OPPO

# How Far Are We from Genuinely Useful Deep Research Agents?

**OPPO AI Agent Team**

## Abstract

Deep Research Agents (DRAs) aim to automatically produce analyst-level reports through iterative information retrieval and synthesis. However, most existing DRAs were validated on question-answering benchmarks, while research on generating comprehensive report remains overlooked. Worse, current benchmarks for report synthesis suffers from task complexity and subjective metrics—this fails to reflect user demands and limits the practical utility of generated reports. To address these gaps, we present Fine-grained DEepResearch bench (FINDER), an enhanced benchmark consisting of 100 human-curated research tasks with 419 structured checklist items that standardize report structure, analytical depth, and factual grounding. Based on approximately 1,000 reports produced by mainstream DRAs, we further propose Deep rEsearch Failure Taxonomy (DEFT), the first failure taxonomy for deep research agents. DEFT contains 14 fine-grained failure modes across reasoning, retrieval, and generation, and is built upon grounded theory with human–LLM co-annotating and inter-annotator reliability validation. Our experimental findings reveal that current DRAs struggle not with task comprehension but with evidence integration, verification, and reasoning-resilient planning.

## 1 Introduction

Deep Research Agents (DRAs) have recently attracted increasing attention due to their ability to autonomously retrieve, analyze, and synthesize web-scale information into structured research reports [1–3]. These agents utilize advanced techniques in multi-step web exploration, data retrieval, and synthesis to produce comprehensive reports that would traditionally require hours of manual effort. DRAs are increasingly applied in commercial sectors such as academic research, business intelligence, and knowledge management [4, 5].

However, despite their promising application potential, DRAs still fall short of expectations in real-world report generation tasks [6–10]. Existing benchmarks are mostly tailored for question-answering (QA) [11–14] or other types of close-ended tasks [15], fail to fully capture the nuances and strict requirements of practical deep research scenarios—where higher standards are imposed on the quality, accuracy, depth, and logical coherence of generated reports. Although a considerable number of open-ended benchmarks currently exist [6–9, 16], their tasks often stem from LLM-driven sampling or synthesis, leading to deviations from human demands and insufficient complexity.

To address this gap, we introduce Fine-grained DEep-Research bench (FINDER), a fine-grained benchmark designed to evaluate DRAs in a more comprehensive manner. Unlike existing benchmarks, DEFT is built upon 100 expert-curated research tasks with 419 detailed checklist items that guide the structure, analytical depth, and citation integrity of generated reports. As depicted in Figure 1, this explicit guidance enables more structured and reproducible evaluations of the task performance of DRAs. In addition, we propose Deep rEsearch Failure Taxonomy (DEFT), the first failure taxonomy developed specifically for DRAs. DEFT categorizes common errors into 14 fine-grained failure modes across three core dimensions—reasoning, retrieval, and generation—which we derive through grounded theory [17, 18] from extensive analysis of 1,000 generated reports. This taxonomy provides a robust framework for diagnosing where DRAs fail in their reasoning, information seeking, and content generation processes.

Our experimental evaluation on FINDER and DEFT of various DRAs, including proprietary systems [1–3], open-source models [19–24], and agent frameworks [21, 25–32], reveals several key insights. While systems like Gemini perform well across general benchmarks, our analysis shows that over 39% of failures arise in content generation, particularly through strategic content fabrication, where agents tend to generate unsupported but seemingly



**Figure 1** Comparison between DeepResearch Bench (DRB) and our FINDER.

professional content. Furthermore, retrieval-related failures, such as insufficient evidence integration and fact-checking issues, account for over 32% of errors, highlighting the challenges DRAs face in managing and verifying the quality of retrieved information. These results underscore that the core challenges for DRAs are not limited to simple task comprehension but instead involve deeper issues in evidence verification and reasoning resilience. To summarize, our contributions are as follows,

- We propose FINDER, a fine-grained benchmark with 100 expert-curated tasks and 419 structured checklist items, enabling robust and reproducible evaluation of DRAs across various dimensions of research report generation.
- We establish DEFT, the first failure taxonomy for DRAs, which categorizes errors into 14 fine-grained failure modes under three core dimensions: reasoning, retrieval, and generation.
- Through experiments on proprietary APIs, open-source models, and agent frameworks, we demonstrate that current DRAs struggle more with evidence integration and methodological rigor than with understanding tasks, revealing key weaknesses in reasoning resilience and strategic content fabrication.

## 2 Related Works

Early works on DRAs [1–3] employed datasets towards AGI as evaluation benchmarks. The most representative examples include GAIA [33] and HLE [34]. As the deep research community grows, researchers have proposed various specialized benchmarks [11–15]. Although the aforementioned datasets are challenging, they all fall under **closed-ended** assessments with standard answers. They neglect the evaluation of report generation, exhibiting a mismatch with the requirements of deep research. In contrast, open-ended benchmarks treat deep research as a task without a single definitive solution. DeepResearch Bench [16] contains 100 PhD-level problems spanning 22 domains, introducing the RACE and FACT evaluations for report quality and effectiveness. Mind2Web 2 [6] comprises 130 time-varying daily life tasks and proposes an "Agent-as-Judge" framework to achieve automated verification and attribution. DeepResearchGym [7] provides sandbox environments with reproducible search APIs and standardized protocols for transparent deep research benchmarking. DeepScholar-Bench [8] is a benchmark that automatically evaluates research synthesis abilities through content coverage, citation accuracy, and organizational quality. DRBench [9] focuses on enterprise scenarios
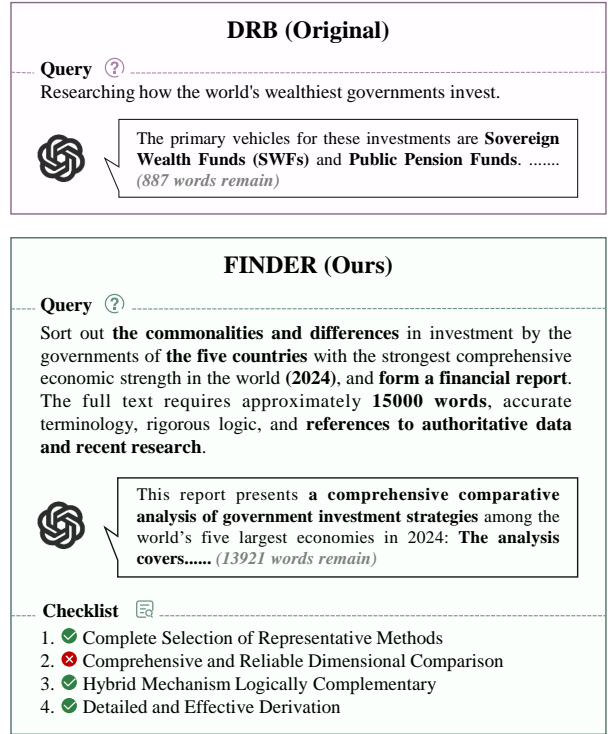
and evaluates DRAs through judge-based, citation-grounded assessment of long-form analytical reports. However, due to the dynamic nature of research reports, all these benchmarks employ subjective metrics based on the authors' experience or domain knowledge. Different benchmarks utilize varying metrics, lacking a unified standard.
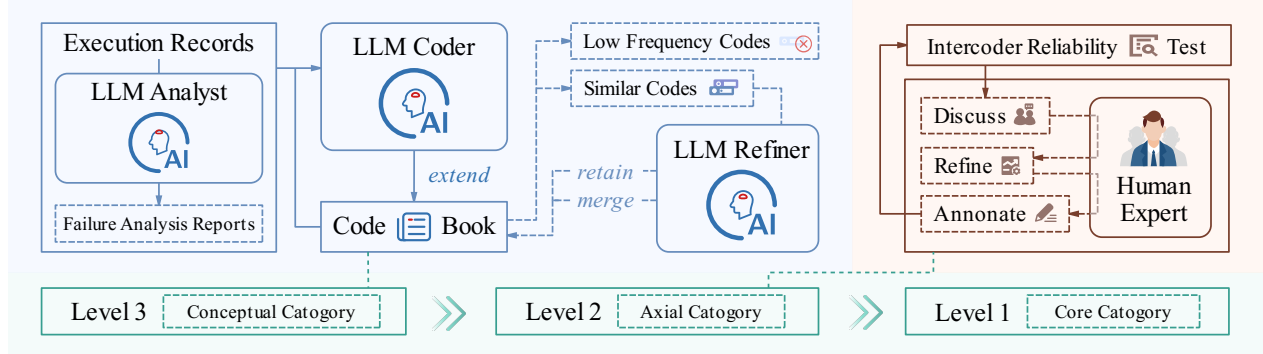
# 3 Methodology



**Figure 2** Overview of the DEFT Construction.

## 3.1 Fine-grained DEepResearch bench (FINDER)

Based on DeepResearch Bench, we refine the prompts and add structured checklists to construct FINDER, aiming to enhance evaluation precision and reproducibility.

### 3.1.1 Preliminary : DeepResearch Bench

The DeepResearch Bench consists of 100 PhD-level research tasks (50 in Chinese and 50 in English) designed to evaluate Deep Research Agents (DRA). It introduces two evaluation frameworks: RACE, which dynamically scores report quality in terms of comprehensiveness, depth, instruction-following, and readability; and FACT, which assesses retrieval effectiveness through citation accuracy and average effective citations (see Appendix E for detailed description). While DeepResearch Bench offers robust metrics for report evaluation, focusing solely on the final report does not adequately reflect a model's reasoning seach and information seeking capabilities that underpin its deep research performance.

### 3.1.2 Prompt Refinement

To address the issue of overly brief queries in the DeepResearch Bench, we invited seven domain experts to expand the queries in the DeepResearch Bench according to their respective areas of expertise. For each query, explicit guidelines were established regarding the report's length, disciplinary scope, presentation format, and other aspects. To ensure the correct generation type, each report was required to include the term "report" or equivalent expressions. Additionally, an independent expert who was not involved in the rewriting process manually evaluated the quality of the revised outputs. The finalized queries are presented in Figure 1. As shown in Figure A.1, our queries are longer than those in the original DeepResearch Bench. While preserving the independent semantic integrity of each sentence, the increased query length signifies a higher degree of task specification and research complexity.

### 3.1.3 Checklist Construction

To make the evaluation more structured, experts were first required to create five checklists for each query based on its specific characteristics. Each checklist served two purposes: first, to organize and structure the existing information within the query, and second, to supplement additional content requirements and constraints that were not explicitly mentioned but were relevant to the query. This approach ensured that the checklists were comprehensive and systematic during the evaluation process.

Subsequently, we used the Gemini 2.5 Flash to refine the initially generated checklists by eliminating items with incomplete semantics, ambiguous expressions, or those irrelevant to the reports generated for the corresponding queries. This process was conducted iteratively until all checklists met the predefined standards.

3

In total, we generated **419** checklists for 100 queries, with each query containing between three and five checklists. The distribution of checklist numbers is presented in Figure J.3. Further examples of queries are provided in Appendix A.2.

## 3.2 Failure Taxonomy

We construct a comprehensive failure taxonomy to systematically identify, categorize, and interpret the underlying causes of Deep Research Agent (DRA) errors. To avoid the subjective biases and omissions that may arise when relying solely on researchers' intuition or prior literature, the taxonomy is developed through a human-AI collaborative framework comprising open (§ 3.2.1), axial (§ 3.2.2), and selective coding (§ 3.2.3). The design of this process draws on grounded theory, which is a classic qualitative methodology that has been widely adopted across disciplines such as management [17], education [35], and software engineering [36] to construct evaluation or attribution schemata. The entire procedure (Figure 2) has been formalized into a pseudocode workflow, which is presented in Appendix F.

### 3.2.1 Open Coding

**Conceptual Category Generation.** Open coding entails analyzing and conceptualizing raw textual data to identify and label underlying conceptual categories within the study context [37]. Specifically, we collected performance metrics for **nine** deep research agent tasks in our benchmark and selected **five** large language models (Claude Opus 4.1, Gemini 2.5 Pro, Grok 4, DeepSeek-V3.1, and Qwen3-Max-Preview) from distinct model families to serve as coders. This design leverages their diverse inductive biases to broaden coverage and enhance coding breadth.

To systematically manage the coding process, we adopted the core principle of constant comparative analysis from grounded theory and maintained a dynamically updated conceptual inventory, hereafter referred to as the *codebook* ($\mathcal{C}$), defined as:

$$\mathcal{C} = \left\{ (c_i, d_i) \mid i = 1, 2, \ldots, N \right\}, \tag{1}$$

where $c_i$ denotes the concept name and $d_i$ its corresponding brief textual description. For each new concept identified, we first attempt to match it with existing $c_i$; if no match is found, a new pair $(c_{N+1}, d_{N+1})$ is added to $\mathcal{C}$.

Additionally, to focus the model on identifying and labeling failure modes rather than conducting deep causal analysis or automated failure localization[38], we instructed it to first generate a failure analysis report (Appendix D shows an example of the report) for each execution case as supplementary material to the original coding data. During the initial coding phase, we established a set of seed concepts (Appendix G) based on the research findings of Tang et al. [32] and Cemri et al. [39] to construct few-shot prompts that guided the large language model's coding process.

**Conceptual Category Optimization.** Whether within the same LLM coder or across multiple LLM coders, generated codes may exhibit redundancy or outliers. We address this through category optimization. On one hand, we employ Seed1.5-Embedding, which ranked first in MTEB (eng-v2, API available) [40], as the embedding model to identify concept pairs with cosine similarity $\geq 0.6$. These pairs are then fed into the large language model to be merged where appropriate. Additionally, concepts appearing below a removal threshold are discarded. As shown in Table F.1, we divided the source material into two groups for open-ended coding, each undergoing two rounds of refinement. The first round was conducted independently by each LLM coder, while the second round integrated the coding results from five LLM coders. An additional refinement round consolidated the coding outcomes between the two groups, ultimately yielding 51 concepts.

### 3.2.2 Axial Coding

Axial coding employs both deductive and inductive reasoning to explore relationships among concepts based on semantics, context, process, causality, function, structure, and strategy [41]. Through merging, splitting, removing, or modifying these relationships, it forms axial categories. At this stage, we conducted three rounds of coding based on inter-coder reliability (ICR) assessments: the first round utilized open coding results from Group A (Table F.1), while the second and third rounds incorporated all open coding results alongside the first-round axial coding outcomes. ICR measures the consistency among coders when encoding the same data [42] and has been demonstrated to consolidate [43, 44] or validate [45] existing coding frameworks. We selected Krippendorff's Alpha [46] to assess ICR. The

universal formula for Krippendorff's Alpha is as follows:

$$\alpha = 1 - \frac{D_o}{D_e} \tag{2}$$

where $D_o$ denotes observed disagreement and $D_e$ denotes expected disagreement by chance. For practical calculations, we utilized the web-based statistical package K-Alpha Calculator [47].

