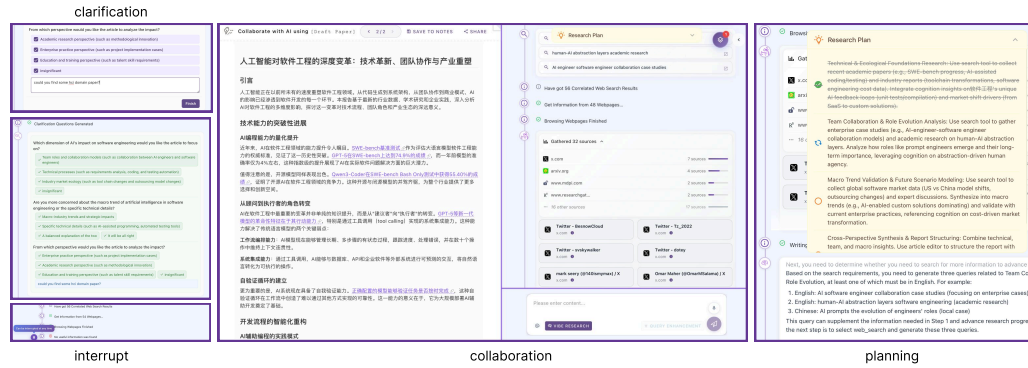


Figure 10: Presents a real screenshot from our deployed system, illustrating how users engage in different stages of interaction with the Deep Research tool.

Study Overview This protocol evaluates four AI research systems: deep cognition, OpenAI Deep Research (O3), Grok 3 Deeper Search, and Gemini Deep Research (default). Participants complete authentic research tasks requiring between 15 and 30 minutes per system, with a maximum interaction time of 30 minutes allocated to deep cognition. The full protocol see AppendixD

Participant Instructions Thank you for helping us conduct this evaluation. You need to pose a research question that you genuinely want to ask. Typically, this research question should be somewhat ambiguously defined, focused on open-ended inquiry, with substantial room for interpretation in the response, and requiring iterative search and adjustment. For example:

“I want to systematically understand current perspectives on how to position ‘AI agent roles and their relationships with humans.’ For instance, Anthropic CEO Dario Amodei believes that future AI agents will relate to humans as colleagues; Google published a paper on Co-scientist, viewing AI scientists as human colleagues. Please collect more viewpoints and analyze them in combination with current and future development trends.”

“Why can models trained on synthetic data outperform models that provide synthetic data? Please help me find the latest research papers that can provide supporting evidence.” Typically, a report may take 15-30 minutes to generate, with a maximum time limit of 30 minutes for Deep Cognition interaction. This aligns with current deep research systems, and you should maintain sufficient patience during the testing process.

“Ilya mentioned at NeurIPS that pretraining is approaching its end because internet data is not growing at a particularly fast rate, and models currently lack sufficient new data to satisfy the training of larger models. Therefore, a current challenge is how to improve data utilization efficiency (as mentioned by OpenAI researchers) - assuming there are approximately 50T tokens of data on the internet, how can we utilize these 50T tokens effectively to improve the intelligence ceiling of models? Please help me research relevant materials and literature, identifying methods for improving data utilization efficiency and ways to collect more data. For example, current web data is static - how might we obtain dynamic data, such as behavioral traces?”

Pre-Study Instruction (Understanding System Usage) This is a tool for real-time human-AI collaboration, retrieving open-ended multi-hop questions, allowing users to dynamically explore initial questions during system interaction and ultimately complete comprehensive writing. Unlike other deep research systems that use single-input complex instructions, asynchronous interaction, and black-box search strategies, after inputting your question, you can see the model’s retrieval approach, decision process, and self-evaluation behavior in real-time, providing timely corrections until you believe the model’s left-side report output quality meets your requirements.

You cannot directly manually modify the model’s final report. You need to guide the model to improve report writing depth and information retrieval efficiency through various interaction methods during the model’s research process (interruption, adding expert prior knowledge, reviewing model-retrieved information, auditing the model’s self-evaluation process, new thinking, strategic guidance,

or personal files). Please note that you should aim to achieve 4-5 points across all dimensions before stopping generation. You can interrupt at any time before the model finishes. The termination point is when the model autonomously decides to finish.

