

Writing As Reasoning: Interleaving Drafting and Deepening for Open-Ended Deep Research

Anonymous ACL submission

Abstract

Generating deep research reports requires large-scale information acquisition and the synthesis of insight-driven analyses, posing a major challenge for current language models. Existing methods largely follow a *plan-then-write* paradigm, whose performance strongly depends on the quality of the initial outline. We propose a **Writing As Reasoning Policy (WARP)** framework, which enables models to dynamically revise outlines during report writing. In this policy, the agent alternates between *Evidence-Based Drafting* and *Reasoning-Driven Deepening* phases, jointly supporting information acquisition, knowledge refinement, and outline updating. We further introduce a **Multi-Stage Agentic Training**—including cold start, atomic skill RL, and holistic pipeline RL—that enables small models to operate the WARP effectively. Experiments on DeepResearch Bench, DeepConsult, and DeepResearch Gym show that our approach allows small models to surpass leading closed-source systems, particularly with substantial improvements in *Insight* metric.

1 Introduction

Open-ended deep research requires artificial agents to navigate vast information landscapes and synthesize their findings into coherent, insightful reports (OpenAI, 2025; Google, 2025; x.AI, 2025; Perplexity, 2025; Kimi, 2025; ByteDance, 2025). In the context of such complex inquiry, *writing* is far more than the mere transcription of retrieved data. Instead, it mirrors the *knowledge-transforming* process in cognitive psychology (Scardamalia and Bereiter, 1987): since the information landscape is initially opaque, researchers rarely execute a rigid, end-to-end plan derived from pre-existing thoughts. Rather, the act of writing serves as a reasoning mechanism itself, helping to uncover what is not yet known. Researchers frequently identify gaps, contradictions,

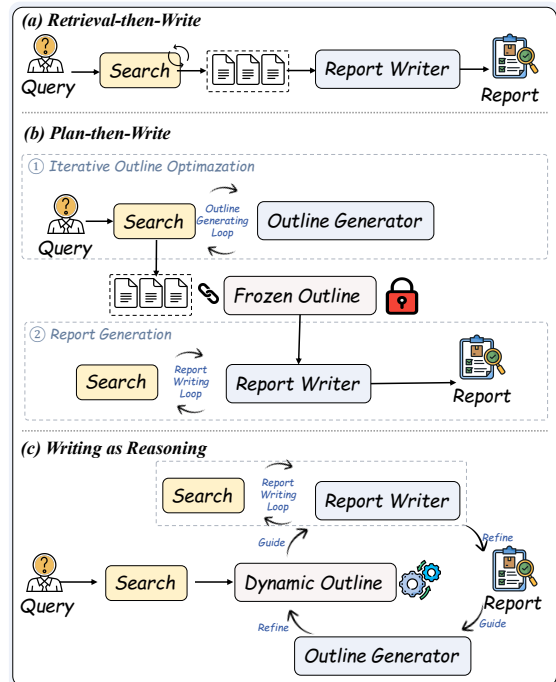


Figure 1: Comparison of different writing paradigms.

and new directions only while drafting, indicating that effective synthesis relies on a tight coupling between planning and writing.

However, existing approaches have struggled to replicate this dynamic. Early methods relied on a *retrieval-then-write* paradigm (Fig. 1a) (Hu et al., 2025), where agents generated content sequentially. While flexible, this unstructured approach often deteriorates into incoherence over long horizons. To address this, recent frameworks (Fig. 1b) (Wang et al., 2024, 2025; Yan et al., 2025) like WebWeaver (Li et al., 2025d) adopted a *plan-then-write* paradigm. By freezing a comprehensive outline before writing, these systems ensure structural stability. However, this approach relies on the *assumption of initial information completeness*—a premise often fallacious in open-ended research. By treating the downstream writer merely as an executor

of a static blueprint, this dichotomy prevents the agent from capturing *emergent insights*—nuanced connections that surface only when articulating specific arguments. Consequently, such baselines encounter an *insight ceiling*, producing reports that are structurally sound but intellectually shallow.

To bridge this gap, we introduce **WARP** (Writing As Reasoning Policy), which transcends the limitations of static planning by modeling deep research as an iterative refinement loop. Instead of adhering to a fixed outline, the agent alternates between two macro-states: *Evidence-based Drafting* and *Reasoning-driven Deepening*.

Crucially, WARP operates as a dynamic policy rather than a static rule-based system. In the *Reasoning-driven Deepening* state, the agent autonomously decides whether to terminate or continue deepening, by evaluating the semantic density and logical coherence of the current draft. If deepening, it decomposes broad sections into granular inquiries and updates the outline based on the actual writing feedback, effectively mirroring the human knowledge-transforming process.

Considering the expansive action space and extended decision horizons inherent to WARP, we design a systematic training pipeline to ensure stable convergence under reasonable resource constraints. Specifically, we introduce a *trajectory pruning* strategy to rigorously filter data for higher quality supervision, and adopt a multi-stage reinforcement learning curriculum that sequentially fine-tunes local atomic actions before optimizing end-to-end performance. This approach instills a robust policy that autonomously balances research depth against computational cost, triggering recursive optimization cycles only when they yield significant informational gain.

Extensive experiments on DeepResearch Bench, Deep Consult, and DeepResearch Gym validate the efficacy of our framework, demonstrating marked improvements in report quality, particularly in the *Insight* metric. Notably, our 8B agent outperforms many closed-source deep research systems, demonstrating that specialized WARP inference diagram can enable small-scale models to rival proprietary large-scale models such as Gemini-2.5-Pro.

2 Method

In this section, we formally present WARP (Writing As Reasoning Policy). Unlike traditional pipelines (Wang et al., 2024; Shi et al., 2025) that

treat planning and writing as temporally segregated stages, we formulate deep research as a unified sequential decision-making process. The core intuition of WARP is that the intermediate draft serves not merely as an output buffer, but as a *dynamic reasoning context* that should actively shape future planning. By maintaining the draft as an evolving state variable, our framework enables an *interleaved* inference loop: the agent alternates between **Evidence-Based Drafting** (filling the report based on the current plan) and **Reasoning-Driven Deepening** (expanding the plan based on the generated report). This mechanism allows the system to pivot from a broad initial skeleton to fine-grained exploration, ensuring the final report achieves both comprehensive breadth and analytical depth.

2.1 Problem Formulation

In specific, we formulate open-ended deep research as an iterative hierarchical decision-making process: At any interaction loop i , the agent observes a global state $S_i = (Q, O_i, D_i, C_i)$, comprising the user query Q , a dynamic outline O_i , the current draft D_i , and the context C_i retrieved at the current loop i . At the j -th step $t_{i,j}$ in a loop i , the agent executes an action $A_{i,j}$ selected from a defined action space: {INITIALIZE, SEARCH, WRITE, EXPAND, TERMINATE}.

This formulation unifies planning and writing: outline adjustments ($O_i \rightarrow O_{i+1}$) and content generation ($D_i \rightarrow D_{i+1}$) are treated as equivalent state transitions driven by the policy.

2.2 The WARP Inference Diagram

WARP begins with a coarse-to-fine initialization strategy designed to establish a comprehensive research scope before diving into details. Starting from the initial state S_0 , the agent analyzes the query Q to generate broad search queries q_0 . Upon retrieving the background context C_0 , it synthesizes an initial Level-1 outline O_0 :

$$O_0 \leftarrow \text{INITIALIZE}(Q, C_0). \quad (1)$$

In contrast to static planners (Li et al., 2025d) that attempt to generate a fully detailed hierarchy, our O_0 is intentionally sparse, consisting only of high-level section titles and brief writing intents. This design mitigates the risk of ungrounded.

To maximize the final report quality, the agent operates under a unified policy π_θ that orchestrates the research trajectory. Crucially, this workflow is

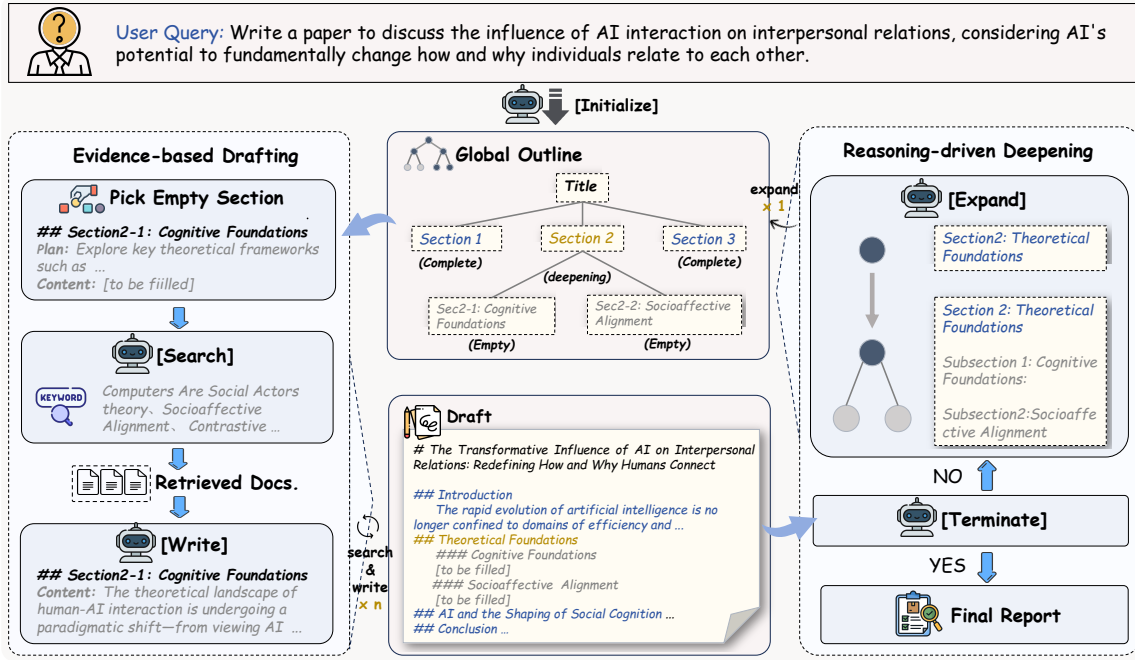


Figure 2: The WARP framework. The agent interleaves **Evidence-Based Drafting** (writing content) and **Reasoning-Driven Deepening** (updating the dynamic outline). This loop allows the agent to discover and bridge logical gaps that emerge only during the writing process.

not linear but iterative, alternating between Drafting and Deepening.