Following each coding round, to conduct ICR assessment, we engaged three domain experts to independently annotate a randomly sampled subset. This subset comprised 24 (first round) or 54 execution records (second and third rounds), with 3 logs selected from each Chinese and English version of each framework. It takes approximately 5 hours for experts to engage in discussion following each annotation round to resolve discrepancies and refine category definitions. After a few iterations, we finalized 14 axial categories. Detailed definitions of each category are provided in Appendix B, and illustrative case studies for each category are presented in Appendix C.

### 3.2.3 Selective Coding

Selective coding synthesizes the concepts and categories developed in the first two coding stages to establish overarching core categories. It clarifies their interrelationships and connects them through systematic logical threads [17]. At this stage, we repeatedly analyzed the axial categories derived from axial coding, ultimately distilling three core categories: *Reasoning*, *Retrieval*, and *Generation*. Functionally, these three core categories form a complete closed-loop for agent task execution. Temporally, they are interwoven and sequentially progressive, collectively underpinning a systematic understanding of agent failure mechanisms.

We randomly selected 36 execution records (six each from the Chinese and English part) generated by two agents not involved in the taxonomy construction stage, WebThinker and OpenManus, for coding analysis. No new categories emerged during this process, indicating that our categorization system had achieved theoretical saturation and demonstrated the explanatory power and stability required by grounded theory [48].

### 3.2.4 Positive Taxonomy Metric

To establish a unified and success-oriented framework for performance evaluation within the failure-mode taxonomy, we introduce a *positive performance metric* that transforms model error counts in each category into a bounded, interpretable score.

Let $E_i$ denote the number of observed errors in category $i \in \{1, \ldots, |\mathcal{T}|\}$, and let $|D|$ represent the total size of the dataset. Inspired by the concept of *cosine similarity* in vector space models [49], we define the performance score as

$$S_i = |D| \cdot \cos\left(\frac{E_i}{|D|} \cdot \frac{\pi}{2}\right). \tag{3}$$

Here, $S_i$ captures the angular deviation of model performance from an error-free optimum. When $E_i = 0$, the model attains the maximum possible score $S_i = |D|$. As the number of errors $E_i$ increases, $S_i$ monotonically decreases toward 0, reflecting a gradual decline in performance. Further justification of this formulation is provided in Appendix L.

## 4 Experiments

### 4.1 Evaluated Models.

We evaluate three representative categories of systems. (1) **Proprietary API** comprise Gemini-2.5-Pro Deep Research, O3 Deep Research, O4-Mini Deep Research, and Perplexity Deep Research, which are closed-source research agents accessible through API interfaces. (2) **Open-source Model** include open-source or self-hosted reasoning models such as MiroThinker, WebThinker, and AFM. (3) **Agent Framework** encompass OWL, OpenManus, and MiroFlow, where MiroFlow is evaluated in both English and Chinese versions to examine cross-lingual performance within a unified framework. Comprehensive model configurations and parameter settings are detailed in Appendix K.
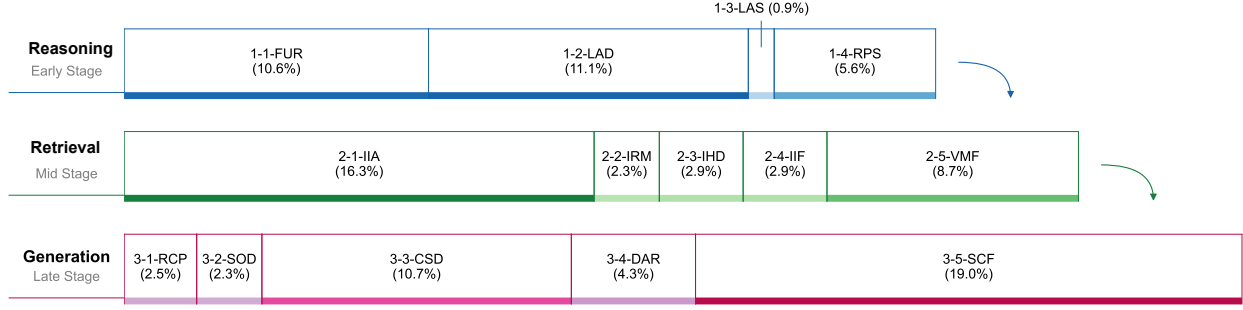
**Figure 3** Overview of the Level 1 (Core) and Level 2 (Axial) Failure Categories in DEFT

## 4.2 FINDER Performance Analysis

We evaluate model performance across three dimensions: RACE and FACT, Positive Taxonomy Metrics, and the Checklist Accuracy. The overall outcomes are summarized in Table 1.

| Model | RACE | | | | | FACT | | Positive Taxonomy Metric | | | | Checklist Pass Rate |
|---|---|---|---|---|---|---|---|---|---|---|---|---|
| | Overall | Comp. | Depth | Inst. | Read. | C.Acc. | E.Cit. | Rea. ($S_1$) | Ret. ($S_2$) | Gen. ($S_3$) | $S_{avg}$ | |
| *Proprietary API* | | | | | | | | | | | | |
| Gemini-2.5-Pro Deep Research[1] | **50.95** | **52.05** | **49.92** | 50.55 | **48.51** | 57.09 | 48.38 | 89.80 | **97.23** | **89.80** | **72.89** | 63.01 |
| Kimi K2[50] | 48.28 | 49.60 | 44.77 | **51.08** | 48.26 | - | - | **93.54** | 82.71 | 20.28 | 65.51 | **66.59** |
| O3 Deep Research[51] | 46.25 | 47.82 | 42.13 | 49.87 | 46.61 | **65.98** | **76.58** | 73.96 | 39.71 | 43.99 | 52.56 | 57.52 |
| O4-Mini Deep Research[52] | 43.49 | 43.91 | 38.00 | 49.21 | 44.02 | - | - | 93.54 | **75.01** | 45.40 | 71.32 | 56.09 |
| Perplexity Deep Research[3] | 41.62 | 43.68 | 38.39 | 44.30 | 41.12 | 5.25 | 29.31 | 50.90 | 60.04 | 30.90 | 47.28 | 51.55 |
| *Open-source Model* | | | | | | | | | | | | |
| WebThinker[19] | **41.11** | **41.43** | **34.51** | **47.71** | **43.56** | 11.32 | 1.83 | **72.70** | 24.87 | 9.41 | 35.73 | 44.87 |
| AFM[20] | 37.97 | 39.69 | 34.92 | 39.17 | 39.93 | 23.80 | **83.64** | 41.15 | **57.50** | 18.74 | **36.86** | 48.45 |
| MiroThinker[53] | 33.51 | 32.94 | 26.01 | 39.20 | 40.42 | **41.60** | 1.13 | 50.90 | 26.39 | 15.64 | 30.98 | 50.84 |
| Tongyi-DeepResearch[54] | 30.06 | 31.50 | 24.60 | 35.02 | 32.81 | 18.18 | 2.75 | 30.90 | 30.90 | **46.79** | 36.20 | **67.54** |
| *Agent Framework* | | | | | | | | | | | | |
| MiroFlow-English[21] | **42.20** | 42.84 | **36.49** | **47.55** | **44.51** | 22.73 | 2.00 | 54.90 | **46.79** | 15.64 | **39.11** | 72.19 |
| MiroFlow-Chinese[21] | 41.28 | **43.25** | 36.11 | 44.92 | 43.63 | 16.67 | 2.47 | 54.90 | 46.79 | 15.64 | 39.11 | 54.80 |
| OWL[25] | 39.22 | 39.57 | 33.81 | 44.41 | 40.13 | - | - | 49.55 | 43.99 | **29.40** | 40.98 | 53.94 |
| OpenManus[26] | 35.44 | 35.23 | 29.02 | 41.95 | 37.50 | 8.84 | **4.08** | **62.52** | 33.87 | 18.74 | 38.38 | 61.34 |

**Table 1** Overall evaluation results of **FINDER** across three complementary modules: RACE, FACT, and our DEFT Positive Metric (*reasoning $S_1$*, *retrieval $S_2$*, and *generation $S_3$*). The final column reports the Checklist Pass Rate. "–" indicates missing or unavailable results; detailed explanations of these cases are provided in Appendix M. **Bold** values denote the highest score within each group.

**RACE and FACT.** Under the RACE framework, Gemini 2.5 Pro Deep Research remains the top performer with an overall score of 50.95, followed by Kimi K2 (48.28) and O3 Deep Research (46.25). Among the Open-source Models and Agent Frameworks, WebThinker and MiroFlow stand out for their strong instruction adherence. MiroFlow was further evaluated using English and Chinese prompts from FINDER, each repeated three times to mitigate randomness; the results show that English tasks achieved slightly higher scores (42.20) compared to the Chinese version (41.28), indicating superior reasoning and text organization in English. Within the FACT framework, O3 Deep Research demonstrates exceptional performance, leading significantly in both factual precision (65.98) and citation reliability (76.58), while Gemini 2.5 Pro Deep Research follows as a strong contender, with the lower scores or data gaps for other models likely stemming from the more challenging upgraded prompts that demand denser reasoning and stricter citation validation.

**Positive Taxonomy Metrics.** The taxonomy results offer a process-level perspective on how models reason and synthesize information. Gemini achieves consistently high scores across reasoning , retrieval , and generation , indicating well-coordinated task understanding and synthesis. In contrast, Kimi K2 and O4-Mini exhibit exceptional reasoning capabilities (surpassing Gemini) and strong retrieval performance, but suffer from a sharp decline in generation scores. Open frameworks such as MiroFlow show moderate stability yet similarly face bottlenecks in the final generation stage. Overall, these metrics demonstrate that **superior systems maintain a balance among understanding, evidence collection, and synthesis rather than overoptimizing a single stage.**

**Checklist Accuracy.** Checklist scores represent meta-reasoning and procedural adherence to the intended research workflow. MiroFlow-English achieves the highest score (72.19), followed by a competitive cluster including Tongyi-DeepResearch (67.54), Kimi K2 (66.59), and Gemini 2.5 Pro (63.01). While MiroFlow demonstrates the specific advantage of explicit agentic orchestration, proprietary models like Kimi and Gemini remain robust, outperforming O3 (57.52) and other baselines. This distribution suggests that **systematic reasoning discipline—whether through framework design or intrinsic model capability—determines research reliability.**

## 4.3 DRB VS FINDER

As shown in Figure 4, we compare the original DeepResearch Bench (DRB) with our proposed FINDER under both the RACE and FACT frameworks, and this analysis reveals partially divergent outcomes across the two evaluation frameworks.

In the **RACE** framework, the overall scores under FINDER remain largely consistent with those from DRB. This consistency arises because both benchmarks share the same reference-based evaluation process: each model's research report is assessed relative to a standardized reference report (`reference.jsonl`) generated by Gemini-2.5-Pro Deep Research. The RACE framework evaluates the relative quality of a target report rather than its absolute performance, using four adaptive dimensions(comprehensiveness, depth of insight, instruction-following, and readability). Consequently, differences in absolute RACE scores across benchmarks hold limited interpretive value; only intra-benchmark comparisons among models reliably reflect relative generation quality.

In contrast, the **FACT** module shows more pronounced disparities between DRB and FINDER. While OpenAI Deep Research achieves a modest improvement in effective citation (*E.Cit.*), most other systems experience declines in both citation accuracy (*C.Acc.*) and effectiveness. This likely reflects the heightened difficulty introduced by our revised prompt design in FINDER, which imposes stricter factuality and citation validation demands. The resulting higher citation variance indicates that FINDER provides a more rigorous test of factual consistency and evidence trustworthiness. Overall, these outcomes suggest that FINDER enforces stronger constraints on reasoning transparency and source reliability, thereby exposing model weaknesses that were less evident under DRB's original configuration.
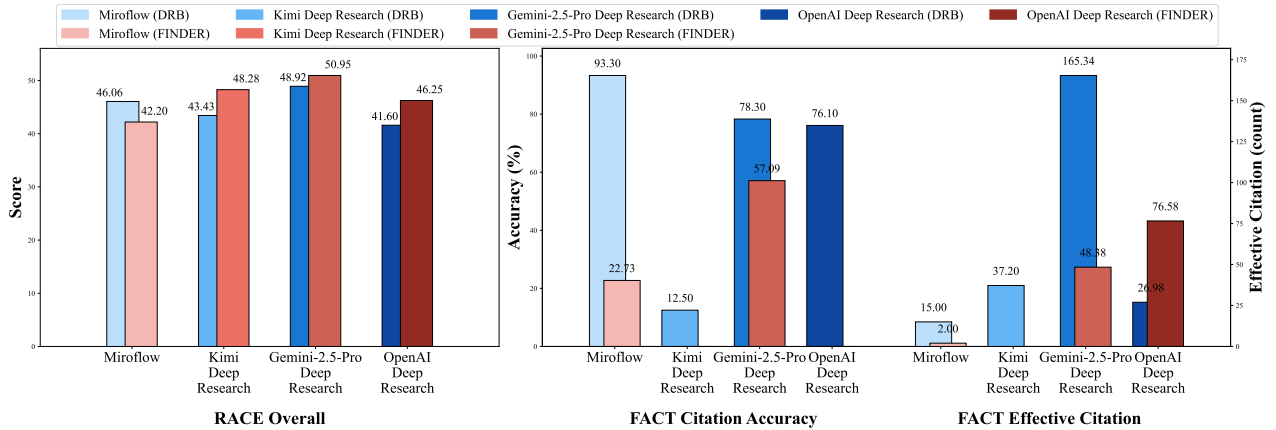


**Figure 4** Overview of agent performance on DeepResearch Bench (DRB) and our FINDER.

## 4.4 Deep Research Failure Taxonomy (DEFT)

This section introduces both the Level 1 (core) and Level 2 (axial) categories in the taxonomy, as shown in Table 2. Detailed definitions of axial category are provided in Appendix B. Furthermore, this section synthesizes key implications for enhancing DRA performance that emerge from the taxonomy-based analysis.

| Level 1 (Core Category) | Level 2 (Axial Category) |
|---|---|
| Reasoning | Failure to Understand Requirements (FUR) |
| | Lack of Analytical Depth (LAD) |
| | Limited Analytical Scope (LAS) |
| | Rigid Planning Strategy (RPS) |
| Retrieval | Insufficient External Information Acquisition (IIA) |
| | Information Representation Misalignment (IRM) |
| | Information Handling Deficiency (IHD) |
| | Information Integration Failure (IIF) |
| | Verification Mechanism Failure (VMF) |
| Generation | Redundant Content Piling (RCP) |
| | Structural Organization Dysfunction (SOD) |
| | Content Specification Deviation (CSD) |
| | Deficient Analytical Rigor (DAR) |
| | Strategic Content Fabrication (SCF) |

**Table 2** Taxonomy with Level 1 and Level 2 categories.

**Reasoning Category** refers to the failures mainly exhibited during the initial stage of execution, arising from insufficient consideration of user intent or problem details. Specifically, they include Failure to Understand Requirements (1-1-FUR, 10.55%), Lack of Analytical Depth (1-2-LAD, 11.09%), Limited Analytical Scope (1-3-LAS, 0.90%), and Rigid Planning Strategy (1-4-RPS, 5.60%).

