**Ill ByteDance | Seed**

# Seed2.0 Model Card:
# Towards Intelligence Frontier for Real-World Complexity

**Bytedance Seed**

## 1 Introduction

Large Language Models (LLMs) now play a central role in modern digital infrastructure. Usage has grown dramatically across both professional and personal contexts [11]. The Seed team has developed a comprehensive model family that includes general-purpose LLMs, multimodal models, open-source releases, code-specialized models, diffusion-based language modeling, formal theorem proving, and generative media systems: Seed1.6/1.8, Seed1.5-VL, Seed-OSS, Seed-Coder, Seed Diffusion, and Seed-Prover. These models currently power a large-scale product ecosystem serving hundreds of millions of daily active users across applications. Meanwhile, the field is moving toward an **agentic** paradigm where LLMs tackle scientific research, complex software development, autonomous documentation learning, and multi-step real-world workflows. This shift motivates **Seed2.0 Series (Pro / Lite / Mini)**, which is designed to deliver optimal user experience in large-scale production environments.

Seed2.0 prioritizes **user experience under large-scale online deployment**, as evidenced by strong results on the LMSYS Chatbot Arena, a public human preference benchmark. [1] In practice, interactive quality is most directly shaped by four factors: the prevalence of visual and multimodal queries, the impact of inference latency on user satisfaction, the need for reliable complex instruction execution, and the demand for seamless coding assistance. Our design reflects these priorities:

- **Robust Visual and Multimodal Understanding.** A substantial fraction of real user queries involve images—screenshots, charts, scanned documents, and mixed-media content. Seed2.0 strengthens visual reasoning with reduced hallucination [20, 35, 117] and improves structured extraction from documents and figures [65, 100].

- **Fast and Flexible Inference.** Inference latency directly impacts user experience. Seed2.0 offers three model sizes (Pro / Lite / Mini), allowing developers to choose the appropriate balance between performance and speed for their specific use case.

- **Reliable Complex Instruction Execution.** In production, we observe that users frequently issue complex, multi-step instructions that require precise execution—tasks where success depends not on factual recall but on structured reasoning and constraint satisfaction. Recent benchmarks such as $DeR^2$ [115] and CL-bench [26] capture exactly this demand. Seed2.0 treats it as a first-class requirement.

Seed2.0 also pursues a broader goal: handling tasks with real-world complexity. The Seed team has focused on raising the intelligence ceiling, moving from Olympiad-style problems toward research-level reasoning tasks. Seed2.0 tackles **Erdos problems** and performs **Scientific Coding**, pushing the boundaries of machine intelligence [4, 8, 70, 76].

Current agent systems, however, show an interesting asymmetry: they solve competition-level problems yet often fail to reliably complete practical tasks end-to-end—like building a well-designed application in one pass [17]. Two factors explain this gap. First, real-world tasks span long horizons and multiple stages, but

---

[1]Seed ranks 6th on the LMSYS Chatbot Arena — Text Arena (Overall) leaderboard and 3rd on the Vision Arena leaderboard as of Feb 16, 2026.
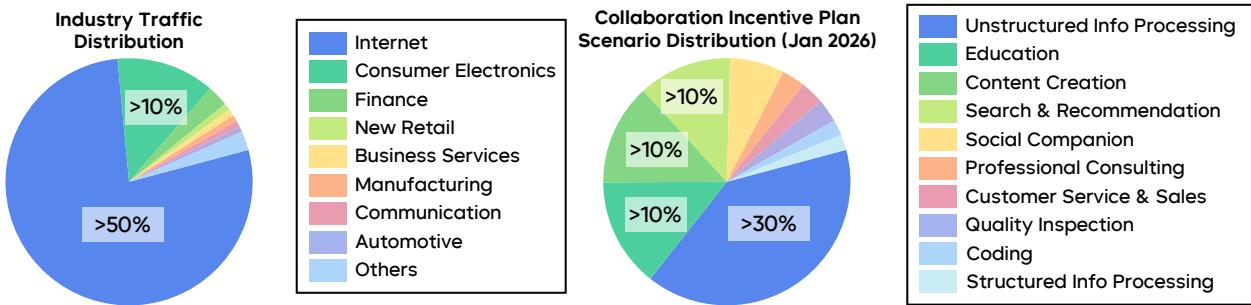
**Industry Traffic Distribution**

>10%
>50%

Internet
Consumer Electronics
Finance
New Retail
Business Services
Manufacturing
Communication
Automotive
Others

**Collaboration Incentive Plan Scenario Distribution (Jan 2026)**

>10%
>10%
>10%
>30%

Unstructured Info Processing
Education
Content Creation
Search & Recommendation
Social Companion
Professional Consulting
Customer Service & Sales
Quality Inspection
Coding
Structured Info Processing

**Figure 1** MaaS usage distribution in mainland China. Left: Industry traffic distribution showing strong dominance of the Internet sector. Right: Business Customer Usage Scenario Distribution. These statistics is named "Doubao Collaboration Incentive Program" , which sourced from the authorization of customers who have signed the Data Authorization Agreement.

existing LLM agents struggle to autonomously construct effective workflows and accumulate experience over extended timescales [5, 37, 60]. Second, real-world knowledge is highly domain-specific and long-tailed; models strong in math and code often provide little value in specialized professional contexts. Seed2.0 addresses this through systematic ingestion of long-tail domain knowledge [40, 132].

This report presents our initial progress toward real-world complexity. To systematically track and guide this effort, we establish an evaluation framework spanning four dimensions: **Science Discovery**, **Vibe Coding**, **Context Learning**, and **Real-World Tasks**—each targeting a core aspect of complex, long-horizon agent performance. This framework serves both as a benchmark suite and an iterative development guide.

Note that **the Seed2.0 Series still have gaps with international frontier LLMs**, while Seed identifies the directions for enhancing the model's capabilities for real-world complexity and makes great efforts to optimize Seed Model Series in this regard. Seed2.0 Series have considerable gaps with Claude in terms of coding, taking SWE-Evo and NL2Repo as examples. Seed2.0 Series have relatively obvious gaps with Gemini in terms of long-tail knowledge closely related to user experience, taking SuperGPQA and SimpleQA-Verified as examples.

In the following sections, we present Seed2.0's performance across standard benchmarks-where it performs on par with leading international frontier models—and showcase representative cases of Seed2.0 solving complex real-world problems.

We invite readers to explore the capabilities of Seed2.0.
Seed2.0 is now accessible on Volcano Engine, under the model id: `Doubao-Seed-2.0-pro`.
The model can be accessed at https://www.volcengine.com/experience/ark.
More details are available on the official page: https://seed.bytedance.com/zh/seed2.

## 2 Seed2.0 Deployment Patterns and Developer Behavior

### 2.1 MaaS Usage in Mainland China

MaaS usage patterns in mainland China concentrate heavily on enterprise-facing digital industries and cognitively intensive applications (Figure 1).

At the industry level, the Internet sector dominates overwhelmingly, accounting for the vast majority of traffic. Consumer electronics, finance, new retail, and business services follow at a considerable distance. Traditional verticals like manufacturing, automotive, and communication each represent less than 1% of total usage, perhaps due to the capability shortcomings of the previous Seed model series. The leading industries share key characteristics: higher information density, faster product iteration cycles, and tighter integration between models and production systems. Seed2.0 operates primarily within large-scale digital infrastructures where models participate directly in core business workflows, not as peripheral productivity tools.
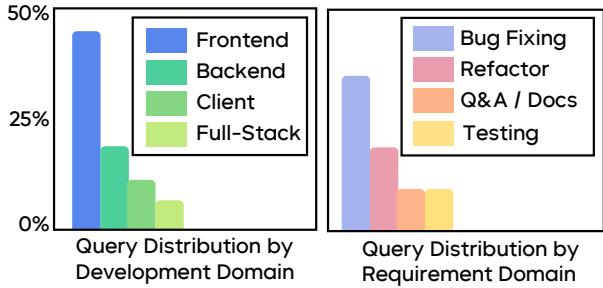
**Figure 2** Query distribution by development and requirement domain. Frontend development and bug fixing substantially dominate agentic coding requests, each far exceeding their respective category alternatives.

**Table 1** API Token Prefill / Decode Price Comparison (USD per 1M tokens). For Seed2.0 models with interval pricing, we report a single representative price.

| Model | Prefill (Input) | Decode (Output) |
|---|---|---|
| GPT-5.2 High | $1.75 | $14.00 |
| Claude-Opus-4.5-thinking | $5.00 | $25.00 |
| Gemini-3-Pro | $2.00–4.00* | $12.00–18.00* |
| Claude-Sonnet-4.5-thinking | $3.00 | $15.00 |
| GPT-5.0-mini High | $0.25 | $2.00 |
| Gemini-3-Flash High | $0.50–1.00* | $3.00* |
| **Seed2.0 Pro** | **$0.47 (¥3.41)** | **$2.37 (¥17.04)** |
| **Seed2.0 Lite** | **$0.09 (¥0.64)** | **$0.53 (¥3.83)** |
| **Seed2.0 Mini** | **$0.03 (¥0.22)** | **$0.31 (¥2.24)** |

*Gemini prices are ranges over context tiers (and input modalities for Flash) from vendor pricing pages.

At the scenario level, we analyze data from the Doubao Collaboration Incentive Program[2][3], an open program that encourages partners and developers to grant authorization for the use of authentic model usage data (Figure 1, right). Unstructured information processing and analysis dominates, representing the largest single share. Education, content creation, and search and recommendation follow as the next tier. Together, these top scenarios account for most deployment cases. Specialized applications—social companion, professional consulting, customer service and sales, quality inspection, coding, and structured information processing—each occupy substantially smaller shares. This pattern reflects how enterprises currently adopt AI: they start with scenarios requiring models to process heterogeneous data at scale, synthesize cross-domain knowledge, and generate actionable insights. More specialized use cases remain in earlier deployment stages.

The dominant scenarios impose specific technical demands: models must handle long contexts, integrate heterogeneous knowledge sources, execute multi-step instructions, and produce structured, high-fidelity outputs for downstream systems. In unstructured information processing—the largest category—enterprises use Seed Models to analyze user feedback, extract insights from multi-source documents, and generate structured reports for decision-making. Education applications power intelligent tutoring systems and personalized learning content. Content creation leverages multimodal capabilities for automated writing, video script generation, and multimedia synthesis. Search and recommendation systems integrate Seed Models for improved semantic understanding and ranking accuracy.

Seed Model functions as a workflow-oriented MaaS foundation rather than a lightweight conversational model, emphasizing multimodal understanding, long-context reasoning, structured generation, and tool-augmented execution for reliable end-to-end enterprise task completion.

Globally, this approach aligns with recent enterprise AI reports from OpenAI, Anthropic, and Google Cloud, which identify software engineering, research, analytics, customer support, and knowledge work as the fastest-growing enterprise AI categories [3, 34, 64].

## 2.2 Query Distribution in Agentic Coding

To understand real-world interaction patterns in agentic coding, we analyze trajectory-level developer usage data. The most striking finding is the dominance of frontend development (Figure 2). Queries related to page layout, styling, and UI logic management far exceed those for backend services, client-side applications, or full-stack integration. This distribution likely reflects both the iterative nature of frontend work—where visual feedback loops encourage frequent model interactions—and the relative accessibility of frontend tasks for AI assistance.

Programming language and framework statistics reinforce this pattern. Frontend-related languages (JavaScript, TypeScript, CSS, HTML) collectively account for the majority of code touched in these sessions. Among

---

[2] https://www.volcengine.com/docs/82379/1391869?lang=zh
[3] https://docs.byteplus.com/en/docs/legal/data_authorization_agreement_modelark

frameworks, Vue.js leads decisively—more than three times the adoption of React—reflecting the developer ecosystem in mainland China where Vue has historically enjoyed stronger community support.

By task type, bug fixing dominates, followed by refactoring and documentation work. This distribution suggests developers primarily turn to AI assistance for reactive maintenance rather than greenfield development. The high proportion of debugging queries indicates that error diagnosis and resolution remain significant pain points where AI provides clear value.

These patterns carry implications for model development. The concentration on frontend work suggests that JavaScript/TypeScript understanding, CSS layout reasoning, and framework-specific knowledge should be prioritized. The prevalence of bug-fixing queries indicates that models benefit from strong debugging capabilities—the ability to trace error messages, understand stack traces, and reason about program state.

## 2.3 Cost Efficiency

A key advantage of Seed2.0 lies in its cost structure. Table 1 compares API pricing across major foundation models. While Seed2.0 achieves comparable performance to frontier models on user experience, its token pricing is roughly an order of magnitude lower.

This cost differential is particularly significant for enterprise MaaS deployments. In scenarios involving high-volume, workflow-integrated usage—such as the unstructured information processing and content generation tasks described above—API costs can become a limiting factor for adoption. Seed2.0's pricing enables use cases that would be economically infeasible with more expensive alternatives, without sacrificing the reasoning and generation quality required for production systems.

The Seed2.0 Series offer tiered options to match different workload requirements. Seed2.0 Pro targets complex reasoning and long-context tasks where capability is paramount. Seed2.0 Lite provides a balanced trade-off for general-purpose applications. Seed2.0 Mini, with decode pricing under $0.50 per million tokens, opens possibilities for high-throughput, latency-sensitive applications where cost per query must remain minimal.

## 3 Comprehensive Evaluation Framework and Methodology

## 3.1 Fundamental Language Capacity

This section evaluates Seed2.0's fundamental capabilities, including reasoning, complex instruction following, and broad knowledge understanding. We compare our results with representative frontier models, including GPT-5.2 High, Claude-Sonnet-4.5, Claude-Opus-4.5, Gemini-3-Pro High, and Gemini-3-Flash High.

Specifically, we evaluate Seed2.0 on AIME 2025 and HMMT 2025 [4], BeyondAIME [8], IMOAnswerBench (no tool) [53], AetherCode [99], LiveCodeBench(v6) [44], Codeforces Elo Rating [74, 127] on the problem set from June to December 2025, GPQA Diamond [76], PhyBench [73], BABE [82], KORBench [54], ARC-AGI-1/2 [72], Inverse IFEval [123], MARS-Bench [108], MultiChallenge [24], COLLIE [113], MMLU-Pro [96], SuperGPQA [28], and LPFQA [132], along with a set of internal benchmarks designed to reflect high-value real-world tasks. For Graphwalks, we use our in-house tokenization pipeline during evaluation, which leads to a mismatch in how tokens are counted compared with the official OpenAI Graphwalks tokenization and scoring setup.

### 3.1.1 Long-tail Professional Knowledge Benchmarks

While benchmarks such as SimpleQA or HLE are valuable probes of factual recall and difficult trivia, they often emphasize rare or idiosyncratic facts that are only occasionally useful in real-world workflows. Inspired by the design philosophy of SuperGPQA [28], we place stronger emphasis on **professionally relevant long-tail knowledge** that arises in practical work settings. To this end, we design two new benchmarks, **LPFQA** and **Encyclo-K**, to directly measure a model's ability to function as a high-end search engine over long-tail professional knowledge.

**LPFQA.** LPFQA (Long-tail Professional Forum-based Question Answering) [132] is constructed from long-tail questions collected from professional forums and expert communities. It evaluates whether a model can correctly answer *realistic, domain-specific questions* encountered in daily work, covering fields such as

programming, finance, engineering, medicine, and applied science. LPFQA therefore measures the model's reliability in retrieving and synthesizing long-tail professional knowledge.

**Encyclo-K.** Encyclo-K [48] evaluates genuine mastery of book-level professional knowledge. It extracts atomic knowledge statements from books and dynamically composes them into evaluation instances, enabling flexible construction of test sets. This design supports both zero-shot and **few-shot in-context learning (ICL)** evaluation, allowing us to probe knowledge acquisition in pre-training and post-training stages. Compared with static QA datasets, Encyclo-K provides a scalable and compositional way to assess whether models internalize structured knowledge from long-form sources.

**HLE-Verified.** HLE-Verified (Humanity's Last Exam—Verified) is a curated subset of Humanity's Last Exam created in response to researchers reporting inaccurate, blurry, or underspecified questions in the original benchmark. We hire domain experts to review and select only questions that are clear, unambiguous, and confidently verifiable, thereby providing a more reliable evaluation of a model's true performance on challenging expert-level problems.

## 3.2 Fundamental Vision Capacity

To comprehensively evaluate the vision capabilities of Seed2.0, we conduct extensive evaluations across 50 public image benchmarks and 24 public video benchmarks. For image understanding, the selected benchmarks cover nine distinct categories: **Math, STEM, Visual Puzzles, Perception & Recognition, General VQA, Pointing & Counting, 2D & 3D Spatial Understanding, Document & Chart Understanding** and **LongContext Understanding**. The following benchmarks are used for evaluating the image understanding capability of Seed2.0:

- **MultiModal Math:** Mathematical reasoning within a visual context constitutes a core challenge for modern multimodal models, requiring rigorous logic and symbol grounding. To evaluate this capability, we employ a suite of benchmarks including MathVista (testmini) [52], MathVision [94], DynaMath [133], MathKangaroo [4], and MathCanvas [81]. Regarding specific metrics, for DynaMath, we report the `worst-case accuracy`, requiring the model to correctly answer all 10 variants of a problem to score. For MathKangaroo, performance is measured by the average accuracy across all bimonthly competitions in 2025, while `complete accuracy` is utilized for MathCanvas.

- **MultiModal STEM:** Mastery of domain-specific knowledge in science and engineering is essential for expert-level assistance. We assess this competence using MMMU [117], MMMU-Pro [118], EMMA [36], SFE [131], HiPhO [116], XLRS-Bench [92], and PhyX [80]. In terms of evaluation settings, for MMMU-Pro, we aggregate scores across the Standard (10 options) and Vision subsets. For HiPhO, the reported metric is the average normalized score from 13 Physics Olympiad competitions. For XLRS-Bench, we calculate the `macro average accuracy` on the lite subset. For PhyX, we report the accuracy over the testmini openended subset.

