



AgentDisCo: Towards Disentanglement and Collaboration in Open-ended Deep Research Agents

Jiarui Jin, Zexuan Yan, Shijian Wang, Wenxiang Jiao, Yuan Lu

Xiaohongshu Inc.

Abstract

Open-ended deep research agents have emerged as a promising paradigm for autonomously performing comprehensive information gathering and synthesis. However, existing approaches typically integrate information **exploration** and **exploitation** into a single unified module—such as an outline generator or a report generator—thereby limiting their flexibility and optimization potential.

In this paper, we introduce **AgentDisCo**, a novel **Disentangled and Collaborative** agentic architecture that formulates deep research as an adversarial optimization problem between information exploration and exploitation. Specifically, a **critic agent** is optimized to evaluate and critique the generated outlines (serving as information exploitation states) and subsequently **refine the search queries (serving as information exploration states)**, while a **generator agent** is optimized to retrieve updated search results based on the refined search queries (serving as information exploration states) and accordingly **update the generated outlines (serving as information exploitation states)**. The resulting outline, progressively refined through iterative adversarial optimization, is subsequently delivered to a downstream report writer module. This module leverages the structured outline alongside the accumulated search results to synthesize a comprehensive, coherent, and well-grounded research report.

The above agentic workflow can be optimized through either handcrafted or automatically discovered design strategies by constructing a **meta-optimization harness** over the adversarial optimization loop, where the generator agent originally tasked with producing target outlines is repurposed as a scoring agent that evaluates and generates quality signals over the critic agent’s outputs, thereby enabling systematic optimization of the search queries. Concretely, powerful code-generation agents—such as Claude-Code or Codex—are employed to systematically explore the space of agent configurations and automatically construct a **policy bank**, a structured repository of reusable and composable design strategies over search query generation across diverse research tasks and search domains, enabling the framework to self-refine its own design strategies without requiring extensive human intervention. We evaluate AgentDisCo on three widely adopted deep research benchmarks—DeepResearchBench, DeepConsult, and DeepResearchGym—with Gemini-2.5-Pro as our base model, demonstrating performance comparable to or surpassing that of leading closed-source deep research agents.

Furthermore, we observe that existing benchmarks predominantly focus on academic or domain-specific consulting queries, which diverge significantly from the breadth and diversity of real-world user needs. To bridge this gap, we introduce **GALA (General AI Life Assistants)**, a novel benchmark constructed via an agentic workflow that automatically mines latent deep research interests from users’ historical browsing behavior, enabling a more faithful reflection of organic, everyday information needs.

As an intuitive and user-friendly interface is essential for bridging the gap between research outputs and end-user consumption, we develop a **rendering agent** capable of transforming structured research reports into visually rich (rednote-style) poster presentations. Building upon this, we further construct a product demonstration—“AutoResearch Your Interest”—which automatically curates and delivers personalized deep research recommendations tailored to individual user profiles derived from their browsing histories. We publicly release our benchmark, code, demo, and evaluation harness to support and accelerate future research in open-ended deep research.

1 Introduction

Open-ended deep research agents—capable of synthesizing vast web-scale information into comprehensive, well-cited reports—have emerged as a critical frontier for large language models (LLMs). On one hand, closed-source commercial offerings—such as GPT Deep Research (OpenAI, 2025a) and Gemini Deep Research (Research, 2025a)—provide neither technical reports nor open implementations. On the other hand, a growing body of open technical reports (Han et al., 2025; Li et al., 2025; Lei et al., 2025) shed light on architectural design choices. Yet, beneath these advances in released architectural designs lies a persistent architectural bottleneck. First, existing deep research agents entangle **information exploitation**—the generation of structured outlines or reports—with **information exploration**—the planning and generation of search queries—into a single, undifferentiated module, fundamentally lacking any guarantee of informational incrementality in the iterative optimization process. Second, current outline-guided iteration loops require LLMs to refine generated outlines without explicit optimization objectives, leaving the model without clear guidance on which parts of the outline are satisfactory and which require further improvement. This absence of structured feedback signals results in directionless and unstable iterative refinement, where the model oscillates between over-revision and under-revision without convergence.

In this paper, we propose **AgentDisCo**, a novel **Disentangled and Collaborative** agentic architecture that formulates deep research as an adversarial optimization problem between information exploration and exploitation. We first argue that iterative optimization should operate on intermediate generated outlines rather than final reports, as outline-level representations offer greater structural flexibility and are more amenable to effective context management. Within the outline optimization loop, we further disentangle the generation of outlines and search queries into two specialized yet interacting agents: a **critic agent** and a **generator agent**. Specifically, the critic agent receives the current **generator state (i.e., the information exploitation state)**—comprising the evolving outlines along with their associated references—and is tasked with evaluating and critiquing the quality and completeness of the generated outlines, upon which it subsequently produces targeted and gap-aware search queries. To further structure this iterative optimization process, we design the critic agent to produce updated **critic state (i.e., the information exploration state)** including a set of **blueprints**, where each blueprint represents a key point to be covered in the final report and is accompanied by a dedicated list of targeted search queries, thereby ensuring that the information retrieval at each iteration is systematically aligned with the intended scope and coverage of the final report. Conversely, the generator agent receives the current **critic state (i.e., the information exploration state)**, and is responsible for retrieving updated search results based on the refined search queries, and accordingly revising the generated outlines along with accompanying references, namely **generator state (i.e., the information exploitation state)**. Through this adversarial yet collaborative interaction, the two agents iteratively drive each other toward more comprehensive information coverage and higher-quality outline generation.

The aforementioned agentic workflow can be optimized through two broad classes of design strategies: handcrafted approaches and **automatically discovered approaches**. Handcrafted optimization relies on domain expert knowledge to manually engineer search heuristics and allocate computational resources across distinct components of the workflow, offering interpretability but limited scalability. To overcome these constraints, recent advances in meta-optimization (Lee et al., 2026) introduce an outer optimization harness that operates over the agentic workflow itself, enabling the agent to systematically explore and refine its own optimization strategies in an automated and adaptive manner—without requiring exhaustive human intervention. Concretely, we employ Claude-Code as our primary code-generation agent and construct a meta-optimization harness around the critic agent. Within this harness, the **generator agent**—originally responsible for producing target outlines—is repurposed as a **scoring agent** that evaluates the critic agent’s outputs and emits structured quality signals. This repurposing enables systematic optimization of search query generation without introducing additional model components. Leveraging the strong code-generation capabilities of state-of-the-art agents such as Claude-Code, the framework systematically explores the space of agent configurations and automatically constructs a **policy bank**—a structured repository of reusable and composable design strategies that govern search query generation across diverse research tasks and retrieval domains. By drawing upon and refining the entries in this policy bank, the framework iteratively self-evolves its design strategies, progressively improving retrieval quality while reducing the need for human intervention.

We evaluate AgentDisCo on three widely adopted deep research benchmarks, namely, DeepResearchBench (Du et al., 2025), DeepConsult (Consult, 2025), and DeepResearchGym (Coelho et al., 2025) with Gemini-2.5-Pro (DeepMind, 2025a) as our base model. Specifically, AgentDisCo w/ Harness achieves a RACE score of **52.11** on DeepResearchBench and **6.86** on DeepConsult, surpassing leading closed-source systems such as Doubao-Research (Research, 2026a), Claude-DeepResearch (anthropic, 2025), and OpenAI-DeepResearch (OpenAI, 2025a).

However, we observe that existing benchmarks predominantly center on academic or domain-specific consulting queries, which diverge substantially from the breadth and diversity of real-world user needs. To bridge this gap, we introduce **GALA** (General AI Life Assistants), a novel benchmark designed to capture authentic, everyday information-seeking behavior. Specifically, we collect over 10,000 highly active users from the Rednote platform¹ along with their browsing and commenting histories, and devise an agentic workflow that automatically mines latent deep research interests and synthesizes personalized queries tailored to each user’s individual preferences. Compared with prior open-sourced benchmarks, the resulting queries exhibit a markedly more lifestyle-oriented character, with the dominant topics shifting from Science & Technology and Finance & Business toward everyday domains such as Home & Hobbies, Fashion & Beauty, and Travel. By grounding evaluation in organic user activity, GALA offers a more faithful reflection of everyday information needs and enables a more realistic assessment of deep research agents in practical deployment scenarios. Alongside the query set, we release a standardized evaluation protocol built upon Gemini-3-Flash (DeepMind, 2025b), in which reports generated by AgentDisCo serve as reference outputs against which competing systems are scored in a pairwise manner. To construct competitive baselines, our human annotation team manually collects reports from the official web interfaces of Doubao-Research (Research, 2026a) and Qwen-Research (Research, 2026c), as well as outputs from OpenAI o3-DeepResearch (Research, 2026b) obtained via its API. Experimental results demonstrate that AgentDisCo consistently outperforms these strong proprietary baselines. Moreover, we find that AgentDisCo achieves stronger performance when relying solely on the Rednote search engine than when relying solely on Google Search, highlighting the superiority of Rednote as a source of community-grounded content for everyday information-seeking tasks.

Motivated by recent advances in AI-generated content (AIGC) and the emergence of paper-to-poster generation systems (Zhang et al., 2025b), we further explore the integration of deep research agents with automated visual presentation. Recognizing that an intuitive and user-friendly interface is essential for bridging the gap between research outputs and end-user consumption, we develop a **rendering agent** that transforms structured research reports into visually rich, Rednote-style poster presentations. Building upon this capability, we construct a product demonstration—AutoResearch Your Interest—which automatically curates and delivers personalized deep research recommendations tailored to individual user profiles inferred from their browsing histories.

Contributions. Our main contributions can be summarized as follows.

- **A novel disentangled and collaborative agentic architecture.** As illustrated in Figure 1, we introduce **AgentDisCo**, which formulates deep research as an adversarial optimization problem between information exploration and exploitation, decoupling outline generation from search query formulation and coordinating them through a dynamic critic-and-generator cycle.
- **A meta-optimization harness for self-evolving search.** We construct a meta-optimization harness over the adversarial optimization loop, in which the generator agent—originally responsible for producing target outlines—is repurposed as a scoring agent that evaluates the critic agent’s outputs and emits structured quality signals. This design enables the systematic and automatic optimization of search query generation without introducing additional model components.
- **A new benchmark for lifestyle deep research needs.** Observing that existing benchmarks predominantly focus on academic or domain-specific consulting queries—diverging substantially from the breadth and diversity of real-world user needs—we introduce **GALA**, a benchmark that captures authentic, lifestyle-oriented information-seeking behavior mined from organic user activity.
- **An open-sourced deep research system with multi-modal render agent.** Recognizing that an intuitive interface is essential for bridging research outputs and end-user consumption, we develop a **rendering agent** that transforms structured reports into visually rich, Rednote-style poster presentations. As depicted in Figure 2, AgentDisCo thus spans the full pipeline—from deep research interest mining (which underpins the GALA benchmark) to the generation of reports and posters. To support and accelerate future research on open-ended deep research, we publicly release our benchmark, code, demo, and evaluation harness.

2 AgentDisCo: A Disentangled and Collaborative Agentic Architecture

2.1 System Overview and Design Philosophy

We consider open-ended deep research questions without ground-truth answers. Given such a question q , the agentic system must perform **information exploration** (searching relevant evidence) and **information exploitation** (synthesizing the evidence into a structured report). We argue that iterative optimization should occur on **intermediate outlines** rather than final reports, and we further disentangle exploration

¹<https://www.xiaohongshu.com/explore>

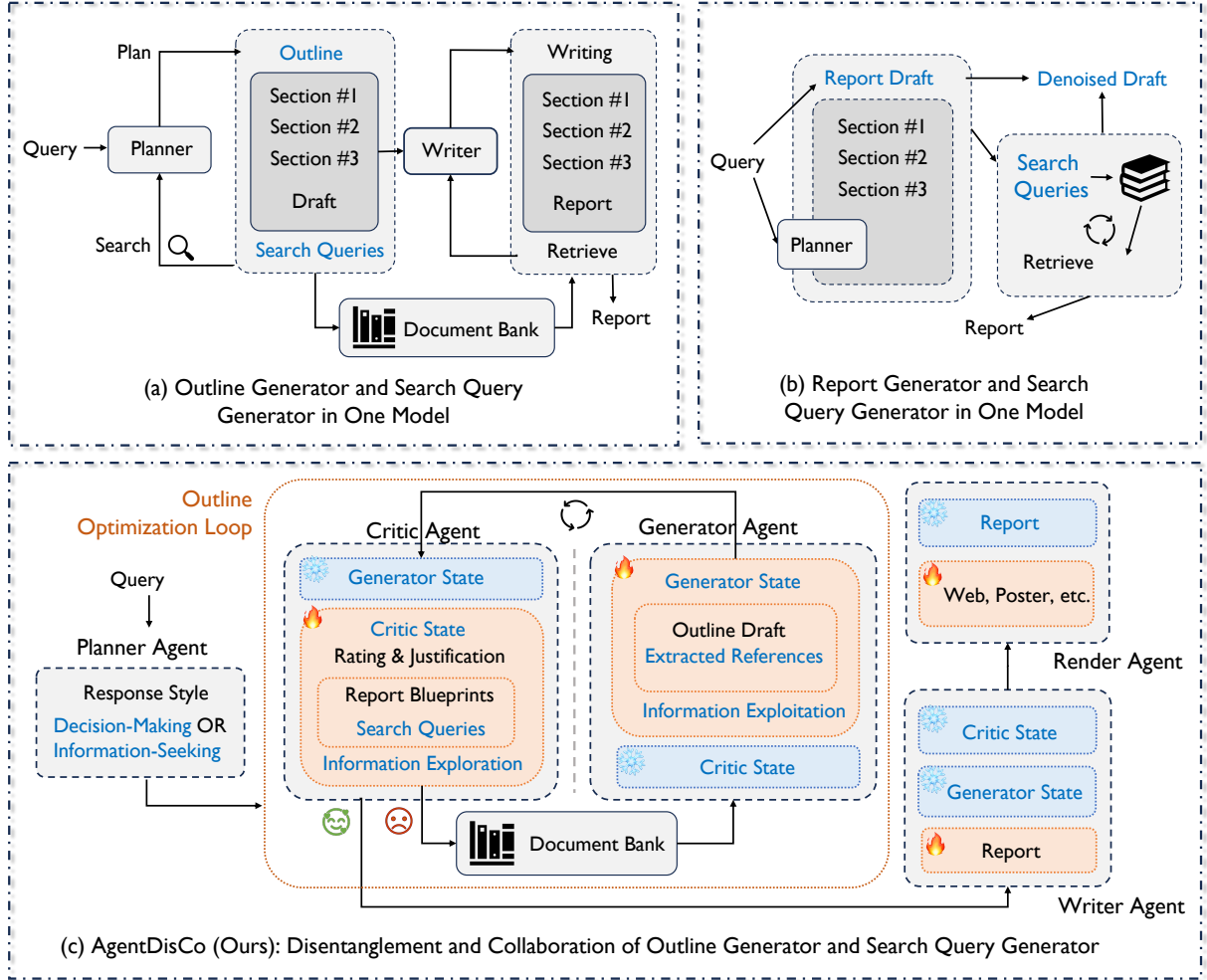


Figure 1: **Comparison of deep research paradigms.** (a) the outline-iterative-optimization paradigm couples outline generation and search query formulation within a single model; (b) the report-iterative-optimization paradigm similarly entangles report generation with search query formulation; (c) in contrast, **AgentDisCo** disentangles the outline generator and the search query generator into separate models, and further coordinates them through a dynamic critic-and-generator research cycle.

and exploitation into two specialized but interacting agents — a **critic agent** π^c and a **generator agent** π^g — whose adversarial yet collaborative interplay drives convergence toward a comprehensive, well-grounded outline.

We formalize this disentangled and collaborative interaction as a dual-agent cooperative MDP (Markov Decision Process):

$$\mathcal{M} = \langle \mathcal{S}^c, \mathcal{S}^g, \mathcal{A}^c, \mathcal{A}^g, \mathcal{P}, \mathcal{R}, T \rangle, \quad (1)$$

where the joint environment state at timestep t is the pair $s_t = (s_t^c, s_t^g)$. The **generator state (i.e., the information-exploitation state)**, denoted $s_t^g = (O_t, R_t)$, consists of the current outline O_t and the set of references R_t attached to it. The **critic state (i.e., the information-exploration state)**, denoted $s_t^c = (B_t, Q_t)$, comprises a set of blueprints B_t together with their associated search queries Q_t . Each blueprint in B_t specifies a key point that the final report should cover, and each such key point is paired with a list of targeted search queries in Q_t dedicated to filling its information gap. Together, the blueprints align retrieval with the intended scope and coverage of the report.

Given a user query q , the system initializes the generator state as $s_0^g = (\emptyset, \emptyset)$, indicating an empty outline and an empty document pool. Conditioned solely on q , the critic agent then performs a coarse decomposition and proposes an initial blueprint set $s_0^c \sim \pi^c(\cdot | q, \emptyset, \emptyset)$, which seeds the subsequent iterative refinement. At each iteration t , the two agents act sequentially rather than simultaneously, with the critic moving first. At the critic’s step, conditioned on the question q , the previous generator state s_{t-1}^g and its own previous state s_{t-1}^c , the critic agent assesses the completeness and quality of the current outline and produces an updated blueprint set: $s_t^c \sim \pi^c(\cdot | q, s_{t-1}^g, s_{t-1}^c)$. That is, the critic’s action is to



Figure 2: Overview of the architecture and applications of AgentDisCo. AgentDisCo spans the full pipeline from mining latent deep research queries in user interaction histories to producing structured reports and rendering visually rich posters. This end-to-end design realizes the vision of “AutoResearch Your Interest”—automatically tracking evolving user interests and delivering personalized deep research recommendations tailored to individual user profiles.

instantiate the next exploration state. At the generator’s step, conditioned on the freshly updated critic state s_t^c together with its previous state s_{t-1}^c , the generator agent executes the queries prescribed by each blueprint, retrieves new evidence via the search tool, and accordingly revises both the outline and its references: $s_t^g \sim \pi^g(\cdot | q, s_{t-1}^g, s_t^c)$. The generator’s action thus realizes the exploitation counterpart of the critic’s exploration. As for environment dynamics, given the actions of both agents, the joint state transition is deterministic: each agent’s output directly instantiates the corresponding component of the next joint state, i.e., $s_{t+1}^c = a_t^c$ and $s_{t+1}^g = a_t^g$. Consequently, all stochasticity of the trajectory originates from the agents’ policies themselves.