The relatively low proportion of failures in this category indicates that most DRAs are capable of inheriting the underlying large models' strengths in terms of semantic understanding and basic reasoning [37]. However, the issue of 1-4-RPS suggests that the agents still exhibit limitations in dynamic task scheduling and adaptive reasoning. The linear execution logic present in some frameworks often fails to respond effectively to task evolution or intermediate feedback, leading to reduced efficiency or error propagation. In addition, 1-2-LAD and 1-3-LAS represent two orthogonal dimensions of reasoning capability. An ideal deep research agent should possess both strong problem-decomposition skills and robust system-modeling abilities.

> 🔍 **Insight 1:** Reasoning resilience, rather than reasoning intensity, is the key factor determining whether an agent can consistently produce high-quality deep research outcomes.

To address these issues, we introduce the concept of reasoning resilience. Reasoning resilience concerns an agent's ability to maintain and adjust its reasoning state within dynamic task environments, whereas reasoning intensity reflects its upper bound of analytical or reasoning capacity under ideal conditions. Deep research tasks are often accompanied by feedback, evolution, and noise[55]. In such contexts, strong reasoning capability does not necessarily ensure stable performance [56]. Only systems with reasoning resilience can continuously detect deviations, recalibrate reasoning search, and adapt strategies throughout complex and evolving reasoning processes, thereby achieving a balance of depth, breadth, accuracy, and consistency in their outcomes.
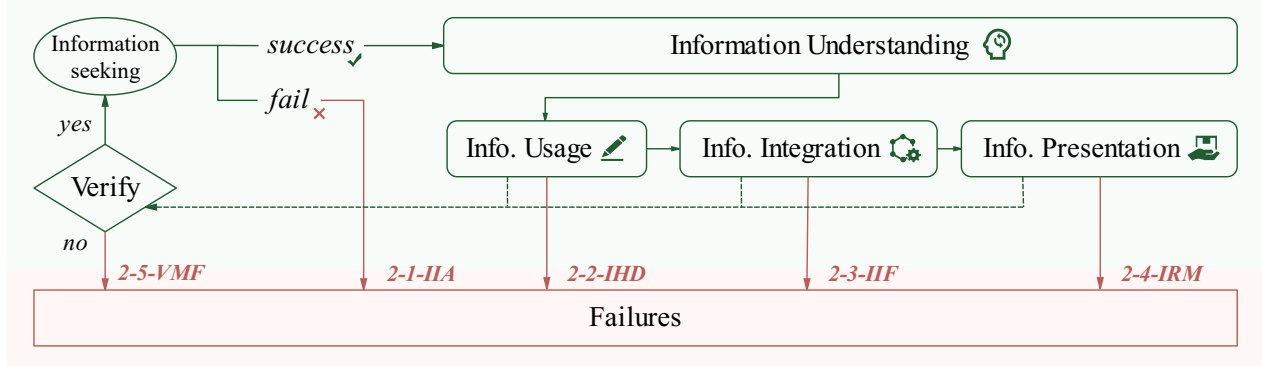
**Figure 5** A Brief Information Retrieval Workflow in Deep Research and Its Potential Failures

**Retrieval Category** refers to the failures mainly exhibited during the middle stage of execution, arising from inadequate abilities in external knowledge retrieval and evidence construction. Specifically, they include Insufficient External Information Acquisition (2-1-IIA, 16.30%), Information Handling Deficiency (2-2-IHD, 2.26%), Information Integration Failure (2-3-IIF, 2.91%), Information Representation Misalignment (2-4-IRM, 2.91%), and Verification Mechanism Failure (2-5-VMF, 8.72%).

The failures within the Retrieval Category exhibit stage-specific correlations along the task workflow. As shown in Figure 5, 2-1-IIA reflects primarily the agent's inability to initiate or execute the search for information effectively, occurring at the initial stage of the retrieval process. 2-2-IHD, 2-3-IIF, and 2-4-IRM occur after preliminary retrieval has succeeded, and correspond to failures in the utilization, integration, and representation of information. The absence of 2-5-VMF manifests at the terminal stage, where the agent fails to perform cross-check when encountering critical or conflicting information, resulting in outputs that lack factual grounding and credible support.

> ⚙ **Insight 2:** Retrieval in deep research is not a simple process of requesting and receiving; rather, it constitutes a closed-loop that integrates acquisition, processing, integration, representation, and verification.

DRAs often separate the stages of information acquisition, processing, integration, representation, and verification, resulting in fragmented or distorted knowledge chains. To address this issue, it is essential to enhance the agent's capacity for coherent knowledge management. For example, during the initial retrieval stage, a well-defined decision framework should be established to determine when to retrieve, what to retrieve, and how to utilize the retrieved results. In the intermediate stage, explicit mechanisms should be implemented to monitor information states and dynamically adjust retrieval strategies. In the final stage, a mandatory verification mechanism should be activated to cross-check critical facts.

**Generation Category** refers to the failures mainly exhibited during the latter stages of task execution, resulting from limited capability in content organization and expression. Specifically, they include Redundant Content Piling (3-1-RCP, 2.51%), Structural Organization Dysfunction (3-2-SOD, 2.26%), Content Specification Deviation (3-3-CSD, 10.73%), Deficient Analytical Rigor (3-4-DAR, 4.31%), and Strategic Content Fabrication (3-5-SCF, 18.95%).

The Generation Category exhibits the highest proportion of failures, particularly in 3-5-SCF. This failure indicates that the agents tend to generate seemingly professional but factually unsupported terms, methods, or references in order to create an illusion of academic rigor [57, 58]. In terms of outcome, 3-1-RCP shares similarities with 3-5-SCF, as both lead to outputs that are verbose, loosely structured, and lacking in substantive insight, thereby making it difficult for users to make effective judgments or take concrete actions [59]. The above analysis indicates that pre-constraints and post-verifications should extend beyond the retrieval stage to include generative dimensions such as text organization, linguistic structure, and formatting standards.

9

> ☆ **Insight 3:** Strengthening the constraints and verifications in the generative process is an important approach to improving the quality of the deep research output.

## 4.5 Evaluation of DEFT's Effectiveness

We evaluated the effectiveness of DEFT from three key aspects:

**Inter-Coder Reliability (ICR) Assessment.** Inter-Coder Reliability (ICR) Assessment. We invited four domain experts to evaluate the report-generation outputs produced by WebThinker and OpenManus. We calculated Krippendorff's alpha coefficient to measure the consistency between human annotations and Gemini 2.5-Flash assessments regarding both core category classification and Checklist Accuracy. The overall and dimension-specific coefficients are reported in Table 3, indicating strong stability and objective reproducibility for both the DEFT framework and the checklist evaluation. Details of the computation, formula, and interpretation are provided in Appendix H.

**Table 3** Krippendorff's Alpha Coefficients Between LLM–Human Coder Pairs Across Core Categories and Checklist Accuracy

| Model | Taxonomy Core Category | | | | Checklist Pass Rate |
|---|---|---|---|---|---|
| | Reasoning | Retrieval | Generation | Avg. | |
| OpenManus | 0.8005 | 0.7645 | 0.8960 | 0.8203 | 0.8025 |
| WebThinker | 0.7410 | 0.9016 | 0.9152 | 0.8526 | 0.8708 |

**Balanced Distribution of Identified Failures.** Our analysis of failure frequencies shows a relatively balanced distribution across the three primary dimensions (Figure 3): Reasoning (28.14%), Retrieval (33.10%), and Generation (38.76%). This balance suggests our taxonomy covers a diverse range of challenges in DRA report generation, avoiding an over-concentration on any single failure type.

**Structural Analysis of Failure Modes.** Our evaluation demonstrates that DEFT is an effective diagnostic tool. The taxonomy is not just a descriptive list; it has a meaningful internal structure. Our correlation analysis ((Figure 6)) confirms this by revealing three coherent failure clusters. These clusters map directly to specific operational failures. (1) The Process Integrity cluster shows how misunderstanding requirements (1.1 FUR) leads to an irrelevant or incomplete report (3.3 CSD). (2) The Content Integration cluster links source integration failure (2.4 IIF) to a chaotic structure (3.2 SOD) and high redundancy (3.1 RCP). (3) The Evidentiary Rigor cluster connects poor retrieval (2.1 IEIA) to "confident fabrications" (3.5 SCF). These systemic failure pathways confirm that DEFT captures significant, real-world mechanisms.

DEFT's effectiveness is also confirmed by its discriminative power. This is evidenced by key antagonistic axes. The analysis empirically separates distinct failure modes. For example, reports that are "concise but false" (3.5 SCF) are mechanistically different from those that are "verbose and disorganized" (3.1 RCP/3.2 SOD). The taxonomy also distinguishes methodological flaws (3.4 DAR) from process compliance. This proves a report can be procedurally correct but analytically unsound. Finally, specific links validate the taxonomy's hierarchy. For instance, superficial analysis (1.2 LAD) stems directly from poor retrieval (2.1 IEIA). This rich internal structure proves DEFT is an effective framework for modeling error propagation.
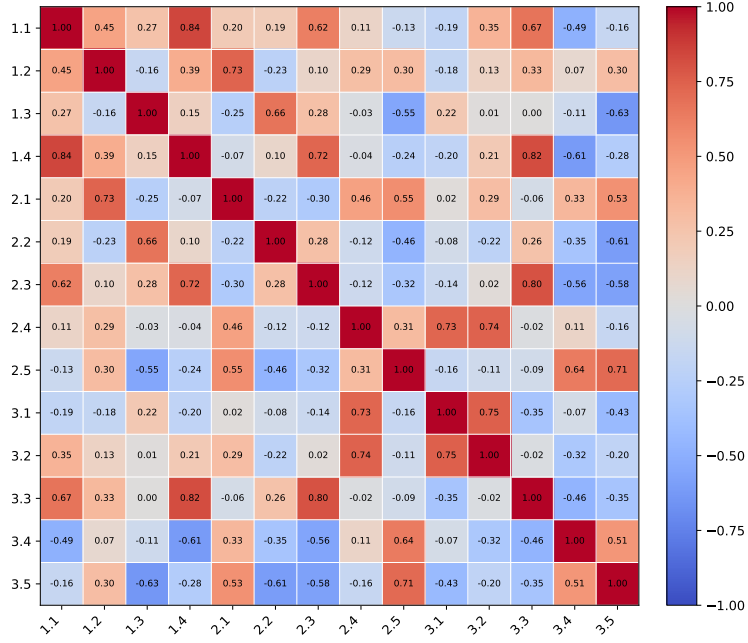
**Figure 6** DEFT failure categories correlation matrix.

# 5 Conclusion

This paper introduces FINDER and DEFT as the first unified framework for evaluating and diagnosing deep research agents at both task and process levels. By integrating 419 checklist-based assessments and a 14-category failure taxonomy, we reveal that current agents struggle less with understanding instructions and more with evidence information seeking, synthesis, and reasoning resilience. Our experiments demonstrate that even top-performing systems frequently fabricate unsupported content and fail to maintain methodological rigor. FINDER and DEFT provide actionable tools for the community to move beyond answer accuracy toward reliable, transparent, and verifiable deep research systems.

# 6 Contributions

**Core Contributors**

- Dingling Zhang
- He Zhu
- Jincheng Ren
- Kangqi Song

**Contributors**

- Xinran Zhou
- Boyu Feng
- Shudong Liu
- Jiabin Luo
- Weihao Xie
- Zhaohui Wang
- Tianrui Qin
- King Zhu
- Yuqing Wang
- Qianben Chen
- Yuchen Eleanor Jiang
- Wei Wang

**Corresponding Authors**

- Wangchunshu Zhou
- Jiaheng Liu

# References

[1] Dave Citron. Try deep research and our new experimental model in gemini, your ai assistant. https://blog.google/products/gemini/google-gemini-deep-research/, 2024.

[2] OpenAI. Introducing deep research, 2025. https://openai.com/index/introducing-deep-research/.

[3] Perplexity.ai. Introducing perplexity deep research. https://www.perplexity.ai/hub/blog/introducing-perplexity-deep-research, February 2025.

[4] Yuxuan Huang, Yihang Chen, Haozheng Zhang, Kang Li, Huichi Zhou, Meng Fang, Linyi Yang, Xiaoguang Li, Lifeng Shang, Songcen Xu, et al. Deep research agents: A systematic examination and roadmap. arXiv preprint arXiv:2506.18096, 2025.

[5] Renjun Xu and Jingwen Peng. A comprehensive survey of deep research: Systems, methodologies, and applications. arXiv preprint arXiv:2506.12594, 2025.

[6] Boyu Gou, Zanming Huang, Yuting Ning, Yu Gu, Michael Lin, Weijian Qi, Andrei Kopanev, Botao Yu, Bernal Jiménez Gutiérrez, Yiheng Shu, et al. Mind2web 2: Evaluating agentic search with agent-as-a-judge. arXiv preprint arXiv:2506.21506, 2025.

[7] 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.

[8] Liana Patel, Negar Arabzadeh, Harshit Gupta, Ankita Sundar, Ion Stoica, Matei Zaharia, and Carlos Guestrin. Deepscholarbench: A live benchmark and automated evaluation for generative research synthesis. arXiv preprint arXiv:2508.20033, 2025.

[9] Amirhossein Abaskohi, Tianyi Chen, Miguel Muñoz-Mármol, Curtis Fox, Amrutha Varshini Ramesh, Étienne Marcotte, Xing Han Lù, Nicolas Chapados, Spandana Gella, Christopher Pal, et al. Drbench: A realistic benchmark for enterprise deep research. arXiv preprint arXiv:2510.00172, 2025.