- **Visual Puzzles:** Abstract reasoning and pattern recognition capabilities are tested through puzzle-solving tasks, which serve as a proxy for general intelligence. Our evaluation includes LogicVista [106], VPCT [7], ZEROBench [77], ArcAGI (Image) [72], and VisuLogic [107]. Specifically, for ZEROBench, we assess accuracy on both main questions and sub-questions. In the case of ArcAGI, models are provided with both text matrices and rendered images. We notice that the additional visual input significantly improves the performance of Seed2.0 on ArcAGI.

- **Perception & Cognition:** Ensuring reliability involves minimizing hallucinations and mitigating biases in visual interpretation. We investigate these fundamental perceptual traits using VLMsAreBiased [90], VLMsAreBlind [75], VisFactor [41], RealWorldQA [105], and BabyVision [15]. For VisFactor, we report the `macro average accuracy` over all constituent tasks.

- **General VQA:** The model's versatility in handling open-ended queries and following instructions is reflected in General Visual Question Answering. This broad capability is gauged via SimpleVQA [20], HallusionBench [35], MME-CC [122], MMStar [16], MUIRBench [91], MTVQA [83], VibeEval [66], and ViVerBench [124]. Distinctively, for MTVQA, we deploy an LLM-judge (Deepseek-V3-0324) rather than rule-based matching to evaluate predictions. For VibeEval, raw scores on the 1–5 scale are normalized to a 0–100 range for reporting.

- **Point & Counting:** Fine-grained visual grounding and precise object enumeration are critical for tasks requiring high spatial fidelity. We measure these skills using CountBench [67], FSC-147 [2], and PointBench [19]. With respect to evaluation metrics, we report the `Mean Absolute Error (MAE)` for FSC-147.

- **2D & 3D Spatial Understanding:** Comprehending geometric relationships and depth is vital for embodied agents and 3D-aware applications. We utilize a comprehensive set of benchmarks including BLINK [33], MMSIBench [112], TreeBench [93], RefSpatialBench [128], DA-2K [111], All-Angles [114], and ERQA [85]. To ensure a robust assessment of spatial consistency in MMSIBench, we apply the `circular evaluation` strategy [49].
- **Document & Chart Understanding:** The ability to extract information from dense texts and interpret complex infographics is key for professional workflows. This domain is covered by ChartQAPro [58], OCRBenchv2 [31], OmniDocBench [65], and CharXiv [100]. For scoring specifics, OCRBenchv2 results represent the average of overall English and Chinese scores. We use `Normalized Edit Distance (NED)` for OmniDocBench 1.5, and for CharXiv, we average accuracy across both descriptive questions (DQ) and reasoning questions (RQ).
- **LongContext Understanding:** Processing extensive visual inputs, such as multi-page documents or long-form videos, tests the model's memory and temporal reasoning. To benchmark this capacity, we select DUDE [89], MMLongBench [98], LongDocURL [22], and MMLongBench-Doc [56].

For video understanding, we conducted an extensive evaluation across six dimensions: **video knowledge**, **video reasoning**, fundamental **video perception and motion understanding**, **long-video understanding**, **multi-video understanding**, and **streaming video understanding**, covering a total of 24 open benchmarks.

- **Video Knowledge:** We mainly adopt VideoMMMU [39], MMVU [125], and VideoSimpleQA [10] to evaluate the model's video knowledge capability. These benchmarks evaluate not only the model's mastery of world knowledge expressed in videos, but also its ability to acquire and internalize knowledge from video content.
- **Video Reasoning:** Multi-hop reasoning and video state tracking are core competencies for video reasoning. To this end, we conduct an in-depth evaluation using VideoReasonBench [51], Morse-500 [9], VideoHolmes [18], and Minerva [62]. Given that Morse-500 places a stronger emphasis on reasoning over physical dynamics, we increase the input frame rate to 5 FPS for this benchmark (the same setting is applied to all compared models).
- **Motion & Perception:** Foundational video perception, particularly motion perception, is central to video understanding. We therefore evaluate this capability using TVBench [21], ContPhy [126], TempCompass [50], EgoTempo [71], TOMATO [78], and MotionBench [38]. Given the prevalence of fast motion and rapid temporal state changes in these benchmarks, we increase the input frame rate to 2 FPS across all evaluations (the same setting is applied to all compared models).
- **Long-Video Understanding:** Long videos remain a major challenge in video understanding and a stringent test of a model's multimodal long-context reasoning. To accurately assess performance under long-video settings, we evaluate on VideoMME [30], CGBench [12], LongVideoBench [103], VideoEval-Pro [55], and LVBench [95], which include many hour-scale videos.
- **Multi-Video Understanding:** Multi-video understanding introduces new challenges for cross-context reasoning and is prevalent in real-world applications. Accordingly, we use CrossVid [46] as the primary benchmark to evaluate the model's performance on cross-video reasoning.
- **Streaming:** To assess real-time perception and interaction capabilities, we employ a comprehensive benchmark suite. OVBench [43] and OVOBench [63] are utilized to evaluate online reasoning and temporal awareness; LiveSports3K [14] emphasizes fine-grained sports event perception; ODVBench [120] tests generalization in autonomous driving scenarios; and ViSpeak [32] focuses on real-time visual referring capabilities.

## 3.3 Fundamental Agentic Capacity

With the transition from passive assistants to agentic systems, LLMs must go beyond single-turn responses and demonstrate the ability to **plan, invoke tools, interact with environments, and complete multi-step tasks**. We define this foundational layer as *Fundamental Agentic Capacity*. To avoid underestimating competing products, the final score is defined as the **maximum** of the score reported in the official documentation and the score obtained in our tests.

For better agentic evaluation, we conduct systematic refactoring of test scripts to optimize execution stability and reproducibility. We eliminate task-level entrypoint configurations to reduce redundant environment provisioning, consolidating execution environments into pre-built images. Where reference-script environments have degraded or contained errors, we perform targeted repairs to restore functionality. Additionally, we

6

replace external package repositories with internal mirrors to ensure consistent dependency resolution.

We also implement quality-based filtering to remove problematic test cases. This exclusion criterion targets several categories: multi-container Docker Compose scenarios that introduced unnecessary complexity; test cases where reference solutions fail to pass their own validation; cases exhibiting non-deterministic behavior across runs; tasks causing abnormal disk consumption; network-dependent problems with inconsistent outcomes; and scenarios requiring extended downloads or complex validation procedures that undermine reproducibility.

Specifically, we evaluate Seed2.0 across five representative dimensions: **Coding Agents**, **Search Agents**, **Tool Use**, **GUI Agents**, and **Deep Research**. These benchmarks cover repository-level software engineering (e.g., Terminal-Bench [59], SWE-Lancer [61], SWE-Bench [45], Multi-SWE-Bench [119], SWE-Bench Pro [23], SWE Multilingual [110], Scicode [87], SWE-Evo [86], Aider Polyglot, ArtifactsBench [121], CodeSimpleQA [109] and SpreadsheetBench Verified [57]), and Trae In-House Bench covering frontend and backend production scenarios), broad and deep information seeking (e.g. BrowseComp [101, 130], HLE [70], WideSearch [102], FinSearchComp [40] and seal-0[69]), tool invocation and orchestration (e.g., $\tau^2$-Bench[5], BFCL-v4 [68], MCP-Mark [104], VitaBench [37]), agentic visual tasks(e.g. Minedojo-Verified [29], HLE-VL and MM-BrowseComp [47] ) and long-horizon research reasoning and synthesis (e.g., DeepConsult [84], Deep Research [27] and ResearchRubrics [79]). Unless explicitly stated, evaluations are conducted *without external tools* and follow each benchmark's official protocol. For Terminal-Bench 2.0, we exclude three cases (`extract-moves-from-video`, `mailman`, and `install-windows-3.11`) due to network access restrictions and security considerations in our evaluation environment. The remaining tasks are adapted to our internal agent framework while preserving the original evaluation criteria.

## 3.4 Advanced Economically & Scientifically Valuable Tasks

As highlighted in the Introduction, the arrival of the **Agent era** fundamentally shifts the role of LLMs from answering isolated prompts to **driving long-horizon, economically and scientifically valuable workflows**. In this paradigm, models are expected to support scientific research, autonomously construct software systems, learn from user-provided context and documentation, and execute complex real-world tasks with tangible economic impact. Motivated by this shift, we build a set of **scenario-grounded evaluations** that directly reflect such real-world agentic workloads.

Specifically, we organize our advanced evaluations into four dimensions: **Scientific Discovery**, **Vibe Coding**, **Context Learning**, and **Real-World Tasks**. Each dimension is anchored by Seed-designed benchmarks targeting concrete failure modes observed in practice.

**Scientific Discovery.** We introduce **Ainstain Bench** [27] and **BABE** [129] to evaluate research-oriented capability. Ainstain Bench emphasizes **scientific coding**, measuring whether models can implement and manipulate computational procedures used in scientific workflows. BABE focuses on reasoning over **interleaved textual and visual scientific information** in the biological domain, assessing whether models can perform research-style inference grounded in multimodal evidence [27, 84].

**Vibe Coding.** We build **NL2Repo-Bench** to measure whether a model can complete an entire software repository from a natural-language specification in a single end-to-end process. This benchmark targets **long-horizon repository construction**, cross-file consistency, and dependency management, reflecting the emerging demand for extreme "vibe coding" scenarios [44, 99, 127].

**Economically Valuable Fields.** Beyond fundamental capabilities such as reasoning and knowledge, we prioritize high-value real-world applications to ensure that Seed2.0's development aligns with practical economic utility. To this end, we have developed a suite of **specialized in-house benchmarks** [1] including:

- **Education:** Evaluates performance in teaching-oriented scenarios, including problem solving, grading, explanation, and question generation, covering core subjects across K–12 levels.

- **Text Classification:** Evaluates the model's ability to analyze text and generate structured outputs or labels. This integrates compositional tasks—where multiple elements like intent and slots are identified in a single inference—with broader information processing, such as sentiment analysis and the synthesis of core viewpoints from unstructured data.

**Table 2** Evaluation on 2025 Olympiad-level Mathematical Competitions. The Gold Medal thresholds are $\geq 35$ for IMO 2025 and $\geq 87$ for CMO 2025.

| Competition | P1 | P2 | P3 | P4 | P5 | P6 | Overall | Medal |
|---|---|---|---|---|---|---|---|---|
| IMO 2025 | 7 | 7 | 7 | 7 | 7 | 0 | 35/42 | Gold |
| CMO 2025 | 21 | 21 | 9 | 21 | 21 | 21 | 114/126 | Gold |

- **Information Extraction:** Assesses structured extraction of relevant elements (e.g., words, sentences, or fields) from heterogeneous documents, including meeting records, legal texts, contracts, and corporate knowledge bases.

**Context Learning.** In enterprise and developer-facing settings, users require strict execution based on supplied context, such as long manuals or internal documents. Beyond existing evaluations like CL-Bench [26] and KOR-Bench [54], we incorporate **DeR**$^2$, which evaluates whether models can extract and utilize useful information from **noisy long-form technical documents** to perform reasoning and problem solving [115]. Furthermore, we introduce two **in-house scenarios** [1] to reflect agentic workloads:

- **Customer Support Q&A:** Assesses the ability to recognize user intent and synthesize answers after retrieving information from enterprise knowledge bases, specifically handling cases where recalled information is highly noisy.
- **Complex Workflow:** Validates the model's capacity to complete sophisticated, continuous tasks by synthesizing complex information and instructions provided within the context.

**Real-World Tasks.** We construct fine-grained internal evaluations for end-to-end task fulfillment. In particular, we curate **GDPVal-Verified**, a reliable subset of GDPVal with rubric-based automatic evaluation, and build the comparable **XpertBench** [1]. In addition, we introduce **WorldTravel** [97] to measure the model's ability to decompose goals and produce executable multi-step plans in real-world scenarios.

Overall, this evaluation suite operationalizes the vision described in our Introduction: assessing whether LLMs can function as **agentic systems that complete real tasks**, rather than merely answering questions. By grounding evaluation in realistic workflows and long-horizon completion, we obtain a more faithful measurement of advanced economically and scientifically valuable capability.

## 4 Results

### 4.1 Fundamental Language Evaluation

We evaluate Seed2.0 (Pro/Lite/Mini) on a comprehensive suite of fundamental language benchmarks spanning knowledge and science, mathematics, code and STEM reasoning, long-context understanding, multilinguality, instruction following, and hallucination robustness. Tables 3 and 4 report the detailed results.

**Overall**, Seed2.0 Pro sits comfortably within the international leading group across core language capabilities. Our optimization is shaped by large-scale product feedback: for user-facing deployments such as Doubao, we prioritize instruction-following robustness, long-tail knowledge coverage, and long-context stability; for coding-oriented products such as Trae, code reasoning and front-end generation quality take precedence. The benchmark results reflect these choices.

**Knowledge and Science.** Seed2.0 Pro leads on HealthBench and remains neck-and-neck with GPT-5.2 and Gemini-3-Pro on SuperGPQA and Encyclo-K, surpassing them in several cases. These gains matter for production: real user queries often probe obscure facts and domain-specific knowledge that generic training underserves.

**Mathematics, Code, and STEM Reasoning.** Seed2.0 Pro demonstrates advanced proficiency in mathematical reasoning. It exhibits highly competitive performance on popular evaluation suites (AIME, HMMT) and

**Table 3** Evaluation on Fundamental Language Capacity Benchmarks (**Large** Models). The highest score is marked in bold, and the second is underlined.

| Capability | Benchmark | GPT-5.2 High | Claude-Sonnet-4.5 | Claude-Opus-4.5 | Gemini-3-Pro High | Seed2.0 Pro |
|---|---|---|---|---|---|---|
| Science | MMLU-Pro | 85.9 | 88.0 | 89.3 | **90.1** | 87.0 |
| | HLE (no tool, text only) | 29.9 | 14.5 | 23.7 | **33.3** | 32.4 |
| | SimpleQA Verified | 36.8 | 29.3 | 48.6 | **72.1** | 36.0 |
| | HealthBench | **63.3** | 28.7 | 36.3 | 37.9 | 57.7 |
| | HealthBench - Hard | **42.0** | 10.9 | 11.0 | 15.0 | 29.1 |
| | SuperGPQA | 67.9 | 65.5 | 70.6 | **73.8** | 68.7 |
| | LPFQA | 54.4 | **54.9** | 52.6 | 51.2 | 52.6 |
| | Encyclo-K | 61.0 | 58.0 | 63.3 | 64.9 | **65.7** |
| Math | AIME 2026 | **97.5** | 82.5 | 92.5 | 93.3 | 94.2 |
| | AIME 2025 | **99.0** | 87.0 | 91.3 | 95.0 | 98.3 |
| | HMMT Feb 2025 | **100.0** | 79.2 | 92.9 | 97.3 | 97.3 |
| | HMMT Nov 2025 | **100.0** | 81.7 | 93.3 | 93.3 | 93.3 |
| | MathArenaApex | 18.2 | 1.0 | 1.6 | **24.5** | 20.3 |
| | MathArenaApex (shortlist) | 80.1 | 26.0 | 47.4 | 71.4 | **82.1** |
| | BeyondAIME | 86.0 | 57.0 | 69.0 | 83.0 | **86.5** |
| | IMOAnswerBench (no tool) | 86.6 | 60.7 | 72.6 | 83.3 | **89.3** |
| Code | Codeforces (no tool) | **3148** | 1485 | 1701 | 2726 | 3020 |
| | AetherCode | **73.8** | 16.4 | 31.6 | 57.8 | 60.6 |
| | LiveCodeBench (v6) | 87.7 | 64.0 | 84.8 | **90.7** | 87.8 |
| STEM | GPQA Diamond | **92.4** | 84.3 | 86.9 | 91.9 | 88.9 |
| | Superchem (text-only) | 58.0 | 32.4 | 43.2 | **63.2** | 51.6 |
| | BABE | **58.1** | 44.7 | 49.3 | 51.30 | 50.0 |
| | Phybench | 74.0 | 48.0 | 69.0 | **80.0** | 74.0 |
| | FrontierSci-research | **25.0** | 16.7 | 21.7 | 15.0 | **25.0** |
| | FrontierSci-olympiad | **75.0** | 60.0 | 71.0 | 73.0 | 74.0 |
| General Reasoning | ARC-AGI-1 | **89.9** | 70.9 | 84.0 | 85.0 | 85.4 |
| | ARC-AGI-2 | **57.5** | 13.6 | 29.1 | 31.1 | 37.5 |
| | KORBench | **79.2** | 73.0 | 77.4 | 73.9 | 77.5 |
| | ProcBench | 95.0 | 87.5 | 92.5 | 90.0 | **96.6** |
| Long Context Performance | MRCR v2 (8-needle) | **89.4** | 47.1 | 56.2 | 79.7 | 54.0 |
| | Graphwalks Bfs (<128k) | **98.0** | 80.5 | 92.0 | 79.9 | 68.9 |
| | Graphwalks Parents (<128k) | **99.7** | 99.0 | 96.2 | **99.7** | 97.6 |
| | LongBench v2 (128k) | 63.2 | 62.0 | 65.0 | **67.4** | 63.8 |
| | Frames | 84.0 | 78.7 | 84.7 | 81.9 | **84.5** |
| | DeR$^2$ Bench | **69.0** | 58.9 | 60.4 | 66.1 | 58.2 |
| | CL-Bench | **23.9** | 18.1 | 22.6 | 15.6 | 20.8 |
| Multilingual | Global PIQA | 93.2 | 93.9 | 93.9 | **95.0** | 92.3 |
| | MMMLU | 90.3 | 89.9 | 91.0 | **91.8** | 88.1 |
| | Disco-X | 76.3 | 70.3 | 78.6 | 76.8 | **82.0** |
| Instruction Following | MultiChallenge | 59.5 | 57.3 | 59.0 | **68.7** | 68.3 |
| | COLLIE | **96.9** | 77.3 | 79.8 | 95.0 | 93.9 |
| | MARS-Bench | **87.9** | 72.9 | 87.7 | 85.6 | 85.6 |
| | Inverse IFEval | 72.3 | 69.3 | 72.4 | **79.6** | 78.9 |
| Hallucination | LongFact-Objects | **99.2** | 98.8 | 99.0 | 98.1 | 92.9 |
| | LongFact-Concepts | **99.7** | 98.5 | 98.8 | 98.5 | 92.8 |
| | FactScore | 91.9 | 90.6 | 91.1 | **92.6** | 71.2 |