Under the above sequential protocol, an interaction trajectory of AgentDisCo unfolds as the alternating sequence $\tau = \{q, s_0^c, s_0^g, s_1^c, s_1^g, \dots\}$, in which the critic and the generator alternately update their respective state components. The likelihood of sampling τ admits the following equivalent factorizations:

$$\begin{aligned}
 p(\tau) &= \prod_{t=0}^{T-1} \underbrace{\pi^c(s_{t+1}^c | q, s_t^g, s_t^c)}_{(1) \text{ Critic agent policy}} \underbrace{\pi^g(s_{t+1}^g | q, s_t^g, s_{t+1}^c)}_{(2) \text{ Generator agent policy}} \\
 &= \prod_{t=0}^{T-1} \underbrace{\pi^c(a_t^c | q, s_t^g, s_t^c)}_{(1) \text{ Critic agent policy}} \underbrace{\mathbb{I}[s_{t+1}^c = a_t^c] \mathcal{P}^c(s_{t+1}^g | q, s_t^g, a_t^c)}_{(2) \text{ Environment dynamic of critic agent}} \\
 &= \prod_{t=0}^{T-1} \underbrace{\pi^g(a_t^g | q, s_t^g, s_{t+1}^c)}_{(1) \text{ Generator agent policy}} \underbrace{\mathbb{I}[s_{t+1}^g = a_t^g] \mathcal{P}^g(s_{t+2}^c | q, a_t^g, s_{t+1}^c)}_{(2) \text{ Environment dynamic of generator agent}}
 \end{aligned} \tag{2}$$

where T denotes the maximum number of optimization rounds and $\mathbb{I}(\cdot)$ is the indicator function. The first line presents the most compact policy-only factorization, while the latter two further decouple each agent’s stochastic decision from its (degenerate) environment transition, thereby making explicit the loci at which learning signals can be injected.

Although both agents are optimized toward a common objective, their per-round roles are functionally adversarial, jointly forming a minimax-style yet cooperative loop. The critic agent π^c adversarially probes the generator’s current outline, surfacing missing evidence and uncovered key points through newly proposed blueprints — thereby pushing exploration outward. The generator agent π^g defensively expands and grounds the outline using the freshly retrieved evidence — thereby consolidating exploitation inward. This adversarial-yet-aligned interaction progressively refines the outline along the axes of coverage, factual grounding, and structural coherence. To reconcile the two locally adversarial roles under a single global objective, both agents share a common cooperative reward that quantifies the

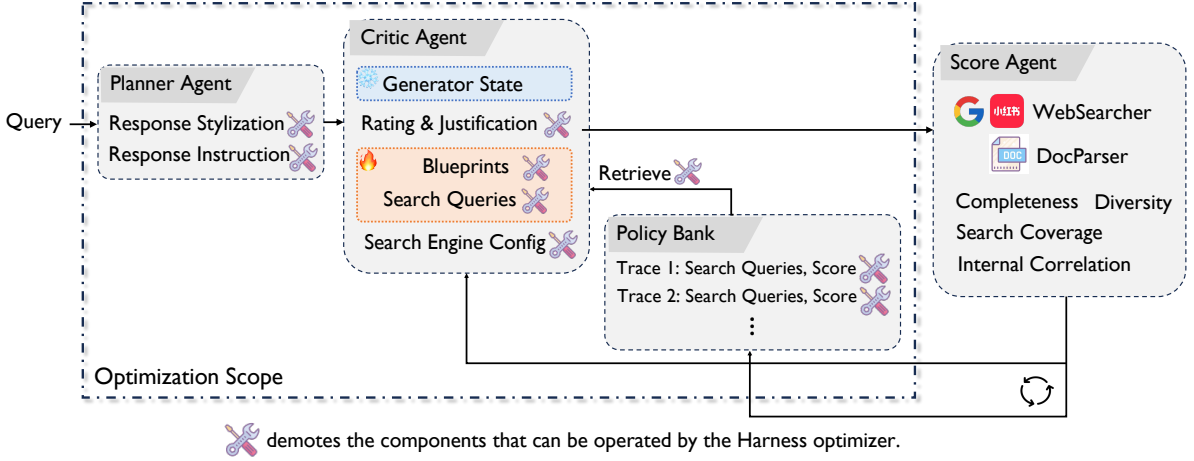


Figure 3: **Overview of the harness optimization in AgentDisCo.** AgentDisCo can automatically discover design strategies by constructing a meta-optimization harness around the adversarial optimization loop. Specifically, the generator agent—originally tasked with producing target outlines—is repurposed as a scoring agent that evaluates the critic agent’s outputs and generates quality signals, thereby enabling systematic optimization of the search queries.

quality of the updated outline for the subsequent writer agent and render agent.

Below, we describe each component of AgentDisCo in detail (as depicted in Figure 1), including the aforementioned (outline) critic and (outline) generator agents.

2.2 Planner Agent

Handcrafted Design. To accommodate the heterogeneous nature of deep research queries encountered on the industrial platform, we introduce a planner agent that classifies each incoming query into one of two top-level categories, each further subdivided into fine-grained intents, namely **information seeking** category (including fact query, status & progress, news & information, deep exploration, and resource locating) and **decision making** category (comparison & selection, recommendations & suggestions, how-to guide, travel planning, and purchase decision). Beyond the categorical label, the planner also infers an expected response style tailored to the predicted intent, which subsequently conditions the downstream critic and generator agents. We instantiate the planner with Gemini-2.5-Pro (DeepMind, 2025a) as the backbone, and provide the full prompt template in Appendix A.3.

Harness Optimization. During harness optimization, we grant the planner agent access to generate instructions for the subsequent agentic workflows. Concretely, we observe that our code-generation agent (i.e., the Claude-Code agent) specifies instructions such as tips for formulating search queries and selecting hyperparameters for the search engines, which are then passed along to the downstream (outline) critic and generator agents. Detailed prompts are provided in Appendix B.2.

Formally, we introduce the notation $\pi_{\text{planner}}(\cdot)$ to denote the planner agent, whose function can be written as $P \sim \pi_{\text{planner}}(q)$. In the handcrafted design, P represents the intent type and response style conditioned on the input query q , whereas during harness optimization, P additionally includes specified instructions. For ease of representation, we omit P , since it can be regarded as a complementary explanation attached to the query and thus integrated into q , namely $q = [P; q]$.

2.3 (Outline) Critic Agent and (Outline) Generator Agent

Handcrafted Design. As described in Section 2.1, the core of AgentDisCo lies in the dual optimization between the critic and generator agents. Eq. (2) reveals that, although both agents assume functionally adversarial roles in each round, they jointly constitute a minimax-style yet cooperative loop: both are optimized toward a shared objective, namely, the production of a high-quality outline that can attain a sufficiently high score under the critic’s evaluation. The resulting artifact, denoted as O_t and R_t , serves as the foundation for the subsequent writing and rendering stages.

To assess outline quality, our handcrafted design employs the critic agent to produce both a numerical score and an accompanying justification for each candidate outline. The optimization loop is further governed by three control parameters: an exit threshold that determines when the outline is deemed

acceptable by the critic, a minimum number of optimization rounds to ensure sufficient refinement, and a maximum number of rounds to bound the overall computational cost. Formally, the reward function for (O_t, R_t) can be expressed as $\mathcal{R}(s_t^g, a_t^g)$, where a_t^g denotes the generator’s action, i.e., the produced outline O_t and its rendering R_t . At iteration $t + 1$, the critic agent receives (O_t, R_t) as input and outputs an evaluation score that constitutes the reward signal. Accordingly, the reward is sampled as $r_t \sim \pi^c(a_{t+1}^c | s_{t+1}^c)$, with the critic’s state defined as $s_{t+1}^c = a_t^g$, thereby coupling the generator’s output directly to the critic’s input.

A persistent challenge in multi-round outline optimization is reference management: as the optimization unfolds, the agent must continually decide which documents retrieved in earlier rounds are worth retaining for downstream use and which can be safely discarded. Naive strategies risk either prematurely dropping high-quality evidence or overloading the context window with stale and redundant content—both of which degrade the quality of the final outline. To address this, we introduce the **document bank**, a lightweight recorder and tracker that maintains a persistent view of retrieved references across rounds. Once the critic agent produces a set of search queries, the document bank parses each retrieved document into fine-grained evidence snippets, scores documents in parallel for relevance, summarizes their content, and extracts key evidence triples; low-scoring documents are filtered out before reaching the generator agent. In this way, the document bank not only compresses raw search results into a structured, citation-ready memory, but also shields the generator agent from contextual noise and redundancy.

To ensure stable progression and prevent information loss across rounds, each new optimization round enforces the following continuity constraints. Firstly, all documents contained in the reference set (i.e., R_t in the outline O_t) are carried over as input to the next round (i.e., $t + 1$), guaranteeing that previously validated evidence remains accessible to both agents. The document bank correspondingly updates its document indices to align with the new round. Secondly, while outlines and search queries are fully re-generated at each round, **the underlying blueprint is only permitted to be modified or expanded**; deletion of existing key points is discouraged, thereby preserving the structural backbone established in earlier rounds. Moreover, each invocation of the critic agent (for refining search queries and blueprints) and the generator agent (for refining outlines) is explicitly conditioned on the concrete content produced in the preceding iteration, ensuring that optimization proceeds incrementally rather than restarting from scratch.

Harness Optimization. Our central insight is that LLMs excel at extracting and summarizing information, but are comparatively weaker at generating effective search queries across heterogeneous retrieval sources and at organizing coherent outlines. In this paper, we therefore focus on harness-based optimization for search query generation, leaving outline organization to future work. Concretely, we develop a harness that optimizes the (outline) critic agent, whose responsibility is to formulate effective search queries and retrieve high-quality information from web-scale sources. As illustrated in Figure 3, we further repurpose the generator agent as a score agent that analyzes the returned search results and provides feedback signals to guide the critic’s query refinement. In practice, since the search queries are organized under a set of blueprints, evaluating individual results in isolation is insufficient. We therefore design several criteria—namely completeness, diversity, search coverage, and internal correlation—and apply them to each search result, aggregating the per-result scores into statistics and distributions over the entire result set. These aggregated signals provide the critic with a holistic view of retrieval quality, enabling more targeted query refinement in subsequent rounds. The detailed prompt used by the score agent is provided in the Appendix B.2.

A particularly noteworthy emergent behavior arises during harness optimization: prompted only by a minimal cue—“you are allowed to store and retrieve traces to evolve”—the coding agent (i.e., the Claude-Code agent) autonomously builds a policy bank that records and reuses relevant historical traces, including critic states, generator states, and the aforementioned criterion scores, to guide subsequent optimization. Formally, the likelihood of sampling τ can be formulated as follows:

$$p(\tau) = \prod_{t=0}^{T-1} \underbrace{\mu^x(m_t | q, s_t^g, s_t^c, M_t)}_{(1) \text{ Retrieve from policy bank}} \underbrace{\pi^c(a_t^c | q, s_t^g, s_t^c, m_t)}_{(2) \text{ Critic agent policy}} \underbrace{\mu^w(M_{t+1} | \cdot)}_{(3) \text{ Update policy bank}} \underbrace{\mathbb{I}[s_{t+1}^c = a_t^c] \mathcal{P}^c(s_{t+1}^g | q, s_t^g, a_t^c)}_{(4) \text{ Environment dynamic of critic agent}} \quad (3)$$

Here, m_t denote the trace(s) sampled at step t from the policy bank M_t . The retrieval process is formulated as $m_t \sim \mu^x(\cdot | q, s_t^g, s_t^c, M_t)$, where μ^x denotes the retrieval function. In practice, μ^x is instantiated as a simple BM25-based retriever (Robertson & Zaragoza, 2009), which is autonomously implemented by the Claude-Code agent during harness optimization. After each step, the newly produced trace is written back into the policy bank, yielding the updated bank $M_{t+1} \sim \mu^w(\cdot | q, s_t^g, s_t^c, m_t, a_t^c, M_t)$, where μ^w is a lightweight writing function. Together, μ^x and μ^w constitute a simple yet effective read-write interface

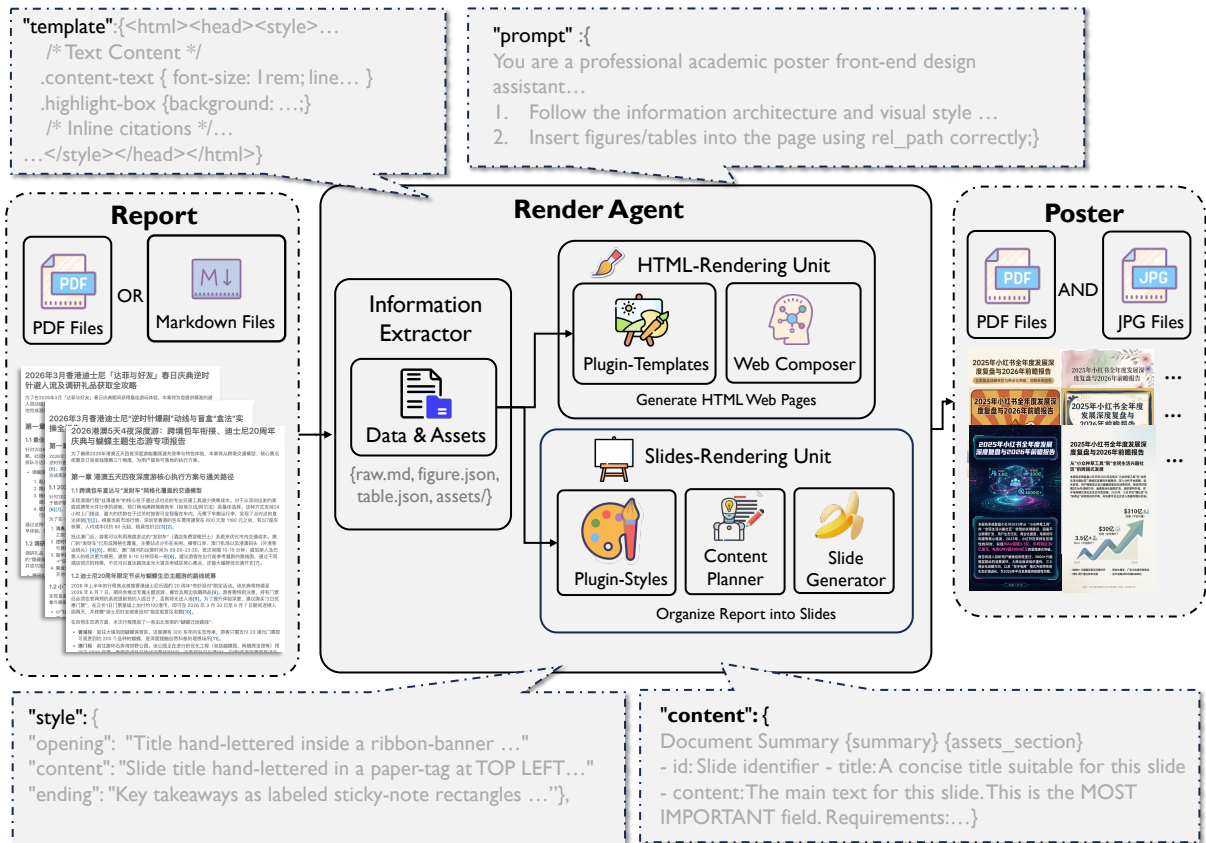


Figure 4: **Overview of the render agent in AgentDisCo.** Our render agent accepts as input a report in either PDF or Markdown format. It first extracts the salient features and structural elements from the report, and then reorganizes the content into one of two presentation modalities: an HTML-based layout or a slide-style layout. Notably, both modalities support pluggable templates and styling components, enabling flexible visual customization. The final output is rendered as a PDF document or a sequence of images, depending on the chosen modality.

that enables the policy bank to evolve continuously throughout the optimization process.

2.4 Writer Agent

As introduced in Section 2.3, the document bank effectively filters out irrelevant content, thereby alleviating the burden on the model’s attentional capacity. Since the generated outline is inherently structured, it is straightforward to partition the outline together with its associated references into a sequence of self-contained chunks, following the strategy of Li et al. (2025). This decomposition reduces the complex task of long-context writing into a series of manageable, attention-focused subtasks, each operating on only the relevant evidence. Because each document has already been assigned a unique index within the outline, the relevant evidence for any given section can be retrieved directly from the document bank. The composition of each section is therefore not a single monolithic action, but rather a deliberate intra-sectional reasoning cycle: at each step, the writer conditions on the previously generated chunks and continues the narrative in a coherent, context-aware manner. This internal monologue is critical for moving beyond shallow summarization toward genuine synthesis across evidence. Finally, the system outputs a Markdown-formatted report in which every cited reference is accompanied by its corresponding URL, ensuring full source traceability. Notably, the writer agent is additionally conditioned on the response style produced by the planner agent (Section 2.2), ensuring that the generated report remains coherent with the input query and faithfully aligned with the user’s intent.

As discussed in Section 2.3, our key premise is that extracting and summarizing information does not constitute the primary bottleneck for LLMs, in contrast to the more challenging task of generating effective search queries across heterogeneous retrieval sources. Accordingly, this paper focuses on harness-based optimization for search query generation, leaving the optimization of the writer agent to future work.