Evidence-Based Drafting Given a tentative outline O_i , the agent executes a *retrieve-then-write* strategy to convert the structural plan into substantiated content. Unlike independent parallel generation, which often leads to fragmentation or redundancy, we enforce contextual consistency by conditioning retrieval queries on the accumulating narrative. For a specific section k , the agent first formulates query $q_{i,k}$ based on user query Q , the section's intent O_i^k , and the draft context D_i :

$$q_{i,k} \leftarrow \text{SEARCH}(Q, O_i^k, D_i). \quad (2)$$

Then, the retrieval tools will acquire new content C_i^k based on the query $q_{i,k}$. This ensures that new information strictly extends the logical flow of previous sections. The agent then synthesizes the section content c_k by grounding the text in retrieved evidence to guarantee faithfulness:

$$D_i^k \leftarrow D_i^{k-1} \oplus \text{WRITE}(Q, O_i^k, D_i^{k-1}, C_i^k). \quad (3)$$

The objective is to achieve information integration—synthesizing disparate sources into a coherent argument—rather than mere aggregation. This phase focuses on writing, iteratively populating the outline to produce a new draft D_{i+1} that serves as the foundation for deeper reasoning.

Reasoning-Driven Deepening Initial outlines are inevitably constrained by the model's pre-retrieval knowledge, which often creates an *insight ceiling*: the structure may cover the breadth of the topic but fail to capture its nuanced depth. To break this ceiling, the policy π_θ periodically shifts from local drafting to global planning, treating the newly generated draft D_{i+1} as a fresh observation for reasoning and diagnosis.

Since D_{i+1} provides a concrete reasoning context, the agent can detect logical gaps or superficial arguments that were invisible during initial planning. If section k^* lacks depth, the agent generates a **local sub-sections** to decompose the topic, updating O_i and triggering a targeted drafting cycle:

$$O_{i+1} \leftarrow O_i \oplus \text{EXPAND}\{k^*\}(Q, O_i, D_{i+1}). \quad (4)$$

The process concludes only when the agent verifies that the logical chain is complete and the content depth aligns with the query's complexity:

$$\text{End} \leftarrow \text{TERMINATE}(Q, O_i, D_{i+1}). \quad (5)$$

2.3 Multi-Stage Agentic Training

While large-scale models have demonstrated strong capabilities within our WARP framework (see §3.3.1), training small-scale models (such as 8B) for open-ended research is non-trivial. The

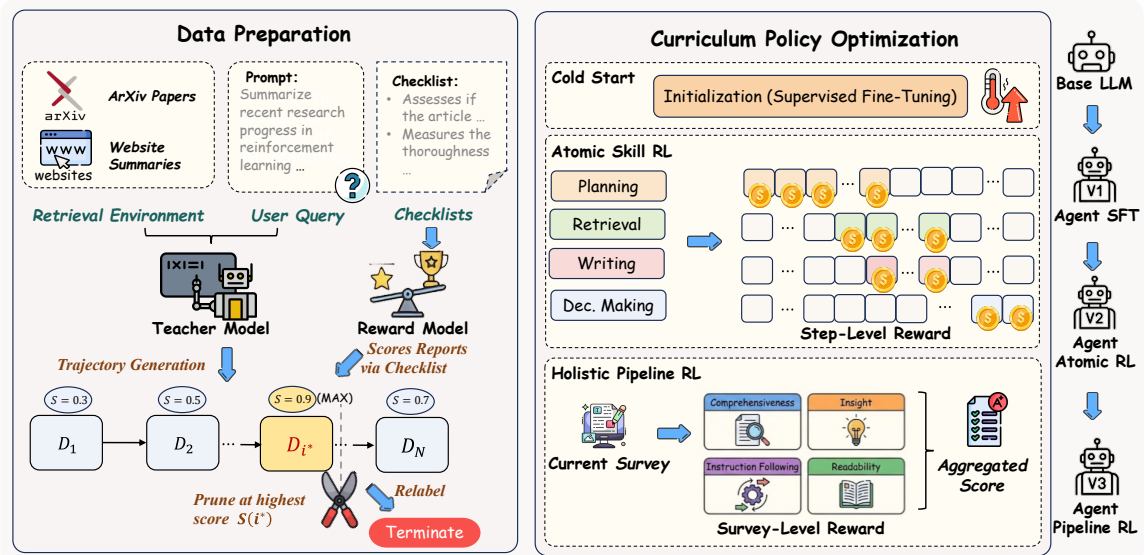


Figure 3: Overview of our multi-stage agentic training process.

Table 1: Reward definition for **Atomic Skill RL**. For each metric, we indicate whether it requires references (**Ref.**) or relies on an **LLM-Judge**. "Dec. Making" indicates the decision-making ability.

Atomic Ability	Metric Name	Need Ref.?	LLM-Judge?	Description & Reward Objective
Planning	Basic Properties	✗	✗	Section number and language consistency.
	Holistic Quality	✓	✓	Quality metrics such as guidance, logic, clarity, etc.
	Faithfulness	✗	✓	Whether the content in the plan is real and reliable.
Retrieval	Relevance Recall	✓	✗	Overlap score between retrieved docs and golden ones.
	Basic Properties	✗	✗	Content length, citation number, and language consistency.
Writing	Holistic Quality	✓	✓	Quality metrics such as relevance, coverage, depth, etc.
	Faithfulness	✗	✓	Penalizes unsupported claims.
	Citation Precision	✓	✗	Rewards citations that overlaps with golden ones.
Dec. Making	Accuracy	✓	✗	Whether terminate at the suitable time.

challenges are twofold: (1) *Ambiguous Termination*: even teacher models struggle to determine the optimal stopping point for research; (2) *Sparse Rewards*: long horizons make reward assignment difficult. To address these, we first introduce a trajectory pruning strategy in §2.3.1. Then, we propose a curriculum learning pipeline in §2.3.2.

2.3.1 Data Preparation

A critical challenge in training is the scarcity of expert trajectories that exhibit efficient decision-making. Teacher models often either expand indefinitely or terminate arbitrarily, giving rise to what we term the **optimal stopping problem**. To solve this, we introduce **trajectory pruning** strategy. Instead of cloning the teacher’s termination behavior, we force the teacher to "over-expand" recursively. This generates a sequence of drafts with varying granularity $\{D_1, \dots, D_N\}$. We then retroactively identify the optimal point i^* where the report draft

D_{i^*} has the highest score. We prune the trajectory at i^* , relabeling the action to **TERMINATE**. This provides a supervision signal for *information saturation*, teaching the agent to stop based on report quality rather than arbitrary imitation. Details on our query construction and retrieval environment are provided in App. A.2 and App. A.3.

2.3.2 Curriculum Policy Optimization

Our optimization includes three stages, as shown in Figure 3: (1) **SFT for Cold Start**: Establishes basic instruction following and format adherence. (2) **Atomic Skill RL**: Uses teacher trajectories as anchors to master local execution and stabilize exploration. (3) **Holistic Pipeline RL**: Optimizes global report quality, enabling the agent to refine its strategy beyond the teacher’s limitations.

Atomic Skill RL To tackle the reward assignment problem, we first decompose the global objective to atomic abilities: **planning** (*Initialize, Ex-*

248 *pand*), **retrieval** (*Search*), **writing** (*Write*), and
249 **decision-making** (*Terminate*). Then, we design
250 different reward functions for them, which com-
251 bine execution results (e.g., basic properties, holis-
252 tic quality, and faithfulness) with *reference align-*
253 *ment*, as shown in Table 1. This stage ensures the
254 agent masters the "how"—producing valid plans,
255 precise searches, and coherent paragraphs—before
256 attempting to optimize global strategy.

257 **Holistic Pipeline RL** Local correctness (e.g., a
258 valid paragraph) does not guarantee global coher-
259 ence. Thus, the final stage shifts to end-to-end
260 optimization to evaluate the final report quality,
261 such as *Comprehensiveness* and *Insight*. Crucially,
262 this stage empowers the agent to deviate from the
263 teacher’s path. By propagating the holistic report
264 score backward, the agent learns to trigger the *deep-*
265 *ening* only when it yields significant informational
266 gain. This effectively refines the quality-efficiency
267 frontier, suppressing redundant expansions that the
268 teacher might have made.

269 3 Experiments

270 3.1 Settings

271 **Implementation Details.** We implement WARP
272 using *MiniCPM4.1-8B* (Team et al., 2025) as the
273 backbone for our deep research agent system.
274 Training follows the curriculum described in §2.3,
275 progressively scaling from atomic skills to holistic
276 reporting. During training and inference phase, we
277 cap the report structure at three levels and limit
278 the number of deepening steps to 12 to ensure effi-
279 ciency. Detailed hyperparameters and data statis-
280 tics per stage are provided in App. B.1.