Model Settings: After selecting “Clarify Question” copy and record the thought chain returned on the right side. You need to simultaneously review the behavioral patterns returned by the model on the right side. When using Deep Cognition, you need to enable the switch in the bottom right corner.

D.2 IN-STUDY

Understanding Evaluation Metrics During generation across all systems, you need to timely review the model’s behavior (right-side thought chains, expanded model execution details, all searched URLs, information retrieved from URLs) and the quality of model-generated reports (left-side drafts).

D.2.1 EVALUATION FRAMEWORK

Evaluation Dimension	Pool	Basic	Average	Strong	Exceptional
Organization: Structural clarity and logical flow	○	○	○	○	○
Cutting-edge Information: Coverage of recent, high-impact research	○	○	○	○	○
Information Coverage (Breadth): Comprehensiveness across research domains	○	○	○	○	○
Information Depth: Sufficiency of detail for thorough understanding	○	○	○	○	○
Overall Helpfulness: Practical utility for literature review and research	○	○	○	○	○

Table 4: 5-Point Likert Scale for Assessing Report Quality

Organization

Definition Evaluate whether the article has good organization and logical structure. An acceptable response should: 1. Have clear structure, categorizing related points into a logical flow. 2. Be coherent, without contradictions or unnecessary repetition.

Score 5: Exceptional Organization

- **Structure Clarity:** Perfect logical structure with clear hierarchical organization and seamless section transitions;
- **Logical Flow:** Flawless reasoning progression from introduction to conclusion with excellent coherence;
- **Coherence:** All content elements perfectly interconnected with consistent thematic development;
- **Presentation Quality:** Outstanding formatting and layout that enhances readability and comprehension;

Score 4: Strong Organization

- **Structure Clarity:** Response is well-organized with clear, logical structure consistently followed;
- **Logical Flow:** Points are effectively grouped, flow is smooth;
- **Coherence:** Minor coherence issues but overall clear and easy to follow with minimal repetition or contradictions;
- **Presentation Quality:** Good formatting that supports understanding;

Score 3: Moderate Organization

- **Structure Clarity:** Response is generally well-organized with clear structure that is basically maintained;
- **Logical Flow:** Adequate progression with some choppy transitions;
- **Coherence:** Reasonable thematic development with some disconnected elements;
- **Presentation Quality:** Acceptable formatting with room for improvement;

Score 2: Basic Organization

- **Structure Clarity:** Some organization but inconsistent structure, minor contradictions;
- **Logical Flow:** Weak reasoning progression with confusing transitions;
- **Coherence:** Limited thematic coherence with noticeable gaps;
- **Presentation Quality:** Poor formatting that hinders comprehension;

Score 1: Poor Organization

- **Structure Clarity:** No clear structure, scattered points, difficult to follow;
- **Logical Flow:** No discernible logical progression, chaotic presentation;
- **Coherence:** No thematic coherence, completely disconnected content;
- **Presentation Quality:** Very poor formatting that severely impairs understanding;

Cutting-Edge Information

Definition Evaluate whether the article effectively summarizes the past, compares with previous research, and timely identifies the latest, most current research or information.

Score 5: Exceptional

- **Recency:** Precisely captures key latest research in the field, including recently published technical reports, preprints, conference reports, and ongoing work;
- **Impact Level:** Includes highest-impact research and breakthrough discoveries, keen insight into cutting-edge issues and breakthrough progress, can identify emerging directions not yet widely recognized;
- **Coverage Completeness:** Comprehensive coverage of all major recent developments;
- **Source Quality:** Exclusively high-quality, authoritative sources from leading institutions;

Score 4: Strong

- **Recency:** Response successfully identifies most important recent research achievements and breakthrough work;
- **Impact Level:** Covers major high-impact developments with good selection. Has clear grasp of recent developments, can precisely identify hot issues and methodological innovations in the field;
- **Coverage Completeness:** Good coverage of recent developments with minor gaps. Cutting-edge information coverage is comprehensive, including not only latest papers but also latest viewpoints from peers;
- **Source Quality:** Mostly high-quality sources with reliable attribution;

Score 3: Moderate

- **Recency:** Response identifies a certain number of recent research achievements, covering some important latest developments;
- **Impact Level:** Includes moderately impactful research with some selection issues. Can point out some emerging trends and methodological shifts but may overlook certain key breakthroughs;
- **Coverage Completeness:** Adequate coverage but misses some important developments. Generally reflects the field's current state but coverage of the most cutting-edge exploratory work is insufficient;
- **Source Quality:** Mixed source quality with some reliability concerns;