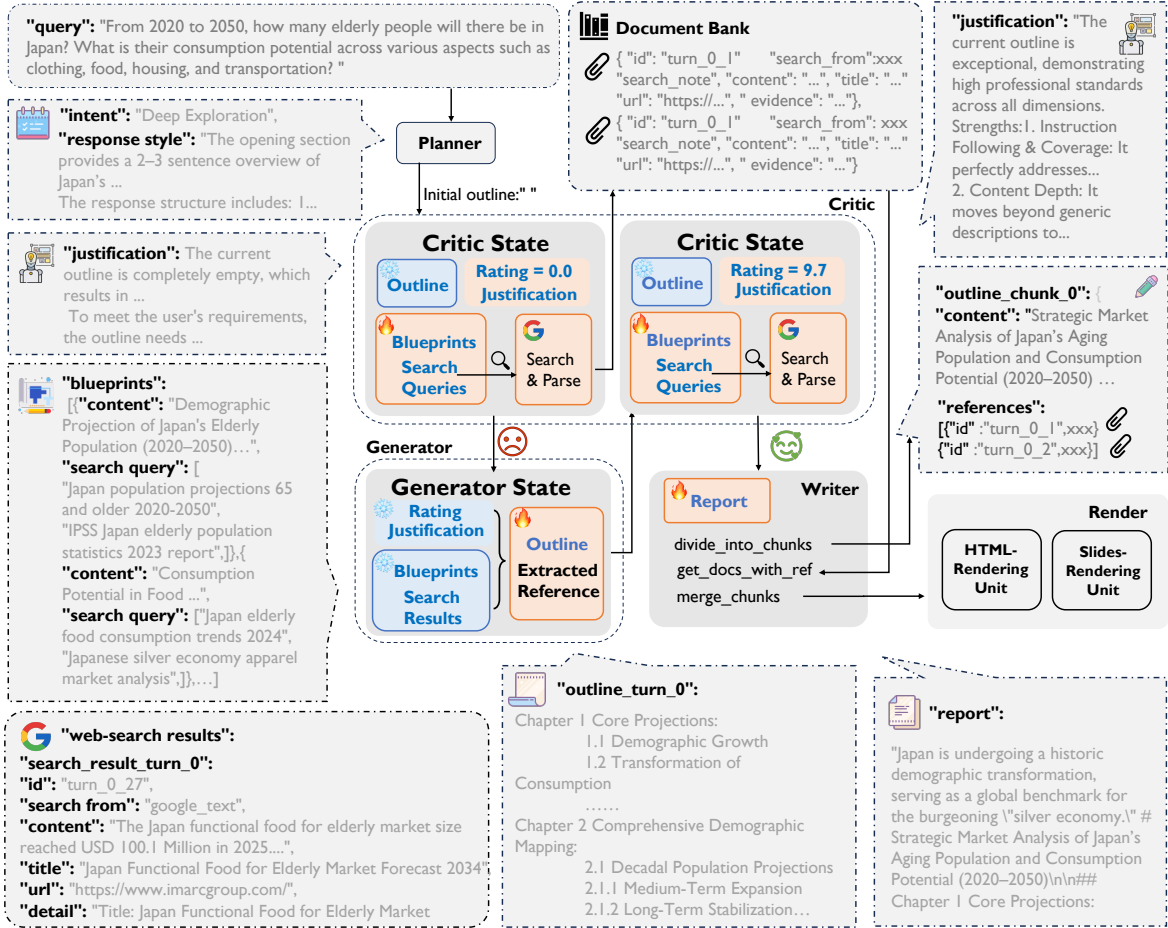


Figure 5: A showcase of AgentDisCo. The figure illustrates the end-to-end processing pipeline of an input query: starting from the planner agent, proceeding through the iterative optimization loop between the (outline) critic agent and the (outline) generator agent, and finally passing to the writer agent and the render agent.

2.5 Render Agent

Considering that a clear, intuitive, and visually engaging interface is essential for bridging the gap between raw research outputs and end-user consumption, we develop a render agent that transforms structured research reports into visually rich presentations—ranging from RedNote-style posters to slide decks and HTML pages—thereby allowing users to readily digest, share, and act upon the generated content. Existing open-source render agents (Sun et al., 2025; Zhang et al., 2025a; Ma et al., 2025; Pang et al., 2025; Yan et al., 2026) rely heavily on complex, tightly coupled pipelines that target a single output modality, making it difficult to accommodate heterogeneous user preferences over presentation forms. Moreover, their intricate iterative design and rigid workflow orchestration limit the extensibility of packaging multimodal pipelines as reusable tools or skills, while incurring additional deployment and debugging overhead. In contrast, we present a clean and flexible render agent that generalizes across diverse rendering tasks within a unified framework.

As depicted in Figure 4, our render agent incorporates an information extractor to extract key information points from the report. To enhance its performance, we also feed the blueprints and response style into the information extractor as auxiliary inputs. To ensure flexibility in use, we construct a webpage template set \mathcal{T}^w and a slide template set \mathcal{T}^s . Given an input document, the information extractor first produces structured multimodal assets conditioned on either \mathcal{T}^w or \mathcal{T}^s . We then apply Gemini-2.5-Pro (DeepMind, 2025a) as the web composer to generate webpages in HTML format, or Gemini-3-Pro-Image (DeepMind, 2026b) as the slide generator to produce multiple images for constructing posters or slides. In addition, we provide an option to generate Rednote-style posters, in which the textual content is produced by Gemini-2.5-Pro (DeepMind, 2025a) acting as the content planner. We present multiple showcases of our rendered posters in the Appendix C.1, along with interactive demos.

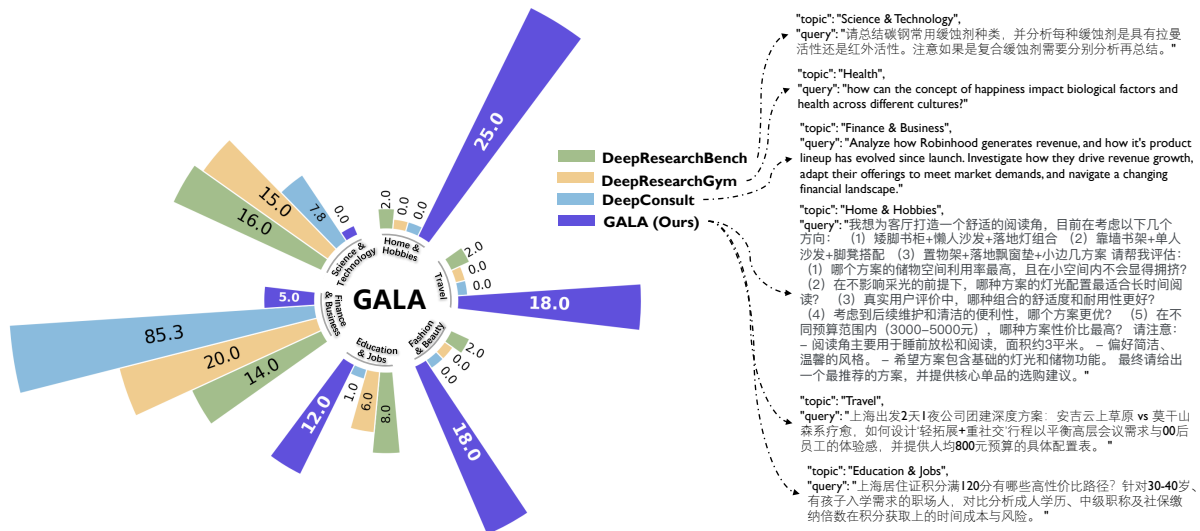


Figure 6: Comparison between our proposed GALA benchmark with existing benchmarks DeepResearchBench, DeepResearchGym, DeepConsult. .

2.6 A Running Example

To illustrate the end-to-end workflow of AgentDisCo, we present a real execution trace in Figure 5 for the example query: "From 2020 to 2050, how many elderly people will there be in Japan, and what is their consumption potential across various aspects such as clothing, food, housing, and transportation?".

First, the planner agent generates a response style based on the input query, which is then passed, together with the query, into the optimization loop between the outline critic agent and the outline generator agent. At round 0, the critic agent observes an empty outline and assigns a rating of 0. At round 1, the critic agent evaluates the draft and identifies that the housing section is underdeveloped. Accordingly, it emits an updated set of blueprints with refined search queries specifically targeting elderly housing consumption. The generator then revises the outline in response, while the document bank preserves the valid citations accumulated from turn 0, ensuring incremental knowledge accumulation across turns. This trace illustrates how the outline critic and generator agents co-evolve through the exchange of blueprints, rather than one dominating the other. Finally, the outline that receives the highest score is forwarded to the downstream agents: the writer agent composes the report by sequentially elaborating each section of the outline, drawing on the associated references retrieved from the document bank, and the render agent subsequently transforms the report into the desired presentation format.

3 GALA: A Benchmark for General AI Life Assistants

3.1 Data Collection and Synthesis on Deep Research Queries

Early deep research agents primarily focused on isolated tasks such as question answering and translation, and later advanced through tool integration to enable autonomous information retrieval and synthesis. To evaluate such systems, a variety of deep research benchmarks (Consult, 2025; Coelho et al., 2025; Du et al., 2025) have been developed, with their queries typically sourced from in-house datasets of raw user interactions with web-search-enabled LLM chatbots. However, these collected queries predominantly fall into a narrow set of categories, such as scientific reports or consulting-style solutions, and thus fail to reflect the diversity of real-world information needs. In contrast, Rednote has emerged as a popular search platform for daily-life information needs, covering a broad spectrum of topics ranging from travel and lifestyle to consumption and entertainment. Leveraging the massive user-interaction and content data on the Rednote platform, we collect the top 10,000 highly active users together with their historical interactions—including clicks, comments, and browsing records—to mine their latent deep research needs. For example, as shown in Figure 2, if a user repeatedly browses or comments on notes related to "Hong Kong Disneyland Duffy and Friends Springtime Festival", we synthesize a corresponding deep research query that captures the user's underlying intent, such as "During the 'Duffy and Friends' Springtime limited-edition event at Hong Kong Disneyland in March 2026, how can one plan a counter-clockwise touring route to avoid crowds, and obtain the latest question bank and practical procedures for the staff survey gifts?" Concretely, we employ a strong LLM (i.e., Gemini-3-Flash (DeepMind, 2025b)) to

extract latent deep research interests from user interactions and reformulate them into well-structured deep research queries. Detailed prompts are provided in Appendix A.1.

To ensure rigor and high quality, we further apply an inspection pipeline that distills 100 high-quality queries from an initial pool of 260,000 generated candidates. The pipeline combines automated LLM-based screening with human verification, together with an optional difficulty-expansion step in between. Specifically, we form a review committee powered by Gemini-3-Pro (DeepMind, 2026a) to automatically evaluate each query along the following criteria: (i) naturalness and clarity: whether the query is fluent, unambiguous, and faithfully reflects a plausible user intent; (ii) indispensability of Rednote-specific knowledge: whether answering the query genuinely relies on Rednote’s user-generated content or community insights, rather than being trivially solvable via generic web search, while ensuring that no user privacy is compromised; (iii) real-world plausibility: whether the query type aligns with realistic information-seeking behaviors observed in daily life, rather than representing a contrived or synthetic use case. Queries that pass this automatic review are then forwarded to human annotators for final verification, ensuring that the resulting benchmark is both authentic and challenging.

3.2 Statistics and Comparisons with Existing Benchmarks

Our GALA benchmark consists of 100 carefully curated deep research queries. Following the taxonomy proposed in DeepResearchBench (Du et al., 2025), each query is categorized into one of the following 22 topics: “Finance & Business”, “Science & Technology”, “Software Development”, “Education & Job”, “Health”, “Literature”, “History”, “Hardware”, “Industrial”, “Art & Design”, “Games”, “Crime & Law”, “Entertainment”, “Sports & Fitness”, “Software”, “Transportation”, “Religion”, “Home & Hobbies”, “Travel”, “Food & Dining”, “Fashion & Beauty”, “Social Life”. We adopt Gemini-3-Pro (DeepMind, 2026a) as the query classifier; the detailed prompts are provided in Appendix A.2. The comparison reveals a clear divergence in topical focus across benchmarks. Existing benchmarks are heavily skewed toward professional and technical domains: DeepResearchBench (Du et al., 2025) is dominated by “Science & Technology” (26.8%), “Finance & Business” (17.0%), and “Education & Jobs” (10.0%); DeepConsult (Consult, 2025) is overwhelmingly concentrated in “Finance & Business” (85.3%), followed by “Science & Technology” (7.8%); and DeepResearchGym (Coelho et al., 2025) is led by “Finance & Business” (20.0%), “Science & Technology” (15.0%), “History” (14.0%), “Social Life” (13.0%), and “Health” (11.0%). In contrast, GALA exhibits a markedly different distribution centered on everyday-life domains, with “Home & Hobbies” (25.0%), “Travel” (18.0%), “Fashion & Beauty” (18.0%), and “Education & Jobs” (12.0%) as its dominant categories. This contrast highlights GALA’s unique role in complementing existing benchmarks: rather than re-emphasizing professional research scenarios, GALA targets authentic, daily-life information needs that have been largely underrepresented in prior evaluations.

For clarity, Figure 6 illustrates the top six taxonomic categories of GALA along with their corresponding distributions across DeepResearchBench, DeepConsult, and DeepResearchGym. As can be observed, the proposed GALA serves as a critical complement to existing benchmarks, offering essential coverage for the “Home & Hobbies”, “Travel”, and “Fashion & Beauty” domains that are notably underrepresented in the aforementioned datasets.

3.3 Evaluation Protocols

One prevailing challenge in evaluating open-ended deep research reports is the absence of an exact ground truth for each query. To address this, we adopt the RACE metric proposed in DeepResearch Bench (Du et al., 2025). The evaluation proceeds in two stages: (i) **dynamic dimension weight allocation**: An LLM (i.e., Gemini-3-Flash (DeepMind, 2025b) in practice), acting as a meta-evaluator, analyzes the input query to determine the relative importance of four evaluation dimensions: Comprehensiveness, Insight, Instruction-Following, and Readability. (ii) **reference-based pair-wise scoring**: A separate LLM (i.e., Gemini-3-Flash (DeepMind, 2025b) in practice) then scores the target report against criteria derived for each dimension. The final score is computed as a weighted summation over the per-dimension scores. Since RACE evaluation requires a reference report to produce pairwise scores, we release—alongside the open-sourced queries—reports generated by our AgentDisCo to serve as the corresponding reference reports.

We do not adopt the FACT metric, which is also developed in DeepResearch Bench (Du et al., 2025). FACT verifies the factual accuracy of references by fetching their corresponding web content; however, as time passes, many referenced web pages become inaccessible (e.g., returning 404 errors), making the metric unreliable for reproducible evaluation.

Table 1: Performance of agents on DeepResearch Bench in terms of comprehensiveness (Comp.), insight, instruction-following (Inst.), readability (Read.), effective citations (Eff. c.), and citation accuracy (C. acc.). The best results are highlighted with purple color, and the second-best results are highlighted with underlines.

Agent systems	RACE				FACT		
	Overall	Comp.	Insight	Inst.	Read.	Eff. c.	C. acc.
Langchain-Open-Deep-Research	43.44	42.97	39.17	48.09	45.22	-	-
Doubao-Research	44.34	44.84	40.56	47.95	44.69	52.62	52.86
Kimi-Research	44.64	44.96	41.97	47.14	45.59	-	-
Claude-Research	45.00	45.34	42.79	47.58	44.66	-	-
Openai-Deepresearch	46.45	46.46	43.73	49.39	47.22	39.79	75.01
Gemini-2.5-Pro-Deepresearch	49.71	49.51	49.45	50.12	50.00	165.34	78.30
AgentDisCo (Gemini-2.5-Pro)	51.44	51.23	52.49	51.57	50.39	63.94	89.06
AgentDisCo w/ Harness (Gemini-2.5-Pro)	52.11	51.89	53.43	51.87	50.45	69.65	89.55
AgentDisCo (Claude-Opus-4.6)	54.02	53.38	56.65	53.11	51.53	89.88	93.56

4 Experiments

4.1 Setups

Benchmarks. We evaluate AgentDisCo on three publicly available benchmarks, together with our proposed GALA benchmark, as detailed below.

- **DeepResearch Bench** (Du et al., 2025) comprises 100 PhD-level complex research tasks meticulously formulated by domain experts across 22 distinct fields, including Science & Technology, Finance & Business, Software Engineering, and Art & Design.
- **DeepConsult** (Consult, 2025) is a specialized collection of prompts tailored for in-depth research within the business and consulting domains. Its queries span a wide range of topics, such as marketing strategy, financial analysis, emerging technology trends, and business planning.
- **DeepResearchGym** (Coelho et al., 2025) is used to assess performance on real-world, complex queries. It contains 100 queries sampled from the large-scale Researchy Questions dataset (Rosset et al., 2024), which comprises approximately 96,000 authentic information-seeking queries.
- **GALA** is a deep research benchmark for general AI life assistants, introduced in Section 3.

Metrics. We adopt the official evaluation metrics and recommended judge LLMs for each benchmark.

- **DeepResearch Bench** (Du et al., 2025) employs two suites of metrics to evaluate different aspects of the system’s output: (i) RACE (Report Quality) assesses the quality of the generated report against a reference report along four dimensions—Comprehensiveness (Comp.), Insight/Depth (Insight), Instruction-Following (Inst.), and Readability (Read.)—with an overall score computed as a weighted sum of these components. (ii) FACT (Web Retrieval via Google Search) measures the effectiveness and reliability of the information retrieval process, including Citation Accuracy (C. Acc.) and the Average Effective Citations per Task (Eff. c.). Following the benchmark’s protocol against Gemini-2.5-Pro-DeepResearch, we adopt Gemini-2.5-Pro (DeepMind, 2025a) as the judge model.
- **DeepConsult** (Consult, 2025) evaluates performance via pairwise comparison against the OpenAI-DeepResearch baseline. The primary metrics are win rate, tie rate, and loss rate, supplemented by an average quality score. The judge model is GPT-4.1-20250414 (OpenAI, 2025c).
- **DeepResearchGym** (Coelho et al., 2025) employs an LLM judge to assess the generated report along several quality dimensions, including clarity, insightfulness, depth, balance, breadth, and support, as well as an overall average quality score. The judge model is GPT-4.1-mini-20250414 (OpenAI, 2025b).
- **GALA** evaluates report quality using the RACE metric with Gemini-3-Flash (DeepMind, 2025b) as the judge model.