281 **Benchmarks and Metrics.** To ensure compre-
282 hensive evaluation, we test on three diverse bench-
283 marks: (1) **DeepResearch Bench** (Du et al., 2025)
284 (100 PhD-level scientific tasks); (2) **DeepCon-**
285 **sult** (Dee, 2025) (102 business and financial analy-
286 sis queries); and (3) **DeepResearchGym** (Coelho
287 et al., 2025) (100 general-purpose information-
288 seeking tasks). We adhere to the standard eval-
289 uation protocols of each benchmark, employing
290 *Gemini-2.5-Pro*, *o3-mini*, and *GPT-4.1-mini* respec-
291 tively as impartial judges.

292 **Baselines.** We compare WARP against three dis-
293 tinct categories of latest systems: (1) **Proprietary**
294 **Systems:** Leading commercial deep research sys-
295 tems including OpenAI (OpenAI, 2025), Gem-
296 ini (Google, 2025), and Claude (claude, 2025), and

Doubao (ByteDance, 2025). (2) **Prompt-Based**
297 **Frameworks:** WebWeaver (Li et al., 2025d), En-
298 terprise DR (Prabhakar et al., 2025), and RhinoIn-
299 sigh (Lei et al., 2025). (3) **Trained Open Models:**
300 Recent open-source research agents including Web-
301 Shaper (Tao et al., 2025), WebThinker (Li et al.,
302 2025c) and DR Tulu (Shao et al., 2025). 303

304 3.2 Main Results

305 Our results on three benchmarks are summarized
306 in Table 2 and Figure 4.

307 **(1) Our WARP framework has strong per-**
308 **formance on *Insight* and *Comprehensiveness*.**
309 Across these benchmarks, our method achieves
310 nearly the best performance in both *Insight* and
311 *Comprehensiveness* metrics despite using the small-
312 est model. Specifically, on the DeepResearch
313 Bench, it achieves an *Insight* score of 52.64 and
314 a *Comprehensiveness* score of 50.54, surpassing
315 Gemini-2.5-Pro-deepresearch (49.45 and 49.51, re-
316 spectively). On the DeepResearch Gym, it gets
317 the highest 100.0 score in the *Depth*, *Breadth*,
318 and *Insightfulness* metrics. These gains stem di-
319 rectly from our *reasoning-driven deepening*. On
320 one hand, the agent continuously **extracts insights**
321 **from condensed intermediate drafts**, enabling
322 deeper reasoning and synthesis. On the other hand,
323 by revisiting intermediate outputs, it can **identify**
324 **missing topics and globally assess which sections**
325 **require further expansion**, resulting in broader
326 and more balanced coverage.

327 **(2) The *Multi-Stage Agentic Training* brings sta-**
328 **ble and comprehensive improvement.** The per-
329 formance of WARP-8B steadily improves from
330 SFT to Atomic RL and finally to Pipeline RL across
331 all metrics on these benchmarks. On DeepResearch
332 Bench, the metric *comprehensiveness* rises from
333 46.24 to 50.54, *Insight* from 48.10 to 52.64, and
334 *Readability* from 41.79 to 44.17. On DeepCon-
335 sult, average score grows from 6.04 to 6.60, win
336 rate increase from 54.17% to 57.60%, and loss rate
337 drop from 35.54% to 28.68%. These consistent
338 gains demonstrate that each stage of the curricu-
339 lum contributes to mastering the full deep research
340 workflow, yielding a more stable and capable agent.

341 **(3) Small-scale agent systems can rival large-**
342 **scale ones.** Averaged across benchmarks, our
343 deep research system demonstrates excellent per-
344 formance. Our WARP-8B (Pipeline RL) achieves
345 an *Overall* score of 50.11 on DeepResearch Bench,

Table 2: Performance of agent systems on DeepResearch Bench in terms of comprehensiveness (Comp.), insight, instruction-following (Inst.), readability (Read.) and DeepConsult (Avg., Win, Tie, Lose).

Agent systems	DeepResearch Bench					DeepConsult			
	Overall	Comp.	Insight	Inst.	Read.	Avg.	Win	Tie	Lose
<i>Proprietary Deep Research Systems</i>									
Doubao-research	44.34	44.84	40.56	47.95	44.69	5.42	29.95	40.35	29.70
Claude-research	45.00	45.34	42.79	47.58	44.66	4.60	25.00	38.89	36.11
OpenAI-deepresearch	46.45	46.46	43.73	49.39	47.22	5.00	0.00	100.00	0.00
Gemini-2.5-Pro-deepresearch	49.71	49.51	49.45	50.12	50.00	6.70	61.27	31.13	7.60
<i>Prompt-Based Frameworks</i>									
WebWeaver (Qwen3-30B-A3B)	46.77	45.15	45.78	49.21	47.34	4.57	28.65	34.90	36.46
WebWeaver (Claude-Sonnet-4)	50.58	51.45	50.02	50.81	49.79	6.96	66.86	10.47	22.67
Enterprise DR (Gemini-2.5-Pro)	49.86	49.01	50.28	50.03	49.98	6.82	71.57	19.12	9.31
RhinoInsigh (Gemini-2.5-Pro)	50.92	50.51	51.45	51.72	50.00	6.82	68.51	11.02	20.47
<i>Trained Open Models</i>									
WebShaper-32B	34.93	31.58	26.17	44.81	40.38	1.63	3.25	3.75	93.00
WebThinker-32B-DPO	–	39.40	35.40	46.00	43.50	–	–	–	–
DR Tulu-8B	–	41.70	41.80	48.20	41.30	–	–	–	–
<i>Our Deep Research Systems</i>									
WARP-8B (SFT)	46.73	46.24	48.10	47.61	41.79	6.04	54.17	10.29	35.54
WARP-8B (Atomic RL)	48.81	48.70	51.36	48.64	42.25	6.06	56.13	11.03	32.84
WARP-8B (Pipeline RL)	50.11	50.54	52.64	48.87	44.17	6.60	57.60	13.73	28.68

surpassing Gemini-2.5-Pro-deepresearch (49.71). It also attains state-of-the-art results on DeepResearch Gym, with an average score of 98.48. These results show that integrating WARP with *multi-stage agentic training* enables small models to reach the performance level of leading proprietary research systems. These findings suggest that, for deep research tasks, the primary bottleneck lies not in model size, but in the design of effective cognitive and planning processes that fully leverage a model’s inherent capabilities.

3.3 Analysis

3.3.1 Does Reasoning As Writing Policy remain effective without training?

To assess whether our framework is intrinsically effective without training, we conduct a prompt-based comparison on DeepResearch Bench using a larger model, *Qwen3-235B-A22B-Instruct-2507* (Yang et al., 2025). We compare two policies: (1) *Plan-then-write*, where the model first constructs a detailed outline through retrieval and then generates the report from this fixed plan; and (2) WARP, which starts from a simple outline and interleaves writing with iterative deepening.

As shown in Table 3, WARP consistently outperforms the *Plan-then-write* paradigm across all metrics, with a notable gain in *Insight* (+1.19) and *Comprehensiveness* (+0.98). By using the evolving draft as a reasoning context, WARP can de-

Table 3: Evaluation for different generation paradigms.

Paradigm	Overall	Comp.	Insight	Inst.	Read.
Plan-then-write	49.90	49.35	51.60	50.13	46.46
WARP (ours)	50.72	50.33	52.79	50.32	47.20

tect underdeveloped or ambiguous content during writing and trigger targeted *Deepening* with additional evidence, whereas *Plan-then-write* remains constrained by a static outline. This confirms that draft-aware deepening is a key driver of insight.

3.3.2 How Multi-Stage Training shapes agent actions and report structure?

In this section, we analyze how agent behavior evolves across training stages by examining statistics of its actions and report sections. Specifically, we focus on the *Write* and *Expand* actions, which directly determine the report structure.

Table 4: Evolution of action usage and hierarchical sectioning across training stages on DeepResearch Bench.

Stages	Actions		Sections		
	Write	Expand	Level-1	Level-2	Level-3
SFT	21.24	4.44	6.27	10.11	4.86
Atomic RL	36.89	8.88	6.49	14.17	16.50
Pipeline RL	39.51	8.63	6.52	15.75	17.32

As shown in Table 4, a clear behavioral shift emerges when moving from SFT to the RL-based

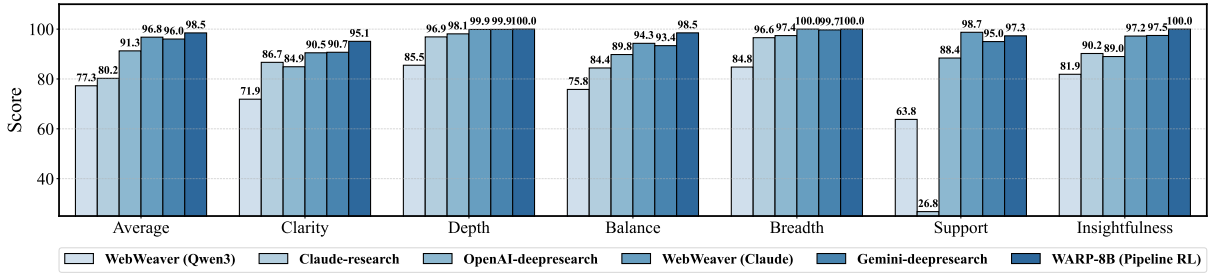


Figure 4: Performance of agent systems on DeepResearch Gym.