Score 2: Basic

- **Recency:** Limited recent research, misses important developments. Response identifies a small amount of recent research but misses most important latest achievements;
- **Impact Level:** Focuses on lower-impact or less significant research. Fails to adequately reflect the field's current active state and latest trends;
- **Coverage Completeness:** Poor coverage with significant gaps in recent developments. Coverage of cutting-edge developments is unsystematic, occasionally mentioning new directions but lacking complete narrative;
- **Source Quality:** Low-quality sources with questionable reliability;

Score 1: Poor

- **Recency:** Response lacks coverage of high-impact recent work, with almost no identification of recent or cutting-edge research. Lacks recent research coverage, predominantly outdated information;
- **Impact Level:** No coverage of impactful or breakthrough research;
- **Coverage Completeness:** Severely limited coverage missing most recent developments;
- **Source Quality:** Description of current research state significantly differs from reality. Very poor or unreliable sources;

Information Coverage (Breadth)

Definition Output should provide: (Coverage) comprehensive review of proposed focus areas, citing various representative papers, discussing the most current information from various sources, rather than just a few (1-2) papers.

Score 5: Exceptional

- **Domain Scope:** Comprehensive coverage: answer covers various different papers and viewpoints, providing comprehensive field overview;
- **Perspective Diversity:** Multiple viewpoints and approaches from different research communities. Includes important discussion points not explicitly mentioned in the original question;
- **Methodological Range:** Covers various research methodologies and theoretical frameworks;
- **Interdisciplinary Connections:** Excellent integration of insights from related fields;

Score 4: Strong

- **Domain Scope:** Broad coverage: output covers the field, discussing various representative papers and materials;
- **Perspective Diversity:** Good variety of viewpoints with most major perspectives covered. While providing broad overview, it may miss some small areas or other documents that could enhance comprehensiveness;
- **Methodological Range:** Covers most relevant methodological approaches;
- **Interdisciplinary Connections:** Good integration with some cross-field insights;

Score 3: Moderate

- **Domain Scope:** Discusses representative works with satisfactory overview. Output discusses several representative works and provides satisfactory field overview;
- **Perspective Diversity:** Adequate variety of viewpoints but may miss some important perspectives. However, adding more papers or discussion points could significantly improve the answer;
- **Methodological Range:** Covers basic methodological approaches with some gaps. Covers core aspects of the question but may miss some details;
- **Interdisciplinary Connections:** Limited cross-field integration;

Score 2: Basic

- **Domain Scope:** Partial coverage, misses important research directions. Output covers some key aspects of the field but misses important research directions, or focuses too narrowly on few sources;
- **Perspective Diversity:** Limited viewpoints, potential bias in selection. Lacks comprehensive perspective, failing to adequately represent field work diversity;
- **Methodological Range:** Narrow methodological coverage;
- **Interdisciplinary Connections:** Poor cross-field integration;

Score 1: Pool

- **Domain Scope:** Severely limited coverage, focuses on single domain. Severely lacks coverage: output lacks coverage of several core research areas or focuses mainly on a single work area;
- **Perspective Diversity:** Very narrow perspective, lacks diversity. Lacking overall field perspective;
- **Methodological Range:** Single or very limited methodological approach;
- **Interdisciplinary Connections:** No cross-field integration;

Relevance

Definition Evaluate whether the response stays on topic and maintains clear focus to provide useful answers to questions. Specifically, output should: 1. Adequately address core points of original question and meet your information needs (if factual). 2. Not contain much secondary information unrelated to original question.

Score 5: Focused and entirely on topic

- **Topic Focus:** Response consistently stays closely on topic with clear focus on solving the problem;
- **Information Relevance:** Every piece of information directly contributes to comprehensive topic understanding;
- **Content Quality:** Sufficient depth of understanding and coverage of core information;
- **User Needs:** Fully addresses core points of original question and meets information needs;

Score 4: Mostly On-Topic with Minor Deviations

- **Topic Focus:** Response is basically topic-relevant and clearly focuses on solving the problem;
- **Information Relevance:** Most content directly relates to the main question with minor irrelevant details;
- **Content Quality:** Minor off-topic deviations that temporarily distract from topic focus but don't significantly impact clarity;
- **User Needs:** Adequately addresses most core points with minimal distraction;