**Table 4** Evaluation on Fundamental Language Capacity Benchmarks (**Efficient** Models). The highest score is marked in bold, and the second is underlined.

| Capability | Benchmark | GPT–5–mini High | Gemini–3–Flash High | Seed2.0 Mini | Seed2.0 Lite |
|---|---|---|---|---|---|
| Science | MMLU-Pro | 84.1 | **87.8** | 83.6 | <u>87.7</u> |
| | HLE (no tool, text only) | 17.6 | **31.7** | 13.3 | <u>28.2</u> |
| | SimpleQA Verified | <u>26.0</u> | **65.4** | 18.9 | 24.0 |
| | HealthBench | **62.5** | <u>51.6</u> | 30.0 | 51.2 |
| | HealthBench - Hard | **38.6** | <u>21.5</u> | 15.3 | 20.0 |
| | SuperGPQA | 60.5 | **72.7** | 61.6 | <u>67.5</u> |
| | LPFQA | 50.7 | **51.6** | 47.2 | <u>50.9</u> |
| | Encyclo-K | 53.0 | <u>60.0</u> | 52.1 | **64.5** |
| Math | AIME 2026 | <u>92.5</u> | **93.3** | 86.7 | 88.3 |
| | AIME 2025 | 90.3 | **95.2** | 87.0 | <u>93.0</u> |
| | HMMT Feb 2025 | <u>93.3</u> | **100** | 70.0 | 90.0 |
| | HMMT Nov 2025 | **96.7** | 96.7 | 80.0 | <u>86.7</u> |
| | MathArenaApex | 2.1 | **17.7** | 4.2 | <u>4.7</u> |
| | MathArenaApex (shortlist) | 43.4 | **71.9** | 31.1 | <u>52.6</u> |
| | BeyondAIME | 72.0 | **82.0** | 69.0 | <u>76.0</u> |
| | IMOAnswerBench (no tool) | 72.1 | **84.4** | 71.6 | <u>81.6</u> |
| Code | Codeforces | 1985 | **2727** | 1644 | <u>2233</u> |
| | AetherCode | <u>42.6</u> | **56.1** | 29.8 | 41.5 |
| | LiveCodeBench (v6) | 62.6 | **84.7** | 64.1 | <u>81.7</u> |
| STEM | GPQA Diamond | 82.1 | **90.7** | 79.0 | <u>85.1</u> |
| | Superchem (text-only) | 34.8 | **54.4** | 16.2 | <u>48.0</u> |
| | BABE | 49.2 | **55.2** | 40.4 | <u>50.2</u> |
| | Phybench | 60.0 | **77.0** | 56.0 | <u>73.0</u> |
| | FrontierSci-research | **18.3** | <u>11.7</u> | 3.3 | **18.3** |
| | FrontierSci-olympiad | 69.0 | **73.0** | 44.0 | <u>70.0</u> |
| General Reasoning | ARC-AGI-1 | 54.5 | **86.9** | 43.3 | <u>75.7</u> |
| | ARC-AGI-2 | 3.5 | **34.3** | 2.3 | <u>14.8</u> |
| | KORBench | 74.2 | <u>76.0</u> | 72.8 | **77.0** |
| | ProcBench | 87.5 | <u>90.0</u> | 80.1 | **92.4** |
| Long Context Performance | MRCR v2 (8-needle) | <u>50.1</u> | **79.0** | 51.4 | 33.6 |
| | Graphwalks Bfs (<128K) | **85.5** | <u>84.2</u> | 64.1 | 82.5 |
| | Graphwalks Parents (<128K) | 96.6 | <u>99.7</u> | 93.0 | **100.0** |
| | LongBench v2 (128K) | 56.7 | **64.0** | 52.3 | <u>59.6</u> |
| | Frames | <u>82.9</u> | **83.7** | 80.5 | 83.4 |
| | DeR$^2$ Bench | 50.3 | **66.0** | 46.6 | <u>57.3</u> |
| | CL-Bench | **25.2** | 16.1 | 14.8 | <u>20.0</u> |
| Multilingual | Global PIQA | 91.6 | **95.6** | 89.2 | <u>92.1</u> |
| | MMMLU | 86.3 | **91.8** | 81.6 | <u>87.7</u> |
| | Disco-X | 67.7 | 71.9 | <u>73.0</u> | **80.3** |
| Instruction Following | MultiChallenge | 59.0 | **69.3** | 61.1 | <u>63.2</u> |
| | COLLIE | **97.4** | <u>96.5</u> | 91.2 | 94.0 |
| | MARS-Bench | 66.1 | **84.6** | 62.4 | <u>80.5</u> |
| | Inverse IFEval | 74.8 | **80.9** | 69.3 | <u>77.1</u> |
| Hallucination | LongFact-Objects | **99.2** | <u>97.9</u> | 87.0 | 92.2 |
| | LongFact-Concepts | **99.5** | <u>98.6</u> | 91.4 | 92.4 |
| | FactScore | **96.1** | <u>92.0</u> | 50.4 | 62.4 |

**Table 5** Evaluation on Putnam-200 [88], Pass@8. Agent-based multi-turn setup with Lean, Python, and Lean search tools. The highest score is marked in **bold**, and the second is <u>underlined</u>.

| Benchmark | Deepseek-Prover-V2 | Seed-1.5 Prover | Gemini-3-Pro | Seed2.0 Lite | Seed2.0 Pro |
|---|---|---|---|---|---|
| Putnam-200 | <4.0 | 26.5 | 26.5 | <u>30.5</u> | **35.5** |

**Table 6** Complex instruction-following benchmark across multiple test sets and instruction categories.

| Test Set | Description | Example |
|---|---|---|
| ● Format | Provide format-related instructions and require the model to produce outputs in the specified format. | You are an expository-instruction rewriter designed for young children. You can transform user-provided instructions for different objects into introductions that are easy to understand and rich in imagination.... You must append a cute, positive emoji to the end of every sentence throughout the entire text 😊. |
| ● Conditional | Provide instructions with conditional rules and require the model to generate outputs that follow the rules. | **Character**<br>◆ You are a customer-service bot for a cross-border e-commerce platform responsible for men's and women's...<br>◆ You must ensure that the user's information is complete. The required information includes "gender," "foot length"... |
| ● Content | Require the response to include specified content. | **Skill 1: Standard interaction**<br>1. When user greets you, identify yourself as a food calorie calculator and invite the user to ask a question …<br>**Skill 2: Provide calorie data per 100 g of food...** |
| ● Phrasing | Require the response to include specified phrasing. | **Role:** You are a question-loving bot. Your "head" is filled with different questions. Although you love asking questions, you consistently offer users distinctive reflections and experiences.<br>**Requirements:**<br>◆ Your task is to generate questions corresponding to the topic proposed by the user.<br>◆ Ensure diversity in the questions; each time you ask questions, you must provide three... |
| ● Tone | Require LLMs to generate responses that match a specified tone or style (e.g., humorous, formal, sarcastic). | You particularly dislike people around you commenting on your current life situation, especially when they ask why you still have not had children or bought a house; you respond with sarcastic, cutting remarks... |
| ● Emoji | Require the response to include emojis, or a specified type of emojis. | … When a pedestrian and a car are involved in a traffic accident, you represent it as "🚗💥🚶". When multiple pedestrians and a single car are involved, represent it as "🚗💥🚶🚶". When vehicles are involved in a chain-reaction collision, represent it as "🚗💥🚙💥🚗💥🚙"... |
| ● Few-shot | Require the model to reference or follow the exemplar format and/or content when generating outputs. | You need to revise the text I provide into an ordered-list format, using the following rules: 1. Convert the text into an ordered list; you do not need to consider whether the revised sentences are fluent...<br>**Fewshot:** … |
| ● Chinese | Evaluate whether the model can generate text that satisfies specified Chinese length constraints. | Please help me come up with a title based on the following content. The title must be exactly 10 Chinese characters (no more, no less). Provide only one title—do not give multiple options... |
| ● English | Evaluate whether the model can generate text that satisfies specified English length constraints. | Please explore the coordination and impact of monetary policy in the context of globalization, for economists and policy makers, combined with international trade theory and exchange rate mechanism, to distinguish the difference in policy response between developed and developing countries... |

performs on par with SOTA models on extremely challenging benchmarks such as IMOAnswerBench and MathApex. Notably, Seed2.0 Pro achieves gold-medal level performance in Olympiad-level mathematical competitions: 2025 International Mathematical Olympiad (IMO) and 2025 China Mathematical Olympiad (CMO). Beyond natural language reasoning, Seed2.0 also excels in formal theorem proving, demonstrating superior performance on the Putnam-200 (a random subset of 200 problems from PutnamBench [88]) and making strides in open Erdős problems. Table 2 and Table 5 report the results. Implementation details are provided in Appendix E. In coding, it reaches a Codeforces Elo of 3020 and posts strong LiveCodeBench numbers. STEM-oriented benchmarks such as FrontierSci-research tell a similar story, with Seed2.0 Pro matching or edging out Gemini-3-Pro in multiple cases. This reflects our continued investment in deep reasoning and code-centric training.

**Long Context Understanding.** Seed2.0 Pro ranks first on the Frames leaderboard. While retrieval-heavy tasks such as MRCR and Graphwalks still show some headroom compared to other leading models, these gaps have limited impact on real-world user experience—and similar weaknesses are observed across other frontier models.

**Instruction Following, Multilinguality, and Hallucination.** Instruction-following capability sees clear gains over Seed1.8, which directly benefits complex user requests in deployment. Multilingual performance and hallucination robustness remain competitive, with room for further improvement in future iterations.

**Complex Instruction Following.** We run a fine-grained analysis using an in-house benchmark (see Table 6) tailored for Chinese-language production scenarios: 912 test cases across 17 weighted dimensions reflecting practical deployment importance.

**Table 7** Complex instruction following performance on in-house Chinese benchmark. Bold indicates the better result.

| Model | Overall | Format | Conditional | Content | Phrasing | Tone | Emoji | Few-shot | Chinese | English |
|---|---|---|---|---|---|---|---|---|---|---|
| Seed-1.8 | 72.89 | 45.33 | 81.25 | **90.48** | 75.00 | 61.29 | 92.14 | 55.78 | 77.67 | 62.05 |
| Seed2.0 Pro | **75.26** | **46.00** | **88.19** | 87.76 | **85.31** | **76.45** | **93.62** | **65.31** | **80.00** | **67.41** |

As shown in Table 7, Seed2.0 Pro scores 75.26%, a +2.37% absolute gain over Seed1.8. The largest jumps come in tone control (+15.16%), phrasing adherence (+10.31%), and few-shot learning (+9.53%). In practice, this means more precise Chinese pragmatic effects—including subtle stylistic modes like tsundere-like affect and sarcastic or ironic ("yin-yang") expression—and more reliable multi-constraint prompting. At the few-shot level, the model better infers structural and stylistic requirements implied by exemplars, including quantitative constraints such as paragraph counts or emoji quotas, and shows improved compliance with strict length limits in both Chinese and English.

**Lite and Mini Variants.** Seed2.0 Lite and Seed2.0 Mini offer strong efficiency–quality trade-offs. Across a wide range of benchmarks, our small models hold their own against counterparts from OpenAI and Google. Seed2.0 Lite in particular posts strong math and reasoning numbers while maintaining solid instruction-following performance— well-suited for latency-sensitive and cost-constrained deployments.

## 4.2   Vision Task Evaluation

We evaluate the performance of Seed2.0 (including Pro, Lite, and Mini variants) on a comprehensive suite of public visual-language benchmarks, comparing it against its predecessor Seed1.8 and several state-of-the-art models such as Gemini-3-Pro, GPT-5.2, and Claude-Opus-4.5. The evaluation covers a broad spectrum of capabilities, ranging from mathematical and STEM reasoning to visual puzzles, spatial understanding, and long-context document processing. Table 8 presents the detailed results, where Seed2.0 Pro demonstrates superior performance, achieving the highest scores across the majority of benchmarks.

**Math and STEM Reasoning.** In mathematical reasoning, Seed2.0 Pro demonstrates exceptional capability, achieving state-of-the-art results on MathVision (88.8), MathKangaroo (90.5), and MathCanvas (61.9), while tying for the top score on MathVista (89.8). Notably, Seed2.0 Lite also excels, securing the top spot on DynaMath (70.5). In STEM tasks, Seed2.0 Pro leads in EMMA (72.0), XLRS-Bench (54.6), and PhyX (72.1). On the remaining benchmarks, Seed2.0 Pro also delivers transformative gains compared to Seed1.8: scores on HiPhO and MMMU-Pro improved by 15.8 and 5.0, respectively, rapidly narrowed the gap with best methods.

**Visual Puzzles and Logic.** Seed2.0 Pro exhibits significant advancements in logical reasoning and visual puzzle solving. It achieves the highest scores on LogicVista (81.4) and ZeroBench (Main 12.0, Sub 47.6), demonstrating robust problem-solving abilities. On the challenging VisuLogic benchmark, Seed2.0 Pro scores 47.4, outperforming all counterparts by a notable margin. Although GPT-5.2 leads on ArcAGI-Image, Seed2.0 Pro secures a competitive second place (88.8 on ArcAGI1 and 43.3 on ArcAGI2), showing substantial improvement over Seed-1.8.

**Perception and General VQA.** For perception and recognition, Seed2.0 Pro demonstrates exceptional capabilities. It achieves SOTA results on VLMsAreBiased (77.4), VLMsAreBlind (98.6), and BabyVision (60.6) . In general VQA benchmarks, Seed2.0 Pro leads in SimpleVQA (71.4), MUIRBench (81.8), and VibeEval (81.4). Seed2.0 Pro consistently outperforms Seed1.8 across all general VQA tasks, highlighting its refined visual understanding.

**Spatial Understanding and Counting.** Seed2.0 Pro excels in spatial understanding and fine-grained localization. It sets new state-of-the-arts on DA-2K (92.3), RefSpatialBench (72.6), and BLINK (79.5), surpassing the strong baseline of Gemini-3-Pro. In counting tasks, Seed2.0 Pro achieves best performance on FSC-147 with a mean absolute error of 11.3 (lower is better), significantly improving upon Seed1.8 and competitor models.

**Document and Long-Context Understanding.** A standout feature of Seed2.0 Pro is its dominance in long-

**Table 8** Performance of Seed2.0 on public visual-language benchmarks compared to previous models. We report Pass@1 in these benchmarks. The best score for each benchmark is marked in **bold**, and the second best is underlined. Results marked with an * are sourced from the technical report.

| Capability | Benchmark | Claude-Opus-4.5 | GPT-5.2 High | Gemini-3-Pro High | Seed1.8 | Seed2.0 Mini | Seed2.0 Lite | Seed2.0 Pro |
|---|---|---|---|---|---|---|---|---|
| Math | MathVista | 80.6 | 83.1 | **89.8** | 87.7 | 85.5 | 89.0 | **89.8** |
| | MathVision | 74.3 | 86.8 | 86.1 | 81.3 | 78.1 | 86.4 | **88.8** |
| | DynaMath | 52.5 | 70.1 | 63.3 | 61.5 | 58.9 | **70.5** | 68.9 |
| | MathKangaroo | 69.6 | 86.9 | 84.4* | 73.8 | 79.8 | 86.3 | **90.5** |
| | MathCanvas | 52.9 | 55.3 | 58.8 | 53.6 | 53.2 | 61.1 | **61.9** |
| STEM | MMMU | 81.6 | 83.7 | **87.0** | 83.4 | 79.7 | 83.7 | 85.4 |
| | MMMU-Pro | 70.8 | 79.5* | **81.0*** | 73.2 | 71.4 | 76.0 | 78.2 |
| | EMMA | 60.4 | 69.4 | 66.5 | 60.9 | 57.0 | 65.5 | **72.0** |
| | SFE | 55.8 | 50.1 | **61.9** | 51.2 | 48.4 | 53.4 | 55.6 |
| | HiPhO | **81.8** | 77.7 | 79.1 | 58.3 | 55.8 | 72.5 | 74.1 |
| | XLRS-Bench (macro) | 50.4 | 49.9 | 51.7 | 39.9 | 49.9 | 53.7 | **54.6** |
| | PhyX (openended) | 61.3 | 71.5 | 71.0 | 65.9 | 65.0 | 62.8 | **72.1** |
| Visual Puzzles | LogicVista | 68.9 | 81.0 | 80.8 | 78.3 | 73.8 | 79.6 | **81.9** |
| | VPCT | 29.0 | 56.0 | **90.0** | 61.0 | 48.0 | 73.0 | 76.0 |
| | ZeroBench (main) | 4.0 | 11.0 | 10.0 | 11.0 | 7.0 | 8.0 | **12.0** |
| | ZeroBench (sub) | 30.8 | 38.9 | 42.2 | 37.7 | 36.2 | 42.2 | **47.6** |
| | ArcAGI1-Image | 75.8 | **93.1** | 69.4 | 31.4 | 29.8 | 80.9 | 88.8 |
| | ArcAGI2-Image | 26.1 | **54.4** | 21.5 | 1.3 | 1.5 | 28.3 | 43.3 |
| | VisuLogic | 27.6 | 37.0 | 39.0 | 35.8 | 40.4 | 47.3 | **47.4** |
| Perception & Recognition | VLMsAreBiased | 21.4 | 28.0 | 50.6* | 62.0 | 58.4 | 74.8 | **77.4** |
| | VLMsAreBlind | 77.2 | 84.2 | 97.5* | 93.0 | 93.1 | 97.0 | **98.6** |
| | VisFactor | 24.5 | 33.6 | **45.8** | 20.4 | 23.6 | 33.4 | 36.8 |
| | RealWorldQA | 75.9 | 82.1 | 84.7 | 78.0 | 81.6 | 81.7 | **86.0** |
| | BabyVision | 16.2 | 37.4 | 49.7* | 30.2 | 38.7 | 57.5 | **60.6** |
| General VQA | SimpleVQA | 57.9 | 54.1 | 69.7 | 65.4 | 68.7 | 67.2 | **71.4** |
| | HallusionBench | 65.3 | 67.7 | **69.9** | 63.9 | 65.1 | 66.0 | 68.0 |
| | MME-CC | 25.2 | 44.4 | 56.9 | 43.4 | 40.8 | 50.2 | **57.0** |
| | MMStar | 73.9 | 78.2 | **83.1** | 79.9 | 79.1 | 80.7 | 83.0 |
| | MUIRBench | 78.9 | 77.4 | 78.2 | 78.7 | 78.0 | 76.2 | **81.8** |
| | MTVQA | **53.1** | 48.5 | 50.8 | 47.3 | 50.6 | 51.1 | 51.1 |
| | WorldVQA | 36.6 | 26.3 | 47.5 | 40.4 | 47.6 | 44.0 | **49.9** |
| | VibeEval | 70.3 | 73.1 | 77.7 | 74.0 | 76.5 | 76.5 | **81.4** |
| | ViVerBench | 72.4 | 74.8 | 75.9 | 74.6 | 73.9 | **80.0** | 75.9 |
| Pointing & Counting | CountBench | 90.3 | 91.2 | **97.3** | 96.3 | 95.5 | 97.1 | 95.5 |
| | FSC-147↓ | 20.9 | 21.1 | 12.1 | 13.6 | 17.3 | 11.9 | **11.3** |
| | Point-Bench | - | - | 85.5* | 76.5 | 77.0 | 79.0 | 81.4 |
| 2D & 3D Spatial Understanding | BLINK | 68.1 | 70.3 | 77.1 | 74.3 | 73.4 | 75.6 | **79.5** |
| | MMSIBench (circular) | 20.2 | 26.1 | 25.4 | 25.8 | 19.7 | 28.3 | **32.5** |
| | TreeBench | 53.6 | 58.8 | 62.7 | 58.5 | 57.3 | 64.2 | **64.7** |
| | RefSpatialBench | - | 25.5 | 65.5* | 56.3 | 55.6 | 66.4 | **72.6** |
| | DA-2K | 70.3 | 78.9 | 82.1 | 90.7 | 86.4 | 90.3 | **92.3** |
| | All-Angles | 63.1 | 71.5 | **73.5** | 61.6 | 61.3 | 65.2 | 72.1 |
| | ERQA | 48.3 | 59.8 | **70.5*** | 58.8 | 56.3 | 65.8 | 68.5 |
| Document & Chart Understanding | ChartQAPro | - | 67.6 | 69.0 | 63.0 | 65.2 | 70.3 | **71.2** |
| | OCRBenchv2 | 55.5 | 55.6 | **63.3** | 52.6 | 58.5 | 62.4 | 62.5 |
| | OmniDocBench 1.5 ↓ | 0.153 | 0.143* | 0.115* | 0.106 | 0.110 | 0.102 | **0.099** |
| | CharXiv-DQ | 92.7 | 93.8 | **94.4** | 88.0 | 91.9 | 93.3 | 93.5 |
| | CharXiv-RQ | 65.5 | **82.1*** | 81.4* | 71.4 | 70.8 | 79.9 | 80.5 |
| LongContext Understanding | DUDE | 55.6 | 68.2 | 70.1 | 69.4 | 68.9 | 72.1 | **72.4** |
| | MMLongBench | - | - | 73.6 | 72.4 | 66.7 | 70.8 | **74.8** |
| | LongDocURL | - | - | 72.0 | 74.5 | 71.3 | **75.1** | 74.7 |
| | MMLongBench-Doc | - | - | 59.5 | 57.0 | 48.9 | 55.1 | **61.4** |