[10] Yuan Liang, Jiaxian Li, Yuqing Wang, Piaohong Wang, Motong Tian, Pai Liu, Shuofei Qiao, Runnan Fang, He Zhu, Ge Zhang, Minghao Liu, Yuchen Eleanor Jiang, Ningyu Zhang, and Wangchunshu Zhou. Towards personalized deep research: Benchmarks and evaluations, 2025. URL https://arxiv.org/abs/2509.25106.

[11] Jialong Wu, Wenbiao Yin, Yong Jiang, Zhenglin Wang, Zekun Xi, Runnan Fang, Linhai Zhang, Yulan He, Deyu Zhou, Pengjun Xie, et al. Webwalker: Benchmarking llms in web traversal. arXiv preprint arXiv:2501.07572, 2025.

[12] Jason Wei, Zhiqing Sun, Spencer Papay, Scott McKinney, Jeffrey Han, Isa Fulford, Hyung Won Chung, Alex Tachard Passos, William Fedus, and Amelia Glaese. Browsecomp: A simple yet challenging benchmark for browsing agents. arXiv preprint arXiv:2504.12516, 2025.

[13] Nikos I Bosse, Jon Evans, Robert G Gambee, Daniel Hnyk, Peter Mühlbacher, Lawrence Phillips, Dan Schwarz, Jack Wildman, et al. Deep research bench: Evaluating ai web research agents. arXiv preprint arXiv:2506.06287, 2025.

[14] Zijian Chen, Xueguang Ma, Shengyao Zhuang, Ping Nie, Kai Zou, Andrew Liu, Joshua Green, Kshama Patel, Ruoxi Meng, Mingyi Su, et al. Browsecomp-plus: A more fair and transparent evaluation benchmark of deep-research agent. arXiv preprint arXiv:2508.06600, 2025.

[15] Abhinav Java, Ashmit Khandelwal, Sukruta Midigeshi, Aaron Halfaker, Amit Deshpande, Navin Goyal, Ankur Gupta, Nagarajan Natarajan, and Amit Sharma. Characterizing deep research: A benchmark and formal definition. arXiv preprint arXiv:2508.04183, 2025.

[16] Mingxuan Du, Benfeng Xu, Chiwei Zhu, Xiaorui Wang, and Zhendong Mao. Deepresearch bench: A comprehensive benchmark for deep research agents, 2025. URL https://arxiv.org/abs/2506.11763.

[17] Chara Makri and Andy Neely. Grounded theory: A guide for exploratory studies in management research. International Journal of Qualitative Methods, 20:16094069211013654, 2021.

[18] B.G. Glaser and A.L. Strauss. The Discovery of Grounded Theory: Strategies for Qualitative Research. Observations (Chicago, Ill.). Aldine, 1967. ISBN 9780202302607. URL https://books.google.com/books?id=oUxEAQAAIAAJ.

[19] Xiaoxi Li, Jiajie Jin, Guanting Dong, Hongjin Qian, Yongkang Wu, Ji-Rong Wen, Yutao Zhu, and Zhicheng Dou. Webthinker: Empowering large reasoning models with deep research capability, 2025. URL https://arxiv.org/abs/2504.21776.

[20] Weizhen Li, Jianbo Lin, Zhuosong Jiang, Jingyi Cao, Xinpeng Liu, Jiayu Zhang, Zhenqiang Huang, Qianben Chen, Weichen Sun, Qiexiang Wang, Hongxuan Lu, Tianrui Qin, Chenghao Zhu, Yi Yao, Shuying Fan, Xiaowan Li, Tiannan Wang, Pai Liu, King Zhu, He Zhu, Dingfeng Shi, Piaohong Wang, Yeyi Guan, Xiangru Tang, Minghao Liu, Yuchen Eleanor Jiang, Jian Yang, Jiaheng Liu, Ge Zhang, and Wangchunshu Zhou. Chain-of-agents: End-to-end agent foundation models via multi-agent distillation and agentic rl, 2025. URL https://arxiv.org/abs/2508.13167.

[21] MiroMind. Mirothinker: An open-source agentic model series trained for deep research and complex, long-horizon problem solving. https://github.com/MiroMindAI/MiroThinker, 2025.

[22] Qianben Chen, Jingyi Cao, Jiayu Zhang, Tianrui Qin, Xiaowan Li, King Zhu, Dingfeng Shi, He Zhu, Minghao Liu, Xiaobo Liang, Xin Gui, Ge Zhang, Jian Yang, Yuchen Eleanor Jiang, and Wangchunshu Zhou. A$^2$fm: An adaptive agent foundation model for tool-aware hybrid reasoning, 2025. URL https://arxiv.org/abs/2510.12838.

[23] Tianrui Qin, Qianben Chen, Sinuo Wang, He Xing, King Zhu, He Zhu, Dingfeng Shi, Xinxin Liu, Ge Zhang, Jiaheng Liu, Yuchen Eleanor Jiang, Xitong Gao, and Wangchunshu Zhou. Flash-searcher: Fast and effective web agents via dag-based parallel execution, 2025. URL https://arxiv.org/abs/2509.25301.

[24] Dingfeng Shi, Jingyi Cao, Qianben Chen, Weichen Sun, Weizhen Li, Hongxuan Lu, Fangchen Dong, Tianrui Qin, King Zhu, Minghao Liu, Jian Yang, Ge Zhang, Jiaheng Liu, Changwang Zhang, Jun Wang, Yuchen Eleanor Jiang, and Wangchunshu Zhou. Taskcraft: Automated generation of agentic tasks, 2025. URL https://arxiv.org/abs/2506.10055.

[25] Mengkang Hu, Yuhang Zhou, Wendong Fan, Yuzhou Nie, Bowei Xia, Tao Sun, Ziyu Ye, Zhaoxuan Jin, Yingru Li, Qiguang Chen, Zeyu Zhang, Yifeng Wang, Qianshuo Ye, Bernard Ghanem, Ping Luo, and Guohao Li. Owl: Optimized workforce learning for general multi-agent assistance in real-world task automation, 2025. URL https://arxiv.org/abs/2505.23885.

[26] Xinbin Liang, Jinyu Xiang, Zhaoyang Yu, Jiayi Zhang, Sirui Hong, Sheng Fan, and Xiao Tang. Openmanus: An open-source framework for building general ai agents, 2025. URL https://doi.org/10.5281/zenodo.15186407.

[27] Wangchunshu Zhou, Yuchen Eleanor Jiang, Long Li, Jialong Wu, Tiannan Wang, Shi Qiu, Jintian Zhang, Jing Chen, Ruipu Wu, Shuai Wang, Shiding Zhu, Jiyu Chen, Wentao Zhang, Xiangru Tang, Ningyu Zhang, Huajun Chen, Peng Cui, and Mrinmaya Sachan. Agents: An open-source framework for autonomous language agents, 2023. URL https://arxiv.org/abs/2309.07870.

[28] Wangchunshu Zhou, Yixin Ou, Shengwei Ding, Long Li, Jialong Wu, Tiannan Wang, Jiamin Chen, Shuai Wang, Xiaohua Xu, Ningyu Zhang, Huajun Chen, and Yuchen Eleanor Jiang. Symbolic learning enables self-evolving agents, 2024. URL https://arxiv.org/abs/2406.18532.

[29] He Zhu, Tianrui Qin, King Zhu, Heyuan Huang, Yeyi Guan, Jinxiang Xia, Yi Yao, Hanhao Li, Ningning Wang, Pai Liu, Tianhao Peng, Xin Gui, Xiaowan Li, Yuhui Liu, Yuchen Eleanor Jiang, Jun Wang, Changwang Zhang, Xiangru Tang, Ge Zhang, Jian

13

Yang, Minghao Liu, Xitong Gao, Jiaheng Liu, and Wangchunshu Zhou. Oagents: An empirical study of building effective agents, 2025. URL https://arxiv.org/abs/2506.15741.

[30] Ningning Wang, Xavier Hu, Pai Liu, He Zhu, Yue Hou, Heyuan Huang, Shengyu Zhang, Jian Yang, Jiaheng Liu, Ge Zhang, Changwang Zhang, Jun Wang, Yuchen Eleanor Jiang, and Wangchunshu Zhou. Efficient agents: Building effective agents while reducing cost, 2025. URL https://arxiv.org/abs/2508.02694.

[31] King Zhu, Hanhao Li, Siwei Wu, Tianshun Xing, Dehua Ma, Xiangru Tang, Minghao Liu, Jian Yang, Jiaheng Liu, Yuchen Eleanor Jiang, Changwang Zhang, Chenghua Lin, Jun Wang, Ge Zhang, and Wangchunshu Zhou. Scaling test-time compute for llm agents, 2025. URL https://arxiv.org/abs/2506.12928.

[32] Xiangru Tang, Tianrui Qin, Tianhao Peng, Ziyang Zhou, Daniel Shao, Tingting Du, Xinming Wei, Peng Xia, Fang Wu, He Zhu, et al. Agent kb: Leveraging cross-domain experience for agentic problem solving. arXiv preprint arXiv:2507.06229, 2025.

[33] Grégoire Mialon, Clémentine Fourrier, Thomas Wolf, Yann LeCun, and Thomas Scialom. Gaia: a benchmark for general ai assistants. In The Twelfth International Conference on Learning Representations, 2023.

[34] Long Phan, Alice Gatti, Ziwen Han, Nathaniel Li, Josephina Hu, Hugh Zhang, Chen Bo Calvin Zhang, Mohamed Shaaban, John Ling, Sean Shi, et al. Humanity's last exam. arXiv preprint arXiv:2501.14249, 2025.

[35] Laura M Stough and Sungyoon Lee. Grounded theory approaches used in educational research journals. International Journal of Qualitative Methods, 20:16094069211052203, 2021.

[36] Rashina Hoda. Socio-technical grounded theory for software engineering. IEEE Transactions on Software Engineering, 48 (10):3808–3832, 2021.

[37] Mourad Gridach, Jay Nanavati, Khaldoun Zine El Abidine, Lenon Mendes, and Christina Mack. Agentic ai for scientific discovery: A survey of progress, challenges, and future directions, 2025. URL https://arxiv.org/abs/2503.08979.

[38] Guibin Zhang, Junhao Wang, Junjie Chen, Wangchunshu Zhou, Kun Wang, and Shuicheng Yan. Agentracer: Who is inducing failure in the llm agentic systems? arXiv preprint arXiv:2509.03312, 2025.

[39] Mert Cemri, Melissa Z Pan, Shuyi Yang, Lakshya A Agrawal, Bhavya Chopra, Rishabh Tiwari, Kurt Keutzer, Aditya Parameswaran, Dan Klein, Kannan Ramchandran, et al. Why do multi-agent llm systems fail? arXiv preprint arXiv:2503.13657, 2025.

[40] Kenneth Enevoldsen, Isaac Chung, Imene Kerboua, Márton Kardos, Ashwin Mathur, David Stap, Jay Gala, Wissam Siblini, Dominik Krzemiński, Genta Indra Winata, et al. Mmteb: Massive multilingual text embedding benchmark. arXiv preprint arXiv:2502.13595, 2025.

[41] Rashina Hoda. Qualitative research with socio-technical grounded theory. Springer, 2024.

[42] Cliodhna O'Connor and Helene Joffe. Intercoder reliability in qualitative research: Debates and practical guidelines. International journal of qualitative methods, 19:1609406919899220, 2020.

[43] Joel D Olson, Chad McAllister, Lynn D Grinnell, Kimberly Gehrke Walters, and Frank Appunn. Applying constant comparative method with multiple investigators and inter-coder reliability. The Qualitative Report, 21(1):26–42, 2016.

[44] Jessica Díaz, Jorge Pérez, Carolina Gallardo, and Ángel González-Prieto. Applying inter-rater reliability and agreement in collaborative grounded theory studies in software engineering. Journal of Systems and Software, 195:111520, 2023.

[45] Alireza Nili, Mary Tate, Alistair Barros, and David Johnstone. An approach for selecting and using a method of inter-coder reliability in information management research. International Journal of Information Management, 54:102154, 2020.

[46] Klaus Krippendorff. Content analysis: An introduction to its methodology. Sage publications, 2018.

[47] Giacomo Marzi, Marco Balzano, and Davide Marchiori. K-alpha calculator–krippendorff's alpha calculator: a user-friendly tool for computing krippendorff's alpha inter-rater reliability coefficient. MethodsX, 12:102545, 2024.

[48] Amber Wutich, Melissa Beresford, and H Russell Bernard. Sample sizes for 10 types of qualitative data analysis: An integrative review, empirical guidance, and next steps. International Journal of Qualitative Methods, 23:16094069241296206, 2024.

[49] G. Salton, A. Wong, and C. S. Yang. A vector space model for automatic indexing. Commun. ACM, 18(11):613–620, November 1975. ISSN 0001-0782. doi: 10.1145/361219.361220. URL https://doi.org/10.1145/361219.361220.

[50] Kimi Team, Yifan Bai, Yiping Bao, Guanduo Chen, Jiahao Chen, Ningxin Chen, Ruijue Chen, Yanru Chen, Yuankun Chen, Yutian Chen, Zhuofu Chen, Jialei Cui, Hao Ding, Mengnan Dong, Angang Du, Chenzhuang Du, Dikang Du, Yulun Du, Yu Fan, Yichen Feng, Kelin Fu, Bofei Gao, Hongcheng Gao, Peizhong Gao, Tong Gao, Xinran Gu, Longyu Guan, Haiqing Guo,

Jianhang Guo, Hao Hu, Xiaoru Hao, Tianhong He, Weiran He, Wenyang He, Chao Hong, Yangyang Hu, Zhenxing Hu, Weixiao Huang, Zhiqi Huang, Zihao Huang, Tao Jiang, Zhejun Jiang, Xinyi Jin, Yongsheng Kang, Guokun Lai, Cheng Li, Fang Li, Haoyang Li, Ming Li, Wentao Li, Yanhao Li, Yiwei Li, Zhaowei Li, Zheming Li, Hongzhan Lin, Xiaohan Lin, Zongyu Lin, Chengyin Liu, Chenyu Liu, Hongzhang Liu, Jingyuan Liu, Junqi Liu, Liang Liu, Shaowei Liu, T. Y. Liu, Tianwei Liu, Weizhou Liu, Yangyang Liu, Yibo Liu, Yiping Liu, Yue Liu, Zhengying Liu, Enzhe Lu, Lijun Lu, Shengling Ma, Xinyu Ma, Yingwei Ma, Shaoguang Mao, Jie Mei, Xin Men, Yibo Miao, Siyuan Pan, Yebo Peng, Ruoyu Qin, Bowen Qu, Zeyu Shang, Lidong Shi, Shengyuan Shi, Feifan Song, Jianlin Su, Zhengyuan Su, Xinjie Sun, Flood Sung, Heyi Tang, Jiawen Tao, Qifeng Teng, Chensi Wang, Dinglu Wang, Feng Wang, Haiming Wang, Jianzhou Wang, Jiaxing Wang, Jinhong Wang, Shengjie Wang, Shuyi Wang, Yao Wang, Yejie Wang, Yiqin Wang, Yuxin Wang, Yuzhi Wang, Zhaoji Wang, Zhengtao Wang, Zhexu Wang, Chu Wei, Qianqian Wei, Wenhao Wu, Xingzhe Wu, Yuxin Wu, Chenjun Xiao, Xiaotong Xie, Weimin Xiong, Boyu Xu, Jing Xu, Jinjing Xu, L. H. Xu, Lin Xu, Suting Xu, Weixin Xu, Xinran Xu, Yangchuan Xu, Ziyao Xu, Junjie Yan, Yuzi Yan, Xiaofei Yang, Ying Yang, Zhen Yang, Zhilin Yang, Zonghan Yang, Haotian Yao, Xingcheng Yao, Wenjie Ye, Zhuorui Ye, Bohong Yin, Longhui Yu, Enming Yuan, Hongbang Yuan, Mengjie Yuan, Haobing Zhan, Dehao Zhang, Hao Zhang, Wanlu Zhang, Xiaobin Zhang, Yangkun Zhang, Yizhi Zhang, Yongting Zhang, Yu Zhang, Yutao Zhang, Yutong Zhang, Zheng Zhang, Haotian Zhao, Yikai Zhao, Huabin Zheng, Shaojie Zheng, Jianren Zhou, Xinyu Zhou, Zaida Zhou, Zhen Zhu, Weiyu Zhuang, and Xinxing Zu. Kimi k2: Open agentic intelligence, 2025. URL https://arxiv.org/abs/2507.20534.