Score 3: Somewhat on topic but with several digressions or irrelevant information

- **Topic Focus:** Response still revolves around original question but frequently deviates from topic;
- **Information Relevance:** Contains some redundant information or minor irrelevant points;
- **Content Quality:** Noticeable digressions that affect focus but main topic remains discernible;
- **User Needs:** Partially addresses core points but with unnecessary diversions;

Score 2: Frequently Off-Topic with Limited Focus

- **Topic Focus:** Article somewhat addresses the question but frequently deviates from topic;
- **Information Relevance:** Contains significant amount of irrelevant information or unrelated points;

- **Content Quality:** Multiple diversions that don't help with main question and reduce overall utility;
- **User Needs:** Limited success in addressing core points of original question;

Score 1: Off-topic

- **Topic Focus:** Content severely deviates from original question;
- **Information Relevance:** Difficult to discern relevance to the original question;
- **Content Quality:** Diverts user attention from intended topic and fails to provide useful answers;
- **User Needs:** Fails to address core points and does not meet information needs;

Information Depth

Definition Evaluate whether the article provides sufficient information. Depth provides sufficient relevant information so readers can thoroughly understand each argument in the article.

Score 5: Excellent Coverage and Amount (depth)

- **Detail Sufficiency:** Provides necessary and sufficient information with selective deep exploration. Can select materials requiring deep exploration for detailed discussion;
- **Technical Accuracy:** Highly accurate technical details with proper context;
- **Analytical Depth:** Deep analytical insights with sophisticated reasoning. Response provides all necessary and sufficient materials;
- **Contextual Understanding:** Excellent understanding of broader implications and context;

Score 4: Good Coverage and Amount (depth)

- **Detail Sufficiency:** Includes most relevant information needed to understand the topic. Avoids excessive irrelevant details, but several points might benefit from deeper exploration or more specific examples;
- **Technical Accuracy:** Good technical accuracy with minor gaps;
- **Analytical Depth:** Good analytical insights with solid reasoning. Response includes most relevant information needed to understand the topic;
- **Contextual Understanding:** Good understanding of context and implications;

Score 3: Acceptable Coverage and Amount (depth)

- **Detail Sufficiency:** Acceptable amount of relevant information, may lack some useful details;
- **Technical Accuracy:** Adequate technical accuracy with some inaccuracies;
- **Analytical Depth:** Output provides reasonable amount of relevant information, though it may lack some useful details.;
- **Contextual Understanding:** Basic understanding of context;

Score 2: Limited Coverage and Amount (depth)

- **Detail Sufficiency:** Provides some relevant information but misses important details;
- **Technical Accuracy:** Poor technical accuracy with significant errors;
- **Analytical Depth:** Response provides some relevant information but misses important details that would aid full topic understanding.;
- **Contextual Understanding:** Poor understanding of broader context;

Score 1: Lack of Coverage and Amount (depth)

- **Detail Sufficiency:** Lacks basic details needed for topic understanding;
- **Technical Accuracy:** Very poor technical accuracy with major errors;
- **Analytical Depth:** Output either lacks basic details needed for adequate topic understanding (e.g., method definitions, relationships between methods);
- **Contextual Understanding:** No understanding of context or implications;

Overall Helpfulness

Definition Do you find the provided answer overall helpful? Does it assist with your literature review? Evaluate the overall utility of the response for research and learning purposes.

Score 5: Super Useful. I can fully trust the answer

- **Question Addressing:** Answer provides comprehensive field overview and fully answers the question;
- **Source Quality:** Provides high-quality, trustworthy sources with comprehensive coverage;
- **Research Utility:** Serves as complete foundation for research without need for independent verification;
- **Information Reliability:** I believe I don't need to independently search for other papers or detailed information;

Score 4: Useful. I may try to verify some details, but overall gives great summary

- **Question Addressing:** Answer provides detailed information and good overview of the area of interest;
- **Source Quality:** Provides high-quality, fresh sources across multiple sources with good diversity;
- **Research Utility:** Requires minimal additional editing, serves as excellent foundation for further work;
- **Information Reliability:** May need to check details of 1-2 specific papers/sources, but overall highly reliable;

Score 3: Provides some useful discussions and papers, though requires independent reading