**Table 9** Performance of Seed2.0 on public video understanding benchmarks compared to previous models. The highest score in each benchmark is marked in bold, and the second is underlined. For benchmarks marked with a ‡, we include subtitles for evaluation. Results marked with an * are sourced from the technical report.

| Capability | Benchmark | Human | Gemini-3-Pro | Gemini-3-Flash | Seed1.8 | Seed2.0 Mini | Seed2.0 Lite | Seed2.0 Pro |
|---|---|---|---|---|---|---|---|---|
| Knowledge | VideoMMMU [39] | 74.4 | 87.6* | **88.1** | 82.7 | 80.6 | 84.1 | 86.9 |
| | MMVU [125] | 49.7 | 76.3 | 77.9 | 73.1 | 69.0 | 75.0 | **78.2** |
| | VideoSimpleQA [10] | - | **71.9** | 70.7 | 67.8 | 67.7 | 66.6 | **71.9** |
| Reasoning | VideoReasonBench [51] | 73.8 | 59.5 | 61.2 | 52.8 | 40.5 | 64.2 | **77.8** |
| | Morse-500 [9] | 55.4 | 33.0 | 32.4 | 29.2 | 32.2 | 32.2 | **37.4** |
| | VideoHolmes‡ [18] | - | 64.2 | 65.6 | 65.5 | 58.6 | 63.8 | **67.4** |
| | Minerva‡ [62] | - | 65.0 | 64.4 | 62.4 | 54.7 | 63.8 | **66.5** |
| Motion & Perception | TVBench [21] | 94.8 | 71.1 | 69.6 | 71.5 | 70.5 | 71.5 | **75.0** |
| | ContPhy [126] | - | 58.0 | 60.5 | 54.9 | 55.9 | 56.1 | **67.4** |
| | TempCompass [50] | 97.3 | 88.0 | 88.3 | 86.9 | 83.7 | 87.0 | **89.6** |
| | EgoTempo [71] | 63.2 | 65.4 | 58.4 | 67.0 | 67.2 | 61.8 | **71.8** |
| | MotionBench [38] | - | 70.3* | 68.9 | 70.6 | 64.4 | 70.9 | **75.2** |
| | TOMATO [78] | 95.2 | 59.6 | **60.8** | **60.8** | 47.4 | 57.3 | 59.9 |
| | + *Thinking with Tracking* | 95.2 | - | 64.0 | 61.0 | 51.3 | 59.2 | **65.3** |
| Long Video | VideoMME‡ [30] | - | 88.4* | 85.2 | 87.8 | 81.2 | 87.7 | **89.5** |
| | CGBench [12] | - | **65.5** | 65.3 | 62.4 | 59.2 | 59.3 | 65.0 |
| | LongVideoBench [103] | - | 76.7 | 74.5 | 77.4 | 74.8 | 77.3 | **80.3** |
| | VideoEval-Pro [55] | - | **52.7** | 51.9 | 45.9 | 43.7 | 44.3 | 48.0 |
| | LVBench [95] | - | - | - | 73.0 | 66.6 | 73.0 | **76.4** |
| Multi Video | CrossVid [46] | 89.2 | 53.0 | 48.7 | 57.3 | 58.6 | 57.7 | **60.3** |
| Streaming | OVBench [43] | - | 62.7 | 59.2 | 65.1 | 60.1 | 65.5 | **69.2** |
| | LiveSports-3K [14] | - | 74.5 | 71.5 | 77.5 | 73.3 | 77.8 | **78.0** |
| | OVOBench [63] | 92.8 | 70.1 | 68.7 | 72.6 | 70.4 | 76.7 | **77.0** |
| | ODVBench [120] | 91.4 | 63.6 | 56.7 | 63.5 | 65.1 | 69.6 | **72.5** |
| | ViSpeak [32] | 96.0 | **89.0** | 86.0 | 79.0 | 77.5 | 84.0 | 78.5 |

context understanding. It achieves SOTA scores on DUDE (72.4), MMLongBench (74.8), and MMLongBench-Doc (61.4). In document and chart understanding, Seed2.0 Pro also leads on ChartQAPro (71.2) and OmniDocBench 1.5, proving its efficacy in processing complex, information-dense visual inputs.

In summary, Seed2.0 Pro delivers state-of-the-art performance across a diverse array of visual-language tasks, particularly excelling in mathematics reasoning, perception proficiency, spatial reasoning, and long-context understanding, while Seed2.0 Lite and Mini offer competitive efficiency-focused alternatives.

## 4.3 Video Task Evaluation

We conduct a comprehensive evaluation of the Seed2.0 family (including Seed2.0 Mini, Seed2.0 Lite, and Seed2.0 Pro) across multiple dimensions of video capabilities, including knowledge, reasoning, perception, motion understanding, as well as long-term, multi-video, and streaming video analysis.

As shown in Table 9, Seed2.0 pushes the performance frontier for video understanding and delivers state-of-the-art results on multiple benchmarks, with exceptional performance in motion perception, reasoning, and streaming video understanding.

**Video Knowledge.** In terms of knowledge, Seed2.0 achieves leading performance on the multidisciplinary benchmark MMVU [125], and delivers results comparable to Gemini-3-Pro on VideoMMMU [39] and VideoSim-

**Table 10** Performance of Seed2.0 with video tool-use on long-form video understanding and reasoning. We compare the performance of Seed2.0 when using the VideoCut tool across different benchmarks.

| Benchmark | Average Duration | Gemini-3-Pro | Seed1.8 | Seed1.8 w/ VideoCut | Seed2.0 | Seed2.0 Pro w/ VideoCut |
|---|---|---|---|---|---|---|
| CGBench [12] | 1624 seconds | 65.5 | 62.4 | <u>65.9</u> | 65.0 | **66.8** |
| LVBench [95] | 4104 seconds | - | 73.0 | <u>78.9</u> | 76.4 | **80.0** |
| ZeroVideo | 1672 seconds | 14.3 | 6.9 | <u>18.8</u> | 14.5 | **27.9** |

pleQA [10], representing a substantial improvement over the previously released Seed1.8 [1]. In addition, it is worth noting that our lightweight model, Seed2.0 Mini, also performs strongly.

**Video Reasoning.** In video reasoning, Seed2.0 Pro continues to push the performance frontier. On VideoReasonBench [51], Seed2.0 surpasses human performance, with particularly strong results in visual state tracking, widely regarded as one of the most fundamental capabilities underlying video reasoning.

On Morse-500 [9], a highly challenging reasoning benchmark covering video abstract reasoning, physical reasoning, and planning reasoning, Seed2.0 further achieves a new state of the art, reaching 37.4% accuracy. Despite this progress, a substantial gap to human-level performance remains, and future iterations of the model will continue to close this gap.

**Video Perception and Motion Understanding.** Foundational video perception is critical to video understanding and reasoning. Among these perceptual competencies, motion understanding is particularly crucial for modeling temporal state transitions and tracking dynamic changes over time. Seed2.0 further strengthens its motion perception and understanding capabilities. As summarized in Table 9, Seed2.0 Pro attains state-of-the-art performance on the majority of benchmarks and substantially outperforms Seed1.8 and Gemini-3-Pro across several evaluations.

For TOMATO [78], we explore a **Thinking with Tracking** strategy in which we prompt to encourage the model to explicitly output per-frame bounding boxes for the moving target during reasoning, and then identifies the motion type based on the target's trajectory. This strategy is effective for both the Gemini and Seed models, and yields meaningful gains on TOMATO. Nevertheless, performance on TOMATO and TVBench suggests that a considerable gap to human-level proficiency remains. We will continue to advance motion perception in subsequent iterations to further narrow this gap.

**Long Video Understanding.** For long-video understanding, Seed2.0 Pro achieves a breakthrough performance of 89.5 on VideoMME [30]. Seed2.0 also delivers strong results on LVBench and LongVideoBench, showing substantial improvements over Seed1.8 across these benchmarks. In addition, the entire Seed2.0 model family is equipped with VideoCut tool-use capabilities by default to improve long video reasoning.

**Multiple Video Understanding.** Multi-video understanding requires models to identify and integrate critical evidence across multiple video contexts, and is a core capability for real-world applications as well as video agents. On CrossVid [46], a benchmark for multi-video understanding, Seed2.0 Pro achieves a new state-of-the-art score of 60.3. Looking ahead, we will prioritize further improvements in multi-video understanding, with an emphasis on strengthening cross-context reasoning.

**Streaming.** Moving beyond static analysis, streaming video understanding and interactive reasoning enable systems to comprehend and respond to visual inputs in real-time, providing key capabilities for interactive products (*e.g.*, Doubao video calling). Across multiple benchmarks, Seed2.0 Pro delivers further improvements, while Seed2.0 Lite surpasses Gemini-3-Pro/Flash in several evaluations with high efficiency.

**Video Tool-Use: VideoCut.** Seed2.0 also further enhances video tool-use capability, *i.e.*, VideoCut. When processing long videos or tasks that require high-frame-rate perception, it can adopt VideoCut to replay relevant segments at a higher FPS, which is an intuitive and effective mechanism for video understanding and reasoning. As shown in Table 10, we evaluate the performance of enabling VideoCut for Seed2.0 Pro.

**Table 11** Evaluation on Fundamental Agentic Capacity Benchmarks (**Large** Models). The highest score is marked in bold, and the second is underlined. Some scores differ greatly from the evaluation results in the tech reports by other organizations. The scores in parentheses represent the results under the aligned settings then.

| Capability | Benchmark | GPT-5.2 High | Claude-Sonnet-4.5 | Claude-Opus-4.5 | Gemini-3-pro High | Seed2.0 Pro |
|---|---|---|---|---|---|---|
| Coding Agent | Terminal Bench 2.0[1] | **62.4** | 45.2 | 60.2 | 56.9 | 55.8 |
| | SWE-Lancer | 48.9 | 45.7 | **56.1** | 44.3 | 49.4 |
| | SWE Bench Verified | 80.0 | 77.2 | **80.9** | 76.2 | 76.5 |
| | Multi-SWE-Bench | 47.7 | 47.7 | **52.8** | 50.2 | 45.2 |
| | SWE-Bench Pro | **55.6** | 48.4 | 55.4 | 49.7 | 46.9 |
| | SWE Multilingual | 68.8 | 64.1 | **74.0** | 72.7 | 71.7 |
| | Scicode | 49.7 | 47.9 | 52.8 | **57.7** | 48.5 |
| | SWE-Evo | 12.5 | 16.7 | **27.1** | 8.9 | 8.5 |
| | Aider Polyglot | 91.1 | 82.2 | 92.4 | **94.2** | 80.0 |
| | ArtifactsBench | **71.1** | 59.1 | 68.5 | 58.4 | 66.6 |
| | CodeSimpleQA | 62.3 | 59.6 | **63.0** | 54.7 | 58.0 |
| | SpreadsheetBench Verified | 69.9 | 75.9 | 78.6 | 70.8 | **79.1** |
| Search Agent | BrowseComp | **77.9 (65.3)** | 43.9 (29.5) | 67.8 (57.2) | 59.2 | 77.3 |
| | BrowseComp-zh | 76.1 | 42.4 | 62.4 | 66.8 | **82.4** |
| | HLE-text | 45.5 | 32.0 | 43.2 | 46.9 | **54.2** |
| | HLE-Verified | 68.5 | 37.6 | 56.6 | 67.5 | **73.6** |
| | WideSearch | **76.8** | 65.1 | 76.2 (71.7) | 67.3 | 74.7 |
| | FinSearchComp | **73.8** | 58.6 | 66.2 | 52.7 | 70.2 |
| | DeepSearchQA | 71.3 (66.4) | 36.3 | 76.1 (41.6) | 63.9 | **77.4** |
| | Seal-0 | 51.4 | **53.4** | 47.7 | 45.5 | 49.5 |
| Tool Use | $\tau^2$-Bench (retail) | 82 | 86.2 | 88.9 | 85.3 | **90.4** |
| | $\tau^2$-Bench (telecom) | **98.7** | 98.0 | 98.2 | 98.0 | 94.2 |
| | MCP-Mark | **57.5** | 32.1 | 42.3 | 53.9 | 54.7 |
| | BFCL-v4 | 65.9 | 72.9 | **76.5** | 71.0 | 73.4 |
| | VitaBench | 41.8 | 40.8 | **55.3** | 48.8 | 47.0 |
| Deep Research | DeepConsult | 54.3 | 55.8 | 61.0 | 48.0 | **61.1** |
| | DeepResearchBench | 52.2 | 47.2 | 50.6 | 49.6 | **53.3** |
| | ResearchRubrics | 42.3 | 38.6 | 45.0 | 37.7 | **50.7** |
| Vision Agent | Minedojo-Verified | 18.3 | - | - | 23.3 | **49.0** |
| | MM-BrowseComp | 26.3 | - | - | 25.0 | **48.8** |
| | HLE-VL | 31.0 | - | - | 36.0 | **39.2** |

[1] Using Terminus2.

With VideoCut, Seed2.0 Pro further raises the ceiling of long-video understanding, yielding substantial improvements on CGBench, LVBench, and the challenging ZeroVideo, which consists of challenging, real-world video reasoning scenarios, including fine-grained high-frame-rate motion perception and long-form, multi-hop reasoning over extended videos.

## 4.4 Fundamental Agentic Capacity

We evaluate Seed2.0 (Pro/Lite) on a diverse set of agentic benchmarks covering search agents, deep research, vision agents, coding agents, and tool use. Tables 11 and 12 report detailed results.

**Overall**, Seed2.0 Pro lands in the top tier across agentic capabilities, with clear advantages on search, deep research, and vision agent tasks—the workloads that matter most for high-frequency user scenarios like information seeking, complex reasoning, and multimodal interaction.