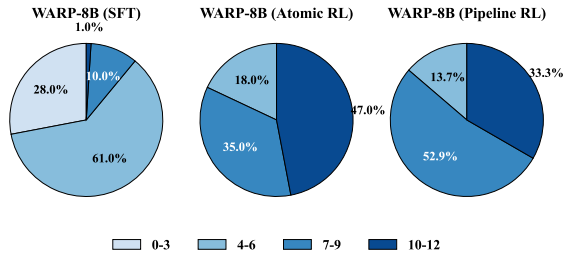


Figure 5: The *Expand* steps on DeepResearch Bench.

training stages. The frequency of *Expand* (Deepening) actions nearly doubles (from 4.44 to around 8.8), which in turn leads to a dramatic growth in fine-grained subsectioning (Level-3: from 4.86 to 17.32). Moreover, Figure 5 illustrates that RL training drives the agent to deepen more compared to SFT, ensuring at least 4 *Expand* steps in all cases. This trend indicates that RL training effectively equips the agent with *Reasoning-Driven Deepening*: Rather than adhering to a shallow outline as in SFT, the agent learns to identify underdeveloped parts of a draft and proactively expand them through iterative refinement, resulting in reports with substantially richer structure.

3.3.3 How does the number of deepening influence report quality?

To examine whether the model has learned an appropriate stopping policy for *deepening* and to quantify how deepening depth affects report quality, we conduct a "Forced Expansion" experiment. Specifically, during inference, we force WARP-8B (Pipeline RL) to apply the *Expand* action exactly k times, where k ranges from 0 to 15. We then compare this forced expansion curve with the actual deepening behaviors of WARP-8B at different training stages (SFT, Atomic RL, and Pipeline RL), overriding their learned termination policies. This allows us to directly evaluate report quality as a function of deepening depth and to compare the model's learned stopping behavior against the opti-

mal expansion point.

The results in Figure 6 reveal three consistent patterns. **First**, performance increases steadily with deeper expansion and begins to plateau at around nine steps, indicating diminishing returns beyond this depth. **Second**, both *Comprehensiveness* and *Insight* rise strongly with deepening, improving by nearly 6 points from shallow to sufficiently deep regimes, confirming the importance of iterative refinement for rich and insightful reports. **Third**, different training stages exhibit distinct stopping behaviors. The SFT agent typically stops within 6 steps and rarely reaches the saturation regime, whereas the Atomic RL and Pipeline RL agents shift their stopping distributions toward 6–15 steps, closely matching the empirically optimal depth.

3.3.4 How does the *Trajectory Pruning* affect the agent training?

To address the optimal stopping problem, we introduce a *trajectory pruning* strategy to construct higher-quality training data. In this section, we isolate its impact by training SFT models using the same teacher-generated trajectories, either with or without pruning. We consider two settings: (1) w/o pruning, which directly uses the raw trajectories produced by the teacher model, and (2) with pruning, which selects the best intermediate draft in one trajectory based on reward scores, retaining only the sub-trajectory before that draft.

Table 5: Effect of trajectory pruning on SFT training.

Trajectory	Overall	Comp.	Insight	Inst.	Read.
w/o pruning	45.80	44.95	47.35	46.71	40.86
with pruning	46.73	46.24	48.10	47.61	41.79

As shown in Table 5, models trained on pruned trajectories consistently outperform those trained on raw teacher ones across all evaluation dimensions. This indicates that trajectory pruning effectively improves the quality of supervision. More

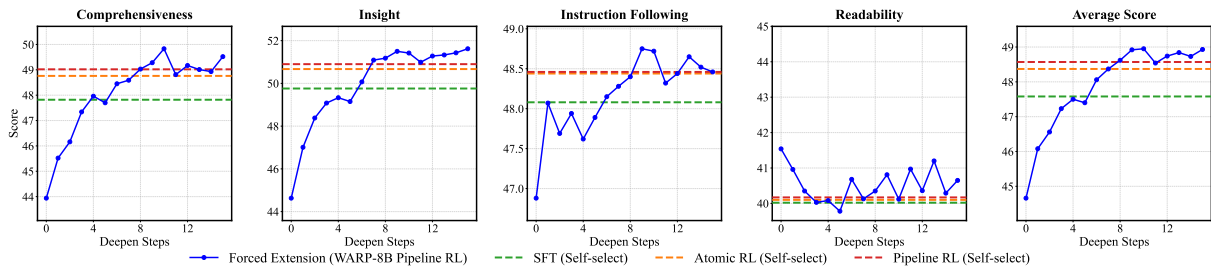


Figure 6: Mean Performance metrics per Deepen Steps on DeepResearch Bench.

453 importantly, these results reveal a key limitation of
 454 large teacher models: although they generate strong
 455 drafts, their termination decisions are often suboptimal. By applying reward-based selection over
 456 intermediate states, trajectory pruning filters out
 457 poorly timed stopping points and provides cleaner
 458 training signals. As a result, the student model
 459 learns a more accurate termination policy, which is
 460 crucial for effective *Reasoning-Driven Deepening*.
 461

4 Related Work

4.1 Deep Information Seeking

462 Retrieval-augmented generation (RAG) meth-
 463 ods (Li et al., 2022; Guu et al., 2020; Lewis
 464 et al., 2020; Wang et al., 2023; Press et al., 2022;
 465 Shao et al., 2023; Trivedi et al., 2022; Yan et al.,
 466 2024; Asai et al., 2023) enhance large models on
 467 knowledge-intensive tasks by incorporating exter-
 468 nal information, mitigating hallucinations, and ex-
 469 tending factual coverage. Beyond passive retrieval,
 470 recent agentic systems enable models to actively
 471 seek information through tool invocation (Li et al.,
 472 2025b; Wu et al., 2025; Schmidgall and Moor,
 473 2025; Wei et al., 2025b; Song et al., 2025; Chen
 474 et al., 2025a) or specialized actions (Tang et al.,
 475 2025; Zheng et al., 2025; Age, 2025). As task com-
 476 plexity increases, effective performance requires
 477 multi-round, reasoning-driven retrieval rather than
 478 isolated queries. We refer to this capability as *deep*
 479 *information seeking*, which has become a core com-
 480 ponent of recent deep research systems (Google,
 481 2025; x.AI, 2025; OpenAI, 2025; Perplexity, 2025;
 482 Li et al., 2025c; Kimi, 2025; Han et al., 2025; Team,
 483 2025; Asai et al., 2024). Existing benchmarks and
 484 systems primarily target either complex question
 485 answering (Mialon et al., 2024; Phan et al., 2025;
 486 Wei et al., 2025a; Zhou et al., 2025) or long-form
 487 research report generation (Du et al., 2025; Dee,
 488 2025; Coelho et al., 2025). Our work focuses on
 489 emphasizing scalable discovery of retrieval direc-
 490 tions to improve information completeness.
 491
 492

4.2 Knowledge-Intensive Long Writing

493 Knowledge-intensive long-form writing aims to
 494 generate reliable and comprehensive articles, such
 495 as Wikipedia-style entries (Shao et al., 2024), aca-
 496 demic surveys (Wang et al., 2024; Hu et al., 2024;
 497 Wang et al., 2025), and research reports (Li et al.,
 498 2025c,d). These tasks require both factual reliabil-
 499 ity—with explicit citations—and coverage of all
 500 major subtopics. Most prior approaches adopt a
 501 *retrieve-then-generate* pipeline combined with a
 502 *plan-then-write* strategy, where an initial outline
 503 guides subsequent hierarchical retrieval and writ-
 504 ing (Yan et al., 2025; Li et al., 2025c; Wang et al.,
 505 2025; Li et al., 2025d). Consequently, the final
 506 quality heavily depends on the completeness of
 507 the initial outline. To address this issue, existing
 508 works improve outline generation through multi-
 509 perspective discussion (Shao et al., 2024; Jiang
 510 et al., 2024), concept-pool expansion (Li et al.,
 511 2025a), or iterative refinement (Yan et al., 2025;
 512 Li et al., 2025d). In contrast, we adopt a *writing*
 513 *as reasoning* policy that integrates outline expan-
 514 sion into the writing process, enabling progressive
 515 refinement and targeted post-hoc revision.
 516

5 Conclusion

517 We present a new paradigm for open-ended deep
 518 research, **Writing As Reasoning Policy (WARP)**,
 519 which tightly integrates outline refinement with
 520 long-form writing. By iteratively alternating be-
 521 tween content generation and structural expansion,
 522 WARP treats intermediate drafts as compact knowl-
 523 edge that continuously updates the plan, overcom-
 524 ing the rigidity of fixed outlines while reducing
 525 reliance on large-scale models. Combined with a
 526 multi-stage agentic training strategy, our compact
 527 8B agent achieves strong performance across multi-
 528 ple benchmarks, surpassing leading closed-source
 529 systems. These results demonstrate that dynamic,
 530 draft-driven planning with small models is suffi-
 531 cient for high-quality deep research.
 532

533 Limitations

534 **Report presentation.** In most existing deep re-
535 search systems, including ours, tables and fig-
536 ures are generated inline with paragraph-level text.
537 However, constructing tabular layouts requires a
538 reasoning process fundamentally different from
539 writing prose, placing heavy demands on a model’s
540 structural and formatting abilities. This coupling
541 partly explains why agent systems based on smaller
542 models often underperform large ones in presenta-
543 tion quality. A promising direction is to decouple
544 presentation from content generation and assign it
545 to a dedicated rendering agent, which could enable
546 small models to achieve comparable or even supe-
547 rior layout quality. Moreover, current readability
548 evaluation remains largely text-based and weakly
549 reflects the true visual structure of rendered reports,
550 suggesting the need for visual-modality evaluation
551 in future work.

552 **Information sources.** Our system relies on a lo-
553 cally deployed textual knowledge base (e.g., arXiv
554 abstracts and web summaries), which ensures sta-
555 bility and reproducibility but limits coverage and
556 timeliness. It also lacks access to images, videos,
557 domain-specific corpora, and personalized data. Fu-
558 ture extensions will expand the knowledge base
559 to support multi-modal content, local and person-
560 alized sources, and continuous updates, enabling
561 richer and more realistic research scenarios.