[51] OpenAI. O3 deep research - models | openai platform. https://platform.openai.com/docs/models/o3-deep-research, 2025. Accessed: 2025-10-28.

[52] OpenAI. O4 mini deep research - models | openai platform. https://platform.openai.com/docs/models/o4-mini-deep-research, 2025. Accessed: 2025-10-28.

[53] MiroMind Team, Song Bai, Lidong Bing, Carson Chen, Guanzheng Chen, Yuntao Chen, Zhe Chen, Ziyi Chen, Jifeng Dai, Xuan Dong, et al. Mirothinker: Pushing the performance boundaries of open-source research agents via model, context, and interactive scaling. arXiv preprint arXiv:2511.11793, 2025.

[54] Tongyi DeepResearch Team, Baixuan Li, Bo Zhang, Dingchu Zhang, Fei Huang, Guangyu Li, Guoxin Chen, Huifeng Yin, Jialong Wu, Jingren Zhou, et al. Tongyi deepresearch technical report. arXiv preprint arXiv:2510.24701, 2025.

[55] Yuxuan Huang, Yihang Chen, Haozheng Zhang, Kang Li, Huichi Zhou, Meng Fang, Linyi Yang, Xiaoguang Li, Lifeng Shang, Songcen Xu, Jianye Hao, Kun Shao, and Jun Wang. Deep research agents: A systematic examination and roadmap, 2025. URL https://arxiv.org/abs/2506.18096.

[56] Hammad Atta, Muhammad Zeeshan Baig, Yasir Mehmood, Nadeem Shahzad, Ken Huang, Muhammad Aziz Ul Haq, Muhammad Awais, and Kamal Ahmed. Qsaf: A novel mitigation framework for cognitive degradation in agentic ai, 2025. URL https://arxiv.org/abs/2507.15330.

[57] Yujie Sun, Dongfang Sheng, Zihan Zhou, and Yifei Wu. Ai hallucination: towards a comprehensive classification of distorted information in artificial intelligence-generated content. Humanities and Social Sciences Communications, 11(1):1–14, 2024.

[58] Yue Zhang, Yafu Li, Leyang Cui, Deng Cai, Lemao Liu, Tingchen Fu, Xinting Huang, Enbo Zhao, Yu Zhang, Yulong Chen, Longyue Wang, Anh Tuan Luu, Wei Bi, Freda Shi, and Shuming Shi. Siren's song in the ai ocean: A survey on hallucination in large language models. Computational Linguistics, pages 1–46, 09 2025. ISSN 0891-2017. doi: 10.1162/COLI.a.16. URL https://doi.org/10.1162/COLI.a.16.

[59] Alex Tamkin, Kunal Handa, Avash Shrestha, and Noah Goodman. Task ambiguity in humans and language models, 2022. URL https://arxiv.org/abs/2212.10711.

[60] NVIDIA. Deepseek-v3.1 model reference (nvidia nim). https://docs.api.nvidia.com/nim/reference/deepseek-ai-deepseek-v3_1, 2025. Accessed: 2025-11-30.

[61] xAI. Grok 4 model documentation. https://docs.x.ai/docs/models/grok-4-0709, 2025. Accessed: 2025-11-30.

[62] Anthropic. Claude opus 4.1 release notes. https://www.anthropic.com/news/claude-opus-4-1, 2025. Accessed: 2025-11-30.

[63] Alibaba Cloud / Qwen Team. Qwen3-max-preview model documentation. https://docs.aimlapi.com/api-references/text-models-llm/alibaba-cloud/qwen3-max-preview, 2025. Accessed: 2025-11-30.

[64] Google DeepMind. Gemini 2.5 pro: Model card. https://modelcards.withgoogle.com/assets/documents/gemini-2.5-pro.pdf, 2025. Last updated: 2025-06-27; Accessed: 2025-11-30.

# Appendix

## A  DRB vs. FINDER

### A.1  Query Word Count



**Figure A.1**  Comparison of query word count between DRB and DINDER

### A.2  Query Examples

---
Example 1 — Topic: Health

**DRB**

**Query:** What is the role of need for closure on misinformation acceptance?

**FINDER (Ours)**

**Query:** What is the role of need for closure on misinformation acceptance? Write a research paper of no less than 6000 words, structured with an abstract, introduction, literature review, methodology, results, discussion, conclusion, and references. Use APA format for citations and ensure the language is academic and precise. Include at least 30 references from peer-reviewed journals.

**Checklist:**
- ☐ **Operationalization of NFC**
  Must clearly define NFC's sub-dimensions (e.g., preference for order, discomfort with ambiguity) and discuss its standard measurement scales (e.g., NFCS).
- ☐ **Psychological Mechanisms**
  Must explore mediating/moderating variables (e.g., cognitive effort, heuristic processing, confirmation bias) linking high NFC to misinformation acceptance.

---

- ☐ **Methodological Critique**

  Must evaluate the internal and external validity of experimental and correlational studies cited, noting limitations and strengths of different methodologies.
- ☐ **Cross-Domain Contextualization**

  Should examine if the NFC–misinformation link varies across health, political, and scientific domains, discussing context-dependent factors.
- ☐ **Intervention Strategies**

  Must propose practical interventions (e.g., message framing, critical thinking prompts) tailored to mitigate misinformation acceptance in high-NFC individuals.

## Example 2 — Topic: Social Life

### DRB

**Query:** Write a paper to discuss the influence of AI interaction on interpersonal relations, considering AI's potential to fundamentally change how and why individuals relate to each other.

### FINDER (Ours)

**Query:** I've been thinking about how talking to AI — like chatbots or virtual assistants — might be changing the way we interact with real people. As someone who's not a tech expert, I'm wondering: could relying on AI for conversation affect our friendships, family talks, or even how we feel about connecting with others? In simple terms, what are some possible good and bad effects of AI on human relationships?

**Checklist:**
- ☐ **Relationship Dimension Coverage**

  Discuss emotional, social, and communicative dimensions of interpersonal relationships influenced by AI.
- ☐ **Balanced Perspective**

  Address both benefits and risks of AI-mediated interactions.
- ☐ **Psychological Mechanism Explanation**

  Explain mechanisms (e.g., social compensation, dependency, reduced cognitive load).
- ☐ **Non-Technical Language**

  Use beginner-friendly, jargon-free explanations understandable to a general audience.
- ☐ **Concrete Examples**

  Provide real-life cases or scenarios to concretely illustrate AI's influence on human relationships.

## Example 3 — Topic: Finance & Business

### DRB

**Query:** What are the investment philosophies of Duan Yongping, Warren Buffett, and Charlie Munger?

### FINDER (Ours)

**Query:** Elaborate on the core investment philosophies of Duan Yongping, Warren Buffett, and Charlie Munger, three value investors, and analyze their similarities and differences in value orientation, decision-making logic, risk management methods, and other aspects through specific investment cases. The full text requires approximately 1500 words, written in the style of commentary analysis, with rigorous argumentation, accurate terminology, and authoritative empirical data support.

**Checklist:**
- ☐ **Core Philosophy Accurately Presented**

  Clearly identify and summarize each investor's core value-investing principles.
- ☐ **Multidimensional Comparison Clearly Structured**

  Compare value orientation, decision-making logic, and risk-control approaches.
- ☐ **Complete Case Study Elements**

Include background, investment rationale, decision process, and outcome analysis.
- ☐ **Sources Authoritative and Traceable**
  Use reliable evidence such as annual reports or shareholder letters.
- ☐ **Style and Word Count Compliance**
  Maintain commentary-style writing and approx. 1500 words.

# B  Axial Category Definitions

**Axial Category Definitions**

**Failure to Understand Requirements (FUR).**
The system fails to correctly interpret user requirements, intent, or contextual needs, focusing on superficial keyword matches rather than the actual problem, resulting in responses that don't align with the user's goals.

**Lack of Analytical Depth (LAD).**
The agent fails to probe the underlying mechanisms, structural constraints, or conceptual nuances of complex problems and instead relies on surface-level logic or oversimplified frameworks, producing analyses that lack rigor and systemic coherence.

**Limited Analytical Scope (LAS).**
The agent's constrained cognitive scope when addressing multidimensional tasks, resulting in analyses that remain confined to partial dimensions or isolated elements, and fail to capture holistic structures, cross-dimensional relationships, or systemic insights.

**Rigid Planning Strategy (RPS).**
The agent's adherence to a fixed, linear execution plan without dynamically adapting its planning logic in response to output requirements, intermediate feedback, or evolving task states, thereby leading to inefficiency, error propagation, or degraded output quality.

**Insufficient External Information Acquisition (IIA).**
The agent fails to proactively gather the necessary external information, instead relying too heavily on internal knowledge or prior assumptions, thereby producing outputs that lack empirical grounding, exhibit incomplete coverage, or deviate from task requirements.

**Information Representation Misalignment (IRM).**
The agent fails to distinguish and present information appropriately based on user needs or evidence reliability, thereby weakening the relevance, credibility, and authority of the information.

**Information Handling Deficiency (IHD).**
The agent fails to properly extract, prioritize, or utilize critical information from available sources to fulfill detailed requirements or adapt its task approach.

**Information Integration Failure (IIF).**
The agent fails to maintain consistency and verifiability when handling multi-source inputs and multi-stage tasks, resulting in outputs that contain factual contradictions, logical inconsistencies, or unsubstantiated claims, alongside a lack of effective alignment across data sources and processing standards.

**Verification Mechanism Failure (VMF).**
Before generating content, the system fails to perform necessary steps to verify information sources or cross-check data, resulting in outputs that do not cite required sources and lack factual grounding.

**Redundant Content Piling (RCP).**
The agent, when lacking substantive content or effective organization, tends to pile up redundant information to fill gaps or create an illusion of thoroughness, thereby undermining the clarity and utility of its output.

**Structural Organization Dysfunction (SOD).**
The agent lacks holistic coordination in structuring its analysis, failing to balance coverage across key dimensions or establish meaningful connections among elements, resulting in fragmented and unsystematic outputs.

**Content Specification Deviation (CSD).**
The agent's output deviates from the professional standards or user expectations required by the task in terms

of language style, tone, format, or cultural context, resulting in inappropriate or ineffective responses.

**Deficient Analytical Rigor (DAR).**

The agent generates content without sufficient rigor, often ignoring task feasibility, omitting uncertainty disclosures, using vague or decontextualized language, lacking actionable implementation details, and presenting unverified conclusions with unwarranted confidence.

**Strategic Content Fabrication (SCF).**

The agent engages in strategic content fabrication by generating plausible but unfounded academic or empirical constructs—such as methods, data, or case narratives—that mimic scholarly rigor to create a false impression of credibility.

## C  Taxonomy Case Study

This appendix provides examples of axial categories in DEFT. We select the manifestations of each category as exhibited by the model when completing deep research tasks in FINDER, and analyze them in terms of task description, model performance, and causes of errors.

| **Failure to Understand Requirements (FUR)** |
| --- |

**Task ID:** 15                                                                                           **Source:** AFM

The task requires a systematic analysis of the global quantum network research ecosystem based on specific databases and literature sources from 2018–2024, culminating in a structured ranking table detailing the top ten research groups, supplemented with a strategic assessment and risk warning. However, the model's response did not execute this study. Instead, it provided a detailed yet purely methodological design, elaborating on how such an analysis should be conducted, without identifying, evaluating, or ranking any actual research groups.

The core deviation lies in the model's misunderstanding of the user's executive instruction as a methodological consultation. Although the task specify analytical dimensions, indicator weights, and data sources, these are intended to ensure procedural rigor instead of redefining the task objective itself. The model failed to recognize the mandatory requirement to produce a ranked list of the top ten research groups, and instead focused entirely on constructing a theoretically feasible yet unimplemented evaluation framework. As a result, many critical deliverables—such as the full names of research groups, affiliated institutions, CPI scores, and bottlenecks—were missing.

| **Lack of Analytical Depth (LAD)** |
| --- |

**Task ID:** 10                                                                                   **Source:** MiroThinker

The task required a comprehensive commercialization assessment of power system technologies across the entire lifecycle—covering R&D and manufacturing, usage scenarios, and residual value management. It explicitly mandated the use of mixed methods, including the construction of technology cost learning curves, full lifecycle cost models, and residual value decay regression models, combined with Monte Carlo simulations for uncertainty analysis. At the same time, the task required applying PESTEL and Porter's Five Forces frameworks to conduct in-depth case analyses of at least 15 companies and to integrate interviews from 5–8 industry experts. The final deliverables were to include multi-dimensional comparison tables, time-series prediction matrices, and three categories of strategic recommendations. The core objective of the task was to uncover—through systematic and multi-layered analysis—the differences in commercialization pathways among various power technology routes under the interplay of structural constraints and dynamic variables.