**Search, Deep Research, and Vision Agents.** Seed2.0 Pro consistently leads or sits near the top on search and research-oriented benchmarks: HLE, BrowseComp, WideSearch and DeepSearchQA. On deep research tasks,

**Table 12** Evaluation on Fundamental Agentic Capacity Benchmarks (**Efficient** Models). The highest score is marked in bold, and the second is underlined.

| Capability | Benchmark | GPT-5-mini High | Gemini-3-Flash High | Seed2.0 Lite |
|---|---|---|---|---|
| Coding Agent | Terminal Bench 2.0[1] | 36.9 | **60.0** | <u>45.0</u> |
| | SWE-Lancer | 43.1 | **51.7** | <u>47.1</u> |
| | SWE Bench Verified | 67.9 | **78.0** | <u>73.5</u> |
| | Multi-SWE-Bench | <u>49.3</u> | **59.0** | 41.1 |
| | SWE-Bench Pro | **51.7** | 46.7 | <u>46.0</u> |
| | SWE Multilingual | 63.8 | **71.1** | <u>64.4</u> |
| | Scicode | 40.2 | **55.0** | <u>52.4</u> |
| | SWE-Evo | 10.4 | **12.8** | <u>10.6</u> |
| | Aider Polyglot | <u>79.1</u> | **92.0** | 76.0 |
| | ArtifactsBench | **68.9** | 52.5 | <u>62.6</u> |
| | CodeSimpleQA | **57.6** | <u>53.7</u> | 53.1 |
| | SpreadsheetBench Verified | 58.1 | <u>65.7</u> | **82.3** |
| Search Agent | BrowseComp | <u>48.1</u> | 41.5 | **72.1** |
| | BrowseComp-zh | 49.5 | <u>63.0</u> | **82.0** |
| | HLE-text | 35.8 | <u>47.6</u> | **49.5** |
| | HLE-Verified | 56.4 | **71.8** | <u>70.7</u> |
| | WideSearch | 37.7 | <u>64.0</u> | **74.5** |
| | FinSearchComp | 38.1 | <u>54.8</u> | **65.1** |
| | DeepSearchQA | 16.7 | <u>54.7</u> | **67.7** |
| | Seal-0 | 34.2 | <u>37.8</u> | **52.3** |
| Tool Use | $\tau^2$-Bench (retail) | - | <u>88.6</u> | **90.9** |
| | $\tau^2$-Bench (telecom) | - | **94.7** | <u>92.1</u> |
| | MCP-Mark | 30.2 | <u>40.5</u> | **46.7** |
| | BFCL-v4 | 57.9 | <u>65.0</u> | **72.9** |
| | VitaBench | 20.8 | **46.7** | <u>41.8</u> |
| Deep Research | DeepConsult | <u>49.8</u> | 26.0 | **60.3** |
| | DeepResearchBench | <u>50.7</u> | 46.1 | **54.4** |
| | ResearchRubrics | <u>43.6</u> | 36.9 | **50.8** |
| Vision Agent | Minedojo Verified | 20.3 | <u>24.3</u> | **39.7** |
| | MM-BrowseComp | 17.9 | <u>22.8</u> | **45.1** |
| | HLE-VL | 18.7 | **36.8** | <u>35.8</u> |

[1] Using Terminus2.

it posts the best results on Deep Research and Research Rubrics. In vision-agent settings, Seed2.0 Pro pulls ahead of reported baselines by a wide margin on Minedojo-Verified, MM-BrowseComp, and HLE-VL—evidence of strong multimodal grounding for action-oriented tasks. Taken together, these numbers place Seed2.0 Pro at or near state-of-the-art on search, deep research, and vision agents.

**Tool Use and Coding Agents.** Tool use and coding-agent performance is solidly first-tier. Seed2.0 Pro holds its own on SWE-Bench Pro and SWE Bench Verified, and tops SpreadSheetBench for structured artifact manipulation. We also see distinctive gains on SWE-Evo, pointing to robust evolutionary code improvement. Beyond this, Seed2.0 Pro also delivers strong performance on general tool use tasks such as MCP-Mark, BFCL-v4 and $\tau^2$-Bench. That said, complex API orchestration and long-horizon code execution still leave room to grow.

**Efficient Agentic Deployment.** Among small models, Seed2.0 Lite consistently ranks first or second on efficient agentic benchmarks. It performs well on search and deep research tasks and posts excellent SpreadsheetBench numbers. Paired with significantly lower inference cost, Lite offers a practical quality–cost trade-off for large-scale and latency-sensitive deployments.

**Table 13** Evaluation on Advanced Economically and Scientifically Valuable Tasks (**Large** Models). The highest score is marked in bold, and the second is underlined.

| Capability | Benchmark | GPT-5.2 High | Claude-Sonnet-4.5 | Claude-Opus-4.5 | Gemini-3-pro High | Seed2.0 Pro |
|---|---|---|---|---|---|---|
| Science Discovery | Scicode | 49.7 | 47.9 | <u>52.8</u> | **57.7** | 52.1 |
| | FrontierSci-research | **25.0** | 16.7 | <u>21.7</u> | 15.0 | <u>23.3</u> |
| | Superchem (text-only) | <u>58.0</u> | 32.4 | 43.2 | **63.2** | 53.0 |
| | BIObench | **58.1** | 44.7 | 49.3 | 51.3 | <u>53.5</u> |
| | AInstein Bench | 41.3 | 33.7 | <u>44.0</u> | 42.8 | **47.7** |
| Vibe Coding | NL2Repo-Bench | **49.3** | 39.9 | <u>43.2</u> | 34.2 | 27.9 |
| | NL2Repo (Pass@1) | **8.0** | 3.0 | 3.0 | <u>4.0</u> | 3.0 |
| | ArtifactsBench | **71.1** | 59.1 | <u>68.5</u> | 58.4 | 66.6 |
| | SWE-Bench Pro | **55.6** | 48.4 | <u>55.4</u> | 49.7 | 46.9 |
| | Terminal Bench 2.0[1] | **62.4** | 45.2 | <u>60.2</u> | 54.2 | 55.8 |
| Context Learning | KORBench | **79.2** | 73.0 | <u>77.4</u> | 73.9 | 77.2 |
| | DeR$^2$ Bench | **69.0** | 58.9 | 60.4 | <u>66.1</u> | 58.2 |
| | CL-Bench | **23.9** | 18.1 | <u>22.6</u> | 15.6 | 21.5 |
| | ToB-Complex Workflows | 45.0 | 61.0 | <u>64.8</u> | **69.2** | 64.7 |
| | ToB-Reference Q&A | 63.6 | 58.9 | 67.9 | <u>68.3</u> | **72.4** |
| Real World Tasks | HealthBench - Hard | **36.6** | 10.9 | 11.0 | 15.0 | <u>28.3</u> |
| | GDPVal-Diamond | **26.9** | 15.2 | 20.7 | 19.4 | <u>21.3</u> |
| | XPert Bench | <u>53.3</u> | 44.7 | 50.5 | 53.1 | **64.5** |
| | ToB-K12 Education | <u>61.6</u> | 50.1 | 56.2 | 59.4 | **62.8** |
| | ToB-Compositional Tasks | 51.5 | 57.3 | <u>63.6</u> | **64.8** | 59.1 |
| | ToB-Text Classification | 62.1 | 64.5 | 63.9 | <u>67.5</u> | **69.0** |
| | ToB-Information Extraction | 44.7 | 48.3 | <u>50.1</u> | 49.0 | **52.0** |
| | World Travel (VLM) | **19.33** | 2.67 | <u>14.0</u> | 8.0 | 12.0 |
| | World Travel (TEXT) | **32.67** | 10.0 | <u>21.3</u> | 14.7 | 23.3 |

[1] Using Terminus2.

## 4.5 Advanced Economically and Scientifically Valuable Tasks

We further evaluate Seed2.0 on advanced tasks that reflect scientific discovery, economically valuable workflows, context-based learning, and real-world task completion. Tables 13 and 14 report the detailed results.

**Scientific Discovery.** Seed2.0 Pro operates at the frontier. It posts strong numbers on FrontierSci-research and scientific coding benchmarks, and leads on AInstein Bench—signs of robust hypothesis-driven reasoning in research-style scenarios.

**Real-World Economic Value.** Beyond science, Seed2.0 Pro translates well to tasks with direct economic value. It tops XPert Bench and remains highly competitive on GDPVal and multiple ToB benchmarks. Notably, the model excels on user-oriented scenarios we explicitly optimize for: customer-service QA, information extraction, intent recognition, and K12 problem solving. These results reflect reliable instruction-following behavior in production.

**Context Learning and Repository-Level Code.** Clear headroom remains on context-driven learning and end-to-end vibe coding. Seed2.0 Pro trails the strongest baselines on DeR$^2$ Bench and NL2Repo-Bench, suggesting that long-horizon context integration and repository-level code generation are still challenging. We have flagged these as priority directions and are actively pushing targeted research.

**Efficiency at Scale.** Seed2.0 Lite holds up well on advanced tasks among efficient models. It posts strong results on ToB real-world benchmarks, solid scientific discovery numbers, and competitive context-learning outcomes. On several advanced benchmarks, Lite beats GPT-5-Mini outright while offering favorable latency and cost.

**Table 14** Evaluation on Advanced Economically and Scientifically Valuable Tasks (**Efficient** Models). The highest score is marked in bold, and the second is underlined.

| Capability | Benchmark | GPT-5-mini High | Gemini-3-Flash High | Seed2.0 Lite |
|---|---|---|---|---|
| Science Discovery | Scicode | 40.2 | **55.0** | 52.4 |
| | FrontierSci-research | **18.3** | 11.7 | **18.3** |
| | Superchem (text-only) | 34.8 | **54.4** | 48.0 |
| | BIObench | 49.2 | **55.2** | 50.2 |
| | AInstein Bench | 35.0 | **44.0** | 38.3 |
| Vibe Coding | NL2Repo-Bench | 19.5 | **27.6** | 24.6 |
| | NL2Repo (Pass@1) | **2.0** | **2.0** | 1.0 |
| | ArtifactsBench | **68.9** | 52.5 | 62.6 |
| | SWE-Bench Pro | **51.7** | 46.7 | 46.0 |
| | Terminal Bench 2.0[1] | 36.9 | **60.0** | 45.0 |
| Context Learning | KORBench | 74.2 | 76.0 | **77.0** |
| | DeR$^2$ Bench | 50.3 | **66.0** | 57.3 |
| | CL-Bench | **25.2** | 16.1 | 20.0 |
| | ToB-Reference Q&A | 53.4 | 64.9 | **68.2** |
| | ToB-Complex Workflows | 42.5 | 61.3 | **62.0** |
| Real World Tasks | HealthBench - Hard | **38.6** | 21.5 | 20.0 |
| | GDPVal-Diamond | 11.5 | 6.46 | **23.2** |
| | XPert Bench | 47.6 | 50.1 | **63.3** |
| | ToB-K12 Education | 53.7 | 59.0 | **63.8** |
| | ToB-Compositional Tasks | 47.5 | **57.7** | 54.8 |
| | ToB-Text Classification | 52.7 | 61.9 | **64.5** |
| | ToB-Information Extraction | 40.5 | 45.4 | **48.4** |
| | World Travel (VLM) | 2.7 | 7.3 | **14.0** |
| | World Travel (TEXT) | 15.3 | 13.3 | **24.0** |

[1] Using Terminus2.

Across agentic capacity and advanced real-world tasks, Seed2.0 delivers consistent top-tier performance. Seed2.0 Pro excels in search, deep research, vision agents, and scientific discovery; Lite provides a compelling efficiency-focused alternative for large-scale deployment.

# 5 Use Cases of Seed2.0

In this section, we demonstrate the applicability of Seed2.0 across diverse real-world scenarios, illustrating how the model addresses practical challenges in vibe coding, application operation, and even scientific discovery.

## 5.1 Vibe Coding

In the "Vibe Coding" setting, the model is evaluated on end-to-end software engineering behavior under realistic constraints rather than isolated code snippets. Tasks include interpreting underspecified requirements, forming executable plans, producing runnable artifacts, and iteratively refining solutions based on feedback until acceptance criteria are met.

To characterize Seed2.0's capabilities in this regime, we present five representative case studies: algorithmically challenging code synthesis under compute constraints (§ 5.1.1), repository-level construction from requirements documents alone (§ 5.1.2), zero-regression debugging guided by unit tests (§ 5.1.3), cross-version upgrade debugging using release notes with full test suite validation (§ 5.1.4), and competition-level programming under official limits (§ 5.1.5). These settings vary supervision completeness, engineering scope, and failure signal observability, enabling comprehensive assessment of planning, implementation, verification, and iterative correction in realistic development workflows.

**Task Overview**

The file feal.c implements a FEAL-like encryption function, with decrypt.c being the corresponding decryption function, and the file pairs.txt contain 32 plaintext-ciphertext pairs...

Implement a known-plaintext attack that recovers the key of the cipher. Prove you have recovered the key by decrypting all the ciphertexts in the file ciphertexts.txt, and saving your answer in the file plaintexts.txt.

**Assistant:**

[Stage 1: Understanding the Cipher Structure]

```
execute_bash(command="ls −la") ...
str_replace_editor(command="view", path="/app/feal.c")...
```

[Stage 2: Attack Design & Implementation]

```
str_replace_editor(command=" create ", path="/app/attack.py", file_text="
```

*Implementing the F-function Inverse*

```python
...
def f_inv_function(output_val):
    output_val = output_val & 0xFFFFFFFF
    y = [0]*4
    for i in range(4):
        y[i] = output_val & 0xFF
        output_val >>= 8

    # Compute x from y
    temp_y1 = rotr2(y[1])
    x1_xor_x0 = (temp_y1 - y[2]) & 0xFF

    temp_y2 = rotr2(y[2])
    x3_xor_x2 = (temp_y2 - x1_xor_x0 - 1) & 0xFF

    temp_y3 = rotr2(y[3])
    x3 = (temp_y3 - y[2]) & 0xFF

    x2 = (x3_xor_x2 ^ x3) & 0xFF
    temp_y0 = rotr2(y[0])

    x0 = (temp_y0 - y[1] - 1) & 0xFF
    x1 = (x1_xor_x0 ^ x0) & 0xFF

    input_val = 0
    for i in range(4):
        input_val |= ([x0, x1, x2, x3][i] << (8*i))
    return input_val & 0xFFFFFFFF
...
```

*Meet-in-the-Middle Attack*

```python
# Step A: Build delta map for k0
k0_delta_map = defaultdict(list)
for k0 in range(0, 1<<20):
    K0_exp = expand(k0)

    # Compute for pair 0
    L0_0, R0p_0, L4_0, R4_0 = preprocessed[0]
    term0 = f_function(K0_exp ^ R0p_0) ^ \
                    L0_0 ^ L4_0
    A0 = f_inv_function(term0)
```

*Key Recovery Implementation*

```python
    # Compute for pair 1
    L0_1, R0p_1, L4_1, R4_1 = preprocessed[1]
    term1 = f_function(K0_exp ^ R0p_1)
                    ^ L0_1 ^ L4_1
    A1 = f_inv_function(term1)
    delta = A0 ^ A1
    k0_delta_map[delta].append((k0, A0, A1))

print(f"Built k0 delta map with", ...)

# Step B: Iterate k3 and find matches
found = None
for k3 in range(0, 1<<20):
    K3_exp = expand(k3)
    # Compute for pair 0
    L0_0, R0p_0, L4_0, R4_0 = preprocessed[0]
    B0 = R4_0 ^ f_function(L4_0 ^ K3_exp)
    # Compute for pair 1
    L0_1, R0p_1, L4_1, R4_1 = \
            preprocessed[1]
    B1 = R4_1 ^ f_function(L4_1 ^ K3_exp)
    ...
# Verify all pairs
print("Verifying all pairs...")
```

*Verification*

```python
# Implement encrypt function to test
def encrypt(plaintext, k0, k1, k2, k3):
    key = [expand(k0), expand(k1), \
            expand(k2), expand(k3)]
    left = getleft(plaintext)
    right = getright(plaintext)
    right ^= left
    for round_num in range(4):
        mixed = f_function(right ^ \
                key[round_num])
        temp = right
        right = left ^ mixed
        left = temp
        left ^= right
    return merge(right, left)
```

```
")
```

[Stage 3: Verification & Decryption]

```
execute_bash(command="cd /app && python3 attack.py")...
```

Verify all 32 pairs: For each (Pi, Ci) in pairs.txt:
Result: All 32 pairs verified successfully! Keys are correct.

Process `ciphertexts.txt` ... Decrypted 100 ciphertexts successfully! Saved to `plaintexts.txt`.

**Figure 3** The case of Seed2.0 with Complex Code Generation Task on TerminalBench 2.0 [59].

### 5.1.1  Complex Code Generation

To evaluate Seed2.0's capacity for complex code synthesis, we examine its performance on the "feal-linear-cryptanalysis" task from TerminalBench 2.0 [59], classified as Hard. This task provides a simplified FEAL-like block cipher as a Feistel network, a reference decryption program, and multiple plaintext–ciphertext pairs. The objective is to construct a known-plaintext attack that recovers the secret key and decrypts additional ciphertexts. Naïve brute-force enumeration over the effective key space is computationally infeasible within benchmark constraints, requiring algorithmic insight rather than exhaustive search.

As shown in Figure 3, Seed2.0 completes this task in 12 interaction rounds through systematic repository analysis. It inspects the project structure, examines `feal.c` to reconstruct the round function, and studies documentation to align its implementation with expected formats and correctness criteria. This structured approach minimizes specification inconsistencies and informs development of `attack.py`, a cryptanalytic program tailored to the specific FEAL variant.

The attack's effectiveness derives from combining algebraic inversion with time–memory trade-offs. Seed2.0 derives inverse transformations for the cipher's internal operations, including the reversal of rotation and modular-addition steps in relevant components, and implements an explicit inverse routine for the $F$-function (`f_inv_function`). Seed2.0 then precomputes an inverse mapping table for the `expand()` operation, a 20-bit to 32-bit expansion described in the task specification, enabling efficient back-mapping from expanded subkey material to seed candidates and eliminating redundant computation during the search phase.

Crucially, Seed2.0 uses a meet-in-the-middle strategy to avoid exhaustive enumeration over the 80-bit key space. It splits the computation at an intermediate state and matches partial forward and backward results, trading memory and preprocessing for a large reduction in runtime and making the attack feasible under benchmark constraints.

After recovering a candidate key, Seed2.0 thoroughly validates it. It checks the key against all 32 plaintext–ciphertext pairs and cross-verified outputs with the provided `decrypt` program, including early spot checks. Only after consistent agreement finishes, it decrypts the 100 ciphertexts in `ciphertexts.txt`, writes outputs to `plaintexts.txt`, and reports successful completion.

This case study is notable because Seed2.0's approach differs from the benchmark reference implementation. This divergence supports the claim that it performed structural program analysis and task-specific cryptanalytic reasoning rather than matching a known solution. The layered verification further indicates a rigorous engineering workflow in a complex setting.