Compared Systems. We benchmark our AgentDisCo system against a suite of leading deep research agents available on the market: **LangChain-Open-Deep-Research** (LangChain, Inc., 2023), **Doubao-Research** (Research, 2026a), **Kimi-Research** (Research, 2025b), **Claude-Research** (anthropic, 2025), **OpenAI-DeepResearch** (OpenAI, 2025a), and **Gemini-2.5-Pro-DeepResearch** (Research, 2025a). Their results on the three public benchmarks are taken directly from Li et al. (2025); Han et al. (2025).

For the GALA benchmark, to construct competitive baselines, our human annotation team manually collected reports from the official web interfaces of **Doubao-Research** (Research, 2026a) and **Qwen-Research** (Research, 2026c), as well as outputs from **OpenAI o3-DeepResearch** (Research, 2026b) obtained via its API, with all data acquired in April 2026. In our evaluations, AgentDisCo, instantiated with

Table 2: Performance of agents on DeepConsult in terms of win rate and average scores and on DeepResearchGym in terms of clarity (Cla.), depth, balance (Bal.), breadth (Brea.), support (Sup.), and insightfulness (Ins.). The best results are highlighted with purple color, and the second-best results are highlighted with underlines.

Agent systems	DeepConsult				DeepResearchGym						
	Win	Tie	Lose	Overall	Cla.	Depth	Bal.	Brea.	Sup.	Ins.	Overall
Doubao-research	29.95	40.35	29.70	5.42	68.85	93.12	83.96	93.33	84.38	83.12	84.46
Claude-research	25.00	38.89	36.11	4.60	86.67	96.88	84.41	96.56	26.77	90.22	80.25
Openai-deepresearch	0.00	100.00	0.00	5.00	84.90	98.10	89.80	97.40	88.40	89.00	91.27
Gemini-2.5-pro-deepresearch	61.27	31.13	7.60	6.70	90.71	99.90	93.37	99.69	95.00	<u>97.45</u>	96.02
AgentDisCo (Gemini-2.5-Pro)	<u>53.26</u>	<u>37.50</u>	<u>9.23</u>	6.75	<u>90.50</u>	<u>100.00</u>	<u>93.75</u>	<u>100.00</u>	<u>96.25</u>	<u>93.75</u>	<u>95.63</u>
AgentDisCo w/ Harness (Gemini-2.5-Pro)	<u>56.86</u>	<u>32.47</u>	<u>10.67</u>	<u>6.86</u>	<u>90.98</u>	<u>100.00</u>	<u>94.30</u>	<u>100.00</u>	<u>97.73</u>	<u>95.22</u>	<u>96.21</u>
AgentDisCo (Claude Opus 4.6)	<u>65.88</u>	<u>22.47</u>	<u>11.65</u>	<u>7.06</u>	<u>90.85</u>	<u>100.00</u>	<u>97.85</u>	<u>100.00</u>	<u>98.93</u>	<u>98.66</u>	<u>97.54</u>

Gemini-2.5-Pro, serves as the reference system. To analyze the effect of retrieval-source selection, we further introduce two variants. **AgentDisCo w/ Rednote** replaces the default retrieval component with the Rednote Search Engine, enabling the agent to retrieve information from Rednote-specific content. In addition, **AgentDisCo w/ Rednote & Google** performs joint retrieval over both the Rednote Search Engine and Google Search Engine, allowing us to examine whether combining social-media-oriented and general web search sources can provide complementary evidence for report generation.

4.2 Main Results

Results on DeepResearch Bench. As shown in Table 1, AgentDisCo consistently outperforms existing deep-research agent systems on DeepResearch Bench. When using Gemini-2.5-Pro as the backbone model, AgentDisCo achieves an overall RACE score of 51.44, surpassing the prior system based on Gemini-2.5-Pro, i.e., Gemini-2.5-Pro-Deepresearch. The improvement is particularly pronounced in insight, comprehensiveness, and instruction following, where AgentDisCo obtains 52.49, 51.23, 51.57, respectively, compared with 49.45, 49.51, 50.12 from Gemini-2.5-Pro-Deepresearch. In contrast, the gain in readability is relatively moderate, suggesting that the advantage of AgentDisCo does not mainly come from more fluent surface-level writing, but rather from producing more substantive, better-supported, and better-structured research content. This improvement can be attributed to the disentangled yet collaborative design between the critic agent and the generator agent in our AgentDisCo framework. The critic agent iteratively evaluates the intermediate report, identifies missing aspects, weak arguments, and insufficient evidence, and then provides targeted feedback to guide subsequent generation. Meanwhile, the generator agent incorporates this feedback to expand the research scope, refine the argument structure, and strengthen evidence grounding. Such an iterative critic-generator collaboration naturally improves comprehensiveness and insight, as the system is encouraged to go beyond a single-pass synthesis and progressively discover under-explored perspectives. Moreover, although AgentDisCo does not produce the largest number of effective citations, it achieves substantially higher citation accuracy. Specifically, AgentDisCo with Gemini-2.5-Pro obtains a citation accuracy of 89.06, improving over Gemini-2.5-Pro-Deepresearch by over 10 points. This indicates that AgentDisCo favors reliable and relevant evidence usage rather than simply increasing the citation count.

We further evaluate the effect of the harness optimization. AgentDisCo w/ Harness improves the Gemini-2.5-Pro-based AgentDisCo from 51.44 to 52.11 in overall RACE score, with consistent gains across all RACE dimensions, including comprehensiveness, insight, instruction-following, and readability. It also improves factual grounding, increasing effective citations from 63.94 to 69.65 and citation accuracy from 89.06 to 89.55. These results demonstrate that the harness optimization provides a stable additional benefit by better coordinating the interaction process and improving the reliability of evidence integration. Finally, we instantiate AgentDisCo with a stronger frontier backbone model, Claude-Opus-4.6 (anthropic, 2026), to examine the scalability of our framework. AgentDisCo with Claude-Opus-4.6 achieves the best overall performance, reaching 54.02 on RACE. Notably, it obtains a substantial insight score of 56.65, far exceeding all other systems, which suggests that AgentDisCo can effectively leverage stronger reasoning capabilities from advanced LLMs. It also achieves the highest citation accuracy of 93.56, showing that the framework remains highly reliable when scaled to a more capable base model. Overall, these results indicate that AgentDisCo is both effective and scalable: its disentangled yet collaborative critic-generator mechanism improves research depth, evidence reliability, and report quality across different backbone models.

Results on DeepConsult and DeepResearchGym. To examine whether AgentDisCo generalizes beyond our main evaluation setting, we further evaluate it on DeepConsult and DeepResearchGym, as reported in Table 2. Overall, AgentDisCo exhibits strong and robust performance across both benchmarks. With Gemini-2.5-Pro as the backbone, AgentDisCo achieves an overall score of 6.75 on DeepConsult, slightly

Table 3: Performance of agents on our proposed GALA benchmark in terms of comprehensiveness (Comp.), insight, instruction-following (Inst.), and readability (Read.). The best results are highlighted in purple, and the second-best results are underlined.

Agent systems	RACE				
	Overall	Comp.	Insight	Inst.	Read.
Doubao-Research (2026-04)	49.82	50.87	47.42	50.65	50.86
Qwen-Research (2026-04)	46.69	45.38	45.36	47.23	49.56
OpenAI o3-DeepResearch (2026-04)	45.88	45.37	42.88	48.04	47.72
AgentDisCo (Gemini-2.5-Pro)	50.00	50.00	50.00	50.00	50.00
AgentDisCo w/ Harness (Gemini-2.5-Pro)	50.58	50.41	51.24	50.16	49.85
AgentDisCo w/ Rednote (Gemini-2.5-Pro)	<u>51.02</u>	50.88	<u>51.11</u>	<u>51.25</u>	50.95
AgentDisCo w/ Rednote & Harness (Gemini-2.5-Pro)	51.90	51.61	53.44	51.78	50.67
AgentDisCo w/ Rednote & Google (Gemini-2.5-Pro)	50.95	<u>51.21</u>	50.44	50.78	49.79

surpassing the Gemini-2.5-Pro-deepresearch baseline in average score. On DeepResearchGym, it obtains nearly perfect scores in both Depth and Breadth, indicating that the proposed iterative research procedure is effective in expanding both the depth and coverage of the generated reports.

As analyzed in the previous subsection, we further evaluate the effect of the harness optimization. Compared with the vanilla Gemini-based AgentDisCo, AgentDisCo w/ Harness improves the DeepConsult win rate from 53.26% to 56.86% and the overall score from 6.75 to 6.86. On DeepResearchGym, the harness also raises the overall score from 95.63 to 96.21, yielding improvements in clarity, balance, support, and insightfulness while preserving nearly perfect scores in Depth and Breadth. These results suggest that the harness optimization does not merely increase coverage but also improves the reliability and controllability of the agent execution process, allowing the planned research workflow to be more effectively translated into high-quality final reports. Moreover, we instantiate AgentDisCo with a stronger frontier backbone model, Claude-Opus-4.6 (anthropic, 2026), to examine the scalability of our framework. This variant achieves the best performance on both benchmarks: it obtains the highest DeepConsult win rate of 65.88% and the best overall score of 7.06, and also achieves the highest DeepResearchGym overall score of 97.54. In particular, it ranks first in Balance, Support, and Insightfulness, while maintaining perfect scores in Depth and Breadth. These results indicate that AgentDisCo can effectively leverage stronger underlying models, suggesting that the proposed framework scales favorably with frontier model capability. Taken together, the results provide evidence for the effectiveness of AgentDisCo’s core design. The consistently high Depth and Breadth scores reflect the benefit of the planner’s iterative research cycle, which enables the system to progressively expand and refine the information space beyond static one-shot planning.

Results on GALA. Beyond existing benchmarks, we further construct a lifestyle-oriented deep research benchmark, GALA, as introduced in Section 3. Table 3 reports the evaluation results in terms of comprehensiveness, insight, instruction-following, and readability. In this reference-based evaluation, AgentDisCo, instantiated with Gemini-2.5-Pro, serves as the reference system and is therefore assigned a score of 50.00 across all dimensions. Notably, mainstream deep research systems, including Doubao-Research, Qwen-Research, and OpenAI o3-DeepResearch, obtain lower overall scores than this reference. This result highlights the advantage of the AgentDisCo framework itself: rather than relying solely on the capability of a strong backbone model, AgentDisCo benefits from its structured workflow that combines iterative planning, targeted evidence acquisition, and hierarchical synthesis, which is particularly important for lifestyle-oriented research tasks requiring practical relevance, contextual understanding, and user-aligned recommendations. We next examine the effect of harness optimization. AgentDisCo w/ Harness improves the overall score from 50.00 to 50.58, with the most evident gain appearing in insight. This suggests that the harness helps stabilize the execution of the multi-step research workflow and enables the agent to produce more informative and analytically useful responses.

We then study the impact of retrieval-source selection. From the table, we can observe that AgentDisCo w/ Rednote achieves a higher overall score of 51.02 and obtains the best readability score among all systems. This indicates that Rednote provides domain-relevant lifestyle content that aligns well with the information needs in GALA, leading to responses that are more natural and accessible for lifestyle-oriented scenarios. Combining domain-specific retrieval with harness optimization yields the strongest performance. AgentDisCo w/ Rednote & Harness achieves the best overall score of 51.90 and ranks first in comprehensiveness, insight, and instruction-following. In particular, its insight score reaches 53.44, showing that lifestyle-oriented social content, when integrated through a more reliable execution harness, can substantially improve the practical and contextual value of the generated reports. This demonstrates that the gains from Rednote are further amplified when the agent’s research process is better controlled and more consistently executed.

4.3 Analysis

Statistics of Outline Optimizations. One of the core ideas of AgentDisCo is its disentangled yet collaborative framework, in which different agents are assigned specialized roles and interact through an iterative optimization process. To directly assess whether this design improves the quality of the generated outline, we evaluate the effect of our outline optimization module in isolation. Since the optimization process is initiated by the critic agents after an initial outline has been produced, at least two rounds are required to observe the effect of critic-guided refinement. To isolate and quantify the contribution of outline optimization, we conduct an ablation study on the end-to-end benchmarks, as reported in Figures 7 and 8. Specifically, we collect samples from DeepResearch Bench and DeepResearchGym and apply up to three rounds of outline optimization, while keeping the subsequent writing strategy unchanged across all settings. This design ensures that performance differences can be primarily attributed to the quality of the optimized outlines rather than variations in the writing process. The benefits of iterative refinement are consistent across both benchmarks. On DeepResearch Bench, the overall score increases steadily as the number of optimization rounds grows, with particularly notable improvements in comprehensiveness and insight. This supports our hypothesis that each optimization round enables the planner to construct a more detailed, coherent, and logically organized outline. A similar trend is observed on DeepResearchGym, where later optimization rounds achieve substantially stronger scores in depth and breadth, indicating more exhaustive coverage of the target topic.

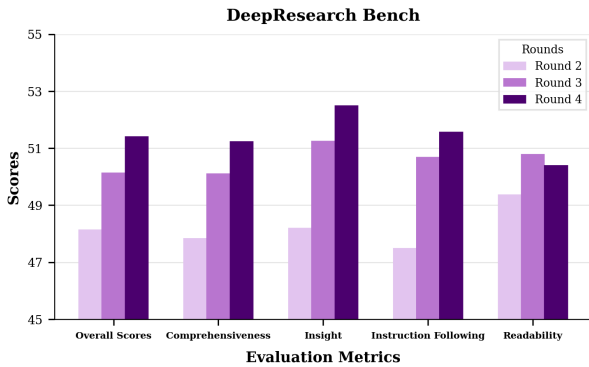


Figure 7: End-to-end scores with varying rounds of outline optimization on Deepresearch Bench.

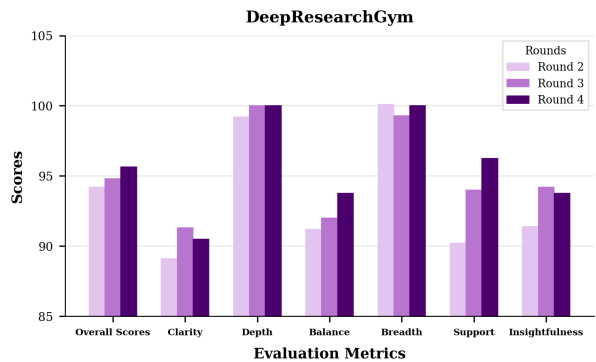


Figure 8: End-to-end scores with varying rounds of outline optimization on DeepresearchGym.

Consistency between Harness-Based Critic Optimization and End-to-End Performance. To evaluate the effectiveness of the proposed harness, we examine whether improvements measured by the harness are consistent with end-to-end performance gains. In the scoring agent of our harness framework, whose detailed prompts are provided in Appendix B.2, we introduce an intermediate metric, named Search Coverage, ranging from 0 to 100, to quantify the quality and coverage of the generated search queries. A higher Search Coverage score indicates that the queries are more likely to capture the key aspects required for answering the research question.

Specifically, we evaluate the Search Coverage scores at optimization rounds 0, 10, and 20, and compare them with the corresponding end-to-end scores on DeepResearch Bench. To reduce evaluation cost while maintaining a representative assessment, we randomly sample 50 examples from DeepResearch Bench as the evaluation pool. As shown in Figure 9, Search Coverage increases from 62.50 at round 0 to 79.25 at round 10 and further to 82.05 at round 20. This trend is accompanied by a consistent improvement in the end-to-end overall score, which increases from 51.41 to 51.82 and then to 52.11. These results indicate that the harness-based optimization signal is well aligned with downstream benchmark performance. In particular, improving the coverage and quality of search queries leads to more effective evidence acquisition, which in turn contributes to better final reports. Therefore, the proposed harness provides a meaningful and practical intermediate optimization objective for improving end-to-end deep research performance.

Superiority of Rednote Search over Life-style Search Queries. To better understand the role of retrieval sources in lifestyle-oriented deep research, we further compare Rednote Search with general web search. Although both sources can be queried with lifestyle-related search queries, they differ substantially in content style and information structure: Google Search primarily retrieves general web pages, while Rednote Search provides user-generated, experience-oriented, and scenario-specific content that is often more aligned with daily-life decision making. This raises an important question: whether the improvement comes merely from using lifestyle-style queries, or from the domain-specific characteristics of Rednote as a retrieval source.

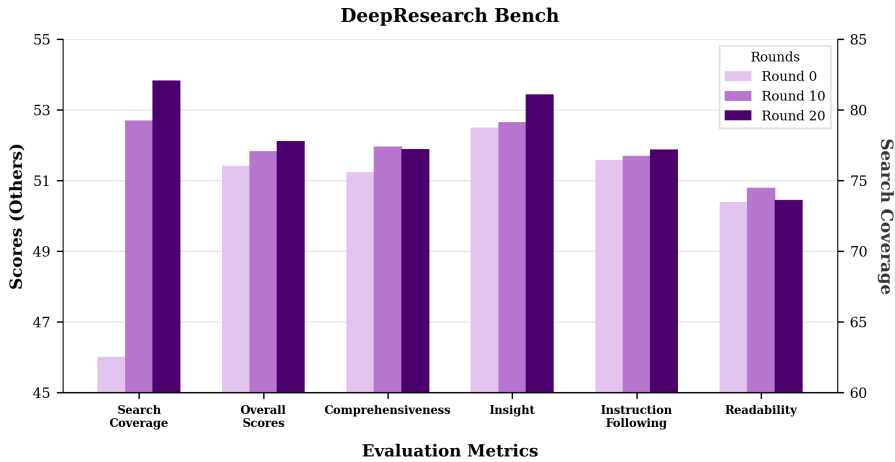


Figure 9: Consistency of harness optimization over critic agent to end-to-end optimization with varying rounds of harness optimization on Deepresearch Bench.

Besides the results reported in Section 4.2, we further evaluate joint retrieval over Rednote and Google and report the results in Table 3. AgentDisCo w/ Rednote & Google obtains an overall score of 50.95, which is substantially higher than using the default Google-based retrieval alone, but slightly lower than using Rednote alone. Notably, joint retrieval achieves the second-best comprehensiveness score of 51.21, indicating that Google can complement Rednote by broadening factual and topical coverage. However, this broader coverage does not translate into better overall performance: compared with AgentDisCo w/ Rednote, the joint-retrieval variant shows lower scores in insight, instruction-following, and readability. This suggests that general web search may introduce less lifestyle-specific or less user-oriented evidence, increasing the burden of evidence filtering and synthesis. Overall, these findings indicate that the advantage of Rednote does not simply come from issuing lifestyle-related queries, but from the nature of the retrieved content itself. Rednote contributes domain-specific, experience-rich, and practically grounded evidence, which is particularly valuable for GALA-style tasks. Meanwhile, adding Google can improve coverage, but without effective filtering and synthesis, increased retrieval breadth may introduce noise and reduce readability. Therefore, for lifestyle-oriented deep research, domain-relevant retrieval quality is more important than retrieval breadth alone.