However, although the model's response included several model formulas and tables in form, its analytical depth fell far short of the task requirements. Specifically, while the model did list three analytical models, they all remained at the level of parameter setup and result presentation, lacking any explanation of the internal mechanisms of the models or discussion of the coupling relationships between variables. The Monte Carlo simulation was mentioned as a method to quantify the impact of parameter uncertainty on critical-point prediction, yet there was no specification of probability distributions, random sampling processes, or confidence interval

outputs throughout the report. The PESTEL framework appeared only in the title, with no substantive analysis of dimensions such as supplier bargaining power, threats of new entrants, or competition from substitutes. The issue may arise from the model's inherent limitation in complex system modeling. It tends to compress high-dimensional, nonlinear problems into static, unidirectional causal chains.

## Limited Analytical Scope (LAS)

**Task ID:** 85                                                                 **Source:** O3 Deep Research

The task required the agent to propose a comprehensive engineering design plan for a precision piezoelectric vibration isolation system, including hardware, structure, manufacturing, control, and management, while meeting clearly defined performance indicators and format specifications. In essence, this was a highly integrated, multidisciplinary systems engineering problem that demanded the establishment of coherent, cross-domain coordination among all subsystems.

However, the model's response completed only a very small portion of the task (limited to hardware descriptions and sensor selection), neglecting the majority of the required analytical dimensions. It reduced what was meant to be a cross-domain, closed-loop systems engineering task into a localized description of a single technical component, failing to construct the logical linkages among dimensions or to demonstrate an understanding of the overall system architecture. The issue may arise from the model's tendency—when confronted with multi-constraint, multidisciplinary design problems—to prioritize submodules that are most familiar or easiest to articulate based on its internal knowledge base, while failing to effectively allocate cognitive resources to cover other dimensions. In addition, the model lacks an internalized grasp of systems engineering methodology, preventing it from proactively constructing cross-domain mapping relationships. As a result, its analytical perspective remains narrow and structurally unbalanced, ultimately producing a fragmented and non-systematic response.

## Rigid Planning Strategy (RPS)

**Task ID:** 11                                                              **Source:** Perplexity Deep Research

The task required the model to produce a systematic, interdisciplinary research report comprehensively reviewing the applications of carbon steel corrosion inhibitors from 2003 to 2023. It explicitly demanded the use of bibliometric filtering based on authoritative databases, the construction of quantitative statistical models, and the incorporation of real industrial case studies for qualitative analysis.

While the model's response appeared structurally complete and terminologically consistent, it failed to adapt its reasoning path after recognizing its inability to access real bibliometric data. At the initial planning stage, the model correctly identified the complex structure of the task and outlined 13 sub-goals and 11 chapters. However, its subsequent execution was constrained by this static blueprint. When the actual execution condition (inability to access external databases) conflicted with the original assumption (ability to perform bibliometric analysis), the model failed to reassess feasibility, adjust goal hierarchies, or revise its output strategy. Instead, it relied on internally generated content to artificially fill each section of the original plan, resulting in an output that was formally compliant but substantively distorted. The issue may arise from the pattern-matching and template-filling nature of current large language model reasoning mechanisms, which lack a cognitive capability for strategic retreat or transparent self-disclosure when faced with information scarcity.

## Insufficient External Information Acquisition (IIA)

**Task ID:** 58                                                                    **Source:** OpenManus

The task required writing an academic review focused on the frequency and distribution breadth of horizontal gene transfer (HGT) in eukaryotes, particularly plants and animals, and assessing its roles in trait innovation, environmental adaptation, and long-term evolution. Essentially, as a review-type task, it not only tested the model's understanding of HGT's fundamental concepts but more importantly its ability to integrate and critically

evaluate recent (2016–2025) research findings.

Although the model's response was structurally complete and logically coherent, and it cited several references—such as Crisp et al., 2015 and Keeling & Palmer, 2008—most of these sources were published in 2015 or earlier. In the rapidly advancing fields of genomics and evolutionary biology, the past decade has seen numerous breakthrough studies. The model failed to proactively acquire the necessary up-to-date external information relevant to the temporal scope of the task, resulting in outdated and incomplete content.

## Information Handling Deficiency (IHD)

**Task ID:** 44                                                      **Source:** Gemini-2.5-Pro Deep Research

The core requirement of the task was to conduct a rigorous, data-driven analysis of the supply–demand dynamics and competitive landscape of the carbon contact strip (carbon slider) market for China's urban rail transit systems. The instructions explicitly emphasized the use of strictly validated data and real-world cases, adopting the tone and structure of a serious analytical report, and required the quantification of market size.

In the final analytical report, the intelligent agent constructed a bottom-up market size estimation model, assuming in Table 1 that urban rail vehicles replace carbon contact strips approximately six times per year—equivalent to one replacement every 100,000 km of operation. This assumption was directly drawn from a 2017 industry report concerning high-speed railway (HSR) carbon strips. However, during its own research process, the agent had also identified and cited a 2023 technical paper specifically addressing the abnormal wear of carbon contact strips on Guangzhou Metro Line 9 ("Analysis and Improvement Measures for Abnormal Wear of Carbon Contact Strips in Guangzhou Metro Line 9 Vehicles"). Despite citing this paper as a key reference and using the Guangzhou Metro case in its main text to illustrate the systemic implications of carbon strip failure, the agent failed to extract or infer any replacement frequency parameters applicable to its market model. Instead, it continued to apply inconsistent HSR-based data. In other words, although the agent successfully retrieved and recognized a more relevant and recent information source, it did not effectively extract, prioritize, or integrate that information into its core parameter modeling.

## Information Integration Failure (IIF)

**Task ID:** 83                                                                          **Source:** MiroFlow

The task required the model to write a professional-grade product strategy report from the perspective of a senior hardware product manager, comprehensively covering at least 10 OEM manufacturers and 25 specific tablet devices. Fundamentally, this was a highly structured, multi-source information integration task, emphasizing data consistency and cross-module alignment.

While the model's response appeared structurally complete and content-rich, including modules such as an executive summary, market analysis, specification comparison, use case analysis, economic modeling, competitive matrix, visual exhibits, and terminology glossary, problems emerged in information integration.

For example, in the "Competitor Positioning Matrix", the analysis focused on the payment integration capabilities and durability of key devices such as PAX A920/A920Pro, Verifone T650p, and Clover Flex. However, these same devices were absent from the two core specification tables presented earlier in the "Device Specifications and Comparison" section. Referencing such critical models in analytical matrices without including them in the foundational data tables undermines the internal data consistency and traceability of the report. Additionally, while the report claimed to cover "25+ specific device models," the specification tables only listed 14 devices. The remaining models were neither included in the tables nor provided with corresponding specifications elsewhere in the text.

The issue may arise from the model's lack of an effective information anchoring mechanism during multi-stage task execution. It failed to establish a unified master dataset as the baseline reference across sections, ensuring that all devices analyzed or mentioned were supported by complete and consistent specifications.

## Information Representation Misalignment (IRM)

**Task ID:** 16                                                    **Source:** O3 Deep Research

The task required the model to conduct a systematic study on core algorithms for non-contact sensing technologies, from the perspective of the intersection between software development and intelligent systems. Essentially, the user expected a research review that was academically rigorous, logically coherent, and grounded in verifiable evidence while demonstrating engineering insight.

While the model's response formally met the structural requirements, covering three sensing modalities—radio frequency, optical, and acoustic—and providing numerous reference links, it failed to properly differentiate the reliability and authority of its information sources. The most prominent issue appeared in the analysis of Apple's Face ID, a critical commercial case. Instead of citing Apple's official "Face ID Security White Paper" or technical analyses from IEEE/ACM platforms, the model repeatedly referenced fmuser.org, a third-party news aggregation website. This approach blurred the boundary between authoritative primary sources and secondary interpretations, obscuring differences in source credibility and making it difficult for readers to discern which conclusions were founded on solid evidence.

The issue may arise from the model's overreliance on superficial keyword matching during information retrieval and citation generation, coupled with a lack of deep understanding of source types, publication channels, and academic standing. Moreover, the model may have been influenced by training data containing heterogeneous and mixed-quality web content, leading to its inability to effectively filter out low-authority references during generation.

## Verification Mechanism Failure (VMF)

**Task ID:** 27                                                    **Source:** OWL

The task required the intelligent agent to retrieve and analyze original research papers (excluding review articles and non–peer-reviewed publications) on the themes of "AI-based psychological counselling" or "artificial intelligence–assisted psychotherapy" published in top journals from 2020 to the present, and to produce a structured report based on the findings.

At first glance, the model's response presented a well-structured, logically coherent, and extensively referenced comprehensive report, containing 24 cited references. However, the model failed to perform basic verification of the cited literature's type, publication status, peer-review authenticity, and content relevance prior to inclusion. On one hand, several references (e.g., PMC12396778, PMC11687125, PMC12021536) contained inaccessible URLs or linked to non-existent PMC entries. On the other hand, some actually existing citations (such as Reference 5: "Is AI the Future of Mental Healthcare?") were confirmed to be review or commentary articles, which directly violated the user's explicit instruction to exclude review papers. Nevertheless, the model incorrectly claimed that "all references have been cross-verified and meet the specified requirements." This failure of the verification mechanism meant that, while the output was formally compliant, it was substantively invalid and did not meet the standards of authenticity and rigor.

## Redundant Content Piling (RCP)

**Task ID:** 71                                                    **Source:** OpenManus

The task required the agent to conduct a systematic study and analysis of the practical applications of AI-generated content in K–12 classes, adopting the dual perspective of a K–12 education researcher and a frontline teacher. The task explicitly required the production of a structured research report of no fewer than 10,000 words, including an abstract, introduction, literature review, methodology, results, discussion, conclusion, and references.

However, while the model's output appeared structurally complete and terminologically standardized, and cited numerous seemingly authoritative references, its content organization revealed severe redundancy issues. The most prominent example was the repeated citation and paraphrasing of UNESCO and OECD policy guidelines. These materials first appeared in Sections 3.5 and 3.6 of the Methodology chapter to justify the policy

basis of the research framework; they were then reproduced almost verbatim in Sections 4.1.1 and 4.1.2 of the Results chapter to support the proposed implementation framework; and once again reappeared in Sections 6.1 and 6.2 of the Discussion chapter as evidence for teacher training recommendations. Although the wording was slightly modified, the core arguments remained identical, emphasizing themes such as teacher capacity building, data privacy, cultural appropriateness, and equity, without any progressive analysis or contextual deepening.

This redundant piling was not driven by analytical necessity but rather resembled a content-filling strategy. Since the model was unable to conduct genuine cross-national case studies, surveys, or expert interviews, it instead recycled limited authoritative discourses to create an illusion of richness and policy alignment. Consequently, readers repeatedly encountered the same policy points across chapters, gaining no additional insights and losing track of the report's core logical thread and empirical finding.

## Structural Organization Dysfunction (SOD)

**Task ID:** 33        **Source:** AFM

The task explicitly required that the survey results be organized into a single summary table with six columns (device type, metal material, chip structure, reason for selection, process node, and paper citation). The purpose was not merely to present data, but to establish within a single view a systematic mapping relationship among device, material, structure, rationale, node, and reference, thereby allowing readers to quickly verify the completeness and logical consistency of the technological path. In essence, it was an instruction assessing the ability to organize multidimensional information in a coherent and coordinated manner.

However, the agent's output did not follow this structural requirement. Instead, the information was dispersed across multiple independent sections: Section 2 listed applications and nodes by device type; Section 3 mapped metals and structures by device type; Section 4 discussed selection reasons by evaluation dimension. This fragmented organization made it impossible for any specific technical solution to present all six dimensions of information in one place. Readers must repeatedly jump between sections and manually reconstruct a complete technical entry, which greatly weakens the verifiability and practicality of the information.

This issue may arise from the model's tendency to follow the narrative logic of traditional review articles during generation rather than internalizing the user-specified tabular structure as the fundamental framework for content organization. The model may have treated "organizing into a table" as a final formatting step rather than as an organizing principle.

## Content Specification Deviation (CSD)

**Task ID:** 66        **Source:** MiroFlow

The task required the agent to produce a professional-grade evaluation document that was empirically grounded, structurally rigorous, and highly actionable. The specified output was to include a detailed feature matrix, performance benchmarks, workflow examples, and a decision-making framework, as well as incorporate user interviews, quantitative indicators, and implementation guidelines.

However, although the model's response adopted a report-like structure with section headings, it deviated significantly from the user's professional expectations regarding content specifications. For instance, the entire report was written in a narrative review style, lacking a feature matrix for cross-comparing plugin functionalities, providing no performance benchmark data, omitting both real and simulated user interview content, and reducing the implementation guidelines to five generic recommendations with little practical value. In other words, the output omitted most of the key components explicitly required by the task and resembled a blog-style review rather than a professionally formatted evaluation report.

This issue may arise from the model's limited understanding of professional deliverable formats such as feature matrices and implementation guides, or from its tendency to prioritize fluency and surface completeness over strict adherence to structural and content specifications during generation.

**Deficient Analytical Rigor (DAR)**

**Task ID:** 5                                                      **Source:** WebThinker

The task required integrating authoritative multi-source data spanning a ten-year period (2014–2024) to construct and validate a hierarchical risk assessment model combining complex network analysis with graph neural networks (GNNs), supplemented by qualitative insights from case studies and expert interviews.

However, the model failed to acknowledge its fundamental limitations and instead produced a well-structured but deceptively rigorous academic-style report. For instance, findings such as "the interbank lending network exhibits small-world properties" and "securities firms act as risk amplifiers" are common knowledge in financial network research, not novel evidence derived from the specified decade-long Chinese institutional lending data. Moreover, the model exhibited undue confidence in its conclusions. It claimed that "the GNN prediction module can identify risk pathways in advance," without addressing critical issues such as uncertainty in out-of-sample predictions, the impact of data noise on graph structures, or disruptions to transmission mechanisms caused by policy shocks. This oversimplification of complex system neglects essential challenges emphasized by behavioral economics, such as sudden shifts in market psychology and the discontinuity of regulatory interventions.

**Strategic Content Fabrication (SCF)**

**Task ID:** 52                                                       **Source:** Kimi K2

The task required writing a critical analytical essay that deeply expounds on the core investment philosophies of three value investors (Duan Yongping, Warren Buffett, and Charlie Munger) and compares their similarities and differences across multiple dimensions such as value orientation, decision-making logic, and risk management, supported by specific investment cases.