### 5.1.2  Project Repository Construction

To evaluate the model's Project Repository Construction capabilities, we use the NL2Repo [25] task, which requires building a library-level Python repository from an empty workspace using only a requirements document. The deliverable is a Decouple-style configuration utility that separates configuration from code while remaining installable, testable, and usable as a standard Python package. Seed2.0 completes the task in 37 interaction rounds. As shown in Figure 4, the workflow follows three standard phases: Specification Analysis, Implementation & Configuration, and Testing & Debugging. Seed2.0 first reads `start.md` and extracts requirements for configuration decoupling, multi-source loading. Implementation centers on `decouple.py`, which defines eight core classes for source precedence, parsing, and type conversion.

Quality assurance relies on a pytest suite (`tests/test_decouple.py`) with 22 model-written test cases. Failures reveal issues that were fixed iteratively: installing a missing pytest dependency; addressing a Python 3.12-related configuration-reading import failure by switching to `parser.read_file()`; adding a guard for boolean casting on empty strings; moving fixtures to global scope to avoid class-scoping conflicts; and correcting a parameter-name mismatch (flat versus choices) to match the intended API. This pattern indicates systematic debugging guided by test output.

After these fixes, all 22 tests pass. The model also runs example code, validates installation via "`pip install -e .`", and confirms cross-directory imports, showing both functional correctness and packaging readiness.

**Task Overview**
… According to the start.md in the workspace, implement the entire project as per the requirements specified in the document, ensuring that the final product can be directly run in the current directory. The running requirements should comply with the <API Usage Guide> section of the document. Please complete this task step by step …

**Assistant:**

**[Stage 1: Specification Analysis]**

Let me first explore the workspace to understand the project structure and read the start.md file.

```
execute_bash(command="pwd && ls –la") …
str_replace_editor(command="view", path="/openhands/workspace/start.md")…
```

**[Stage 2: Implementation & Configuration]**

```
str_replace_editor(command="create", path="…/decouple.py", file_text="…")
str_replace_editor(command="create", path="…/setup.py", file_text="…")
str_replace_editor(command="create", path="…/README.rst", file_text="…")
…
```

**[Stage 3: Testing & Debugging]**

```
execute_bash(command="mkdir -p /openhands/workspace/tests") …
str_replace_editor(command="create", path="…/test_decouple.py", file_text="…")
```
                                                                    **[Create 22 Test Cases]**
```
execute_bash(command=" cd /openhands/workspace && python -m pytest tests/test_decouple.py -v")
```

| Issue | Likely cause | Solution / mitigation |
|---|---|---|
| Pytest missing | Dependency not installed | pip install pytest |
| read_config import error | Python 3.12 API change | Use parser.read_file() |
| Empty string exception | strtobool not handling empty input | if not value: return False |
| Fixture scope error | Class-scoped fixture not shared | Move fixture to module/global scope |
| choices argument confusion | Incorrect example parameter names (flat vs choices) | Align parameter names; use choices correctly |

…
✅ **Project Implementation Complete**
**Core Features Implemented …**
**Files Created …**
**Verification**
✅ All tests pass
✅ Imports work as specified in API Usage Guide
✅ Package installs successfully via pip
✅ Example usage works correctly
✅ All required features are implemented

The Python-Decouple library is now fully functional and ready to use!

**Figure 4** Case study of Seed2.0 performing the Project Repository Construction Task on NL2Repo [25].

Overall, Seed2.0 follows a structured, engineer-like workflow from requirements to implementation, testing, and debugging, producing a complete, convention-compliant repository with documentation.

### 5.1.3 Iterative Code Debugging

To systematically analyze the iterative code debugging capabilities of Seed2.0, we present a detailed case study from SWEBenchPro [23]. This case presents a zero-regression refactoring that relocates the Qt warning suppression logic and its tests from `log.py` to `qtlog.py`, correcting module boundaries without altering external behavior.

**Figure 5** Case study of Seed2.0 performing the Iterative Code Debugging task on SweBenchPro [23].

As shown in Figure 5, the refactoring preserves runtime semantics, the public API, and downstream compatibility across four key stages: Understanding & Location, Edition & Optimization, Local Verification, and Global Verification. The goal is not redesign but stronger verification through test relocation and extension. The core constraint is strict behavioral equivalence, identical behavior to the original, with the same effects on all non-Qt and non-matching logs.

More specifically, during Edition & Optimization, correctness is validated at the handler boundary. Non-matching warnings continue to appear in the console and captured logs, while matching warnings are fully suppressed before reaching handlers. Suppression follows the original semantics: a message is hidden if it equals the filter string or begins with the configured pattern; otherwise, it prints normally. The filtering mechanism, including `hide_qt_warning` and `QtWarningFilter`, is migrated to `qtlog.py`, with tests moved to `test_qtlog.py`. For backward compatibility, `log.py` retains an import layer so existing call sites remain intact, consolidating `qtlog.py` as the single source of truth.

Validation proceeds through five phases: boundary scanning, migration with compatibility wiring, unit

**Task Overview**

... I've uploaded a python code repository in the directory /testbed. Consider the following issue description:

<issue_description>
* GDrive: use lates PyDrive2 version to fix large file uploads (#3592) @shcheklein
* CLI: redirect warnings to logger (#3591) @shcheklein
* push/pull/status/metrics: add support for `--all-commits` (#3587) @efiop ...
</issue_description>

Can you help me implement the necessary changes to the repository so that the requirements specified in the <issue_description> are met? ...
Your task is to make the minimal changes to non-test files in the /testbed directory to ensure the <issue_description> is satisfied.

Follow these phases to resolve the issue:
Phase 1. READING: read the problem and reword it in clearer terms...
Phase 2. RUNNING: install and run the tests on the repository...
Phase 8. FINAL REVIEW: Carefully re-read the problem description and compare your changes with the base commit a338fad036bd9ac8bdd79ecec964f8c5558609c4...

**Assistant:**

**[Stage 1: Environment Exploration]**

```
execute_bash(command="ls –la /testbed") ...
str_replace_editor(command="view", path="/testbed/setup.py")...
```

**[Stage 2: Solving Issue #3581 (~6 rounds)]**

```
str_replace_editor(command="str_replace", path="/testbed/dvc/system.py", old_str ="
if platform.system() == "Windows" and sys.version_info < (3, 8):
    try:
        import speedcopy
        speedcopy.patch_copyfile()
    except ImportError:
        pass
", new_str="
if platform.system() == "Windows" and sys.version_info < (3, 8):
    import sys
    winver = sys.getwindowsversion()
    if winver.major > 6 or (winver.major == 6 and winver.minor >= 1):
        try:
            import speedcopy
            speedcopy.patch_copyfile()
        except ImportError:
            pass
")...
```

```
str_replace_editor(command="str_replace",path="…/dvc/system.py",old_str ="…",new_str="…")...
execute_bash(command="grep -r \"warnings\" /testbed/dvc --include=\"*.py\" | grep -E …")...
str_replace_editor(command="str_replace", path="…/setup.py", old_str="…", new_str="…")...
```

**[Stage 3: Solving Issue #3591 (5 rounds)]**

```
execute_bash(command="mkdir -p /openhands/workspace/tests") ...
str_replace_editor(command="view", path=" /testbed/dvc/main.py") ...
str_replace_editor(command="str_replace", path="…/dvc/main.py", old_str ="…", new_str="…")...
str_replace_editor(command="str_replace", path="…/dvc/main.py", old_str ="…", new_str="…")...
```

**[Stage 4: Solving Issue #3587 (50+ rounds)]**

```
# Add the --all-commits parameter to push/pull/fetch
str_replace_editor(command="str_replace", path="…/ dvc/command/data_sync.py"…) ...
str_replace_editor(command="str_replace", path="…/dvc/command/status.py"…) ...
str_replace_editor(command="str_replace", path="…/dvc/command/ metrics.py"…) ...
# Backend implementation
```

24

```
str_replace_editor(command="str_replace", path=".../dvc/repo/pull.py"...) ...
str_replace_editor(command="str_replace", path=".../dvc/repo/push.py"...) ...
str_replace_editor(command="str_replace", path=".../dvc/repo/fetch.py"...) ...
...
```

**[Stage 5: Solving Issue #3583 (25 rounds)]**

```
str_replace_editor(command="view", path="/testbed/dvc/progress.py")...
# Utilizing `threading. RLock()` to improve thread safety for tqdm
str_replace_editor(command="str_replace", path=".../dvc/progress.py"...) ...
# Multiple tests to verify the correct usage of RLock
...
```

**[Stage 8: Testing & Validation (~10 rounds)]**

```
execute_bash(command="find /testbed/tests/unit -name \"test_pull*\"") ...
execute_bash(command="cd /testbed && python -m dvc --help") ...
execute_bash(command="cd /testbed && python -m dvc status --help") ...
execute_bash(command="cd /testbed && python -m dvc metrics show --help") ...
execute_bash(command="cd /testbed && python -c \"
from dvc.command.metrics import _show_diff
print('Testing _show_diff with empty diff:')
result = _show_diff({})
print(f'Result: [{result}]')
print(f'Length: {len(result)}')
\"")...
...
```

1. ✅ **GDrive: use latest PyDrive2 version to fix large file uploads (#3592)**

   - Updated setup.py to require pydrive2>=1.6.3 (up from 1.4.8)

2. ✅ **CLI: redirect warnings to logger (#3591)** ...

3. ✅ **push/pull/status/metrics: add support for `--all-commits` (#3587)** ...

4. ✅ **progress: better threading (#3583)** ...

5. ✅ **Only import speedcopy on Win > 7 (#3581)** ...

6. ✅ **metrics: diff: don't print "No changes" (#3576)** ...

**Figure 6** The case of Seed2.0 with Complex Iterative Debugging Task on SWE-EVO.

regression testing, integration verification, and final cleanup. Results show 56 passing unit tests, verified import compatibility, accurate handler outputs, and complete reference coverage. A logger-parameter issue found during testing is fixed and revalidated, ensuring precise behavioral equivalence. The refactoring delivers comprehensive test coverage and production-level maintainability.

### 5.1.4 Complex Iterative Debugging

In this Release-Note-only SWE-EVO upgrade case, the agent's objective shifts from resolving a single, well-scoped issue (as in SWE-bench) to upgrading a library across versions using only release notes as guidance. With no explicit file-level directives, the model has to infer which components encoded the promised behavioral changes and treat the full unit test suite as the acceptance oracle, including both historically passing and historically failing tests. As shown in Figure 6, across 79 turns, the workflow converges on an exploration→implementation→verification loop: it first maps the repository and dependency surface (via directory inspection and `setup.py`) and identifies the codebase as DVC, then iteratively implements cross-cutting upgrades spanning OS compatibility, logging semantics, CLI/API consistency, concurrency safety, and user-facing output correctness.

The most involved change is the end-to-end introduction of an `-all-commits` flag across `push`, `pull`, `status`, and `metrics`, because it requires more than CLI plumbing: the agent has to thread a single semantic intent from command entry points through repository operations and into revision traversal, ultimately extending the branch iteration mechanism (`brancher.py`) and ensuring repository initialization invokes the updated traversal path (`init.py`).

**Figure 7** The competitive programming results across five ICPC.

In parallel, it handles platform constraints by gating `speedcopy` imports using `sys.getwindowsversion()` (a robust choice over fragile string parsing) and refines progress reporting by switching tqdm-related locking to `threading.RLock()` to mitigate re-entrancy and multi-thread contention. It also improves UX fidelity by suppressing the "`No changes`" message in `metrics diff` when suppression is requested, demonstrating attention to behavioral edge cases that tests commonly encode. Notably, the `PyDrive2` large-file upload item is resolved as a dependency conformance check, `pydrive2>=1.4.8` is already declared, so the agent avoids unnecessary code churn, aligning with a minimal-change upgrade strategy. Verification proceeds through repeated targeted test runs and final holistic validation, including checking CLI help output to ensure argument exposure matched the intended interface design.

### 5.1.5 Competition-Level Programming

To evaluate capabilities on competition-level programming, we assess model performance across five recent ICPC Official contests in late 2025. These five contests include the ICPC 2025 World Finals Baku (Sept 4, 2025), the ICPC 2025 Xi'an Regional Contest (Oct 19, 2025), the ICPC 2025 Chengdu Regional Contest (Oct 26, 2025), the ICPC 2025 Wuhan Regional Contest (Nov 2, 2025), and the ICPC 2025 Shanghai Regional Contest (Nov 23, 2025). We utilize test cases, time limits, and memory limits consistent with the official competitions. We compare the Pass@8 score for each model, which represents the percentage of problems solved when allowing each model up to eight independent generations per problem. The input prompt for the models contains only the problem statement and a simple instruction requiring the solution in C++, without the use of multi-turn interaction or external tool calls.

The results are shown in Figure 7. Seed2.0 Pro achieves a Pass@8 score of 73.02%, significantly outperforming GPT-5.2 and Gemini-3-Pro. Compared against human team performance, Seed2.0 Pro achieves Gold Medals in all five contests.

## 5.2 Real-world Application Operation

In productivity software, GUI interactions are state-dependent: toolchains vary with workspace or mode changes; parameters interlink across editing rounds; and interaction details (e.g., dialog prompts, focus management, double-clicking, drag-to-align) can trigger cascading failures. To examine whether Seed2.0 exhibits semantic understanding and long-horizon execution for complex interfaces, we evaluate it in two settings: (i) FreeCAD, a CAD environment emphasizing modeling context and constraint logic, and (ii) CapCut, a video-editing environment emphasizing timeline alignment and multi-track editing. These case studies illustrate how the model executes multi-step procedures while tracking interface state and self-corrects when UI errors arise.

### 5.2.1 FreeCAD Operation

As shown in Figure 8, we evaluate whether Seed2.0 can perform semantic GUI understanding in a stateful CAD environment beyond template-based interactions. The model successfully tracks contextual information, including active workbenches, bodies, operational modes, and attachment planes. When encountering errors

**Task Overview**

**Scenario:** You are working in FreeCAD 1.0.2, using the Part Design workbench to create a parametric solid. After completing the modeling, you need to read the solid's volume and area from the FreeCAD ...

**Task Objective:** First, create a circular sketch with a diameter of 80mm and extrude it by 40mm to create a boss. Then, create a rectangular sketch with a length of 50mm and a width of 30mm on the top surface of this boss and extrude it by 20mm to create a secondary boss. Read the final solid's volume and surface area values.

**Assistant:**

[Stage 1: Initial Start]

<think>Now we're on FreeCAD's start page. In the project tile area at the top left of the start page, click the New File icon. The "Parametric Part" tile is the New File icon for a parametric part—click it to create a new parametric solid model project.





The "Select attachment" dialog has popped up. According to the operating procedure, I need to choose the XY-Plane (reference plane). In the list, find "XY-Plane001 (Reference plane)", click to select it, then click OK to enter Sketcher.

[Stage 2: Self-Reflection to Correct Mis-Operation]

The OK button didn't respond just now. This may be because the click wasn't accurate enough, or the system was slow to respond. Look carefully for the OK button: it's the green checkmark icon with the text "OK" next to it. Click this button accurately again to enter the sketch editing interface.





We accidentally clicked the Pocket tool and an error dialog popped up. We need to close this dialog first by clicking OK, then return to the correct step and locate the tool for creating a circle. First, handle the error prompt by clicking the OK button in the dialog box.

The Sketcher menu is now expanded. We need to find the "Circle" option under "Sketcher geometries." Move the mouse to the "Sketcher geometries" submenu, expand it, then find and click the "Circle" tool to start drawing a circle sketch. This helps us select the correct circle tool and avoid clicking the wrong icon again.





We have selected the "Circle by center" tool. Next, We need to place the circle's center at the origin on the canvas (the intersection of the two axes). Click that intersection to set the center, then click another point to define the circle's edge and draw the circle. First, ...





The "Insert diameter" dialog has popped up. Change the diameter value to 80 mm. First select the existing value in the input box, then type 80 to replace the original 47.47 (this correctly updates the value).

[Stage 3: Correct Operation]

27

We're back in the correct interface and need to find the Pad tool again. On the top toolbar in the Part Design workbench, the Pad icon looks like an upward extrusion; it should be the button labeled "Pad," or we can locate it via the menu for better accuracy. Click "Part Design" on the top menu bar, then find and click "Pad" to avoid selecting the wrong icon.

We've placed the first corner point. Now move the mouse to the lower-right and click once to place the second corner and create the rectangle. The initial rectangle size can later be constrained to 50mm (length) and 30mm (width) with dimensions. Choose a suitable lower-right position and click to complete the rectangle.

The vertical length has been set to 30 mm. Click the OK button in the dialog to confirm this dimension constraint and finish constraining the rectangle's vertical edge. Now both rectangle dimensions (50 mm by 30 mm) are fully defined; next we can exit sketch editing and perform the ...

The Python console has printed the volume and surface area values: volume is 231061.92982974678 mm3, and surface area is 23306.192982974673 mm2. The task is complete; submit the results, stating the retrieved volume and surface area values </think>

The modeling is complete, and the final solid's volume and surface area have been read: Volume = 231061.93 mm$^3$, Surface area = 23306.19 mm$^2$ (values rounded to two decimal places). (Successfully calculated the surface area and volume of the drawn object.)

**Figure 8** The case of Seed2.0 with FreeCAD Operation Task.

such as incorrect tool selection or unresponsive dialogs, the model demonstrates expert-like recovery by formulating micro-plans (e.g., "create circle, constrain diameter, exit sketch") and executing them. The model's self-reflection mechanism functions as an online control loop under noisy UI feedback. Rather than simply documenting errors, the model diagnoses failure modes (misclicks, mode errors, or UI lag), applies targeted corrective actions (dismissing dialogs, reselecting tools, or reconfirming constraints), and resumes the parametric workflow. This behavior suggests a generalizable framework for robust GUI agents: maintain system invariants, restore operational context, and proceed with the intended task.

Based on this, correctness is anchored by redundant verification rather than "looks right." The model validates the final solid via both the Properties panel (Data → Shape → Volume/Area) and an independent Python readout (obj.Shape.Volume, obj.Shape.Area), turning a GUI procedure into a checkable outcome and reducing common silent failures like wrong selection or wrong document state.