- **Question Addressing:** Answer is generally helpful and provides good overview with diverse perspectives;
- **Source Quality:** Provides at least 2-3 useful information sources previously unknown to reader;
- **Research Utility:** Can base further reading on recommended papers, good starting point for deeper research;
- **Information Reliability:** May need to independently verify some details or consult other core research papers;

Score 2: Better than searching from scratch but limited utility

- **Question Addressing:** Answer provides at least one useful starting point but discussions are somewhat irrelevant;
- **Source Quality:** Provides at least one useful paper that can be read carefully;
- **Research Utility:** Limited utility for research purposes, requires significant additional work;
- **Information Reliability:** Overall discussions don't provide sufficiently useful information for the topic;

Score 1: Unhelpful

- **Question Addressing:** Answer doesn't address the question or provides confusing information;
- **Source Quality:** Hasn't conducted effective retrieval, still generating using pre-trained knowledge;
- **Research Utility:** Cannot serve as useful starting point for learning or writing relevant content;
- **Information Reliability:** Fails to provide understanding of literature in this field;

Evaluation Dimension	-2	-1	0	+1	+2
Transparency: Decision-making process visibility	○	○	○	○	○
Interruptibility: Real-time intervention capability	○	○	○	○	○
Fine-grained Interaction: Interaction granularity level	○	○	○	○	○
Inspiration: Unexpected discoveries and insights	○	○	○	○	○
Collaboration: Collaborative partnership quality	○	○	○	○	○

Table 5: System Design Assessment Rubric

D.2.2 SYSTEM DESIGN EVALUATION (-2 TO +2 SCALE)

System Design Evaluation Definition

Question: Does the system design provide sufficient transparency in decision-making processes?

Interruptibility (Interruptible at any time): To what extent do you think interruptibility can help correct the model’s research approach and reduce model errors?

Fine-grained and Bidirectional Interaction: How fine-grained do you think the current system’s interaction is? (Interaction refers to nodes where users can provide input to the model)

Inspirational Perspectives (Shared cognitive context as exploration space): How much information in the model’s decision and search process exceeded your expectations? Did it help inspire you?

Inspirational Perspectives (Shared cognitive context as exploration space): How much information in the model’s decision and search process exceeded your expectations? Did it help inspire you?

Long-term Collaboration Willingness: Deep research systems can all interact (Deep Cognition during process, other 3 systems after research process). Research is a dynamic, multi-round complex long-term task. To what extent do these systems’ interaction methods (including input methods and system feedback output methods) make you willing to engage in long-term, multi-round communication and collaboration with the system?

Long-term Collaboration Willingness: Deep research systems can all interact (Deep Cognition during process, other 3 systems after research process). Research is a dynamic, multi-round complex long-term task. To what extent do these systems’ interaction methods (including input methods and system feedback output methods) make you willing to engage in long-term, multi-round communication and collaboration with the system?

+2 points - Excellent:

- **Process Visibility:** Complete visibility of thinking, actions, and browsed content;
- **Decision Rationale:** Clear explanation of all decision-making processes;
- **Source Verification:** Full source verification and citation transparency;
- **Strategy Disclosure:** Complete disclosure of search and analysis strategies;

+1 points - Good:

- **Process Visibility:** Good transparency with some decision process visibility;
- **Decision Rationale:** Adequate explanation of major decisions;
- **Source Verification:** Good source transparency with minor gaps;
- **Strategy Disclosure:** Partial disclosure of strategies and approaches;

0 points - Neutral:

- **Process Visibility:** Neutral/adequate transparency level;

- **Decision Rationale:** Basic explanation of some decisions;
 - **Source Verification:** Adequate source information;
 - **Strategy Disclosure:** Limited strategy disclosure;
- 1 points - Poor:**
- **Process Visibility:** Limited transparency, unclear decision processes;
 - **Decision Rationale:** Poor explanation of decision-making;
 - **Source Verification:** Limited source transparency;
 - **Strategy Disclosure:** Minimal strategy disclosure;
- 2 points - Extremely Poor:**
- **Process Visibility:** Black box operation with no process visibility;
 - **Decision Rationale:** No explanation of decision-making processes;
 - **Source Verification:** No source transparency or verification;
 - **Strategy Disclosure:** No disclosure of strategies or methods;

D.2.3 DEEP COGNITION SPECIFIC EVALUATION

Qualitative indicator: When comparing the Deep Cognition system with other deep research systems, do the system’s functional designs (interruptibility, transparent thinking process, transparent behavioral paths, presenting search queries, displaying retrieved content) enhance this system’s collaborative attributes?