5 Conclusion and Future Work

In this paper, we presented **AgentDisCo**, a disentangled and collaborative agentic framework for open-ended deep research. By separating information exploration from information exploitation and formulating their interaction as an iterative adversarial optimization process, AgentDisCo enables critic and generator agents to progressively refine search queries and research outlines before final report synthesis. We further introduced a meta-optimization harness that automatically discovers reusable design strategies for improving the critic agent, allowing the framework to enhance its own search and planning behavior with limited human intervention.

Extensive experiments on DeepResearchBench, DeepConsult, and DeepResearchGym demonstrate that AgentDisCo achieves competitive or superior performance compared with leading closed-source deep research systems. To better reflect real-world user needs, we also introduced **GALA**, a lifestyle-oriented deep research benchmark mined from users’ browsing histories, and showed that AgentDisCo is effective in this more practical setting. Finally, we developed a rendering agent and a product demonstration, “AutoResearch Your Interest”, to make deep research outputs more accessible and personalized for end users. We release our benchmark, code, demo, and evaluation harness to facilitate future research on open-ended, user-centered deep research agents.

References

- anthropic. Meet claude, 2025. URL <https://www.anthropic.com/claude>.
- anthropic. Introducing claude opus 4.6, 2026. URL <https://www.anthropic.com/news/claude-opus-4-6>.
- João Coelho, Jingjie Ning, Jingyuan He, Kangrui Mao, Abhijay Paladugu, Pranav Setlur, Jiahe Jin, Jamie Callan, João Magalhães, Bruno Martins, et al. Deepresearchgym: A free, transparent, and reproducible evaluation sandbox for deep research. *arXiv preprint arXiv:2505.19253*, 2025.
- Deep Consult. Deep consult. 2025. URL <https://github.com/Su-Sea/ydc-deep-research-evals>.
- Google DeepMind. Gemini 2.5, 2025a. URL <https://blog.google/technology/google-deepmind/gemini-model-thinking-updates-march-2025/>.
- Google DeepMind. Gemini-3-flash, 2025b. URL <https://deepmind.google/models/gemini/flash/>.
- Google DeepMind. Gemini 3 pro, 2026a. URL https://aistudio.google.com/models/gemini-3-pro-image?utm_source=google&utm_medium=cpc&utm_campaign=Cloud-SS-DR-AIS-FY26-global-gsem-1713578&utm_content=text-ad&utm_term=KW_image%20api&gad_source=1&gad_campaignid=23417416052&gbraid=0AAAAACn9t671c5F9JLzrVf8S8bh39Ky16&gclid=Cj0KCQjw2MbPBhCSARISAP3jP9wnU46TLFWDeB3AV4ZZbHsMzkkviWFO03EyKvxaqqB3S1L1wMVSbTEaAi7vEALw_wcB.
- Google DeepMind. Gemini 3 pro image, 2026b. URL <https://deepmind.google/models/gemini/pro/>.
- Mingxuan Du, Benfeng Xu, Chiwei Zhu, Xiaorui Wang, and Zhendong Mao. Deepresearch bench: A comprehensive benchmark for deep research agents. *arXiv preprint arXiv:2506.11763*, 2025.
- Rujun Han, Yanfei Chen, Zoey CuiZhu, Lesly Miculicich, Guan Sun, Yuanjun Bi, Weiming Wen, Hui Wan, Chunfeng Wen, Solène Maître, George Lee, Vishy Tirumalashetty, Emily Xue, Zizhao Zhang, Salem Haykal, Burak Gokturk, Tomas Pfister, and Chen-Yu Lee. Deep researcher with test-time diffusion, 2025. URL <https://arxiv.org/abs/2507.16075>.
- LangChain, Inc. LangChain: Building applications with LLMs through composability, 2023. URL <https://python.langchain.com/>.
- Yoonho Lee, Roshen Nair, Qizheng Zhang, Kangwook Lee, Omar Khattab, and Chelsea Finn. Meta-harness: End-to-end optimization of model harnesses. *arXiv preprint arXiv:2603.28052*, 2026.
- Yu Lei, Shuzheng Si, Wei Wang, Yifei Wu, Gang Chen, Fanchao Qi, and Maosong Sun. Rhinoinstight: Improving deep research through control mechanisms for model behavior and context. *arXiv preprint arXiv:2511.18743*, 2025.
- Zijian Li, Xin Guan, Bo Zhang, Shen Huang, Houquan Zhou, Shaopeng Lai, Ming Yan, Yong Jiang, Pengjun Xie, Fei Huang, Jun Zhang, and Jingren Zhou. Webweaver: Structuring web-scale evidence with dynamic outlines for open-ended deep research, 2025. URL <https://arxiv.org/abs/2509.13312>.
- Qianli Ma, Siyu Wang, Yilin Chen, Yinhao Tang, Yixiang Yang, Chang Guo, Bingjie Gao, Zhening Xing, Yanan Sun, and Zhipeng Zhang. Human-agent collaborative paper-to-page crafting for under \$0.1. In *arXiv preprint arXiv:2510.19600*, 2025. URL <https://arxiv.org/abs/2510.19600>.
- OpenAI. Deep research system card, 2025a. URL <https://cdn.openai.com/deep-research-system-card.pdf>.
- OpenAI. Gpt-4.1-mini-20250414., 2025b. URL <https://openai.com/index/gpt-4-1/>.
- OpenAI. Gpt-4.1-20250414., 2025c. URL <https://openai.com/index/gpt-4-1/>.
- Wei Pang, Kevin Qinghong Lin, Xiangru Jian, Xi He, and Philip Torr. Paper2poster: Towards multimodal poster automation from scientific papers. In *arXiv preprint arXiv:2505.21497*, 2025. URL <https://arxiv.org/abs/2505.21497>.
- Doubao Deep Research. Doubao deep research. 2026a. URL <https://www.doubao.com/chat/>.
- Gemini Research. Gemini research. 2025a. URL <https://gemini.google/overview/deep-research/>.
- Kimi Deep Research. Kimi deep research. 2025b. URL <https://www.kimi.com/>.
- O3 Deep Research. O3 deep research. 2026b. URL <https://developers.openai.com/api/docs/models/o3-deep-research>.

-
- Qwen Deep Research. Qwen deep research. 2026c. URL https://chat.qwen.ai/?inputFeature=deep_research.
- Stephen Robertson and Hugo Zaragoza. *The probabilistic relevance framework: BM25 and beyond*, volume 4. Now Publishers Inc, 2009.
- Corby Rosset, Ho-Lam Chung, Guanghui Qin, Ethan C Chau, Zhuo Feng, Ahmed Awadallah, Jennifer Neville, and Nikhil Rao. Researchy questions: A dataset of multi-perspective, decompositional questions for llm web agents. *arXiv preprint arXiv:2402.17896*, 2024.
- Tao Sun, Enhao Pan, Zhengkai Yang, Kaixin Sui, Jiajun Shi, Xianfu Cheng, Tongliang Li, Wenhao Huang, Ge Zhang, Jian Yang, and Zhoujun Li. P2p: Automated paper-to-poster generation and fine-grained benchmark. In *arXiv preprint arXiv:2505.17104*, 2025. URL <https://arxiv.org/abs/2505.17104>.
- Zexuan Yan, Jiarui Jin, Yue Ma, Shijian Wang, Jiahui Hu, Wenxiang Jiao, Yuan Lu, and Linfeng Zhang. Glyphbanana: Advancing precise text rendering through agentic workflows. In *arXiv preprint arXiv:2603.12155*, 2026. URL <https://arxiv.org/abs/2603.12155>.
- Zhilin Zhang, Xiang Zhang, Jiaqi Wei, Yiwei Xu, and Chenyu You. Postergen: Aesthetic-aware paper-to-poster generation via multi-agent llms. In *arXiv preprint arXiv:2508.17188*, 2025a.
- Zhilin Zhang, Xiang Zhang, Jiaqi Wei, Yiwei Xu, and Chenyu You. Postergen: Aesthetic-aware paper-to-poster generation via multi-agent llms. *arXiv preprint arXiv:2508.17188*, 2025b.

Appendix

A Prompt Design

A.1 Prompt for Deep Research Query Miner

As depicted in Figure 2, the query miner is designed to elicit users' latent interests by systematically analyzing their browsing history. The detailed prompts are listed below.

Prompt: Research Query Generator

Role Definition

You are a **user behavior analysis + research task design** expert.

Your task is: based on the user's Xiaohongshu (RedNote) browsing history, generate 4–6 complex Deep Research queries.

These queries will be submitted to a Research Agent with the following capabilities:

- Can search for information online
- Can access multiple data sources
- Can perform multi-step reasoning
- Can synthesize multiple information sources to draw conclusions

Therefore, queries must:

- Be clear, structured, and actionable
- Focus on decision-making actions, driving users to make specific decisions
- Require support from multiple information sources and cannot be answered directly by common sense

Step 1 — User Interest Modeling (*internal reasoning, do not output*)

Extract three layers of interests from the Xiaohongshu browsing history:

[Long-term Preferences]

Identify ≥ 3 sustained interest directions, based on:

- High-frequency likes / saves / time-on-post over 60 seconds
- Repeated appearances across different time periods
- Involving multiple brands or content formats

Common directions include: fashion & styling / beauty & skincare / home renovation / fitness & sports / food exploration / travel

[Short-term Interests]

Identify ≥ 3 recent concentrated browsing hotspots, based on:

- Highly concentrated recent browsing
- High frequency of appearance, high content similarity
- Pointing to a specific product / brand / event / scenario

[Potential Decision Needs]

Based on long-term preferences and short-term interests, infer actual decisions the user may currently face, such as:

- Whether to purchase a certain product
- How to choose among several options
- How to plan a trip or renovation
- Whether to try a new brand or new solution

When generating queries, must ensure:

- At least 2 queries come from long-term preferences
- At least 1 query comes from short-term interests
- Different queries cover different decision types

Step 2 — Research Query Structure (*internal reasoning, do not output*)

Each query must contain the following five components during construction, and be presented in natural language at the end:

[1] Background Context (1–2 sentences)

Describe a real user scenario, clearly stating the user's goal or confusion, such as:

- I'm planning to purchase ...
- I'm planning to ...
- I'm considering whether to ...
- I'm torn between the following options ...

[2] Research Subjects

Clearly list the specific subjects to be researched, using numbered format:

- (1) ...
- (2) ...
- (3) ...

Research subjects can be: products / brands / solutions / locations / platforms / strategies

[3] Research Questions (3–5)

Use numbered format to list questions that require research to answer, such as:

1. What are real user reviews and reputation like
2. Does actual experience match the advertised claims
3. How is the price and value for money
4. What are common issues or usage risks
5. How high is the long-term usage cost or maintenance difficulty

Research questions must be: specific, searchable, comparable, and analyzable

[4] Constraints (2–4)

Clearly state the user's limiting conditions, such as:

- Budget constraint: total budget not exceeding XXX
- Usage scenario: mainly for commuting / suitable for southern climate / no drilling allowed in rental
- Personal preference: prioritizes durability / dislikes complex maintenance
- Time condition: needs to make a decision soon

[5] Research Goal

Clearly state the action conclusion the user hopes to obtain, such as:

- Provide a ranked recommendation
- Determine whether it is worth purchasing
- Choose the optimal solution
- Formulate a specific action plan

Step 3 — Query Depth Requirements

Each query must satisfy:

- Requires at least 3 different information sources to answer
- Requires comparison of multiple subjects
- Requires synthesis of real user reviews, product data, or usage experience

Strictly prohibited: generating questions that can be answered directly by common sense.

Step 4 — Query Type Diversity

All queries must belong to the **Decision Making** type, must cover the following five categories, and 4–6 queries must not all concentrate on the same type:

Type	Core Feature	Example Phrasing
Comparison & Selection	Weighing pros and cons among multiple options, ultimately choosing one	"Which is more suitable for me, XX or XX? Compare from XX dimensions"
Recommendations	Seeking personalized recommendations with no specific candidates	"In XX scenario, what XX products are worth buying?"
How-to Guide	Seeking specific actionable step-by-step methods	"How to complete XX step by step? Including pitfall-avoidance tips"
Travel Planning	Developing travel plans or itinerary arrangements	"How to plan the best itinerary for X days in XX?"
Purchase Decision	Deciding whether to buy, when to buy, which one to buy	"Is it a good time to buy XX now? Key purchase considerations"

Warning: Strictly prohibited from generating the following types of queries:

-
- Pure information queries: “What is XX”, “History of XX”
 - Status tracking: “Latest developments of XX”, “What’s happening with XX now”
 - Trend analysis, mechanism research, systematic reviews, or other academically-oriented content

Step 5 — Complexity Requirements

Each query must:

- Length: **120–220 words**
- Include structured numbered lists
- Include **3–5 research questions**
- Include explicit constraints
- End with a **specific action goal**

Step 6 — Expression Style Requirements

[Authentic User Language]

Queries must read like real users commissioning a research assistant, using question or imperative sentences.

Wrong example: “Comparative study on anti-aging mechanisms of retinol vs. Proxylane”

Correct example: “I want to start anti-aging skincare but don’t know where to begin — should I go with retinol or Proxylane?”

[Avoid Academic Language]

Prohibited: mechanism research / trend analysis / in-depth analysis / systematic review / paper-title-style stacking

[Avoid Content Stacking]

Do not compress multiple research dimensions into a single sentence — must be expressed in structured, separate lines.

Step 7 — Example Query

I've been seeing a lot of rental apartment renovation content on Xiaohongshu lately, and I want to transform my living room into a creamy minimalist style without drilling or making major structural changes. I'm currently torn between the following directions:

- (1) Multi-functional storage sofa + cream-colored curtains combination
- (2) Rattan storage cabinet + arched floor lamp pairing
- (3) Modular bookshelf wall + minimalist rug solution

Please help me evaluate:

- (1) Which solution has higher practical feasibility under rental restrictions
- (2) In real renovation cases on Xiaohongshu, which direction has a lower failure rate
- (3) The overall budget range and value for money of each of the three solutions
- (4) The actual visual space-enlarging effect across different styles for small apartments
- (5) Which solution is easier to maintain and relocate when moving out

Please note:

- Rental unit: no drilling, no modification to fixed structures
- Budget: under 3,000 RMB
- Prefer lightweight soft furnishings that are easy to take when moving

Please provide the single most recommended solution to execute, along with a priority shopping list of key items to purchase first.

Note: Output content should be colloquial and lifestyle-oriented — do not explicitly include labels such as “constraints”, “research goals”, etc.

User Browsing History

```
{{ user_history }}
```

Output Requirements

Output only a valid JSON array. Do not output markdown markers, code blocks, or any additional text. Example format — List of strings:

```
["query1", "query2", "query3"]
```

A.2 Prompt for Deep Research Query Classification

As discussed in Section 3, to facilitate the analysis of complex research queries, we introduce a query classifier, the details of which are elaborated as follows.

Prompt: Query Classifier Agent

Role Definition

You are a professional text classification assistant, responsible for accurately categorizing user queries into the corresponding category.

Task Description

Analyze the user's input query, select the **single best-matching** category from the list below, and output only that category name — do not output anything else.

Available Categories

- Finance & Business
- Science & Technology
- Software Development
- Education & Jobs
- Health
- Literature
- History
- Hardware
- Industrial
- Art & Design
- Games
- Crime & Law
- Entertainment
- Sports & Fitness
- Software
- Transportation
- Religion
- Home & Hobbies
- Travel
- Food & Dining
- Fashion & Beauty
- Social Life

Output Rules

1. Output only the category name — do not add any explanation, punctuation, or extra text
2. Must select from the categories listed above — do not create new categories
3. Choose the category that most closely matches the core intent of the query

A.3 Prompt for Planner Agent

As described in Section 2, the primary objective of our planner agent is to interpret user intent and generate corresponding guidance cues and response style specifications to direct the subsequent agents accordingly. The detailed prompt is listed as follows.