### 5.2.2 CapCut Operation

Further, we demonstrate Seed2.0's capabilities on operational CapCut tasks. As shown in Figure 9, we present a case where Seed2.0 functions as an execution-oriented editor, tasked with producing a frame-accurate CapCut (Jianying Pro) project. The system maintains coherent goal states across multiple operations, video splitting, transition placement, audio splicing, and global effects application, while employing reflective self-correction to recover from UI-level errors. These errors include missed double-clicks, incorrect effect selection, and misaligned audio placement.

The trace reveals Seed2.0's core strength: adaptive control under partial failure. Upon detecting actions that fail to yield the intended UI state, the system re-identifies the correct control sequence. For instance, it re-targets the "Audio" panel, closes obstructive pop-ups, or switches from "click add" to "drag-to-join" for

## Task Overview

Please help me create a video using CapCut. The source materials are the video file `1.mp4` and the audio file `21.mp3`, both located on my desktop. Strictly follow the timing parameters provided below. The specific requirements are as follows:

**Step 1: Create a new project**

◆ Import the files `1.mp4`.

**Step 2: Video Editing**

◆ Set the first video split point at `00:00:01:24`.

◆ Set the second video split point at `00:00:03:28`. to these timestamps...

## Assistant:

**Think:** The editing task requires launching CapCut (Professional Edition) as a prerequisite. Since the desktop contains a CapCut Professional icon, the first step is to initiate the application by double-clicking this icon...
**Action:** left_double[0.025,0.268]

**Think:** CapCut did not appear to start successfully, which may indicate that the prior double-click was not registered accurately or that the startup latency is longer than expected... (Failure of left-button double-click; reflective retry)
**Action:** left_double[0.025,0.268]

**Think:** The CapCut interface has now loaded, and the "Start Creating" button is visible. Clicking this button enters the editing project workspace, which is required before importing assets...
**Action:** click[0.577,0.166]

**Think:** A yellow notification prompt is present. First, click "Got it" to dismiss the prompt. Next, click the "Import" control in the media panel (the blue plus sign) to open the file selection dialog and import 1.mp4...
**Action:** click[0.142,0.100]

**Think:** The asset 1.mp4 is now visible in the media panel. Click the "+" button at the lower-right corner of the video thumbnail to add the clip to the timeline. The "+" icon should be located at the bottom-right of the 1.mp4 thumbnail...
**Action:** click[0.157,0.260]

**Think:** The playhead was positioned inaccurately; it must be set precisely to 00:00:01:24. First, address any pop-up by clicking "Got it" to close the prompt, and then readjust the playhead location.
**Action:** click[0.973,0.507]

**Think:** The transition panel has finished loading. In the left-side category list, there is a "Dissolve" category. Click "Dissolve" to filter dissolve-related transitions, then locate the dissolve...

**Think:** The previous attempt selected the "Dissolve" effect itself rather than its plus icon, resulting in the effect being highlighted but not applied to the track...
**Action:** drag[0.366,0.220]

**Think:** Clicking "Audio" produced no response, which may be attributable to incomplete UI loading or inaccurate click targeting...
**Action:** click[0.050,0.065]

**[Stage 3: Self-Reflection to Correct Mis-Operation]**

**Think:** A key issue is observed: the audio track begins at approximately 00:00:03:00 rather than 00:00:00. Consequently, when the playhead is at 00:00:05:00, ... At this stage, the playhead should be dragged to the far-left of the timeline (00:00:00) to...
**Action:** drag[0.289 0.588]

**Think:** After the prior drag, the start position of the latter audio segment remains incorrect and requires further adjustment. Move the playhead to 00:00:02:29, then ...
**Action:** drag[0.208,0.843]

**[Stage 4: Failed Local Fixes and Refinement]**

**Think:** The "Snowflake Sparkle" effect clip is currently selected. Drag the right endpoint of this effect segment to 00:00:08:26 so that the effect covers the full duration of the video... The next step is to extend the right endpoint of the "Snowflake Sparkle" effect on the effects track to 00:00:08:26...
**Action:** drag[0.219,0.671]

**Figure 9** The case of Seed2.0 with CapCut Operation Task.

dissolve transitions. Crucially, the system re-anchors the workflow to canonical timeline invariants, project starts at 00:00:00:00 and splice boundary at 00:00:02:29, before proceeding. This re-anchoring mechanism prevents error propagation in multi-stage editing workflows.

The specification's HH:MM:SS:FR timecodes are meaningful only under a consistent project frame rate. Therefore, Seed2.0's emphasis on precise alignment and playhead repositioning serves as a critical control mechanism that guards against silent drift at split and splice boundaries. Within this frame-accurate regime, the system applies the requested "Dissolve" transitions by locating them in CapCut's transitions library and placing them at each clip junction created by the three splits. Similarly, the system extends a snow-style atmospheric overlay to cover the full 00:00:00:00–00:00:08:26 program window by applying an effect clip on the timeline and adjusting its duration to match the target interval.

**Task Overview**

**Input:** Recreate this website and add the following features: number rolling count animation, XP bar animation, 3D tilt effect + highlight sweep animation for Mission cards, and dynamic shimmer texture for progress bars.

**输入：** 重新创建这个网站，并添加以下功能：数字滚动计数动画、经验值（XP）进度条动画、任务卡片的 3D 倾斜效果 + 高光扫掠动画，以及进度条的动态闪烁纹理。

[Input]

[Generated Content]

这是一个完整、可独立运行的落地页实现，包含所有你要求的功能：数字滚动动画、XP 进度条动画、任务卡片的 3D 倾斜 + 高光扫过效果，以及带流光的进度条。

```html
<!DOCTYPE html>（此处折叠具体内容）
```

功能说明：

数字滚动动画：所有统计数字（顶部指标、XP 数值、任务/徽章/成功率）会在滚动进入视口时从 0 平滑增长到目标值，使用 Intersection Observer，确保只有用户看到内容时才触发动画。

XP 条动画：当个人资料卡进入视口时，XP 进度条会从 0 动画到目标宽度（2450/3000 对应 81.6%），并使用平滑的 ease-out 过渡效果。

任务卡片 3D 倾斜 + 高光扫过：
倾斜效果：悬停时卡片会根据鼠标位置自然倾斜，带透视深度和轻微放大。
高光扫过：悬停时会有一道明亮渐变扫过卡片，同时叠加与卡片主题色匹配的柔和发光效果。

动态流光进度条：XP 条的填充部分带有无限循环的移动渐变流光，使其更有质感和动态感。

页面是完全响应式的，视觉上与原始深色设计一致，并且只使用 CDN 资源，因此你可以将其保存为 HTML 文件后直接在任意浏览器中打开。

[Render website]

**Figure 10** The case of Seed2.0 with Image to code website HTML generation task.

### 5.2.3 Image-to-Code Generation

Here are two examples of code generation in real-world application scenario. As shown in Figure 10, the task is to recreate a website according to the given image. Some features are clearly required in the query prompt, including number rolling count animation, XP bar animation, 3D tilt effect + highlight sweep animation for Mission cards, and dynamic shimmer texture for progress bars. By saving the prediction code part as an HTML file and opening it directly in any browser, we could get a website almost the same as the input with all required features. As shown in Figure 11, the task is to recreate a 3D plot as a matplotlib expert, in which x-y projection, x-z projection, and y-z projection relationships are provided. By compiling the generated Python code, we could get a 3D phase space which follows the projection as input.

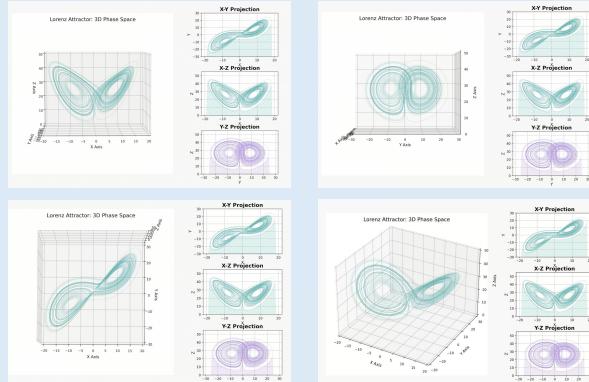**Figure 11** The case of Seed2.0 with Image to code 3D distribution generation task.

## 5.3 Multidisciplinary Scientific Research

To comprehensively examine Seed2.0's utility in cross-disciplinary scientific research, we evaluate its capacity to facilitate scientific discovery through two complementary task scenarios: code generation across scientific domains and cross-domain analytical reasoning.

## Task Overview

Quantum computation represents programs as unitary transformations acting on quantum states. In practice, arbitrary unitaries must be compiled into sequences of hardware-supported gates. Qiskit implements this compilation pipeline by combining quantum physics, numerical methods, and group-theoretic representations.

A core component of this pipeline is the Solovay–Kitaev (SK) algorithm, which approximates arbitrary single-qubit unitaries using a finite universal gate set. Mathematically, single-qubit gates are elements of the Lie group SU(2). However, Qiskit's implementation formulates the recursive core of the SK algorithm in terms of SO(3) rotation matrices, corresponding to rotations on the Bloch sphere.
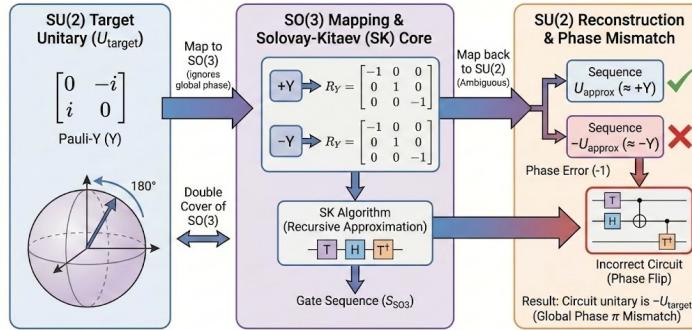


**Figure 1.** Global phase ambiguity in Qiskit's Solovay-Kitaev implementation, where the SU(2) → SO(3) mapping can lead to a -1 phase error upon reconstruction.

This design choice introduces a subtle but critical ambiguity: SU(2) is a double cover of SO(3). Two unitary matrices that differ only by a global phase of −1 correspond to the same SO(3) rotation. As a result, computations performed purely in SO(3) discard global phase information, which must be carefully restored when mapping results back to SU(2).

In the reported issue, applying the Solovay–Kitaev pass to a simple Pauli-Y gate produced a circuit whose unitary matrix was exactly the negative of the target operator. While physically equivalent up to global phase, the circuit reported an incorrect global phase value, leading to a mismatch at the operator level and violating user expectations.

Resolving this issue requires reasoning simultaneously about:
the mathematical relationship between SU(2) and SO(3),
the architectural design of Qiskit's SK implementation, and
the precise location where lost phase information must be reconstructed.

## Assistant:

### [Stage 1: Symptom Identification and Reproduction]

The agent begins by reproducing the reported behavior using a minimal example: decomposing a single-qubit Pauli-Y gate with the Solovay–Kitaev pass. It observes that the resulting circuit implements an operator numerically equal to **−Y**, while reporting a global phase of $\frac{3\pi}{4}$.

Crucially, the agent does not treat this as a numerical approximation error. Instead, it explicitly checks operator equality and verifies that the discrepancy is a **pure sign flip**, corresponding to a phase difference of $\pi$.

### [Stage 2: Mathematical Interpretation of the Error]

Let's interpret the sign flip through **group representation theory.** We can recognize that:
◆ the SK recursion operates in **SO(3)**,
◆ the mapping from **SU(2) → SO(3)** is two-to-one, and
◆ Therefore the recursion cannot distinguish U from -U.

At this stage, the error is identified as **structural rather than accidental**: global phase information is systematically lost during the SO(3)-based recursion and must be recovered afterward.

This reasoning step is independent of any specific code and reflects domain knowledge in quantum mechanics and Lie groups.

### [Stage 3: Error Localization in the Codebase]

The agent then localizes the error to a specific boundary in the implementation:
◆ the recursive SK algorithm itself is mathematically correct,
◆ the loss occurs at the **interface between the SO(3) recursion and the final SU(2) circuit construction**.

Concretely, the agent identifies the output phase assignment logic in Qiskit's solovay_kitaev.py

as the point where the missing ± sign ambiguity must be resolved. This is the *true root-cause localization* step: the error is neither in the basis generation nor in the commutator recursion, but in the **post-processing stage**.

**[Stage 4: Failed Local Fixes and Refinement]**

The agent explores naive fixes, such as adjusting the reported global phase by constant offsets. These attempts fix individual gates (e.g., Y) but break others (e.g., X or parameterized rotations), revealing that ad-hoc phase adjustments are insufficient.

This stage is critical: the agent learns that the solution must be **gate-agnostic** and derived from first principles, not from special-case heuristics.

**[Stage 5: Constructive Phase Recovery]**

The agent derives a general solution based on operator alignment. It reconstructs the SU(2) matrix implemented by the decomposed gate sequence and computes the unique global phase $\phi$ such that

$$e^{i\phi}U_{decomp} = U_{target}$$

This directly compensates for the SU(2)/SO(3) sign ambiguity and restores exact operator equality. The fix is mathematically principled and independent of the specific gate being decomposed.

**[Stage 6: Validation and Generalization]**

Finally, the agent validates the fix across multiple single-qubit gates and confirms that the decomposed circuits now match the original operators exactly (up to numerical precision), without introducing regressions.

The solution generalizes to arbitrary single-qubit unitaries and respects Qiskit's architectural constraint of performing the SK recursion in SO(3).

The solution generalizes to arbitrary single-qubit unitaries and respects Qiskit's architectural constraint of performing the SK recursion in SO(3).

**Figure 12** The case of Seed2.0 with Quantum Computing Task on AInsteinBench [27].

### 5.3.1 Multidisciplinary Scientific Research Coding

In real-world development of multidisciplinary scientific software, problems typically involve coupling among mathematical representations, physical conventions, and large-scale codebases, leading to subtle yet high-cost errors. Seed2.0 demonstrates robust generalization and debugging capabilities on such cross-domain tasks by identifying structural causes of anomalous behavior and producing verifiable fixes through dependency tracing and execution context analysis. The following sections present representative case studies from quantum computing, numerical general relativity, and computational chemistry.

**Quantum Computing.** On specialized and emerging tasks, Seed2.0 demonstrates remarkable generalization capabilities. In the quantum software engineering track of AInstein Bench [27], Seed2.0 successfully resolves a subtle bug in Qiskit's Solovay–Kitaev (SK) compiler stemming from a representation mismatch. The bug arises because Qiskit performs SK recursion in SO(3) (Bloch-sphere rotations), whereas target single-qubit gates belong to SU(2). Because SU(2) double-covers SO(3), the recursion cannot distinguish $U$ from $-U$, so global phase information can be dropped and must be restored when mapping the result back to an SU(2) circuit. The reported symptom, decomposing Pauli-$Y$ yields an operator exactly equal to $-Y$ with an inconsistent global phase, reflects this structural phase ambiguity rather than approximation error.

Seed2.0's approach to this problem exemplifies systematic debugging methodology. The system first reproduces the failure using a minimal test case, verifies that the discrepancy manifested as a pure sign flip (phase difference of $\pi$), and localizes the root cause to the post-processing boundary in `solovay_kitaev.py` where phase assignment occurs. Critically, Seed2.0 avoids modifying the mathematically correct recursion itself. After rejecting gate-specific constant-offset patches that would lack generalizability, it implements a principled solution: aligning operators by computing the phase $\phi$ such that $e^{i\phi}U_{\text{decomp}} = U_{\text{target}}$, then validating that this fix generalizes across all single-qubit gates without introducing regressions. This trajectory demonstrates agentic debugging through domain-theoretic diagnosis, precise code-level intervention at the appropriate
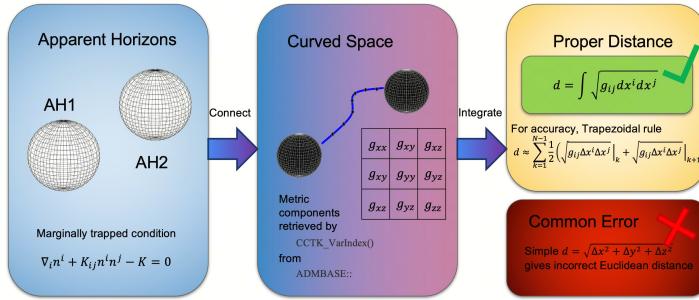
**Task Overview**

A core requirement when tracking multiple horizons is computing the proper distance between them. Unlike coordinate distance (simple Euclidean distance), proper distance in general relativity accounts for the curvature of spacetime. Mathematically, for a path connecting two points, the proper distance is given by integrating the line element:

$$s = \int \sqrt{\{g_{\{ij\}}, dx^i, dx^j\}}$$

where $g_{\{ij\}}$ is the spatial metric tensor. This integral needs to be evaluated numerically by:

1. Sampling points along a path between the two horizons
2. Interpolating the metric components at each point
3. Computing the infinitesimal proper length $ds$ at each step
4. Summing to obtain the total proper distance



The task was to implement AHFinder_dis.F, the subroutine responsible for computing the proper distance between two AH centers, using the methods above.

**Assistant:**

**[Stage 1: Codebase Exploration and Task Identification]**

... Ah there it is! It's referenced in line 5 of the makefile as AHFinder_dis.F, but it doesn't exist in the src directory! So I need to create AHFinder_dis.F. ... Line 1507 is where it calls the subroutine AHFinder_dis... so AHFinder_dis is called when find3 is true, meaning we have to find 3 horizons. This establishes the functional requirement: the subroutine must calculate proper distances between detected horizons when the find3 flag is enabled (for binary black hole simulations).

**[Stage 2: Data Structure Analysis]**

... From line 71-73 of AHFinder_dat I see that we have c0_0, c0_1, c0_2 as well as cc_0, cc_1, cc_2 and cs_0, cs_1, cs_2. Those are coefficients for 3 different horizons!
We identify that each horizon is characterized by:
◆ Center coordinates: xc_0, yc_0, zc_0 (and similarly for horizons 1 and 2)
◆ Shape coefficients: spherical harmonic expansion coefficients c0, cc, cs
◆ Status flags: status_old_0, status_old_1, status_old_2 indicating whether each horizon was found

**[Stage 3: Physical Interpretation and Algorithm Design]**

Proper distance between two events... is the length of the shortest spatial geodesic between them. But wait, do we really want the shortest proper distance? Or just the Euclidean distance? ...Perfect! Now I understand! Proper distance between two points in a curved 3D space is the minimum length of the path between them (the geodesic). To calculate it: 1. Choose a parametric path from point A to point B, 2. Interpolate the 3D metric along the path, 3. Integrate the length of the path segments using the metric.