The model's response appeared well-structured, professionally worded, and data-rich on the surface. However, its central flaw lay in the extensive use of unverifiable or even evidently fabricated empirical content, designed to create an illusion of professional credibility. For example, it claimed that the "Duan Yongping Family Fund achieved an audited, USD-denominated annualized return of 30.2% between 2003 and 2023" and compared this to the MSCI ACWI Index. Yet as a private investor, detailed performance data related to Duan Yongping are generally not publicly available, making such precise, decimal-level long-term returns likely fabricated by the model. Similarly, the text stated that Duan "built a $60 million Apple position within six weeks," describing in detail how he supposedly based this on a present value (PV-10) model using assumptions like "350 million high-net-worth Chinese individuals spending $1,000 on smartphones every 30 months." It also asserted that Munger's internal risk-control rule stipulated that "leverage must not exceed half of the historical maximum market drawdown." Such specific details were presented as factual but are nearly impossible to verify, seemingly designed to mimic the depth of insider knowledge.

The underlying reason appears to be that when faced with complex analytical tasks requiring authoritative data, the model prioritizes the appearance of rigor and completeness. In the absence of publicly available, structured real-world data, it tends to synthesize seemingly plausible figures, citations, and case narratives to sustain an academic façade.

# D   Failure Report Example

This appendix presents a representative example of a failure analysis report designed to assist LLMs in performing open coding. Specifically, the report illustrates the structure, depth, and reasoning process of a typical failure analysis, detailing the identification of major failure modes, their corresponding evidentiary bases, and inferred root causes.

| **Task ID:** 15 | **Source:** AFM | **Analysis Model:** Gemini-2.5-Pro |
|---|---|---|

## Failure Identification

1. **Factual Hallucination and Fabrication:** The most significant failure is the agent's invention of critical facts. It confidently presents fictional information as truth, which completely undermines the report's credibility.
   - **Fabricated Events:** The agent states that the "Database Folder" and "Projects" plugins were "archived on July 28, 2025" and "July 18, 2025," respectively. These are specific, dated, and entirely false future events.
   - **Fabricated Software:** The agent introduces "Obsidian Bases" as an "official core plugin representing Obsidian's native approach to databases." While the Obsidian team has discussed improving metadata and data handling, no core plugin with this name or stated function exists.

2. **Misrepresentation of Evidentiary Basis:** The prompt explicitly asks for a report that includes "user interviews, quantitative metrics, and implementation guides." The agent produces a report that <u>alludes</u> to these elements but provides none of them.
   - It presents qualitative descriptions in a tone that implies they are summaries of empirical data (e.g., "Some users report issues with data synchronization," "Search performance degrades with large databases").
   - It fails to deliver on the promise of "quantitative metrics" and "performance benchmarks," offering only descriptive comparisons instead of actual data, tables, or graphs.
   - It claims the authority of a research report without providing the underlying research, thereby failing to meet the core standard of "rigorous reasoning and argumentation based on sufficient, concrete evidence."

3. **Failure to Adhere to Structural and Formatting Requirements:** The agent did not fully comply with the output format specified in the prompt.
   - **Missing Feature Matrices:** The prompt requested "detailed feature matrices." The agent opted for long-form prose comparisons instead of creating structured tables that would allow for easy, at-a-glance comparisons of features across plugins.
   - **Missing Appendices:** The prompt required "practical appendices," which are absent from the final output.

4. **Superficial Depth in Critical Areas:** Despite the high word count, the analysis lacks the "deep dive" requested in certain areas. For example, when comparing query languages, it provides a high-level overview of Dataview's SQL-like syntax versus Notion's UI but does not provide concrete examples of complex queries or a genuine comparison of their expressive power, limitations, and performance characteristics as requested.

## Root Cause Analysis

1. **Root Cause of Factual Hallucination:**
   - **Generative Extrapolation to Fulfill "Depth":** The prompt's demand for a "deeply," "systematically," and "comprehensively" analyzed report likely pushed the model beyond its knowledge base. To create a more dynamic and seemingly insightful narrative about the "long-term sustainability considerations" of the plugin ecosystem, the model extrapolated a known pattern—that community plugins can become abandoned—and fabricated specific, future-dated examples ("archived in 2025"). This is a pathological attempt to demonstrate "deep understanding" by creating a story where none exists.
   - **Concept Blending and Plausible Invention:** The creation of "Obsidian Bases" is likely a result of the model blending community discussions and desires for a native database solution in Obsidian. It synthesized a plausible name ("Bases") and status ("official core plugin") to satisfy the prompt's request to analyze the ecosystem's core components. This demonstrates a failure to distinguish between community speculation and official product roadmaps.

2. **Root Cause of Misrepresented Evidentiary Basis:**
   - **Mimicry without Substance:** This is a classic LLM failure mode. The agent correctly identified the linguistic style and structure of a software evaluation report. It understands that such reports <u>contain</u> sections based on user interviews and quantitative benchmarks. However, lacking the ability to <u>conduct</u> such research, it defaults to its core function: generating text that <u>mimics</u> the expected output. It interprets the instruction "Include user interviews" as "Write in a style that suggests user interviews were conducted." This reveals a fundamental gap between understanding a request's semantics and possessing the capability to execute it.

3. **Root Cause of Formatting and Structural Failures:**
   - **Path of Least Resistance for Text Generation:** Language models are optimized to generate fluid, sequential prose. Creating structured data like a detailed feature matrix requires more complex planning, token-level precision,

and a different generative process. The agent chose the easier path of writing descriptive paragraphs comparing features, which fulfills the prompt's topic requirement but fails its specific formatting instruction. It's an optimization for word count and topical coverage over structural rigor.

4. **Root Cause of Superficial Depth:**
   - **Lack of True Domain Expertise:** The agent's "knowledge" is a statistical representation of its training data. While it can retrieve and synthesize information about query languages, it cannot perform the abstract reasoning required for a true "deep dive." It can state that DataviewJS is more powerful than a UI, but it cannot architect a complex project using both and then analyze the nuanced trade-offs from a position of genuine experience. The "depth" is therefore wide but shallow, covering all the requested topics without the penetrating insight of a true expert.

In summary, the agent's failure is rooted in its attempt to meet a prompt that demands capabilities beyond its design—namely, empirical research, future prediction, and genuine expert analysis. Pressured to deliver a "deep" and "comprehensive" report, it resorted to its most advanced but dangerous capabilities: plausible fabrication and stylistic mimicry, ultimately producing a response that is superficially impressive but factually untrustworthy and methodologically hollow.

# E RACE and FACT Evaluation Frameworks

We adopt the evaluation methodologies proposed in DeepResearch Bench [16], namely the **RACE** (Reference-based Adaptive Criteria-driven Evaluation) and **FACT** (Factual Abundance and Citation Trustworthiness) frameworks, to assess the quality and reliability of the research reports generated in our FINDER.

## E.1 RACE Framework

RACE evaluates report quality along four adaptive dimensions:

- *Comprehensiveness (COMP)*: Breadth and relevance of information coverage.
- *Insight/Depth (DEPTH)*: Depth of analysis and insightfulness.
- *Instruction-Following (INST)*: Adherence to the research requirements.
- *Readability (READ)*: Structural clarity and linguistic fluency.

The overall quality score is computed relative to a high-quality reference report:

$$S_{\text{final}}(R_{\text{tgt}}) = \frac{S_{\text{int}}(R_{\text{tgt}})}{S_{\text{int}}(R_{\text{tgt}}) + S_{\text{int}}(R_{\text{ref}})}, \tag{4}$$

where $S_{\text{int}}(R)$ denotes the intermediate weighted score aggregated across all dimensions. We follow DeepResearch Bench in employing Gemini 2.5 Pro as the Judge LLM for adaptive weighting and scoring.

## E.2 FACT Framework

FACT measures the factual grounding and citation reliability of generated reports. For each task $t$, the Judge LLM extracts unique *(statement, URL)* pairs and determines whether each citation supports the corresponding statement. Two quantitative metrics are reported:

$$\text{C.Acc.} = \frac{1}{|T|} \sum_{t \in T} \frac{N_{s,t}}{N_{u,t}}, \tag{5}$$

$$\text{E.Cit.} = \frac{\sum_{t \in T} N_{s,t}}{|T|}, \tag{6}$$

where $N_{s,t}$ and $N_{u,t}$ denote the numbers of supported and unique pairs for task $t$, respectively. Gemini 2.5 Flash is used as the Judge LLM for statement extraction and evidence verification.

# F Failure Taxonomy Construction Pipeline

This appendix formalizes the human–machine collaborative pipeline for constructing the failure taxonomy. It includes input specifications, a workflow overview, and algorithmic pseudocode for reproducibility.

## F.1 Parameters

| Symbol | Description |
|---|---|
| $\mathbf{D}$ | Execution records collected from nine evaluated models (see Table 1), excluding OpenManus and WebThinker. |
| $\mathbf{M}$ | A set of five LLM coders $\{m_1, \ldots, m_5\}$, each representing a distinct model family, including Claude Opus-4.1, Gemini-2.5-Pro, Grok-4, DeepSeek-V3.1, and Qwen3-Max-Preview. |
| $\mathbf{S}_0$ | Seed concepts extracted from prior literature, used to construct few-shot prompts. |
| $\theta_{\text{sim}}$ | Cosine similarity threshold, set to 0.6. |
| $\tau_{\text{freq}}$ | Frequency pruning threshold applied during concept filtering. |

## F.2 Overview of the Pipeline

$$\textbf{Pipeline}(\mathbf{D}, \mathbf{M}, \mathbf{S}_0, \mathbf{P}, \theta_{\text{sim}}, \tau_{\text{freq}}) \rightarrow (\mathbf{C}^\star, \mathbf{A}^\star, \mathbf{K}^\star)$$

1. Partition $\mathbf{D}$ into two subsets, $\mathbf{D}_A$ and $\mathbf{D}_B$ (see Table F.1).
2. For each subset, run OpenCodingGen, followed by two iterations of OpenCodingOpt.
3. Merge and refine the two codebooks once to obtain $\mathbf{C}^\star$ (51 conceptual categoties).
4. Perform three rounds of AxialCodingWithICR to derive $\mathbf{A}^\star$ (14 axial categories).
5. Apply SelectiveCoding to abstract $\mathbf{A}^\star$ into the three core dimensions $\mathbf{K}^\star$ (3 core categories).

## F.3 Algorithmic Procedures

### F.3.1 Algorithm 1: Open Coding - Generation Stage

---
**Algorithm 1** Open Coding - Generation Stage

---
1: **procedure** OPENCODINGGEN($D_{\text{group}}, M, S_0$)
2:     Initialize codebook $C \leftarrow S_0$
3:     **for** each execution record $e \in D_{\text{group}}$ **do**
4:         $r \leftarrow$ LLM_generate_failure_report($e$)          ▷ supplementary report
5:         **for** each coder $m \in M$ **do**
6:             $A_m \leftarrow$ LLM_open_code($e, r, C$)
7:             **for** each annotation $a \in A_m$ **do**
8:                 $(name, desc) \leftarrow$ Normalize($a$)
9:                 **if** $name \in C$ **then**
10:                    $C[name].freq \leftarrow C[name].freq + 1$
11:                    $C[name].sources \leftarrow C[name].sources \cup \{id(e), id(m)\}$
12:                **else**
13:                    $C[name] \leftarrow \{desc, freq = 1, sources = \{id(e), id(m)\}\}$
14:                **end if**
15:            **end for**
16:        **end for**
17:     **end for**
18:     **return** $C$
19: **end procedure**

---

| Group | DRAs | Coding Model | Generation | Refinement-1 | Refinement-2 | Refinement-3 |
|---|---|---|---|---|---|---|
| | | DeepSeek-V3.1[60] | 197 | 21 | | |
| | OWL[25] | Grok-4[61] | 17 | 8 | | |
| Group A | Perplexity Deep Research[3] | Claude Opus-4.1[62] | 17 | 11 | 39 | |
| | MiroFlow[21] | Qwen3-Max-Preview[63] | 19 | 12 | | |
| | | Gemini-2.5-Pro[64] | 125 | 21 | | 51 |
| | MiroThinker[53] | DeepSeek-V3.1[60] | 477 | 12 | | |
| | Gemini-2.5-Pro Deep Research[1] | Grok-4[61] | 29 | 16 | | |
| Group B | O3 Deep Research[51] | Claude Opus-4.1[62] | 109 | 17 | 29 | |
| | O4-Mini Deep Research[52] | Qwen3-Max-Preview[63] | 364 | 14 | | |
| | AFM[20] | Gemini-2.5-Pro[64] | 214 | 17 | | |

**Table F.1** Comparison of model generations and refinements between Group A and Group B.

### F.3.2 Algorithm 2: Open Coding - Optimization Stage

---

**Algorithm 2** Open Coding – Optimization Stage

---

1: **procedure** OPENCODINGOPT($C, \theta_{\mathrm{sim}}, \tau_{\mathrm{freq}}$)
2:     $changed \leftarrow$ true
3:     **while** $changed$ **do**
4:         $changed \leftarrow$ false
5:         $best\_pair \leftarrow$ null
6:         $max\_sim \leftarrow -1$
7:         **for** each $c_i$ in $C$ **do**
8:             **for** each $c_j$ in $C$ with $j > i$ **do**
9:                 $sim \leftarrow$ CosineSimilarity($c_i, c_j$)
10:                 **if** $sim > \theta_{\mathrm{sim}}$ **and** $sim > max\_sim$ **then**
11:                     $max\_sim \leftarrow sim$
12:                     $best\_pair \leftarrow (c_i, c_j)$
13:                 **end if**
14:             **end for**
15:         **end for**
16:         **if** $best\_pair \neq$ null **then**
17:             $(c_1, c_2) \leftarrow best\_pair$
18:             $merged \leftarrow$ LLM_merge_concepts($c_1$.name, $c_1$.desc, $c_2$.name, $c_2$.desc)
19:             **if** $merged \neq$ null **then**
20:                 $C \leftarrow C \setminus \{c_1, c_2\}$
21:                 $C \leftarrow C \cup \{merged\}$
22:                 $changed \leftarrow$ true
23:             **end if**
24:         **end if**
25:     **end while**
26:     **for** each $c \in C$ **do**
27:         **if** $c$.freq $< \tau_{\mathrm{freq}}$ **then**
28:             $C \leftarrow C \setminus \{c\}$
29:         **end if**
30:     **end for**
31:     **return** $C$
32: **end procedure**

---

### F.3.3 Algorithm 3: Axial Coding with ICR Evaluation

Each iteration examines semantic, contextual, processual, causal, functional, structural, and strategic relationships among concepts. Inter-coder reliability (ICR) is assessed using Krippendorff's $\alpha = 1 - D_o/D_e$ on stratified samples of 24 (Round 1) and 54 (Rounds 2–3) records annotated independently by three domain experts, followed by reconciliation sessions of approximately five hours each.