Follow-up questions: A. If enhanced, can you provide specific examples? Which functions enhanced collaborative attributes? B. During model behavior review, could the model provide new insights/unexpected search information?

Feature	Description
Text Input	Basic text communication capability
Question Clarification	System’s ability to clarify ambiguous queries
Expert Information Integration	Incorporating domain expertise
Thinking Process Visibility	Transparency of reasoning steps
Decision Process	Clarity of decision-making rationale
Interruptibility	Effectiveness of real-time intervention
Content Summary Reading	Quality of information synthesis
Search Query Visibility	Transparency of search strategies

Table 6: Deep Cognition Feature-Specific Ratings (1-5 Scale)

D.3 POST-STUDY

Deep Cognition Evaluation: -2 for strongly negative, 0 for neutral, 2 for strongly positive

1. Enhanced Effectiveness (Enhance cognitive efficiency or not)

To what extent do you think this collaborative approach can improve final report generation quality (organization and consistency/information coverage/information density (depth)/relevance/overall helpfulness)?

Dimension	Score (-2/-1/0/1/2)	Reason
Organization and consistency		
Information coverage		
Information density (depth)		
Relevance		
Overall helpfulness		

2. Results-worth-effort Interacting with these systems costs your time and energy. Do you think it's worth it? How worthwhile?

System	Score (-2/-1/0/1/2)	Reason
Deep Cognition		
OpenAI		
Gemini		
Grok 3		

3. Research Stage Evaluation

At which stages do you think interrupting the model's operation can effectively improve subsequent report quality? Which stage can enhance your real-time collaboration willingness with the model?

Current model nodes include: evaluating research status, generating search queries, filtering webpage URLs, browsing webpages, extracting summaries from webpages and determining usefulness, prioritizing information retrieved from webpages and organizing arguments.

You may define research stages according to your own understanding when asking this question.

Follow-up questions:

- At which stage of model research development is your collaboration willingness higher?
- Can the model's research process provide you with insights? Can you give an example (screenshot or text)?
- At which stages do you think interrupting the model's operation can more effectively improve subsequent report quality? Which stage can enhance your real-time collaboration willingness with the model?

4. Usage Willingness and Learning Cost (Interaction Willingness)

Quantitative indicators: To what extent are you willing to use this system? How are the learning costs and operational burden?

Aspect	Score (-2/-1/0/1/2)	Reason
Usage willingness		
Ease of operation		

5. Feature Evaluation

How helpful are these features for your research process? Rate (1-5) and explain reasons.

Feature Number	Feature Name	Score	Comments
1	Send text		
2	Clarify questions		
3	Add expert information		
4	Thinking process		
6	Decision		
7	Interruptible		
8	Read summaries		
9	Search queries		

E CASE STUDY

A Dynamic Research Task Case	
Domain: Interdisciplinary Writing	
Initial query: Please assist in investigating the latest innovations in interactive and scalable interfaces designed to enhance the interpretability of large language models in writing.	
Clarify query: What aspect of interpretability are you most interested in exploring? Who is the primary target audience for these interfaces? What writing contexts are you most interested in? What aspects of interface scalability are most relevant to your research? Are you interested in any specific emerging technologies related to LLM interpretability?	
Domain knowledge input: Add Jeff Rzeszutarski’s PhD dissertation, and research in PAIR (People + AI Research Initiative) team.	
Initial goal: Development trend of interpretability of Interpretable Machine Learning Interface	
Last goal: Investigate which research fields the scholars who previously worked in this direction have migrated to.	

F LLM USAGE

Large Language Models (LLMs) were used to aid in the writing and polishing of the manuscript. Specifically, we used an LLM to assist in refining the language of the paper. The model helped with tasks such as sentence rephrasing, grammar checking, and enhancing the overall flow of the text. It is important to note that the LLM was not involved in the identification, research methodology, or experimental design. All research concepts, ideas, and analyses were developed and conducted by the authors. The contributions of the LLM were solely focused on improving the linguistic quality of the paper, with no involvement in the scientific content or data analysis. The authors take full responsibility for the content of the manuscript, including any text generated or polished by the LLM. We have ensured that the LLM-generated text adheres to ethical guidelines and does not contribute to plagiarism or scientific misconduct.