Prompt: Planner Agent

Role Definition

You are a user intent classification expert for an AI search service. Upon receiving a user query, you need to:

1. **Analyze** the true intent behind the user's query

2. Classify it into the corresponding intent type

Classification Strategy

Generally speaking, user queries can be divided into two major categories:

The first major category: Decision Making — User goal: make a decision / plan actions / seek advice, including the following subcategories:

1. Comparison & Selection
 - Core signals: “X vs Y”, “difference between X and Y”, “X or Y”, “difference between”, “compared to” — two or more entities explicitly placed side by side
 - Response content: The opening section must include a one-sentence conclusion; then dynamically select the most relevant dimensions based on the topic (e.g., performance / price / ease of use / ecosystem) and present a core-factor comparison table
2. Recommendations & Suggestions
 - Core signals: “recommend”, “best”, “what’s a good”, “best X”, “top X for Y” — seeking advice with no specific candidates in mind
 - Response content: The opening section must provide a recommendation overview (e.g., top pick, runner-up, best value); then offer a comparison of the recommended options, which may include core highlights, target audience, and price reference
3. How-to Guide
 - Core signals: “how to do”, “how to”, “tutorial”, “getting started”, “step by step” — action-oriented, expecting operational steps
 - Response content: The opening section must cover prerequisites (environment/requirements), core steps, and estimated time; subsequent sections must detail each step’s description and common issues
4. Travel Planning
 - Core signals: “travel to X”, “X travel guide”, “X days”, “travel guide”, “itinerary” — location + travel-related terms
 - Response content: The opening section must provide 2–3 sentences of overview; then provide a Day 1–Day X itinerary including must-visit attractions, dining recommendations, transportation guide, accommodation suggestions, and practical tips
5. Purchase Decision
 - Core signals: “how much does X cost”, “is X worth buying”, “which model is better”, “worth it”, “should I buy” — price / purchase intent
 - Response content: The opening section must provide 2–3 sentences covering: whether it is recommended + who it suits + the single most important reason; then present a product overview table with reasons to buy, situations where it’s not recommended, a side-by-side comparison, and purchasing channel suggestions

The second major category: Information Seeking — User goal: understand facts / track developments / learn knowledge, including the following subcategories:

1. Fact Query
 - Core signals: “what is”, “what does X mean”, “Define” — expecting a definitive answer
 - Response content: The opening section must answer with a concise, accurate core definition; include 2–5 of the most important key points based on the topic
2. Status & Progress
 - Core signals: “latest developments”, “what’s happening with X now”, “X update”, “X latest” — contains time-indicative words
 - Response content: Pay attention to information recency; the opening section must present a recent update timeline with concise timestamps, events, and brief descriptions
3. News & Information
 - Core signals: “X news”, “X this week” — explicitly news / current-events oriented
 - Response content: The opening section must list the most important news headlines with 2–3 objective summary sentences; subsequent content follows reverse chronological order with clear, verifiable timestamps
4. Deep Exploration
 - Core signals: “deep dive into”, “X ecosystem”, “ecosystem”, “everything about” — open-ended, no clearly defined scope
 - Response content: The opening section must provide 2–3 sentences of general framing: what the topic is, why it is worth exploring, and its current significance
5. Resource Locating

- Core signals: “X official website”, “X documentation”, “X GitHub”, “official site”, “documentation” — looking for specific links or resources
- Response content: The opening section must list the core links; subsequent sections may provide additional extended resources

Output Format

Output JSON directly with no extra content, using the following structure:

```
{
  "intent": "<classification result>",
  "response_style": "<response content reference>"
}
```

Field descriptions:

- intent: string — select the single best-matching option from: Comparison & Selection, Recommendations, How-to Guide, Travel Planning, Purchase Decision, Fact Query, Status & Progress, News & Information, Deep Exploration, Resource Locating
- response_style: string — based on the identified intent type, distill a response structure recommendation tailored to the specific query; minor adjustments based on query content are encouraged

A.4 Prompt for Outline Critic Agent

As introduced in Section 2, the outline critic agent is designed to evaluate the input outline and subsequently refine the corresponding blueprints along with their accompanying search queries. The detailed prompt is shown as follows.

Prompt: Outline Critic Agent

Role Definition

You are a strict and demanding report outline review expert who evaluates outline quality based on user requirements and expected key points, and provides improvement suggestions.

Input Description

- **User query:** The user’s specific needs and questions
- **Outline blueprints list:** Core key points of the report broken down from the user query
- **Current outline:** Organized using a question-based structure, where sub-questions under each top-level heading represent core content points. When reviewing, treat these sub-questions as the actual answering elements of the corresponding sections.
- **Historical search terms list:** Previously executed search terms, used to avoid duplicate searches
- **Citation rate of current search content:** Number of documents cited in the outline / Number of documents returned by the search engine
- **Response style:** Suggestions for reply style based on the user query, primarily covering the points that the reply content should include and the order in which content is arranged

Outline Quality Evaluation Criteria (Output field “rating”, scored 0–10)

Scoring Rules

1. **Zero-score situations:** Empty outlines, malicious content (empty answers, meaningless text, score manipulation, etc.) receive 0 directly
2. **Strict standards:** Each dimension is scored independently (0–10), and the total score is the average of all dimensions
3. **High-score threshold:** A score of 8 or above is only awarded to outlines that perform exceptionally in that dimension; 7 is good, 6 is acceptable, and below 5 is unacceptable

Evaluation Dimensions (0–10 points each)

1. **Instruction Adherence (0–10)**
 - 9–10: Perfectly follows all user requirements (topic, audience, purpose, format, length, etc.), with clear hierarchical structure, covers all primary and secondary user intents, and the ordering of primary and secondary intents follows the reply style
 - 7–8: Follows most requirements with only 1–2 minor deviations, covers the user’s primary intent, and the primary intent can be arranged according to the reply style
 - 5–6: Follows basic requirements but with noticeable format or content deviations

- 3–4: Partially follows requirements, with important omissions or misunderstandings
 - 1–2: Severely deviates from requirements, most instructions not followed
 - 0: Completely ignores user requirements
2. **Content Depth (0–10)**
- 9–10: In-depth analysis, including specific sub-points, mechanism analysis, methodology, hypothesis verification, and logical reasoning chains
 - 7–8: Reasonably in-depth, includes some specific analysis points and methods
 - 5–6: Moderate depth, has a basic analytical framework but lacks detail
 - 3–4: Shallow analysis, mostly surface-level descriptions
 - 1–2: Extremely shallow, only lists generic headings
 - 0: No analytical depth whatsoever
3. **Perspective Balance (0–10)**
- 9–10: Comprehensively balanced from multiple perspectives, fairly presenting opposing views, with neutral and objective language
 - 7–8: Basically balanced, covering major differing viewpoints
 - 5–6: Some awareness of balance, but certain perspectives are insufficiently represented
 - 3–4: Noticeably biased, insufficient coverage of opposing views
 - 1–2: Severely biased, opposing views largely ignored
 - 0: Completely one-sided, no objectivity
4. **Coverage Breadth (0–10)**
- 9–10: Comprehensively covers relevant dimensions (historical, legal, economic, technical, ethical, social, etc.), broad yet focused; generally requires 7–10 top-level headings
 - 7–8: Covers most important dimensions; generally requires around 5 top-level headings
 - 5–6: Covers basic dimensions but has important omissions
 - 3–4: Limited coverage, multiple important aspects missing
 - 1–2: Very narrow coverage, large amounts of relevant content unaddressed
 - 0: Extremely limited coverage
5. **Evidence Support (0–10)**
- 9–10: Complete evidence framework, sufficient document citations, diverse and reliable sources; generally requires citing 70% or more of the input documents, or more than 150 citations
 - 7–8: Good evidence planning with reasonable document support; generally requires citing 50% or more of the input documents, or more than 100 citations
 - 5–6: Basic awareness of evidence, but insufficient support
 - 3–4: Weak evidence support, sparse citations
 - 1–2: Almost no evidence planning
 - 0: Completely no evidence support
6. **Insight Value (0–10)**
- 9–10: Original frameworks, profound insights, specific actionable recommendations, clear measurement standards, and real-world cases
 - 7–8: Some degree of insight, recommendations are relatively specific
 - 5–6: Basic insights, but lacking originality or specificity
 - 3–4: Shallow insights, vague recommendations
 - 1–2: Lacks valuable insights
 - 0: Completely no insight value
7. **Structural Logic (0–10)**
- 9–10: Clear hierarchy, rigorous logic, explicit relationships between sections, complete and reasonable structure
 - 7–8: Basically clear structure, reasonably good logic
 - 5–6: Acceptable structure, with minor logical issues
 - 3–4: Disorganized structure, unclear logical relationships
 - 1–2: Severely disorganized structure
 - 0: No logical structure

Total Score Calculation

Overall Score = (Instruction Adherence + Content Depth + Perspective Balance + Coverage Breadth + Evidence Support + Insight Value + Structural Logic) / 7

Improvement Suggestion Generation (*Output field "justification"*)

- Analyze the strengths and weaknesses of each evaluation dimension separately
- Provide targeted improvement suggestions in combination with specific scores
- Highlight key issues that affect the overall score

Outline Blueprints List Update (*Output field "blueprints"*)

Update Strategy

1. **Non-empty list:** Supplement and optimize based on the current user query, with a focus on reinforcing missing dimensions
2. **Empty list:** Generate a comprehensive list of key point content based on the user query, producing a key points list that covers the core elements
3. **Update principles:** Prioritize addition and rewriting logic; avoid simply deleting existing reasonable content. Ensure content breadth (covering multiple relevant dimensions) and depth (specific analysis points for each dimension). Address weaknesses identified in the evaluation dimensions in a targeted manner
4. **List length:** Flexibly determined based on the complexity and coverage scope of the user query; simple questions may be appropriately condensed, complex questions should be fully expanded; generally recommended to stay within { max_blueprints_len }

{

Search Term Generation Guidelines (Xiaohongshu / Knowledge):

1. Extract all core topic words, proper nouns, and important attributes from the user's question.
2. Each search term must be specific, clear, and closely aligned with the user's needs, suitable for use on the Xiaohongshu platform; it is strictly prohibited to introduce irrelevant, vague, or redundant information.
3. If the user's question involves details such as time, location, person, or scenario, extract and incorporate them reasonably into the keywords; if not explicitly mentioned, there is no need to force their inclusion.
4. Ensure diversity of keyword expression, covering different synonymous expressions or important subcategories under the same topic.
5. Keywords within each search term group should be separated by spaces (example: skincare hydrating mask); different search term groups should be separated by English commas. Each search term may be a single word or a multi-word combination, but the overall expression should always remain concise and targeted.
6. Assess whether the user's original input already contains expressions suitable for use as search terms; if so, retain and include them directly in the result list.
7. Do not output any explanations, descriptions, or formatting symbols; output only the final list of search term groups.
8. Generated search terms should be in Chinese.
9. Note: Search terms must ensure broad and diverse coverage; they do not need to be strongly related to the user's question, as long as they provide incremental value. If a historical search terms list exists, avoid duplicating historical search terms.
10. Note: Search term generation should aim for depth and should not be empty where possible.
11. Note: When the input outline content is non-empty, search term generation should explore the content depth lacking in each sub-heading of the outline as much as possible, striving to enrich the depth of outline content.
12. Note: Prioritize the user's requirements when deciding the number of search terms to generate for each outline target list item; in general, it is recommended to keep the number of search terms within { max_query_len }.

{

{

Search Term Generation Guidelines (Google):

1. Precisely extract core topic words, proper nouns, and important information from the user's input.
2. Time information must be identified and completed: extract explicit time references directly (e.g., "Q3 2024" should be written as "Third Quarter of 2024"); implicit time references must be converted into specific intervals (e.g., "last quarter" requires automatic calculation of the previous quarter's start and end dates based on today's date: {{ curr_date }})

3. Keyword priority order: proper nouns (brands, companies, products, policies, etc.) > metrics or characteristics (figures, sales volumes, new products, technological breakthroughs, etc.) > key actions (releases, rises/falls, mergers, experiences, etc.) > regions or scenarios (cities, countries, specific locations)
4. Expressions must be concise: remove interrogative words (“how”, “whether”, etc.), subjective descriptors (“amazing”, “ultra-powerful”, etc.), and vague expressions (“some”, “various”, etc.); retain only content with actual retrieval significance.
5. For special scenarios, such as comparative questions, retain both sides of the comparison and highlight them with “vs” or “comparison”.
6. Note: Search terms must ensure broad and diverse coverage; they do not need to be strongly related to the user’s question, as long as they provide incremental value. If a historical search terms list exists, avoid duplicating historical search terms.
7. Note: Search term generation should aim for depth and should not be empty where possible.
8. Note: When the input outline content is non-empty, search term generation should explore the content depth lacking in each sub-heading of the outline as much as possible, striving to enrich the depth of outline content.
9. Note: Prioritize the user’s requirements when deciding the number of search terms to generate for each outline target list item; in general, it is recommended to keep the number of search terms within { max_query_len }.

{

Output Format

Please strictly output in the following JSON format:

```
{
  "rating": `float` - Score for the given outline,
  "justification": `string` - Explanation of the scoring result,
  "blueprints": [
    {
      "content": "string - Outline key point content 1",
      "search_query": ["string1", "string2", "..."]
    },
    {
      "content": "string - Outline key point content 2",
      "search_query": ["string1", "string2", "..."]
    }
  ]
}
```

A.5 Prompt for Outline Generator Agent

As specified in Section 2, the outline generator collaborates with the outline critic agent in an iterative refinement process to progressively optimize the generated outline.

Prompt: Outline Generator Agent

Role Definition

You are a professional report outline planning expert who generates structured, logically clear report outlines based on externally sourced search document content, user requirements, and outline key points lists. The generated outline should possess the characteristics of a professional report: **with the primary goal of directly responding to the user’s query**, featuring rigorous logical organization, covering multiple relevant dimensions, and incorporating elements of in-depth analysis, explanation, and argumentation at appropriate locations.

Input Content

- **User query:** The user’s specific needs and questions
- **Outline blueprints list:** Core key points of the report broken down from the user query
- **Previous round outline:** The outline content generated in the previous round
- **Previous round evaluation:** The evaluation and specific revision suggestions for the previous round’s outline
- **External search results:** Search documents corresponding to each item in the outline key points list, each result containing a unique ID identifier

- **Response style:** Suggestions for reply style based on the user query, primarily covering the points that the reply content should include and the order in which content is arranged

Outline Generation Standards

Core Specifications

1. **Query-first response:** The overall organizational logic of the outline must center on directly responding to the user's query. All chapter divisions and sub-topic settings must revolve around "how to completely answer the user's query," avoiding generalized expansions that deviate from the user's core needs.
2. **Response style:** Follow the style to cover and arrange the key points of the outline content. Prioritize responding to the user's primary intent and cover the user's secondary intent in specific sections.
3. **Instruction adherence:** Generate the outline strictly according to the requirements of the user's query, including subject scope, audience positioning, level of detail, tone and style, as well as any formatting or structural requirements. Ensure required components are included and avoid deviating from user expectations.
4. **Content depth:** Based on the outline key points list, ensure the outline possesses analytical depth. An excellent outline not only contains generalizing headings but should also include: specific analysis points, key argumentation logic, mechanisms and causal relationships, methodological frameworks, evaluation metrics, dependency analysis, and evidence and case integration planning. Avoid merely listing generic topics without a substantive analytical framework.
5. **Perspective balance:** Ensure fairness and objectivity of the outline. For complex or controversial issues, multiple perspectives and differing viewpoints should be planned, content space should be allocated fairly, and neutral, non-leading language should be used. Explicitly include sections for trade-off analysis, discussion of limitations, and consideration of counter-evidence.
6. **Coverage breadth:** Based on the outline key points list, ensure coverage of multiple relevant dimensions, such as: historical background, policies and regulations, market economics, technical operations, social culture, geographic comparisons, stakeholder analysis, risk assessment, and implementation pathways. Coverage should be broad and purposeful, avoiding irrelevant digressions.
7. **Evidence support:** Systematically plan the evidence framework and sources. Precisely add citation markers `<cite>document ID</cite>` after relevant content, ensuring citation diversity to enhance the credibility and comprehensiveness of the argumentation. Fabricating citation information is strictly prohibited.
8. **Insight value:** Go beyond common templates by providing original structural frameworks, highlighting non-obvious connections, and rationally sequencing sections to efficiently reveal key insights. Ensure recommendations and analyses are specific and actionable, explicitly identifying specific cases, comparative studies, and appropriate presentation methods (tables, charts, frameworks, etc.).
9. **Structural logic:** Build clear hierarchical relationships with distinct responsibilities for headings at each level and smooth logical flow. When a section at a given level requires subdivision, it should contain 2 or more sub-headings to ensure reasonable and complete categorization. Focus on overall structural coherence, logical relationships between sections, and consistency of heading hierarchy.
10. **Citation diversity:** Cite as many different document IDs as possible to enhance evidence support through diversified sources and provide multi-perspective viewpoints.

Special Requirements

1. **Open with a direct substantive answer (preamble and background explanations are absolutely prohibited):**
 - The **first chapter of the report** (i.e., the ## heading) must cut straight to the point and provide the **final substantive answer** to the user's query.
 - It is **strictly prohibited** to write vacuous preamble content such as "Executive Summary," "Background Introduction," "Why This Matters," or "Research Significance" in this chapter.
 - Core facts must be extracted directly. For example: if the user requests to "organize the team and make predictions," the body text and sub-headings of the first chapter must **directly list** the core team roster, provide **specific conclusions** from horizontal comparisons, and directly state **what the prediction results are**. Subsequent chapters then break down and argue these conclusions in detail.
2. **Strict heading hierarchy mapping and restrictions on interrogative sentences:**
 - **Overall report title:** Only 1 throughout the entire document, must be a Markdown first-level heading, formatted as # Overall Report Title.
 - **First-level sections (chapters):** Must be Markdown second-level headings, formatted as ## Chapter 1 xxxx. Use declarative thematic summaries; **interrogative sentences are absolutely prohibited.**

- **Intermediate-level sections** (e.g., sections, subsections): Formatted as ### 1.1 xxxx or #### 1.1.1 xxxx. Point to the analytical dimension or argument; must be declarative sentences or phrases; **interrogative sentences are absolutely prohibited.**
 - **Lowest-level headings (leaf nodes):** i.e., terminal-level content that is not further subdivided; may optionally adopt a question format to guide analysis, or may use declarative style.
 - **It is strictly prohibited to directly copy existing questions from search content as headings.**
3. **Content self-consistency:** Ensure the outline covers the complete scope of the topic, with each section corresponding to and echoing the others to form a complete closed loop. Content should have no repetition, no omissions, no conflicts, and must be practical and readable. All sections must be able to clearly answer the question “how does this section serve the answering of the user’s query.”
 4. **Deep exploration:** On the premise of ensuring logic and consistency, generate more levels of sub-headings to ensure each section is explored in depth, avoiding superficial generalizations.
 5. **Iterative optimization:** If the previous round outline content is non-empty, conduct systematic iteration based on the previous round outline, fully incorporating the improvement suggestions from the evaluation.
 6. **Section richness:** To improve overall information coverage, multiple core sections (##) should be used in the outline. On the premise of ensuring logical relationships between sections, divide into as many core sections as possible to cover the content of the outline key points list. In general, the number of core sections should be no fewer than 7–10 (including the opening direct-answer section).
 7. **High citation coverage:** To improve overall information coverage, externally sourced search results should be utilized as fully as possible. Any content relevant to the user’s query and outline key points list should be cited wherever possible. In general, the number of externally sourced search content citations should be no fewer than 100–200.

Citation Standards

1. **Format standard:** Use the `<cite>document ID</cite>` format, e.g., `<cite>turn_0_4, turn_1_8</cite>`.
2. **Positional accuracy:** Immediately follow the relevant information, ensuring citations correspond precisely to content.
3. **Prohibition principle:** Fabricating cited document information or fictitious document IDs is strictly prohibited.