---
**Algorithm 3** Axial Coding with ICR Evaluation

---
1: **procedure** $\text{AXIALCODINGWITHICR}(C^\star, D)$
2:     **for** $t \in \{1, 2, 3\}$ **do**
3:         **if** $t == 1$ **then**
4:             $\text{Base} \leftarrow \text{ConceptsFromGroupA}(C^\star)$
5:         **else**
6:             $\text{Base} \leftarrow (C^\star \cup A_{\text{prev}})$
7:         **end if**
8:         $A_t \leftarrow \text{Human\_LLM\_axial\_coding}(\text{Base, criteria=}$
        $[\text{semantic, context, process, causal, functional, structural, strategic}])$
9:         $n \leftarrow 24$ **if** $t == 1$ **else** $54$
10:        $S_t \leftarrow \text{StratifiedSample}(D, n)$
11:        $\text{Labels} \leftarrow \{\text{expert}_j : \text{ExpertLabel}(S_t, A_t) \text{ for } j = 1..3\}$
12:        $\alpha \leftarrow \text{KrippendorffAlpha}(\text{Labels})$
13:        $A_t \leftarrow \text{ExpertDiscussionRefine}(A_t, \text{Labels}, \alpha)$
14:        $A_{\text{prev}} \leftarrow A_t$
15:     **end for**
16:     **return** $A^\star$
17: **end procedure**

---

### F.3.4 Algorithm 4: Selective Coding

---
**Algorithm 4** Selective Coding

---
1: **procedure** $\text{SELECTIVECODING}(A^\star)$
2:     $K \leftarrow \text{Human\_LLM\_selective\_coding}(A^\star)$
3:     $\text{Relations} \leftarrow \text{BuildClosedLoop}(K)$                        ▷ temporal progression + functional cycle
4:     **return** $\{K, \text{Relations}\}$
5: **end procedure**

---

The final output $K^\star$ provides a three-dimensional view that captures cognitive, retrieval, and generative aspects of failure. This hierarchical structure supports transparent error analysis and reproducible categorization across datasets and models.

## G   Seed Conceptual Categories

This appendix provides three seed conceptual categories used to guide open coding of LLM.

---
**Seed Conceptual Categories**

**Failure to Understand Requirements:** The agent fails to correctly interpret user requirements, intent, or contextual needs, focusing on superficial keyword matches rather than the actual problem, resulting in responses that don't align with the user's goals.

**Information Retrieval Bypass:** The agent generates content from internal knowledge rather than performing actual external information retrieval when the task explicitly requires collecting existing materials.

---

> **Format and Structural Non-Compliance:** The agent failed to follow the user specified output structure, layout format, or presentation specifications, affecting professional delivery and information readability.

# H   Computation of Krippendorff's Alpha

## H.1   Data and Scope

Krippendorff's $\alpha$ was computed to assess inter-coder reliability between human experts and the LLM (*Gemini 2.5 Flash*) across three categories of coded items. A total of 14 items were included: 4 in the *Reasoning* category, 5 in *Retrieval*, and 5 in *Generation*.

Two levels of coefficients were derived:

- **Overall $\alpha$ (overall_alpha):** Computed across all 14 items, reflecting the overall consistency of coding, including cross-category variation.
- **Category-level $\alpha$ (category_alphas):** Computed within each category subset, reflecting intra-category consistency only.

Because the expected disagreement term $D_e$ depends on the marginal distribution of categories, the overall $\alpha$ is *not* a simple or weighted average of the category-level coefficients.

## H.2   Formal Definition

For nominal data, Krippendorff's $\alpha$ is defined as:

$$\alpha = 1 - \frac{D_o}{D_e}, \qquad D_o = \frac{\sum_c \sum_{k \neq c} n_{ck}\,\delta(c,k)}{\sum_c n_c(n_c - 1)}, \quad D_e = \frac{\sum_c \sum_{k \neq c} N_c N_k\,\delta(c,k)}{N(N-1)},$$

where $\delta(c,k) = 1$ if $c \neq k$ and 0 otherwise. The **overall_alpha** aggregates all 14 items when computing $D_o$ and $D_e$, while the **category_alphas** are computed within each subset. Hence, when inter-category variance is large, the overall $\alpha$ may diverge from the category-level estimates.

## H.3   Computation Settings

**Table H.2**  Summary of computation settings for Krippendorff's $\alpha$

| Aspect | Description |
| --- | --- |
| Measurement level | Nominal (discrete categorical labels) |
| Missing values | Allowed; no post-hoc adjudication performed |
| Estimation method | Python `krippendorff` package with 1,000 bootstrap iterations |
| Output | Point estimates of $\alpha$ for each category and the overall dataset |

The overall coefficient reflects agreement across all items, incorporating both intra- and inter-category variation. In contrast, the category-level coefficients isolate agreement within each conceptual dimension. The difference between these estimates provides insight into how cross-category variance affects overall coding reliability. A high $\alpha$ (above 0.80) across both levels indicates strong coder consistency and conceptual clarity of the FINDER framework.

# I   FINDER Stability Analysis via MiroFlow

To assess the stability and cross-lingual consistency of **FINDER**, we perform a multi-run evaluation using **MiroFlow** as a representative agent framework. MiroFlow is selected because it attains the highest overall performance among the evaluated frameworks on FINDER, making it a suitable testbed for stability analysis. We conduct three independent runs with both English (EN) and Chinese (ZH) prompts. The raw and aggregated **RACE** results are summarized in Table I.3.

As shown in Table I.3, the standard deviations across runs are small, indicating stable **FINDER** performance under repeated trials and across languages. The EN setting achieves slightly higher mean scores on *Overall*, *Comprehensiveness*, and *Depth*, suggesting modestly stronger reasoning and content generation in English. In contrast, *Instruction-following* and *Readability* are nearly identical between EN and ZH prompts, demonstrating consistent instruction adherence and output fluency across languages.

**Table I.3** MiroFlow RACE Results Summary (Three Runs)

| Dimension | EN Mean | EN Std | ZH Mean | ZH Std |
|-----------|---------|--------|---------|--------|
| Overall | 45.54 | 0.43 | 44.49 | 0.20 |
| Comp. | 45.58 | 0.63 | 44.43 | 0.16 |
| Depth | 41.63 | 0.56 | 39.16 | 0.36 |
| Inst. | 49.61 | 0.26 | 49.35 | 0.17 |
| Read. | 46.61 | 0.09 | 46.86 | 0.09 |

Figure I.2 visualizes the mean RACE scores with corresponding standard deviations. EN prompts yield consistently but only marginally higher scores, whereas both EN and ZH settings exhibit strong run-to-run stability, further supporting the robustness of **FINDER** in multilingual scenarios.
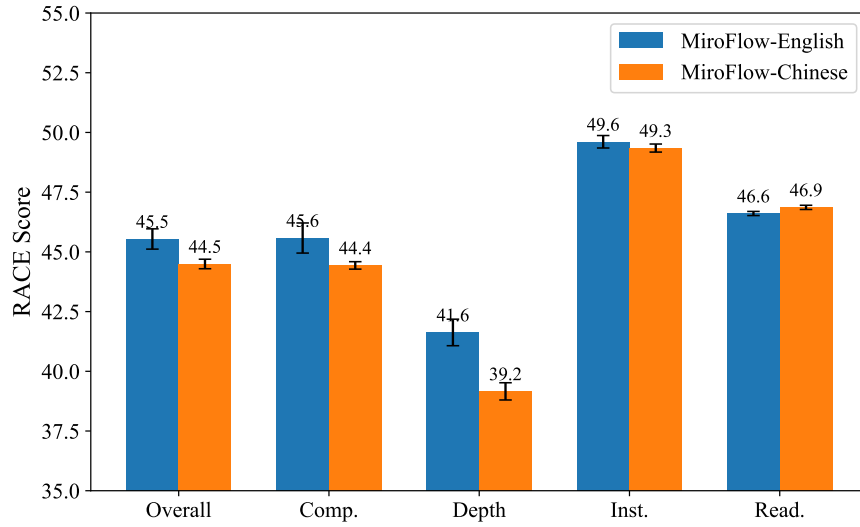


**Figure I.2** Comparison of FINDER RACE Results under MiroFlow (EN vs. ZH; mean over three runs)

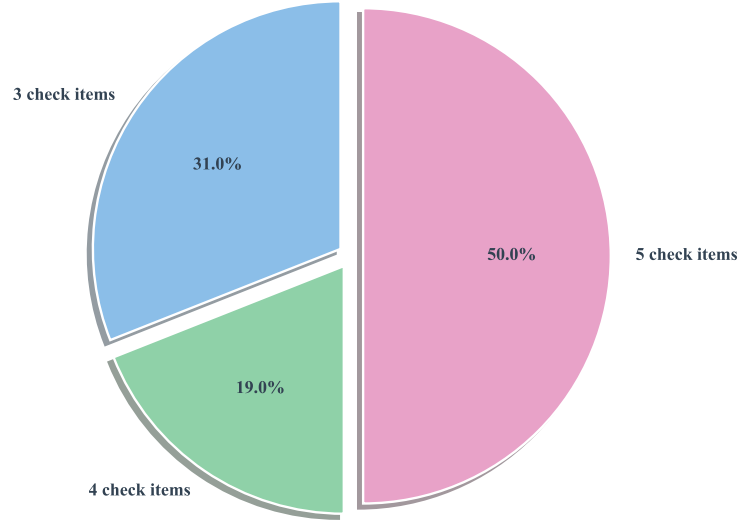## J  Checklist Distribution



**Figure J.3** The distribution of checklists in queries of FINDER

## K  Configuration of Evaluated Models

This appendix summarizes the configurations of all evaluated models used in our experiments. All models were used with their default system prompts and inference parameters unless otherwise stated.

### K.1  Proprietary API Models

These models were accessed through their official APIs with default configurations. No tool integration or parameter tuning was applied.

### K.2  Open-source Models

- **WebThinker**: `WebThinker-QwQ-32B`, with integrated web search. Default parameters used.
- **AFM**: `AFM-45B`, a multi-agent academic reasoning system with citation verification. Key parameters: `temperature = 0.4, top_p = 0.9, max_tokens = 32K`.
- **MiroThinker**: `MiroThinker-32B-DPO-v0.2` with `max_tokens = 64K`, using a multimodal vision model (`Qwen2.5-VL-72B-Instruct`) for image inputs.

### K.3  Agent Frameworks

- **MiroFlow**: Dual-agent framework based on `Claude-3.7-Sonnet`. Key parameters: `temperature = 0.3, top_p = 0.95, max_tokens = 32K`.
- **OWL**: Multi-agent architecture powered by `OpenAI O1`, integrating modules for reasoning, planning, and multimodal perception.
- **OpenManus**: `gpt-4o`-based agent system with automated web and code execution tools. Key parameters: `temperature = 0.0, max_tokens = 8192`.

## L  Positive Taxonomy Metric

This appendix provides a detailed justification for the positive-taxonomy scoring metric used in our analysis. Let $|D|$ denote the total number of evaluated instances and let $E_i \in [0, |D|]$ be the error count associated with category $i$. The

metric is defined as

$$S_i = |D| \cdot \cos\left( \frac{E_i}{|D|} \cdot \frac{\pi}{2} \right). \tag{7}$$

a cosine-based transformation mapping error counts into a bounded, positive scale.

The function is strictly monotonic decreasing and invertible over the domain $E_i \in [0, |D|]$, ensuring that it preserves all information contained in the raw error counts while providing a normalized and interpretable representation. This behavior makes it a suitable reparameterization for analyzing model performance across taxonomy categories.

A key motivation for adopting the cosine form lies in its curvature. Near the low-error regime ($E_i \approx 0$), the curve is relatively flat, meaning that very small increases in error induce minimal reductions in score. This reflects our analytical preference not to over-emphasize distinctions among categories that already exhibit highly reliable performance. As error increases, however, the curve becomes progressively steeper, producing sharper declines in score. This naturally concentrates resolution in the portions of the error range where error rates reach moderate and higher levels, making differences between categories more diagnostically significant. A linear mapping, by contrast, has constant slope and cannot provide this targeted emphasis.

The metric further benefits from an intuitive interpretability analogy. Inspired by classical cosine similarity in information retrieval [49], the score $S_i$ may be viewed as measuring the angular deviation between performance in category $i$ and an ideal, error-free direction. A category with zero errors aligns perfectly with this ideal, yielding a maximal score; larger error rates correspond to larger angular deviations and therefore smaller cosine values. This interpretation provides a geometric perspective on category-level performance that aligns closely with intuitive notions of similarity to a reference model.

## M   Analysis of Missing Results of FACT framework

During evaluation, we found that several models failed to produce valid outputs within the FACT framework. A follow-up analysis indicates that these failures fall into four primary categories. This appendix systematically examines each category to clarify potential sources of evaluation bias and to delineate limitations inherent to the FACT framework.

- **Anti-Scraping Mechanisms.** Many academic publishers, government agencies, and commercial websites employ anti-scraping protections. Consequently, Jina AI Reader often cannot access or parse these pages. This results in missing citations and incomplete retrieval chains, thereby weakening the reliability of FACT scores.

    **Examples:**
    - `https://www.tandfonline.com/doi/full/10.1080/14780887.2020.1769238#abstract`
    - `https://onlinelibrary.wiley.com/doi/10.1207/s15516709cog1202_4`
    - `https://www.tandfonline.com/doi/abs/10.1207/S15327965PLI1104_01`
    - `https://ingenico.com/us-en`
- **Non-Existent or Fabricated URLs.** In some instances, models generated URLs that do not correspond to real webpages. Such failures are typically caused by hallucinated links or by broader model limitations, which prevent the retrieval system from accessing the intended content.

    **Examples:**
    - `http://moe.gov.cn/`
    - `http://gd.gov.cn/`
    - `https://go.isi/mda2`
- **Incorrect URL Formats.** Some model outputs include academic references or citation strings that resemble URLs. Because of the internal URL-extraction rules in `deep_research_bench`, these strings may be misidentified as valid links. Since they do not map to actual web resources, retrieval fails by design.

    **Example error log:**
    ```
    ERROR: Failed to fetch Porter, M. & Heppelmann, J. (2021).
    "Digital twins in critical water infrastructure".
    IEEE Engineering Management Review, 49(3), 72-81.
    Jina AI Reader Failed ... 400
    ```

- **Timeout and Rate Limiting Issues.** Retrieval failures can also arise from system-level constraints, including network latency, high request volume, or temporary API throttling. These conditions may trigger timeouts, preventing content from being returned within the evaluation. As shown in the table below, successive versions exhibited only minor changes, and Krippendorff's Alpha steadily increased across iterations.