562 References

563 2025. Agenticseek: Private, local manus alterna-
564 tive. <https://github.com/Fosowl/agenticSeek>.
565 GitHub repository.

566 2025. Deepconsult: A deep research benchmark for
567 consulting/business queries. <https://github.com/youdotcom-oss/ydc-deep-research-evals>.
568 GitHub repository.

570 Akari Asai, Jacqueline He, Rulin Shao, Weijia Shi,
571 Amanpreet Singh, Joseph Chee Chang, Kyle Lo,
572 Luca Soldaini, Sergey Feldman, Mike D’arcy, and
573 1 others. 2024. Openscholar: Synthesizing scien-
574 tific literature with retrieval-augmented lms. *arXiv*
575 *preprint arXiv:2411.14199*.

576 Akari Asai, Zeqiu Wu, Yizhong Wang, Avirup Sil, and
577 Hannaneh Hajishirzi. 2023. Self-rag: Learning to
578 retrieve, generate, and critique through self-reflection.
579 In *The Twelfth International Conference on Learning*
580 *Representations*.

581 ByteDance. 2025. Doubao deep research. [https://](https://www.doubao.com/chat/)
582 www.doubao.com/chat/.

Mingyang Chen, Tianpeng Li, Haoze Sun, Yijie Zhou,
Chenzheng Zhu, Haofen Wang, Jeff Z. Pan, Wen
Zhang, Huajun Chen, Fan Yang, Zenan Zhou, and
Weipeng Chen. 2025a. Research: Learning to rea-
son with search for llms via reinforcement learning.
Preprint, arXiv:2503.19470. 583 584 585 586 587 588

Yuxuan Chen, Dewen Guo, Sen Mei, Xinze Li, Hao
Chen, Yishan Li, Yixuan Wang, Chaoyue Tang, Ruob-
ing Wang, Dingjun Wu, and 1 others. 2025b. Ul-
trarag: A modular and automated toolkit for adap-
tive retrieval-augmented generation. *arXiv preprint*
arXiv:2504.08761. 589 590 591 592 593 594

Meet claude. 2025. Claude deep research. [https://](https://www.anthropic.com/claude)
www.anthropic.com/claude. 595 596

João Coelho, Jingjie Ning, Jingyuan He, Kangrui Mao,
Abhijay Paladugu, Pranav Setlur, Jiahe Jin, Jamie
Callan, João Magalhães, Bruno Martins, and 1 others.
2025. Deepresearchgym: A free, transparent, and
reproducible evaluation sandbox for deep research.
arXiv preprint arXiv:2505.19253. 597 598 599 600 601 602

Mingxuan Du, Benfeng Xu, Chiwei Zhu, Xiaorui Wang,
and Zhendong Mao. 2025. Deepresearch bench: A
comprehensive benchmark for deep research agents.
arXiv preprint arXiv:2506.11763. 603 604 605 606

Google. 2025. Gemini deep research. [https://](https://gemini.google/overview/deep-research/)
gemini.google/overview/deep-research/. 607 608

Kelvin Guu, Kenton Lee, Zora Tung, Panupong Pasu-
pat, and Mingwei Chang. 2020. Retrieval augmented
language model pre-training. In *International confer-*
ence on machine learning, pages 3929–3938. PMLR. 609 610 611 612

Rujun Han, Yanfei Chen, Zoey CuiZhu, Lesly Miculi-
cich, Guan Sun, Yuanjun Bi, Weiming Wen, Hui Wan,
Chunfeng Wen, Solène Maître, and 1 others. 2025.
Deep researcher with test-time diffusion. *arXiv*
preprint arXiv:2507.16075. 613 614 615 616 617

Chen Hu, Haikuo Du, Heng Wang, Lin Lin, Mingrui
Chen, Peng Liu, Ruihang Miao, Tianchi Yue, Wang
You, Wei Ji, and 1 others. 2025. Step-deepresearch
technical report. *arXiv preprint arXiv:2512.20491*. 618 619 620 621

Yuntong Hu, Zhuofeng Li, Zheng Zhang, Chen Ling,
Raasikh Kanjani, Boxin Zhao, and Liang Zhao. 2024.
Hireview: Hierarchical taxonomy-driven automatic
literature review generation. 622 623 624 625

Yucheng Jiang, Yijia Shao, Dekun Ma, Sina J Sem-
nani, and Monica S Lam. 2024. Into the unknown
unknowns: Engaged human learning through partici-
pation in language model agent conversations. *arXiv*
preprint arXiv:2408.15232. 626 627 628 629 630

Kimi. 2025. Kimi-researcher: End-to-end rl train-
ing for emerging agentic capabilities. [https://](https://moonshotai.github.io/Kimi-Researcher/)
moonshotai.github.io/Kimi-Researcher/. 631 632 633

Yu Lei, Shuzheng Si, Wei Wang, Yifei Wu, Gang
Chen, Fanchao Qi, and Maosong Sun. 2025. Rhi-
noinight: Improving deep research through control 634 635 636

637	mechanisms for model behavior and context. <i>arXiv preprint arXiv:2511.18743</i> .	Akshara Prabhakar, Roshan Ram, Zixiang Chen, Silvio Savarese, Frank Wang, Caiming Xiong, Huan Wang, and Weiran Yao. 2025. Enterprise deep research: Steerable multi-agent deep research for enterprise analytics. <i>arXiv preprint arXiv:2510.17797</i> .	692
638			693
639	Patrick Lewis, Ethan Perez, Aleksandra Piktus, Fabio Petroni, Vladimir Karpukhin, Naman Goyal, Heinrich Küttler, Mike Lewis, Wen-tau Yih, Tim Rocktäschel, and 1 others. 2020. Retrieval-augmented generation for knowledge-intensive nlp tasks. <i>Advances in neural information processing systems</i> , 33:9459–9474.	Ofir Press, Muru Zhang, Sewon Min, Ludwig Schmidt, Noah A Smith, and Mike Lewis. 2022. Measuring and narrowing the compositionality gap in language models. <i>arXiv preprint arXiv:2210.03350</i> .	694
640			695
641			696
642			697
643			698
644			699
645			700
646	Derek Li, Jiaming Zhou, Leo Maxime Brunswic, Abbas Ghaddar, Qianyi Sun, Liheng Ma, Yu Luo, Dong Li, Mark Coates, Jianye Hao, and 1 others. 2025a. Omni-thinker: Scaling multi-task rl in llms with hybrid reward and task scheduling. <i>arXiv preprint arXiv:2507.14783</i> .	Marlene Scardamalia and Carl Bereiter. 1987. Knowledge telling and knowledge transforming in written composition. <i>Advances in applied psycholinguistics</i> , 2:142–175.	701
647			702
648			703
649			704
650			705
651			706
652	Xiaoxi Li, Guanting Dong, Jiajie Jin, Yuyao Zhang, Yujia Zhou, Yutao Zhu, Peitian Zhang, and Zhicheng Dou. 2025b. Search-ol: Agentic search-enhanced large reasoning models . <i>Preprint</i> , arXiv:2501.05366.	Samuel Schmidgall and Michael Moor. 2025. Agentriv: Towards collaborative autonomous research. <i>arXiv preprint arXiv:2503.18102</i> .	707
653			708
654			709
655			710
656	Xiaoxi Li, Jiajie Jin, Guanting Dong, Hongjin Qian, Yutao Zhu, Yongkang Wu, Ji-Rong Wen, and Zhicheng Dou. 2025c. Webthinker: Empowering large reasoning models with deep research capability . <i>Preprint</i> , arXiv:2504.21776.	Rulin Shao, Akari Asai, Shannon Zejiang Shen, Hamish Ivison, Varsha Kishore, Jingming Zhuo, Xinran Zhao, Molly Park, Samuel G Finlayson, David Sontag, and 1 others. 2025. Dr tulur: Reinforcement learning with evolving rubrics for deep research. <i>arXiv preprint arXiv:2511.19399</i> .	711
657			712
658			713
659			714
660			715
661	Yongqi Li, Wenjie Li, and Liqiang Nie. 2022. Dynamic graph reasoning for conversational open-domain question answering. <i>ACM Transactions on Information Systems (TOIS)</i> , 40(4):1–24.	Yijia Shao, Yucheng Jiang, Theodore A Kanell, Peter Xu, Omar Khattab, and Monica S Lam. 2024. Assisting in writing wikipedia-like articles from scratch with large language models. <i>arXiv preprint arXiv:2402.14207</i> .	716
662			717
663			718
664			719
665	Zijian Li, Xin Guan, Bo Zhang, Shen Huang, Houquan Zhou, Shaopeng Lai, Ming Yan, Yong Jiang, Pengjun Xie, Fei Huang, and 1 others. 2025d. Webweaver: Structuring web-scale evidence with dynamic outlines for open-ended deep research. <i>arXiv preprint arXiv:2509.13312</i> .	Zhihong Shao, Yeyun Gong, Yelong Shen, Minlie Huang, Nan Duan, and Weizhu Chen. 2023. Enhancing retrieval-augmented large language models with iterative retrieval-generation synergy. <i>arXiv preprint arXiv:2305.15294</i> .	720
666			721
667			722
668			723
669			724
670			725
671	Grégoire Mialon, Clémentine Fourrier, Thomas Wolf, Yann LeCun, and Thomas Scialom. 2024. GAIA: a benchmark for general AI assistants . In <i>The Twelfth International Conference on Learning Representations</i> .	Xiaofeng Shi, Qian Kou, Yuduo Li, Ning Tang, Jinxin Xie, Longbin Yu, Songjing Wang, and Hua Zhou. 2025. Scisage: A multi-agent framework for high-quality scientific survey generation. <i>arXiv preprint arXiv:2506.12689</i> .	726
672			727
673			728
674			729
675			730
676	Sewon Min, Kalpesh Krishna, Xinxi Lyu, Mike Lewis, Wen-tau Yih, Pang Wei Koh, Mohit Iyyer, Luke Zettlemoyer, and Hannaneh Hajishirzi. 2023. Factscore: Fine-grained atomic evaluation of factual precision in long form text generation. <i>arXiv preprint arXiv:2305.14251</i> .	Huatong Song, Jinhao Jiang, Yingqian Min, Jie Chen, Zhipeng Chen, Wayne Xin Zhao, Lei Fang, and Ji-Rong Wen. 2025. R1-searcher: Incentivizing the search capability in llms via reinforcement learning . <i>Preprint</i> , arXiv:2503.05592.	731
677			732
678			733
679			734
680			735
681			736
682	OpenAI. 2025. Deep research. https://openai.com/index/introducing-deep-research/ .	Jiabin Tang, Tianyu Fan, and Chao Huang. 2025. Autoagent: A fully-automated and zero-code framework for llm agents . <i>Preprint</i> , arXiv:2502.05957.	737
683			738
684	Perplexity. 2025. Perplexity deep research. https://www.perplexity.ai/hub/blog/introducing-perplexity-deep-research .	Zhengwei Tao, Jialong Wu, Wenbiao Yin, Junkai Zhang, Baixuan Li, Haiyang Shen, Kuan Li, Liwen Zhang, Xinyu Wang, Yong Jiang, and 1 others. 2025. Webshaper: Agentic data synthesizing via information-seeking formalization. <i>arXiv preprint arXiv:2507.15061</i> .	739
685			740
686			741
687	Long Phan, Alice Gatti, Ziwen Han, Nathaniel Li, Josephina Hu, Hugh Zhang, Chen Bo Calvin Zhang, Mohamed Shaaban, John Ling, Sean Shi, and 1 others. 2025. Humanity’s last exam. <i>arXiv preprint arXiv:2501.14249</i> .	MiniCPM Team, Chaojun Xiao, Yuxuan Li, Xu Han, Yuzhuo Bai, Jie Cai, Haotian Chen, Wentong Chen, Xin Cong, Ganqu Cui, and 1 others. 2025. Minicpm4: Ultra-efficient llms on end devices. <i>arXiv preprint arXiv:2506.07900</i> .	742
688			743
689			744
690			745
691			746
			747