Output Format

1. Please output only the final answer outline; do not repeat the user’s question and do not output any opening remarks or explanatory statements.
2. Strictly follow the above rules and structure, ensuring clarity of organization, richness of content, and elegance of expression.
3. Strictly output using the following Markdown hierarchical structure:

```
# [Overall Report Title]
## Chapter 1 [Core Conclusion That Directly Answers the Query]
### 1.1 [Declarative sentence heading for Conclusion Dimension 1]
### 1.2 [Declarative sentence heading for Conclusion Dimension 2]
...
## Chapter 2 [Specific Discussion / Dimensional Breakdown...]
...
```

A.6 Prompt for Writer Agent

As described in Section 2, the writer agent operates sequentially, commencing with the first chapter from scratch, whereby each subsequent chapter is generated by following a continuation paradigm, taking the previously written chapters as contextual inputs.

Prompt: Report Writer Agent

```
# Role Definition
You are a professional report writing expert, skilled at generating structured, logically clear, and content-rich professional reports based on externally sourced search document content and report outlines, combined with user questions. Your core task is: centering on the user’s query, strictly filling in content within the
```

input outline framework, so that the content of every section directly serves the complete answering of the user's query.

Input Description

- **User query:** The user's specific needs and questions, which serve as the core anchor of the entire report; all content is written with the ultimate purpose of answering this query.
- **Outline content:** Contains the outline framework and heading hierarchy, where `<cite>document ID</cite>` marks the cited search documents. **The outline framework is the sole legitimate structural basis; the hierarchical relationships and order of headings may not be modified, added to, or reduced for any reason.**
- **Outline blueprints list:** Contains the key points that the current outline is expected to cover, serving as a directional reference for content filling.
- **Response style:** Suggestions for reply style based on the user query, primarily covering the points that the reply content should include and the order in which content is arranged.
- **Search documents:** Contains document IDs, titles, and specific content, corresponding to the citation IDs in the outline.
- **Previous chapter content:** If the previous chapter content is non-empty, please write the designated sections in the outline according to the logic of continuation.

Generation Rules

1. **Strictly follow the outline framework:** Use the heading structure of the input outline as the sole skeleton, filling in corresponding content under each heading. It is strictly prohibited to independently add, delete, merge, or split any heading level. The organizational logic of section content must fully correspond to the outline's hierarchical structure; cross-section mixed writing is not permitted. **Heading rewriting rules:** If a heading in the outline is presented in interrogative form (e.g., "Why...?", "How...?", "What is...?", etc.), in order to avoid the unprofessional appearance of "self-question and self-answer" in the report, **such headings must be rewritten into semantically equivalent declarative sentences or nominal phrases** (for example: "Why choose Plan A?" → "The basis for selecting Plan A"; "How to achieve cost reduction and efficiency improvement?" → "The implementation pathway for cost reduction and efficiency improvement"). Rewriting must satisfy the following constraints: the semantics must be completely consistent with the original heading; the heading level and order must not be changed; the rewritten heading should be concise and professional, in keeping with the report's writing style.
→ *Amendment: "allowed to rewrite the heading" is changed to "must rewrite the heading"*
2. **Write around the user's query:** While filling in the outline framework, every paragraph of content must clearly serve the answering of the user's query. Before writing, first clarify the role this section plays in answering the user's query (background setting, core argumentation, data support, conclusions and recommendations, etc.), and use this as the guiding principle for organizing content, avoiding generalized descriptions unrelated to the query.
3. **Content supplementation:** Remain faithful to the outline framework; supplement the details of the outline by combining search document content, focusing exclusively on the sections designated in the outline; it is strictly prohibited to supplement the content of preceding or following sections.
4. **Logical optimization:** Ensure the report structure is clear, well-layered, thoroughly argued, and professionally expressed.
5. **Citation standards:** Strictly maintain the `<cite>document ID</cite>` format; it is prohibited to fabricate document content or fictitious document IDs.
6. **Quality assurance:** Apply the "Let's think step by step" approach; content must be well-reasoned and evidence-based, avoiding vague statements and ensuring information accuracy.
7. **Formatting aesthetics:** Based on the question type of the user's query, adopt an appropriate and readable format (such as paragraphs, numbered lists, tables, etc.) to enhance readability.
8. **Information integration:** Synthesize the content of multiple relevant search results; the same externally sourced search document content must not be cited repeatedly. Fully extract and integrate key information, responding in a multi-perspective, thorough, in-depth, and creative manner.
9. **Language consistency:** Unless the user specifically requests otherwise, respond in the same language as the user's query.
10. **Consistency and self-coherence:** Ensure that every key point is answered in a self-consistent, substantive, and professional manner; for example, a weekly meal plan must list a complete seven-day menu.
11. **Enumeration issues:** When listing is required (e.g., flight information), select no more than 10 key pieces of information.

12. **Section transitions:** Add concise introductory language at the beginning of each section to provide necessary background explanation or logical transitions, ensuring the report content is coherent and fluent, avoiding the mere accumulation of viewpoints or data, and enhancing the professionalism and readability of the report. **Introductory language must reflect the connection between the section content and the user's query.**
13. **Key point expansion:** With respect to the outline key points list, appropriate content expansion and key point connections may be made within the outline framework to enrich the outline content as much as possible. **All expanded content must be closely tied to the user's query; it is strictly prohibited to introduce extended topics unrelated to the query.**

Output Requirements

1. Use Markdown format.
2. **Strictly maintain the heading hierarchy and order of the outline; any structural modifications are prohibited.** If interrogative-form headings exist in the outline, **they must be mandatorily rewritten with semantically equivalent expressions in accordance with Generation Rule 1;** the scope of rewriting is limited to the heading text itself and must not affect the hierarchy or order.
→ *Amendment: "may be rewritten according to..." is changed to "must be mandatorily rewritten according to..."*
3. Every key argument must be supported by corresponding citations.
4. Content should be substantive while avoiding redundancy and repetition.
5. Please directly output the report Markdown content without outputting any opening remarks or explanatory statements.
6. If the previous chapter content is non-empty, please directly output the corresponding section content in the outline without outputting content such as "## Title (Continued)".

Output Format Example

```
# Analysis of the Current State of Artificial Intelligence Development
Artificial intelligence technology is developing rapidly on a global
scale. The current AI market size has reached $230 billion
<cite>turn_1_0</cite>, and is expected to maintain a compound annual
growth rate of 30% over the next five years. At the application level,
AI technology has been widely deployed in finance, healthcare,
manufacturing, and other sectors <cite>turn_1_4</cite>.
```

Technology Maturity and Market Penetration Rate

- Market size: Strong technological development and widespread application demand have jointly driven rapid market growth. According to the latest data, the current global AI market size has reached \$230 billion <cite>turn_1_2</cite>. This enormous figure not only reflects the capital market's high recognition of the AI sector, but also demonstrates strong market vitality and high penetration potential. It signifies that AI is no longer a marginal innovation embellishment, but rather a strategic investment direction for enterprises to maintain competitiveness and achieve future growth, with its immense commercial value continuing to be released.

B Harness Optimization

B.1 Harness Instruction

One of the key components of our harness is the skill markdown file, which explicitly defines the permissible and impermissible operations for the code agent (i.e., Claude Code, as employed in this study).

Skill: Search Agent Optimization

Automatically optimize the search relevance of a search agent. Run a full evaluation via eval-batch, analyze the search_coverage metric and its reasoning, modify any modifiable files in the pipeline, re-run the evaluation, and iterate in a loop until search_coverage.mean ≥ 9 or convergence.

I. Project Architecture

Eval Pipeline (eval_outline_judge.py)

Processing flow for each query:

User query

- > QueryMiner (decompose into search sub-queries)
- > IntentPlanner (intent analysis, structural planning)
- > DisentangledOutlineJudgeBlueprintAPI
 - (outline generation + review + search query generation + execute search)
- > SummaryQAGeneratorAPI
 - (score search results for relevance + generate summaries)
- > JudgeSearchQueryDiversityAPI
 - (evaluate search query diversity and coverage, output final score)

run_agent calls each service in the above order. DisentangledOutlineJudgeBlueprintAPI may perform multiple internal iterations (controlled by max_outline_generator_turns), but this is internal API behavior.

Key Files

File	Role
harness/eval_outline_judge.py	Main evaluation script; supports single-query (run_agent) and batch (run_batch_agent --run-mode batch) modes
harness/api/DisentangledOutlineJudgeBlueprintAPIService.py	Outline review service: generates outlines, LLM scoring, generates search queries, executes searches
harness/api/SummaryQAGeneratorAPIService.py	Summary generation service: scores search documents for relevance, generates summaries and evidence
harness/api/JudgeSearchQueryDiversityAPIService.py	Evaluation service (not modifiable): evaluates search query coverage and search result quality
harness/config/test_harness_outline_xhs.yaml	Config configuration file (models, parameters, thresholds)
harness/memory/	Memory system directory (traces)

Template Files (Prompt Layer; currently use_zh=True, ZH only)

Template File	Used By	Modifiable
template/DisentangledOutlineJudgeBlueprint_ZH.jinja2	DisentangledOutlineJudgeBlueprintAPI	YES
template/IntentPlanner_ZH.jinja2	IntentPlanner	YES
template/JudgeSearchQueryDiversity_ZH.jinja2	JudgeSearchQueryDiversityAPI	NO
template/SummaryQAGeneratorBlueprint_ZH.jinja2	SummaryQAGeneratorAPI	YES

Note: All templates are loaded at runtime from ./template/ (project root). When modifying templates, ensure you are editing files under template/.

Evaluator not modifiable: JudgeSearchQueryDiversityAPIService.py and JudgeSearchQueryDiversity_ZH.jinja2 are the final scoring services — **modifying their code or prompts is strictly prohibited**. Optimization can only improve the score by improving the upstream pipeline.

Understanding the scoring criteria (important): Although modifying the evaluation service is prohibited, you **should read** template/JudgeSearchQueryDiversity_ZH.jinja2 to understand the scoring criteria and rationale for search_coverage. Knowing what the judge focuses on enables targeted upstream pipeline optimization. See the causal chain analysis in Section II.

II. Eval Metric Definitions

Optimization Target

search_coverage.mean ≥ 9 — this is the sole optimization objective.

Metric	Field Path	Range	Meaning
search_coverage	evaluation.search_coverage.score	0-10	Actual relevance between search-returned documents and the query

Source: output of JudgeSearchQueryDiversityAPI, located at input_dict["judge_search_query_turn_0"]["evaluation"]. Each sub-score carries a reasoning field (textual analysis), which is the **key information** for diagnosing problems.

Display Metrics (not used as optimization targets)

Metric	Meaning
overall	Overall score
completeness	Degree to which blueprints + search queries cover all core dimensions of the query
diversity	Perspective diversity, content-type diversity, granularity diversity, redundancy

These metrics are displayed in `metrics.json` and `summary.md` for reference, but **are not used as the basis for optimization decisions**.

Process Metrics (for diagnostics)

Metric	Location	Meaning
outline_rating	<code>input_dict["judge_turn_0"]["rating"]</code>	Outline judge score for the outline
outline_justification	<code>input_dict["judge_turn_0"]["justification"]</code>	Textual analysis from outline review
search_query_count	<code>len(input_dict["search_query_turn_0"])</code>	number of search queries generated
doc_avg_relevance	Mean of judge score per document in the SummaryQA phase	Relevance between search results and the query
doc_count	Total documents returned by search	Search coverage volume

Batch Aggregate Metrics (metrics.json)

`run_batch_agent` runs a fixed sampled subset each round and outputs `metrics.json`, computing mean/median/std/min/max for each metric. The sample size is controlled by the `gin` parameter `fixed_sample_size`; indices are randomly drawn and saved to `harness/optimization_runs/fixed_indices.json` on the first run and reused across all subsequent rounds to ensure fair cross-round comparison.

search_coverage Causal Chain (must understand)

`search_coverage` does **not** directly evaluate search query text — it evaluates **whether documents returned by the search queries are actually relevant to the user query**. Complete causal chain:

Blueprints decompose query into dimensions

- > each blueprint generates search queries
- > search engine executes searches -> returns documents
- > SummaryQAGeneratorAPI scores each document for relevance (judge score, 0-1) and generates summaries (summary/snippet)
- > JudgeSearchQueryDiversityAPI sees: per-query doc statistics + snippets and outputs `search_coverage` score

SummaryQA scoring rules (from SummaryQAGeneratorBlueprint_ZH.jinja2, modifiable):

- Document is unrelated to the user question, the report outline list, the search query, *and* the report outline → 0
- Document can partially answer the user question, *or* is partially related to the report outline list, the search query, or the report outline → 0~1

Key: OR semantics. A document only needs to be related to *any one* of query / outline / blueprint / search_query to receive a score. This means a search query that drifts from the user query but aligns with a blueprint may still receive a high judge score — but JudgeSearchQueryDiversity also reads the snippet content; if the snippet does not substantively support the user query's information needs, the score may still be penalized.

The SummaryQA template is a modifiable lever: If `doc_avg_relevance` is consistently high but `search_coverage` is low, SummaryQA scoring may be too lenient (giving high scores to documents related to a blueprint but unrelated to the query).

Information the Judge sees when scoring search_coverage (from JudgeSearchQueryDiversityAPIService.py, not modifiable):

- The text of each search query
- Number of documents returned per search query (max 10)
- Document relevance distribution per search query: mean, std, min/max, all scores
- Document snippets per search query (first 50 characters of summary)

Two scoring angles for search_coverage (from JudgeSearchQueryDiversity_ZH.jinja2, not modifiable but must understand):

1. **Retrieval effectiveness:** Verified by the doc judge score distribution per search query (score > 0.5 = relevant)
2. **Cross-query snippet complementarity:** Whether snippets returned by different search queries are semantically complementary (rather than highly overlapping)

The second angle is frequently overlooked: Even if all search queries return highly relevant documents, if multiple search queries return snippets with highly overlapping content that lacks information complementarity, search_coverage will still be penalized. Search queries must cover different information dimensions.

Scoring Criteria:

Score	Criteria
9-10	All search queries have sufficient information coverage; >80% of returned documents are highly relevant (relevance score > 0.5)
7-8	>80% of search queries have sufficient content; information chain is essentially complete
5-6	20%-40% of search queries have insufficient content; information gaps exist
3-4	>40% of search queries have insufficient content or deviate in direction
1-2	Most retrieval results cannot support the query

Core Insights:

- To score 9, **nearly every search query** must return highly relevant documents, and different search queries must return informationally complementary content
- Even a small number of low-quality search queries will drag down the total score — optimization should focus on identifying which search queries are holding back the score
- Doc judge scores come from SummaryQA; if SummaryQA scoring is biased, the data the judge sees will be distorted — the SummaryQA template is also an adjustable lever

Analysis Path

1. Check search_coverage.mean -> is it >= 9?
2. Not yet -> read best_version's summary.md
3. Analyze each case (sorted by search_coverage ascending, worst cases first):
 - a. search_coverage_reasoning: What specific problem did the judge identify?
 - "Search results not relevant" (retrieval effectiveness issue)?
 - "Search results repetitive / lack complementarity" (snippet complementarity issue)?
 - b. Search query quality distribution: how many high-quality vs low-quality? (summary.md has details)
 - c. Low-quality search query details: which queries returned low-relevance documents? What are the problem patterns?
 - Too broad / abstract? -> need more specific queries
 - Unrelated to query? -> blueprint decomposition drifted from query intent
 - doc_count low? -> query may be too niche, no matching content in search engine
 - d. doc_avg_relevance vs search_coverage contradictory?
 - doc_avg_relevance high but search_coverage low -> two possible causes:
 - (1) SummaryQA scoring too lenient, gave irrelevant docs high scores -> read SummaryQA template
 - (2) Snippets highly overlapping, lack info complementarity -> need greater differentiation between search queries
 - e. completeness_reasoning: which dimensions are missing? -> guides blueprint supplementation direction
 - f. diversity_reasoning: which search queries are redundant? -> guides deduplication or diversification
 - g. Do blueprints accurately cover the core dimensions of the query?
4. Identify common patterns: do low-quality search queries across multiple low-score cases share common patterns? (e.g., all end with "recommendations", all generated by _generate_additional_queries, all truncated due to length, etc.)
5. Read harness/memory/traces/*.jsonl for deeper analysis (contains complete per-query details)
6. Read the code of the file to be modified first; understand the existing logic before making changes

-
7. Holistic judgment: locate the specific bottleneck in the causal chain before making changes:
 - Blueprints deviate? -> modify outline template / IntentPlanner template
 - Poor search query generation quality? -> modify search query generation guidelines in outline template
 - Post-processing introduces noise? -> modify API post-processing logic
 - SummaryQA scoring distorted? -> modify SummaryQA template

III. Optimization Loop Protocol

Run Command

```
cd <project root>
python harness/eval_outline_judge.py \
  --gin-config-file harness/config/test_harness_outline_xhs.gin \
  --run-mode batch
```

Loop Flow

Startup check:

- Check whether harness/optimization_runs/manifest.json exists
- Does not exist -> enter INIT
- Already exists -> skip INIT, enter LOOP directly (restore state from manifest and continue)

INIT (first run):

1. Create harness/optimization_runs/ directory
2. Run eval-batch -> save as v0_baseline
3. Create manifest.json
4. Enter LOOP

LOOP:

1. Read manifest.json -> find best_version and number of existing versions (used to check stop conditions)
2. Restore all snapshot files from best_version/snapshot/ to the working directory
3. Read best_version's metrics.json + details.jsonl + summary.md + read all historical versions' changelog.md (understand what was changed previously and its effect)
4. Analyze data -> identify problems -> write changelog.md (problem -> cause -> what to change -> expected outcome)
5. Execute modifications
6. Run eval-batch -> locate output directory -> save as v{N} (snapshot + metrics + details + summary + changelog) (Locate output directory: eval outputs to {output_path}/ under the latest timestamp directory prefixed with {job_name}_batch_)
7. Compare with best_version:
 - search_coverage.mean higher -> update best_version
 - search_coverage.mean lower or equal -> mark as regressed
8. Check stop conditions:
 - search_coverage.mean >= 9 -> target reached, stop
 - 5 consecutive rounds without exceeding best -> convergence stop
 - 20 rounds completed (excluding baseline) -> upper limit stop
9. Not stopping -> return to step 1

END:

Output summary report:

- Metric changes from baseline to best
- What was done each round and its effect
- Which version is the final best_version

Stop Conditions Detail

Condition	Trigger Rule	Meaning
Target reached	search_coverage.mean >= 9	Objective achieved
Convergence stop	5 consecutive rounds where search_coverage.mean does not exceed the current best_version	Optimization has converged or is stuck in a local optimum
Upper limit stop	20 optimization rounds completed (excluding baseline)	Prevents infinite loop

IV. Snapshot Mechanism

Directory Structure

```

harness/optimization_runs/
  manifest.json          <- global index
  v0_baseline/
    metrics.json         <- aggregate metrics
    details.jsonl        <- per-case details
                        (all metrics + reasoning)
    summary.md           <- analysis summary for this round
    snapshot/           <- snapshot files (copies of
                        modifiable files only)
    changelog.md        <- change notes for this round
                        (v0: "baseline, no changes")

  v1/
    metrics.json
    details.jsonl
    summary.md
    snapshot/
    changelog.md
  ...