748	Tongyi DeepResearch Team. 2025. Tongyi deep-	Xiangchao Yan, Shiyang Feng, Jiakang Yuan, Ren-	801
749	research: A new era of open-source ai re-	qiu Xia, Bin Wang, Bo Zhang, and Lei Bai. 2025.	802
750	searchers. https://github.com/Alibaba-NLP/	Surveyforge: On the outline heuristics, memory-	803
751	DeepResearch .	driven generation, and multi-dimensional evalua-	804
		tion for automated survey writing. <i>arXiv preprint</i>	805
752	Harsh Trivedi, Niranjan Balasubramanian, Tushar	<i>arXiv:2503.04629</i> .	806
753	Khot, and Ashish Sabharwal. 2022. Interleav-		
754	ing retrieval with chain-of-thought reasoning for	An Yang, Anfeng Li, Baosong Yang, Beichen Zhang,	807
755	knowledge-intensive multi-step questions. <i>arXiv</i>	Binyuan Hui, Bo Zheng, Bowen Yu, Chang	808
756	<i>preprint arXiv:2212.10509</i> .	Gao, Chengen Huang, Chenxu Lv, and 1 others.	809
		2025. Qwen3 technical report. <i>arXiv preprint</i>	810
757	Haoyu Wang, Yujia Fu, Zhu Zhang, Shuo Wang, Zirui	<i>arXiv:2505.09388</i> .	811
758	Ren, Xiaorong Wang, Zhili Li, Chaoqun He, Bo An,		
759	Zhiyuan Liu, and 1 others. 2025. Llm mapreduce-v2:	Yuxiang Zheng, Dayuan Fu, Xiangkun Hu, Xiaojie Cai,	812
760	Entropy-driven convolutional test-time scaling for	Lyumanshan Ye, Pengrui Lu, and Pengfei Liu. 2025.	813
761	generating long-form articles from extremely long	Deepresearcher: Scaling deep research via reinforce-	814
762	resources. <i>arXiv preprint arXiv:2504.05732</i> .	ment learning in real-world environments . <i>Preprint</i> ,	815
		<i>arXiv:2504.03160</i> .	816
763	Yidong Wang, Qi Guo, Wenjin Yao, Hongbo Zhang,		
764	Xin Zhang, Zhen Wu, Meishan Zhang, Xinyu Dai,	Peilin Zhou, Bruce Leon, Xiang Ying, Can Zhang,	817
765	Qingsong Wen, Wei Ye, and 1 others. 2024. Autosur-	Yifan Shao, Qichen Ye, Dading Chong, Zhiling	818
766	vey: Large language models can automatically write	Jin, Chenxuan Xie, Meng Cao, and 1 others. 2025.	819
767	surveys. <i>Advances in neural information processing</i>	Browsecomp-zh: Benchmarking web browsing abil-	820
768	<i>systems</i> , 37:115119–115145.	ity of large language models in chinese. <i>arXiv</i>	821
		<i>preprint arXiv:2504.19314</i> .	822
769	Yile Wang, Peng Li, Maosong Sun, and Yang Liu.		
770	2023. Self-knowledge guided retrieval augmen-	A Method Details	823
771	tation for large language models. <i>arXiv preprint</i>		
772	<i>arXiv:2310.05002</i> .	A.1 The Prompts in WARP Framework	824
773	Jason Wei, Zhiqing Sun, Spencer Papay, Scott McK-	In our WARP framework, there are five actions	825
774	inney, Jeffrey Han, Isa Fulford, Hyung Won Chung,	in all three stages: <i>Initialize</i> , <i>Search</i> , <i>Write</i> , <i>Ex-</i>	826
775	Alex Tachard Passos, William Fedus, and Amelia	<i>expand(Deepen)</i> , and <i>Terminate</i> . In the Initialization	827
776	Glaese. 2025a. Browsecomp: A simple yet chal-	stage, the agent generates the initial Level-1 out-	828
777	lenging benchmark for browsing agents . <i>Preprint</i> ,	line with writing plans by the <i>search</i> and <i>initialize</i>	829
778	<i>arXiv:2504.12516</i> .	actions. The prompt for the <i>search</i> is shown in	830
		Figure 8, and the prompt for the <i>initialize</i> is shown	831
779	Zhepei Wei, Wenlin Yao, Yao Liu, Weizhi Zhang, Qin	in Figure 7. Then, in the Evidence-Based Draft-	832
780	Lu, Liang Qiu, Changlong Yu, Puyang Xu, Chao	ing stage, the agent writes the paragraphs by the	833
781	Zhang, Bing Yin, and 1 others. 2025b. Webagent-	<i>search</i> and <i>write</i> actions. The prompt for the <i>write</i>	834
782	r1: Training web agents via end-to-end multi-turn	is shown in Figure 9. After that, in the Reasoning-	835
783	reinforcement learning . <i>Preprint</i> , <i>arXiv:2505.16421</i> .	Driven Deepening stage, the agent will decision	836
784	<i>ArXiv preprint</i> .	whether to expand a section for more details by the	837
785	Junde Wu, Jiayuan Zhu, Yuyuan Liu, Min Xu, and	<i>expand</i> action or to end the total process directly by	838
786	Yueming Jin. 2025. Agentic reasoning: A stream-	the <i>terminate</i> action. The prompt for both actions	839
787	lined framework for enhancing LLM reasoning with	is shown in Figure 10.	840
788	agentic tools . In <i>Proceedings of the 63rd Annual</i>		
789	<i>Meeting of the Association for Computational Lin-</i>	A.2 User Query Construction	841
790	<i>guistics (Volume 1: Long Papers)</i> , pages 28489–	We constructed a dataset of approximately 2000	842
791	28503, Vienna, Austria. Association for Computa-	user queries with corresponding scoring checklists	843
792	tional Linguistics.	to support the multi-stage training process. Of	844
793	x.AI. 2025. Grok 3 beta — the age of reasoning agents .	these, around 700 queries are focused on special-	845
794		ized academic survey topics, while the remaining	846
795	Zekun Xi, Wenbiao Yin, Jizhan Fang, Jialong Wu, Run-	1300 address general research reporting topics.	847
796	nan Fang, Ningyu Zhang, Jiang Yong, Pengjun Xie,	For the academic survey queries, we employed	848
797	Fei Huang, and Huajun Chen. 2025. Omnithink: Ex-	a <i>reverse question construction</i> approach: we first	849
798	panding knowledge boundaries in machine writing	selected 700 surveys from ArXiv and then used a	850
	through thinking. <i>arXiv preprint arXiv:2501.09751</i> .	large model to generate user questions based on	851
799	Shi-Qi Yan, Jia-Chen Gu, Yun Zhu, and Zhen-Hua Ling.		
800	2024. Corrective retrieval augmented generation.		

You are a professional report generation expert, skilled at creating high-quality report outlines. Now, you need to analyze the users question and provide a simple article outline structure (only top-level sections).

**** User Query ****
[user query]

**** Latest Retrieved Information ****
[current information]

Notes

1. The outline must be comprehensive, logically sound, and aligned with the users stated preferences and requirements.
2. The output language must match the language of the users query.

**** Available Actions ****

- initialize: Generate the top-level section outline along with an appropriate title.

Action Format:

```
<action> {"name": "initialize", "title": "...", "sections": [{"title": "...", "plan": "..."}, {"title": "...", "plan": "..."}, ...]} </action>
```

**** Output Format ****

```
<thought> Provide a detailed reasoning process </thought>
```

```
<action> Action (in JSON format) </action>
```

Please output strictly according to the specified format.

Figure 7: The LLM prompt for the *initialize* action.

these articles. The prompts used for this process are available in Figure 11. For general research reporting queries, we selected 1300 real questions.

In addition, we constructed a **preference checklist** inspired by DeepResearch Bench (Du et al., 2025) for each user query. A large model was used to generate weighted scores for different evaluation aspects based on the existing questions. The final report score for a query is computed as a weighted sum across these aspects. The method for checklist generation is same as DeepResearch Bench. Out of the 2000 queries, approximately 1500 were used to **construct trajectory data** for SFT and single-step RL, while the remaining 500 were directly used for end-to-end RL training.

A.3 Retrieval Environment Setup

We constructed a local database containing approximately 2.86 million documents to serve as the agent’s retrieval environment. This database supports both trajectory data collection and interactive RL training. Among these documents, roughly 2.71 million are abstracts of papers from ArXiv, sourced via Kaggle¹. The remaining 150k documents come

¹<https://www.kaggle.com/api/v1/datasets/download/Cornell-University/arxiv>

from general web pages, for which we employed *Gemini 2.0-Flash* to generate concise summaries while controlling for document length and quality.

To ensure efficient retrieval, we built a vector database. Specifically, all documents were vectorized using the embedding model MiniCPM-Embedding-Light² and indexed with Faiss³. The pipeline is implemented via UltraRAG(Chen et al., 2025b), an open-source library for constructiong Retrieval-Augmented Generation (RAG) systems.

A.4 Trajectory Data Construction

Base on the user queries in §A.2 and the retrieval environment in §A.3, we collected 1,500 actual execution trajectories within our WARP framework.

We chose *Qwen3-235B-A22B-Instruct-2507* as the teacher model, and use the prompts in §A.1. Despite their scale, current large language models still struggle with high-level decision-making and cannot reliably determine when to stop. To address this, we introduce an **trajectory pruning** strategy. The gathering process here slightly differs from the standard inference phase of the WARP framework:

²<https://huggingface.co/openbmb/MiniCPM-Embedding-Light>

³<https://github.com/facebookresearch/faiss>

You are a searcher within a multi-agent system consisting of "Analyst-Searcher-Writer". You must perform retrieval based on instructions from the "Analyst". Carefully select the most accurate search keywords and strictly adhere to the specified output format. You should focus on the user's query and the current article outline to determine the most relevant keywords for searching. You can give one or less to five keywords. The content should be in the same language as the user's query.

** User Query **
[user query]

** Current Article Outline **
[current outline]

** Analyst's Instruction **
[current instruction]

Action Example:

```
<action> {"name": "search", "keywords": [keyword-1, keyword-2, ..]} </action>
```

** Output Format **
<thought> Your reasoning process </thought>
<action> Action (in JSON format) </action>

Please output strictly according to the specified format.

Figure 8: The LLM prompt for the *search* action.

during each *Reasoning-Driven Deepening* stage, we explicitly force the agent to select a position for outline deepening, instead of allowing the model to autonomously decide whether to expand or terminate. This ensures that the agent continuously expands the outline and revises the report until a maximum of 12 expansions is reached. All the results are scored by the survey-level reward mentioned in §A.6.2, and the highest-scoring result is selected as the endpoint of the trajectory.

For each user query, we collect a single high-quality execution trajectory. In total, 1500 trajectories were obtained corresponding to the 1500 user queries. Among these, 1200 trajectories were used as SFT data for cold-start training, while the remaining 300 were reserved for atomic skill RL.

A.5 Action Data Distribution

For a complete trajectory collected in Section A.4, it typically contains one *initialize* action, one *termination* action, several *expand* actions, and many *search* and *write* actions. These actions correspond to four core agent capabilities: *planning*, *retrieval*, *writing*, and *decision-making*. However, the natural distribution of actions is highly imbalanced with respect to training needs: the more critical and challenging abilities—such as *planning* and *decision-making*—are underrepresented, while easier-to-learn abilities—such as *searching*

and *writing*—dominate the trajectories.

To address this issue, we introduce an **action-level balanced sampling** strategy that increases the sampling probability of more important and difficult actions and decreases that of easier ones, thereby providing more effective supervision for training the agent’s core capabilities.

From the 1,500 collected trajectories, we obtained approximately 100k actions in total. When grouped by user query type, actions from academic review tasks and general report tasks followed an approximate ratio of 3:5. We used about 33k actions for cold-start training (SFT) and about 5k actions for atomic skill RL, with the remaining data discarded.

A.6 Reward System

In Reinforcement Learning (RL) training, the design of the reward functions is crucial, often directly impacting training efficiency and stability. In this section, we introduce our reward system from two aspects. First, in Section A.6.1, we introduce our different ability-specific reward functions for four abilities: **planning**, **retrieval**, **writing**, and **decision-making**, primarily used for single-step RL training. These functions are used for optimizing five actions: *initialize*, *expand-plan*, *search*, *write*, and *terminate*, as shown in Table 6. Second, in Section A.6.2, we introduce report-level

Table 6: The mapping between actions and the four agent abilities.

Agent Ability	Action Name	Action Parameters
planning	initialize	{"title": "...", "sections": [{"title": "...", "plan": "..."}, {"title": "...", "plan": "..."}, ...]}
	expand	{"position": "section-x.y.z", "content": "...", "subsections": [{"title": "...", "plan": "..."}, ...]}
retrieval	search	{"keywords": [keyword-1, keyword-2, ...]}
writing	write	{"position": "section-x.y.z", "title": "...", "content": "..."}
decision-making	terminate	{}

only contains user questions. We will directly evaluate the final result of the entire process (the generated report) from several aspects, instead of scoring the intermediate actions. We hope that in this phase, the agent can explore more freely, generate more diverse results, and ultimately surpass the capabilities of the teacher model. This section will detail our evaluation of the final result from four aspects: comprehensiveness, insight, instruction-following, and readability. The judgment model we use is *Qwen3-32B* (Yang et al., 2025).

B Experiments Details

B.1 Training Details

We adopt a three-stage training pipeline consisting of cold start training (SFT), atomic skill RL (single-step RL), and holistic pipeline RL (end-to-end RL). All our experiments are run on 8 A100 GPUs, and the training settings are shown in Table 7.

Cold-Start Training For cold-start, we collect approximately 33k action-level samples from 1,200 trajectories. The model is trained using SFT with a learning rate of $1.5e-5$ and batch size as 32 for 4 epochs, taking about 2 days to complete.

Atomic Skill RL We further perform single-step RL using approximately 5150 action-level samples from 300 trajectories. We set the learning rate to $2.5e-6$, batch size to 8, rollout number to 8, and train for 200 optimization steps, taking about 2 days to complete.

Holistic Pipeline RL Finally, we conduct end-to-end RL on 500 user queries, optimizing the entire report generation pipeline jointly. The learning rate remains $1e-6$, with a batch size of 8, rollout number of 4, and a total of 50 training steps, taking about 4 days to complete.

B.2 Metrics Details

We conducted evaluations on three benchmarks: DeepResearch Bench (Du et al., 2025), DeepConsult (Dee, 2025), and DeepResearch Gym (Coelho

Table 7: The training settings for different stages.

Parameters	SFT	Single-Step RL	End-to-End RL
user queries	1,200	300	500
train samples	33,292	5150	500
learning rate	$1.5e-5$	$2.5e-6$	$1e-6$
batch size	32	8	8
rollout number	–	8	4
train epochs	4	–	–
train steps	–	200	50

et al., 2025).

DeepResearch Bench It consists of 100 PhD-level research tasks spanning 22 academic domains. It adopts the RACE and FACT evaluation frameworks. RACE assesses Comprehensiveness, Insight/Depth, Instruction Following, and Readability, while FACT evaluates effective citations per report and citation reliability. We evaluate it by *Gemini-2.5-Pro*.

Deep Consult It includes 102 queries from business and consulting scenarios. Evaluation is conducted via pairwise comparisons against an *OpenAI-DeepSearch* baseline, reporting win, tie, and loss rates, together with average quality scores on instruction following, comprehensiveness, completeness, and writing quality. We evaluate it by *o3-mini-2025-01-31*.

DeepResearch Gym It is built on the Researchy Questions dataset. Following WebWeaver (Li et al., 2025d), we sample 100 queries from the top 1,000 test queries and evaluate the report quality in six aspects: clarity, depth, balance, breadth, support, and insightfulness. We evaluate it by *GPT-4.1-mini-20250414*.

C The Usage of LLMs

We use LLMs (e.g., ChatGPT) to help polish the paragraphs of our paper.

You are a writer operating within a multi-agent system consisting of "Analyzer-Searcher-Writer". Based on instructions from the "Analyzer", the current writing status, and the most recently retrieved information, you are to compose a new paragraph while ensuring logical coherence and accurate citation of facts.

You should give a paragraph with breadth and depth, ensuring it is informative and engaging. You are encouraged to incorporate examples, **tables**, code snippets, and other elements to enhance the content. But don't write other sections or chapters that are not assigned to you. You'd better give analytical and comparative content, not just a summary of facts. Please attention the coherence and logical flow of the entire article and the other sections. You can extract the claims from the retrieved information, and design how to write the paragraph based on the claims in the thought process.

BE FAITHFUL! Make sure all your claims, especially the facts and the numbers, can be supported by the retrieved information, with your citations. Don't add any claims can't be supported by your citations. All the facts or data in your claims should can be found in the retrieved information you cited.

You should ensure that the content you write is not redundant with other sections.

And you should strictly follow the citation format like `\cite{bibkey}` or `\cite{bibkey1, bibkey2..}` for any referenced information. The content should be in the same language as the user's query.