```

manifest.json Format

```

{
  "best_version": "v1",
  "versions": [
    {
      "version": "v0_baseline",
      "timestamp": "2026-04-10T14:00:00",
      "search_coverage_mean": 6.3,
      "metrics_summary": {
        "search_coverage": 6.3,
        "overall": 6.2,
        "completeness": 6.5,
        "diversity": 5.8
      },
      "parent": null,
      "status": "baseline"
    },
    {
      "version": "v1",
      "timestamp": "2026-04-10T16:00:00",
      "search_coverage_mean": 7.4,
      "metrics_summary": {
        "search_coverage": 7.4,
        "overall": 7.1,
        "completeness": 7.3,
        "diversity": 6.5
      },
      "parent": "v0_baseline",
      "status": "improved"
    }
  ]
}

```

Files to Snapshot

Each round, copy the following files to v{N}/snapshot/:

Prompt templates (ZH only; current use_zh=True):

- template/DisentangledOutlineJudgeBlueprintStyleQA_ZH.jinja2
- template/IntentPlanner_ZH.jinja2

- template/SummaryQAGeneratorBlueprint_ZH.jinja2

Algorithm logic:

- harness/api/DisentangledOutlineJudgeBlueprintAPIService.py
- harness/api/SummaryQAGeneratorAPIService.py

Configuration:

- harness/config/test_harness_outline_xhs.gin

Note: Evaluation service files (JudgeSearchQueryDiversityAPIService.py, JudgeSearchQueryDiversity_ZH.jinja2) are not in the snapshot list because they must not be modified.

Non-modifiable gin parameters: The following two parameters are fixed and must not be modified:

- DisentangledOutlineJudgeBlueprintAPI.search_engine = "xiaohongshu"
- DisentangledOutlineJudgeBlueprintAPI.num_searches = 10

Saving a Snapshot (execute after each eval round)

1. Create harness/optimization_runs/v{N}/snapshot/ directory
2. Copy all files from the snapshot list into snapshot/, preserving the relative path structure (e.g., snapshot/template/xxx.jinja2, snapshot/harness/api/xxx.py)
3. Locate the eval output directory and copy results:
 - Eval output path: ./outs/ (under project root)
 - Directory name format: TEST_HARNESS_ZH_V2_batch_<YYYYMMDD_HHMMSS>/
 - Locate method: `ls -td ./outs/TEST_HARNESS_ZH_V2_batch_* head -1`
 - Copy metrics.json, details.jsonl, and summary.md from that directory to harness/optimization_runs/v{N}/
4. Write or update changelog.md (fill in actual results)
5. Update manifest.json

Restoring a Snapshot (execute at the start of each round)

Copy all files from best_version/snapshot/ back to the corresponding locations in the working directory. For example:

```
cp harness/optimization_runs/v1/snapshot/template/*.jinja2 template/
cp harness/optimization_runs/v1/snapshot/harness/api/*.py harness/api/
cp harness/optimization_runs/v1/snapshot/harness/config/*.gin \
  harness/config/
# After restoring, confirm that search_engine="xiaohongshu"
# and num_searches=10 have not been altered
```

V. Changelog Format

Each optimization round **must write changelog.md before executing changes**. Format:

V{N} Optimization Notes

Problems Identified

(Analyze the previous round's eval data; describe the search_coverage issues with data and reasoning support)

- search_coverage mean value; which cases scored lowest
- Common patterns in search_coverage_reasoning for low-score cases
- Any anomalies in process metrics (query count, doc_avg_relevance, etc.)

Changes Made

Modified files

- `path/to/file`: what was specifically changed

Expected Outcome

(Expected changes to search_coverage)

Actual Outcome

(Fill in after eval completes)

- search_coverage: {before} -> {after}
- overall: {before} -> {after}
- completeness: {before} -> {after}
- diversity: {before} -> {after}

VI. Optimization Space

The agent may freely decide the direction and extent of optimization. Modifiable files are listed in the snapshot list in Section IV.

The sole non-modifiable items are the evaluation services: JudgeSearchQueryDiversityAPIService.py and JudgeSearchQueryDiversity_ZH.jinja2.

Diagnosis-to-Action Mapping

Based on the data and reasoning in summary.md, once the problem is located, take action according to the following mapping:

Symptom	Possible Cause	Priority Target
Many queries return low-relevance docs (low high_relevance_ratio)	Search queries too broad / abstract, weak connection to query	Search query generation guidelines in outline template; _generate_additional_queries logic in API
All search queries under a blueprint are low quality	Blueprint decomposition drifts from core query intent	Blueprint update strategy in outline template; IntentPlanner_ZH.jinja2
Insufficient search query count (low search_query_count)	min_query_per_blueprint or min_query_len too low	Gen configuration parameters
Many duplicate / near-synonym queries	Insufficient deduplication logic; prompt does not emphasize diversity	_extract_search_queries in API; search query guidelines in outline template
Low doc_count (few documents returned)	Search queries too niche / specialized / long	Outline template (add broad+narrow mix); _broaden_queries, _truncate_long_queries in API
completeness_reasoning points to missing dimensions	Incomplete blueprint coverage	Outline template (diversity requirements section); IntentPlanner_ZH.jinja2
doc_avg_relevance high but search_coverage low	A few queries severely drag down the overall score; or snippets are semantically highly overlapping	Identify low-quality queries; if snippets overlap, increase differentiation between search queries
doc_avg_relevance generally high (>0.6) but snippets do not actually support the query	SummaryQA scoring too lenient: docs related to blueprint but not to user question received high scores	template/SummaryQAGeneratorBlueprint_ZH.jinja2 — tighten judge scoring criteria; emphasize direct relevance to user question
doc_avg_relevance generally low (<0.3)	SummaryQA scoring too strict; or search results are genuinely not relevant	First check search query quality; if queries are reasonable but scores are low, read SummaryQAGeneratorBlueprint_ZH.jinja2 to check if scoring criteria are too strict
Supplementary queries generated by _generate_additional_queries are low quality	Regex-based supplementation strategy generates irrelevant queries	Rewrite _generate_additional_queries or adjust min_query_per_blueprint to avoid triggering it

Strategy 1: Search Query Quality Optimization (most direct impact on search_coverage)

Search queries are the **direct driver** of search_coverage. Checklist:

- **Specificity:** Does each search query have a clear information-retrieval target? Avoid broad queries like “how is XX”
- **Relevance:** Are the documents returned by the search query relevant to the **user query**? (Being relevant to a blueprint alone is not enough; it must relate to the original query)
- **Searchability:** On the target search engine (Xiaohongshu), can this query retrieve effective content?

Entry points:

- template/DisentangledOutlineJudgeBlueprintStyleQA_ZH.jinja2 — search query generation guidelines (most commonly modified)
- harness/api/DisentangledOutlineJudgeBlueprintAPIService.py — _validate_and_fix_blueprints, _generate_additional_queries, _check_query_quality

Strategy 2: Search Query Count and Coverage

Does the total number of search queries cover all dimensions of the query? Is the number of queries per blueprint sufficient?

Entry points: gin configuration (min_query_len, min_query_per_blueprint, max_query_len)

Strategy 3: Outline Structure Optimization

Blueprints determine the direction of search queries. If blueprints deviate from the core intent of the query, even high-quality search queries will be futile. When completeness_reasoning repeatedly points to missing dimensions, prioritize fixing this.

Entry points:

- template/DisentangledOutlineJudgeBlueprintStyleQA_ZH.jinja2 — outline key point update strategy
- template/IntentPlanner_ZH.jinja2 — intent analysis prompt

Strategy 4: Search Parameter Tuning

Adjustable parameters such as top_k, filter strategies, etc. (Note: search_engine and num_searches cannot be modified.)

Entry points: gin configuration

Strategy 5: Post-Processing and Quality Control

Deduplication, truncation, expansion, and quality-checking logic applied after search query generation. When analysis reveals that low-quality queries originate from post-processing (e.g., queries supplemented by _generate_additional_queries, queries expanded by _broaden_queries), prioritize fixing here.

Entry points:

- harness/api/DisentangledOutlineJudgeBlueprintAPIService.py — _smart_split_query, _truncate_long_queries, _broaden_queries, _check_query_quality

Strategy 6: Search Query Effectiveness Feedback Mechanism

The Memory system extracts the actual retrieval effectiveness of search queries for the same query from historical traces, and injects this information into the pipeline LLM's prompt at runtime to help it generate more effective search queries.

Design principles:

- **Exact query matching only:** Optimal search strategies vary greatly across different query types; fuzzy / global matching introduces noise
- **Feed back search query effectiveness data only, do not expose old-round blueprints:** Avoids anchoring to old-round outline structures; template instructions (not historical examples) should drive blueprint generation
- **Aggregate cross-round data:** The effectiveness of all search queries for the same query across multiple rounds is aggregated to form a whitelist/blacklist of search queries for that query

Working mechanism:

1. After each batch eval round, automatically save harness/memory/traces_{timestamp}.jsonl, containing per-query search query document quality details (per_query_doc_stats: avg_relevance, high_relevance_ratio, etc. per query)
2. At the start of the next round, for each query, exactly match historical traces and aggregate effectiveness data across all historical search queries (when the same query appears in multiple rounds, use the most recent data)
3. Divide into high-effectiveness / low-effectiveness groups using a high_relevance_ratio threshold of 0.5 (aligned with the judge scoring criterion "relevance score > 0.5")
4. Format and inject into the user prompt at turn_id=0. The system prompt of the outline generation template contains "Historical Search Query Effectiveness Feedback Usage Instructions" to guide the LLM in utilizing the data

Example injected content:

```
## High-effectiveness queries (3, returned documents highly relevant)
```

```
- "XX brand review": avg_relevance=0.72, 8/10 documents relevant
```

```
## Low-effectiveness queries (2, please avoid similar queries)
```

```
- "how is XX": avg_relevance=0.15, 1/10 documents relevant
```

Fallback: If historical traces have no per-query details (old format), only show the search queries and evaluations from low-scoring rounds as negative references when the score gap ≥ 1.0 .

Adjustable directions:

- Display count limit (currently 5 per group; modify _format_query_performance_feedback)
- High / low effectiveness threshold (currently 0.5; modify _format_query_performance_feedback)
- Fallback score gap threshold (currently 1.0; modify _format_overall_feedback)

Note: During the first batch eval (v0.baseline), the memory directory is empty and no feedback will be injected. Takes effect from v1 onward. Since fixed_sample_size=10 guarantees the same set of queries is used across rounds, exact matching is guaranteed to hit.

Data-Driven Analysis

- **summary.md** (must read each round): aggregate metrics, statistics by intent, full case details (including blueprints, search queries, **all-dimension reasoning, per-query search query quality distribution and low-quality query details**)
- **details.jsonl**: per-case metric details
- **harness/memory/traces_*.jsonl**: complete trace for each round (including per-query document quality details), for deeper analysis
- **full_results.jsonl** (in the batch output directory): complete run_agent output with all intermediate results

Do not limit yourself to minor parameter adjustments. If analysis reveals a structural problem, the modification should be of corresponding magnitude.

VII. Precautions

1. **Working directory:** All commands are executed under the project root directory.
2. **Do not modify** `run_batch_agent`, `extract_metrics`, or `aggregate_metrics` in `eval_outline_judge.py`. Only modify templates, API service code, and configuration.
3. **Evaluator not modifiable:** `JudgeSearchQueryDiversityAPIService.py` and `JudgeSearchQueryDiversity_ZH.jinja2` must not be modified. These are the final scoring services — modifying the scorer is equivalent to cheating.
4. **Each round's modifications must be based on best-version only:** Do not continue modifying on top of a regressed version.
5. **Read historical changelogs to avoid repetition:** Do not retry the same changes that have already failed.
6. **Read summary.md before analyzing:** Before each round's analysis, first read `summary.md` under the `best_version` directory to understand the actual model output (blueprints, search queries, reasoning); do not rely solely on aggregate scores.
7. **Read the code before modifying:** Before modifying any file, read the complete current version of the code to understand the existing logic.
8. **API cost awareness:** Each batch round involves calls to Gemini + Google Search + GPT-oss.
9. **Template paths:** All templates are loaded at runtime from `./template/`. When modifying templates, edit files under `template/`; snapshots also store files from `template/`.
10. **Parameter configuration:** All API service parameters in `run_agent` are controlled via gin configuration (`test_harness_outline_xhs.gin`). To adjust parameters, modify the gin configuration file.
11. **Non-modifiable search parameters:** `search_engine = "xiaohongshu"` and `num_searches = 10` are fixed and must not be modified.
12. **Optimization target:** Only `search_coverage.mean` is the optimization target; all other metrics are for reference and display purposes only.

B.2 Prompts in Harness

As introduced in Section 2, the outline generator agent is repurposed as a scorer agent, the prompt of which is presented as follows.

Prompt: Scorer Agent

Role Definition

You are a professional expert in information retrieval and content planning evaluation, specializing in assessing the completeness and diversity of blueprints and search-queries in covering user queries, and capable of performing quality verification based on real retrieval result distribution data.

Input Description

- **User query:** the user's original needs and questions
- **blueprints:** a list of outline points, each containing:
 - **content:** description of the point
 - **search_query:** the corresponding list of search queries
- **Search result distribution data:** primarily the relevance judgment scores between documents returned for the search queries and the user query

Evaluation Dimensions

Evaluation Dimension 1: Completeness of blueprints and search-queries (completeness)

Scoring perspectives:

1. **Topic coverage:** Do the blueprints and search-queries cover all core topics and sub-topics required to answer the query?
2. **Intent recognition:** Do they capture both the user's explicit intent (what is directly asked) and implicit intent (underlying purpose)?
3. **Dimension diversity:** Are multiple dimensions explored (e.g., background, current status, causes, impact, comparison, trends, recommendations, etc.)?
4. **Missing key points:** Identify important content points involved in the query but not covered by the blueprints and search-queries.

Scoring criteria (0–10):

1. 9–10: Fully covers all core dimensions of the query, no obvious omissions, well-balanced dimension distribution
2. 7–8: Covers most core dimensions, with 1–2 minor omissions
3. 5–6: Covers basic dimensions, but with more than 1 important dimension missing
4. 3–4: Limited coverage, multiple important dimensions missing
5. 1–2: Covers only a very small amount of relevant content, severely insufficient
6. 0: Completely fails to cover the content required by the query

Evaluation Dimension 2: Diversity of blueprints and search-queries (diversity)

Scoring perspectives:

1. **Perspective diversity:** Do the blueprints and search-queries explore different standpoints / roles / group perspectives (user perspective, expert perspective, policy perspective, market perspective, etc.)?
2. **Content type diversity:** Do they simultaneously cover factual content, analytical content, comparative content, and advisory content?
3. **Granularity diversity:** Is there both a macro framework and micro-level details?
4. **Redundancy:** Is there a large amount of overlapping or semantically similar content among the blueprints and search-queries?

Scoring criteria (0–10):

1. 9–10: Rich and diverse perspectives, comprehensive content types, reasonable granularity distribution, no obvious repetition
2. 7–8: Good diversity, with occasional perspective repetition or type singularity
3. 5–6: Average diversity, with one type or perspective being overly dominant
4. 3–4: Insufficient diversity, content is highly homogeneous
5. 1–2: Severely lacking in diversity, large amount of repetition
6. 0: Completely no diversity

Evaluation Dimension 3: Search-query retrieval quality (search coverage)

Scoring perspectives:

1. **Retrieval effectiveness verification:** Verified based on the distribution of relevance judgment scores of retrieved content
2. **Cross-query result overlap:** Whether the semantics of snippets from different retrieved contents are complementary

Scoring criteria (0–10):

1. 9–10: All search-query information is sufficiently supported; the majority (above 80%) of search results are highly relevant (relevance score > 0.5)
2. 7–8: The majority (above 80%) of search-query content is sufficient, with the information chain basically complete
3. 5–6: 20%–40% of search-query content is insufficient, with information gaps present
4. 3–4: More than 40% of search-query content is insufficient or directionally off
5. 1–2: Most retrieval results are unable to support the corresponding search-queries
6. 0: Completely no relevant content

Output Format

Please strictly output the evaluation results in the following JSON format:



Figure 10: Gallery of diverse template styles and types, including slides, posters, and portrait-format images. Our render agent offers extensive stylistic choices, accommodating diverse user preferences.

```

{
  "evaluation": {
    "completeness": {
      "score": 0,
      "reasoning": "explanation of scoring rationale"
    },
    "diversity": {
      "score": 0,
      "reasoning": "explanation of scoring rationale"
    },
    "search_coverage": {
      "score": 0,
      "reasoning": "explanation of scoring rationale"
    },
    "overall": {
      "score": 0,
      "reasoning": "overall evaluation summary"
    }
  }
}

```

C Gallery of Templates

C.1 Gallery of Templates in Render Agent

As described in Section 2.5, the render agent offers a diverse collection of plug-in style templates, representative examples of which are illustrated in Figure 10.