PLEASE JUST OUTPUT THE CONTENT IN ANALYZER'S INSTRUCTION, DO NOT OUTPUT OTHER SECTIONS.

THE OUTPUT SHOULD BE IN THE SAME LANGUAGE AS THE USER'S QUERY.

User Query
[user query]

Current Article Summary
[current survey]

Analyzers Instruction
[current instruction]

Retrieved Information
[current information]

Action Example:
<action> content </action>

Output Format
<thought> Your thought process </thought>
<action> Your Content (in Markdown format) (include BIBKEY for citations within the content) </action>

Please strictly follow the specified output format.

Figure 9: The LLM prompt for the *write* action.

You are a professional report-generation expert skilled at crafting high-quality report outlines. Based on the the users stated preferences, you must now determine whether any section requires expansion into subsections.

Important Notes:

1. Select only the single section or subsection most in need of expansion.
2. If no expansion is needed, output a "terminate" (no operation) action.
3. If you think the expansion is necessary to make the article more comprehensive or insightful, feel free to expand it.
4. Make sure the new subsections aren't redundant or overly detailed with other sections. If it's too detailed or redundant with other sections, just terminate it.
5. Make sure the new subsections are relevant and coherent with other sections.
6. You can only expand the section in 1 level and 2 level, do not expand the section in 3 level or more.
7. Please don't extend the section that is already extended.
8. Just extend one hierarchy level at a time, the subsections you give should not have more than one hierarchy level.
8. The output language must match the language of the users query.

** User Query **
[user query]

** Current Full Report **
[current survey]

** Available Actions **
- extend-plan: Expand a section by adding subsections (e.g., section-1 to section-1.1, section-1.2, section-1.3).
- terminate: No operation.

Action Format:

```
<action> {"name": "expand", "position": "section-x.y.z", "subsections": [{"title": "...", "plan": "..."}, {"title": "...", "plan": "..."}, ...]} </action>  
<action> {"name": "terminate"} </action>
```

** Output Format **
<thought> Provide a detailed reasoning process </thought>
<action> Action (in JSON format) </action>

Please output strictly according to the specified format.

Figure 10: The LLM prompt for the *expand* and *terminate* actions.

You are a Instruction-writing expert. Below is a Survey title. Your task is to infer the Instruction the user might have.

Survey Title:
{title}

Before crafting the Query, analyze the following:

1. **Avoid using exact titles.**
2. Use domain-specific keywords, synonyms.

First, output your analysis starting with "Thought:", simulating the Survey authors thought process.

Then, generate one or more Queries based on the analysis, starting with "Instruction:". **DONT OUTPUT EXPLANATIONS AFTER THE Instruction.**

Figure 11: The LLM prompt for reverse question (user query) construction.

Rating	Description
<i>Guide</i>	
Score 1	The outline fails to guide content generation, omitting significant aspects of the topic or providing insufficient direction.
Score 2	The outline provides limited guidance, covering some key areas but lacking depth or completeness in addressing the topic.
Score 3	The outline provides moderate guidance for content generation, addressing most key areas but leaving some gaps or ambiguities.
Score 4	The outline effectively guides content generation, covering all significant aspects with clear direction, though minor refinements could enhance comprehensiveness.
Score 5	The outline is exemplary in guiding content generation, thoroughly addressing all aspects of the topic with clear, detailed direction and no significant gaps.
<i>Hierarchical</i>	
Score 1	The outline exhibits no discernible hierarchical structure. Topics and subtopics are jumbled together without logical separation or clear levels, making it nearly impossible to follow or identify any organization.
Score 2	The outline attempts to establish a hierarchy but fails to maintain logical consistency. Main topics and subtopics are frequently misclassified, and the structure is overly rigid or disjointed. Subtopics may be missing, misplaced, or redundant, making it hard to grasp the intent of the structure.
Score 3	The outline demonstrates a basic level of logical coherence. Most topics follow a general sequence, but some sections feel forced, with weak or unclear transitions. There are small jumps in logic, causing slight confusion or loss of flow at certain points.
Score 4	The outline displays a clear, logical, and diverse hierarchical structure. Main topics are distinct, and subtopics are properly nested. While most elements are well-placed, there may be minor redundancies or opportunities to introduce more diverse formats for subtopics. Slight adjustments could achieve better precision and variety in style.
Score 5	The outline showcases an exceptional, flawless hierarchical structure. Each main topic is distinct, and subtopics are logically nested with absolute clarity and stylistic diversity. The outline demonstrates flexibility in structure and organization, adapting its style where appropriate for the content and logic. No further refinement is necessary.
<i>Coherence</i>	
Score 1	The outline is highly disjointed and incoherent. Topics and subtopics appear in a random, unordered manner, with no logical flow or sense of progression. Major conceptual gaps and illogical jumps are present throughout the structure.
Score 2	The outline shows some attempt at logical organization, but it contains frequent inconsistencies, abrupt shifts, or logical missteps. Topics and subtopics are misaligned or lack proper transitions, making the reader work hard to follow the structure.
Score 3	The outline demonstrates a basic level of logical coherence. Most topics follow a general sequence, but some sections feel forced, with weak or unclear transitions. There are small jumps in logic, causing slight confusion or loss of flow at certain points.
Score 4	The outline exhibits a strong sense of logical flow, with ideas presented in a mostly smooth and connected manner. Transitions between topics and subtopics are clear, but a few minor adjustments could make the flow more seamless or natural. The logic is sound, but room for refinement exists.
Score 5	The outline achieves exceptional logical coherence. Each topic and subtopic follows a deliberate, thoughtful progression, with clear, natural, and intuitive transitions. The reader experiences a seamless flow of ideas, and no adjustments are required to improve logical consistency or flow.

Table 8: The plan scoring criteria rating scale 1-5.

```

###Task Description:
An instruction (might include an Input inside it), a response to evaluate, a reference answer that gets a score of 5, and a score rubric representing a evaluation criteria are given.
1. Identify the major and minor errors in this Response. Write a detailed list of the errors in the response strictly based on the given score rubric, not evaluating in general.
2. After writing the list of errors, write a score that is an integer between 1 and 5. You should refer to the score rubric.
3. The output format should look as follows: "(write the list of errors for criteria) [RESULT] (an integer number between 1 and 5)"
4. Please do not generate any other opening, closing, and explanations.

5. Please be fair, don't hesitate to give a low score like 1 or 2.
6. Note that Major errors refer to actual errors that affects the task severely, may change the meaning of the output, and Minor errors refer to smaller imperfections, and purely subjective opinions about the output.

###The instruction to evaluate:
{instruction}

###Response to evaluate:
{response}

###Reference Answer (Score 5):
{reference_answer}

###Score Rubrics:
{rubric}

###Feedback:

```

Figure 12: The Outline quality reward prompt template.

```

Here is an academic survey about the topic "[TOPIC]":
---
[SURVEY]
---

<instruction>
Please evaluate this survey about the topic "[TOPIC]" based on the criterion provided below, identify the major and minor errors in this survey, and give a score from 1 to 5 according to the score description:
---

Criterion Description: [Criterion Description]
---
Score 1 Description: [Score 1 Description]
Score 2 Description: [Score 2 Description]
Score 3 Description: [Score 3 Description]
Score 4 Description: [Score 4 Description]
Score 5 Description: [Score 5 Description]
---

Note that Major errors refer to actual errors that affect the task severely, may change the meaning of the output, and Minor errors refer to smaller imperfections, and purely subjective opinions about the output.
There may be multiple errors or no errors in the output.
After listing the errors, then, please score the survey with 1 to 5.
Return the score without any other information at the end of the output.

```

Figure 13: The Content quality reward prompt template.

Rating	Description
<i>Relevance</i>	
Score 1	Very poor focus; discourse diverges significantly from the initial topic and intent with many irrelevant detours.
Score 2	Poor focus; some relevant information, but many sections diverge from the initial topic.
Score 3	Moderate focus; mostly stays on topic with occasional digressions that still provide useful information.
Score 4	Good focus; maintains relevance and focus throughout the discourse with minor divergences that add value.
Score 5	Excellent focus; consistently relevant and focused discourse, even when exploring divergent but highly pertinent aspects.
<i>Coverage</i>	
Score 1	Severely lacking; offers little to no coverage of the topic's primary aspects, resulting in a very narrow perspective.
Score 2	Partial coverage; includes some of the topic's main aspects but misses others, resulting in an incomplete portrayal
Score 3	Acceptable breadth; covers most main aspects, though it may stray into minor unnecessary details or overlook some relevant points.
Score 4	Good coverage; achieves broad coverage of the topic, hitting on all major points with minimal extraneous information.
Score 5	Exemplary in breadth; delivers outstanding coverage, thoroughly detailing all crucial aspects of the topic without including irrelevant information.
<i>Depth</i>	
Score 1	Very superficial; provides only a basic overview with significant gaps in exploration.
Score 2	Superficial; offers some detail but leaves many important aspects unexplored.
Score 3	Moderate depth; covers key aspects but may lack detailed exploration in some areas.
Score 4	Good depth; explores most aspects in detail with minor gaps.
Score 5	Excellent depth; thoroughly explores all relevant aspects with comprehensive detail, reflecting a deep and dynamic discourse.
<i>Novelty</i>	
score 1	Lacks novelty; the report strictly follows the user's initial intent with no additional insights.
score 2	Minimal novelty; includes few new aspects but they are not significantly related to the initial intent.
score 3	Moderate novelty; introduces some new aspects that are somewhat related to the initial intent.
score 4	Good novelty; covers several new aspects that enhance the understanding of the initial intent.
score 5	Excellent novelty; introduces numerous new aspects that are highly relevant and significantly enrich the initial intent.

Table 9: The content scoring criteria rating scale 1-5.