**[Stage 4: Implementation]**

Metric Component Extraction:
call CCTK_VarIndex(v_gxx, "admbase::gxx")
call CCTK_VarIndex(v_gyy, "admbase::gyy")
call CCTK_VarIndex(v_gzz, "admbase::gzz")
call CCTK_VarIndex(v_gxy, "admbase::gxy")
call CCTK_VarIndex(v_gxz, "admbase::gxz")
call CCTK_VarIndex(v_gyz, "admbase::gyz")
...

```
Proper Distance Integration:
do i = 2, npoints
   ds = sqrt( &
      gxx(i) * dx**2 + &
      gyy(i) * dy**2 + &
      gzz(i) * dz**2 + &
      2 * gxy(i) * dx * dy + &
      2 * gxz(i) * dx * dz + &
      2 * gyz(i) * dy * dz &
   )
   total_distance = total_distance + ds
end do
...
```

**[Stage 5: Integration with Existing Infrastructure]**

```
So the pattern is: you have the filename, then you open it and write to it...
Now where is it for disf? It should have a pattern similar to...
disf = filestr(1:nfile) // "/ahf_d12" // trim(ext) // "tl"
inquire(file=disf, exist=disf_exists)
if (disf_exists) then
   open(11, file=disf, form='formatted', status='old', position='append')
else
   open(11, file=disf, form='formatted', status='replace')
   write(11, "(A18)") '" Horizon distance'
end if
write(11, "(2ES14.6)") cctk_time, total_distance
...
```

**[Stage 6: Error Handling and Robustness]**

```
The agent adds comprehensive error handling:
◆ Checks for valid interpolator and coordinate system handles
◆ Validates metric variable indices...
```

The solution generalizes to arbitrary single-qubit unitaries and respects Qiskit's architectural constraint of performing the SK recursion in SO(3).

**Figure 13** The case of Seed2.0 with General Relativity Task on AInsteinBench [27].

interface, and rigorous generalization testing.

This case underscores Seed2.0's capability in quantum compiler debugging, a demanding scientific software engineering task that cannot be resolved through pattern matching or localized edits alone. Success requires understanding the mathematical abstraction (Lie group coverings), mapping this abstraction onto the software architecture, and identifying where theoretical guarantees must be enforced in the implementation. The demonstrated trajectory indicates that Seed2.0 can effectively navigate tasks requiring a deep integration of mathematical reasoning and software engineering principles.

**General Relativity.**  Seed2.0 demonstrates robust performance on numerical relativity tasks from AInstein-Bench [27]. When tasked with implementing numerical-relativity functionality in the Einstein Toolkit's legacy Fortran codebase (Cactus framework), the agent autonomously diagnoses and resolves a critical integration failure: it identifies that `AHFinder_dis.F`, referenced by the build system, is missing from `src`, then traces its call site in `AHFinder.F` to determine the routine's role in multi-horizon tracking during binary black-hole simulations. This systematic debugging approach, locating missing components and verifying execution context, demonstrates practical software engineering capabilities beyond isolated code generation.

The implementation exhibits both physics insight and numerical rigor. Seed2.0 correctly interprets "distance" as proper distance in curved spacetime rather than coordinate separation, implementing $s = \int \sqrt{g_{ij}\,dx^i\,dx^j}$ with the full spatial 3-metric including all off-diagonal components ($g_{xx}, g_{yy}, g_{zz}, g_{xy}, g_{xz}, g_{yz}$). The solution properly accesses metric variables using Einstein Toolkit conventions, integrated with AHFinder's output

**Task Overview**

In Hartree-Fock and density functional theory calculations, the most computationally expensive step is constructing the Coulomb (J) and exchange (K) matrices from two-electron integrals. Density fitting (also called resolution-of-identity) is a technique that dramatically reduces this cost by approximating the four-center integrals using an auxiliary basis:

$$(\mu\nu|\lambda\sigma) \approx \sum_{P,Q} (\mu\nu|P)(P|Q)^{-1}(Q|\lambda\sigma)$$

PySCF's implementation stores the three-center integrals $(\mu\nu|P)$ in a packed triangular format and uses optimized C library routines to perform the integral transformations required for building J and K matrices... The reported bug manifests as a dramatic energy error (~0.9 Hartree) when computing energies using density fitting with complex MO coefficients, despite real and complex non-DF calculations producing identical results.



Resolving this issue requires:
◆ understanding the mathematical structure of J and K matrix construction
◆ recognizing that the low-level AO transformation routines only support real arrays
◆ implementing proper handling that decomposes complex density matrices(DMs) into real/imaginary parts and processes them correctly

**Assistant:**

**[Stage 1: Symptom Identification and Reproduction]**

The agent begins by creating a reproduction script based on the user's reported issue. Running the script confirms the problem:
E(rrhf) = –15.533592806500874
E(crhf) = –15.533592806497765 (complex, non-DF)
E(dfrrhf)= –15.533610777908837 (real, DF – correct)
E(dfchf)= –14.625586087385603 (complex, DF – WRONG!)...

**[Stage 2: Codebase Exploration and Error Localization]**

The agent systematically traces the code path:
scf.RHF(mol).density_fit() leads to pyscf/df/df_jk.py
The _DFHF class's get_jk method delegates to df.get_jk
The get_jk function at line 239 of df_jk.py is identified as the core computation
The agent examines the implementation and notes that it uses ctypes calls to C library functions:
fmmm = _ao2mo.libao2mo.AO2MOmmm_bra_nr_s2
fdrv = _ao2mo.libao2mo.AO2MOnr_e2_drv
ftrans = _ao2mo.libao2mo.AO2MOtranse2_nr_s2
...

**[Stage 3: Diagnostic Instrumentation]**

The agent adds detailed diagnostics to the reproduction script to isolate exactly which component is incorrect:

vj, vk = dfrhf.get_jk(dfrhf.mol, dm)
e_j = np.einsum('ij,ji->', vj, dm).real
e_k = np.einsum('ij,ji->', vk, dm).real

The output reveals:
e_j dfchf = 21.192171853892166
e_k dfchf = 9.277148921064926 <-- WRONG
e_j crhf = 21.19224885012347 <-- reference
...

**[Stage 4: Mathematical Analysis of the Error]**

The agent reasons why the K matrix fails for complex DMs. The K matrix construction involves:

$$K_{\mu\nu} = \sum_{\lambda\sigma} (\mu\lambda|\sigma\nu)D_{\lambda\sigma}$$

In the DF implementation, this becomes a series of matrix multiplications involving the three-center integrals and the density matrix. When the C routines receive a complex array but interpret it as real (taking only the real parts or misinterpreting memory layout), the result is incorrect.The agent verifies this hypothesis by manually splitting the complex DM:
...

**[Stage 5: Implementation of the Fix]**

The agent implements a fix at the beginning of the get_jk function that intercepts complex density matrices and handles them explicitly:

```
dms = numpy.asarray(dm)
dm_is_complex = numpy.iscomplexobj(dms)
if dm_is_complex:
    dms_real = dms.real
    dms_imag = dms.imag
...
```

**[Stage 6: Validation and Generalization]**

The agent validates the fix by running the original reproduction script:
E(dfrrhf)= -15.533610777908851
E(dfchf)= -15.533610777902968
The energies now match to within numerical precision (~$10^{-11}$). The agent also runs the existing test suite (test_df_jk.py) to ensure no regressions.

Linear Algebra Correctness: The fix relies on the linearity of the J and K matrix constructions. Because $J[D_{re} + iD_{im}] = J[D_{re}] + iJ[D_{im}]$ and similarly for K, processing real and imaginary parts separately and recombining gives mathematically correct results.

**Figure 14** The case of Seed2.0 with Computational Chemistry Task on AInsteinBench [27].

patterns, and includes robustness checks for invalid handles and absent horizons. While the production implementation incorporates additional refinements, grid symmetry handling, horizon–line intersections via spherical harmonics, and improved quadrature, the agent's solution captures the essential physics with appropriate simplifying assumptions.

This performance illustrates Seed2.0's capacity to navigate the multidimensional challenge space of scientific computing: translating physical concepts (Riemannian geometry, geodesic integration) into numerical methods (interpolation, quadrature) within complex software architectures (framework APIs, mixed-language codebases with F77/F90 syntax and preprocessor macros). The agent's successful integration of domain knowledge with system-level engineering reflects the core competencies that AInsteinBench is designed to evaluate.

**Computational Chemistry.** Meanwhile, Seed2.0 can identify and resolve a subtle defect in PySCF's density-fitting (DF/RI) implementation that manifests exclusively with complex-valued density matrices. Density fitting accelerates $J/K$ matrix construction by replacing four-center integrals with auxiliary-basis three-center tensors. The diagnostic symptom, a substantial energy deviation of approximately 0.9 Hartree in DF calculations with complex SCF inputs, occurs despite excellent agreement between real and complex non-DF energies (precision to $10^{-11}$ Hartree). This pattern immediately localizes the failure to the DF $J/K$ pipeline rather than general complex-number handling. Seed2.0 reproduces the error and traces the execution path to pyscf/df/df_jk.py, correctly identifying that the underlying computation routes through low-level C AO2MO kernels whose naming and packing conventions assume real-valued, symmetric arrays. By

**Task Overview**

Is it true that if $1 \le a_1 < a_2 < \ldots$ is a sequence of integers with $\lim\limits_{n \to \infty} a_n^{\frac{1}{2^n}} > 1$ then $\sum_{n=1}^{\infty} \frac{1}{a_n a_{n+1}}$ is irrational?

**Assistant:**

**[Stage 1: Rationality Implies Lower Bound on xₘ]**

... Assume $S = P/Q$ for coprime positive integers P, Q. For each $m \ge 1$, split the sum into the partial sum and tail:

$$S = S_m + R_m, \qquad S_m = \sum_{n=1}^{m} \frac{1}{a_n a_{n+1}}, \qquad R_m = \sum_{n=m+1}^{\infty} \frac{1}{a_n a_{n+1}}.$$

Let $D_m = a_1 a_2 \cdots a_{m+1}$. For $n \le m$, the term $\frac{D_m}{a_n a_{n+1}} = \left( \prod_{i=1}^{n-1} a_i \right) \left( \prod_{i=n+2}^{m+1} a_i \right)$ is an integer, so $K_m = D_m S_m \in \mathbb{Z}$. Rearranging:

$$D_m S = K_m + D_m R_m \implies x_m = D_m R_m = \frac{D_m P - K_m Q}{Q} \ldots$$

**[Stage 2: Loose Tail Bound and Recurrence]**

Because $\{a_n\}$ is strictly increasing, $a_{n+1} - a_n \ge 1$, so:

$$\frac{1}{a_n a_{n+1}} = \frac{1}{a_{n+1} - a_n} \left( \frac{1}{a_n} - \frac{1}{a_{n+1}} \right) \le \frac{1}{a_n} - \frac{1}{a_{n+1}}.$$

The liminf condition implies $a_n \to \infty$: there exists $N_1 \in \mathbb{N}$ such that for all $n \ge N_1$, $a_n^{\frac{1}{2^n}} > \frac{L+1}{2} > 1$, so $a_n > \left[ \frac{L+1}{2} \right]^{2^n} \to \infty$ as $n \to \infty$.

Summing the telescoping inequality from n = m+2 to $\infty$:

$$\sum_{n=m+2}^{\infty} \frac{1}{a_n a_{n+1}} \le \sum_{n=m+2}^{\infty} \left( \frac{1}{a_n} - \frac{1}{a_{n+1}} \right) = \lim_{N \to \infty} \left( \frac{1}{a_{m+2}} - \frac{1}{a_N} \right) = \frac{1}{a_{m+2}} \ldots$$

**[Stage 3: Logarithmic Variables and Boundedness of rₖ]**

Take natural logarithms and define:

$$b_k = \log a_k, \qquad Q_k = \log P_k = \sum_{i=1}^{k} b_i, \qquad r_k = \frac{Q_k}{2^k}.$$

The recurrence $a_k \le K \, P_{k-1}$ becomes:

$$b_k \le \log K + Q_{k-1}, \qquad for \ k \ge 3.$$

Since $Q_k = Q_{k-1} + b_k$, substitute to get:

$$Q_k \le Q_{k-1} + \log K + Q_{k-1} = 2Q_{k-1} + C, \qquad where \ C = \log K.$$

Divide by $2^k$:

$$\frac{Q_k}{2^k} \le \frac{Q_{k-1}}{2^{k-1}} + \frac{C}{2^k} \implies r_k \le r_{k-1} + \frac{C}{2^k}, \qquad for \ k \ge 3 \ldots$$

**[Stage 4: Identity Relating $r_k$ and $d_k$]**

Define the normalized log sequence:

$$d_k = \frac{b_k}{2^k} = \log \frac{a_k}{2^k}.$$

The liminf condition translates to:

$$\lim_{k \to \infty} d_k = l = \log L > 0.$$

Using $Q_k = Q_{k-1} + b_k$, divide by $2^k$ to derive the identity:

$$r_k = \frac{Q_{k-1}}{2^k} + \frac{b_k}{2^k} = \frac{r_{k-1}}{2} + d_k,$$

which holds for all $k \ge 1$ with $r_0 = \frac{Q_0}{2^0} = \log 1 = 0$ (since $P_0 = 1$, the empty product)...

**[Stage 5: Convergence of $r_k$ and $d_k$ ]**

We now show $\lim\limits_{k \to \infty} r_k = R$, so $r_k$ converges to $R$. Let $R' = \lim\limits_{k \to \infty} r_k$.

Choose subsequences:

◆ $k_m$ such that $r_{k_m} \to R$ as $m \to \infty$,

◆ $j_n$ such that $r_{j_n} \to R'$ as $n \to \infty$.

For each $j_n$, let k_m be the largest index satisfying $k_m \le j_n$, so $k_m \le j_n < k_{m+1}$. Since $k_{m+1} > j_n$, apply the inequality from Step 3:

$$r_{k_{m+1}} \le r_{j_n} + \frac{C'}{2^{j_n}} \ldots$$

...

**[Stage 7: Final Contradiction]**

For $m \geq N_2 - 1$, use the tight tail bound to write:
$$x_m = D_m R_m \leq \frac{2 D_m}{a_{m+1} a_{m+2}}.$$
Since $D_m = P_{m+1} = P_m a_{m+1}$, substitute to simplify:
$$x_m \leq \frac{2 P_m a_{m+1}}{a_{m+1} a_{m+2}} = \frac{2 P_m}{a_{m+2}} \cdots$$

... The proof proceeds by contradiction, assuming $S = \sum_{n=1}^{\infty} \frac{1}{a_n a_{n+1}} = \frac{P}{Q}$ is rational. We split $S$ into a partial sum $S_m$ and tail $R_m$, defining $D_m = a_1 a_2 \cdots a_{m+1}$ so that $x_m = D_m R_m$ is a positive rational number with denominator dividing $Q$, hence $x_m \geq \frac{1}{Q}$. A telescoping bound on $R_m$ yields the recurrence $a_k \leq K P_{k-1}$ (with $K = 2Q$) for $k \geq 3$, where $P_k$ is the product of the first $k$ terms of the sequence. Converting to logarithmic variables $r_k = \frac{\log P_k}{2^k}$, we show $r_k \leq r_{k-1} + \frac{C}{2^k}$ for some constant $C$, implying $r_n \leq r_m + \frac{C}{2^m}$ for all $n \geq m$. Using this, we prove $\liminf r_k = \limsup r_k = R$, so $r_k$ converges to $R$. Defining $d_k = \frac{\log a_k}{2^k}$, the identity $r_k = \frac{r_{k-1}}{2} + d_k$ implies $d_k \to D = \frac{R}{2}$, with $D \geq \log L > 0$ by the liminf condition. Using the convergence of $d_k$, we bound the tail $R_m$ by a geometric series, giving $R_m \leq \frac{2}{a_{m+1} a_{m+2}}$ for large $m$. Combining this with the growth rates of $P_m = e^{2^m r_m}$ and $a_{m+2} = e^{4 d_{m+2} 2^m}$, we find $x_m \leq 2 e^{2^m (r_m - 4 d_{m+2})} \to 0$ as $m \to \infty$, contradicting $x_m \geq \frac{1}{Q} > 0$. Thus $S$ is irrational.

**Figure 15** The case of Seed2.0 with Complex Mathematics Solution Task on Erdős.

instrumenting the computation to isolate energy contributions, Seed2.0 determines that $J$ remains accurate under the DF approximation while $K$ exhibited severe corruption. This diagnostic pattern is consistent with complex buffers being misinterpreted as real memory, a failure mode that produces plausible but silently incorrect outputs. The implemented fix is minimal yet mathematically rigorous. Seed2.0 intercepts complex density matrices, decomposes $D$ into real and imaginary components, executes the existing optimized real DF pathway twice, and recombines results via linearity: $K(D_{\mathrm{re}} + iD_{\mathrm{im}}) = K(D_{\mathrm{re}}) + iK(D_{\mathrm{im}})$. Post-patch validation confirms that DF complex energies match both DF real and non-DF references to numerical precision ($\sim 10^{-11}$ Hartree). Seed2.0 further verifies robustness by executing the existing test suite and ensuring comprehensive coverage of complex-DF behavior.

This resolution exemplifies Seed2.0's capability to integrate domain expertise, understanding $J/K$ matrix construction under density fitting, with systems-level debugging skills, including analysis of ctypes-to-C boundaries and packed storage conventions. The resulting patch maintains computational efficiency while preserving the stability of the real-valued fast path.

## 5.4 Multidisciplinary Scientific Research Analysis

In this section, we systematically evaluate Seed2.0's scientific research analysis capabilities across tasks that mirror real-world workflows. We design four thematic categories spanning theoretical derivations, computational modeling, and experimental design. Our assessment examines whether Seed2.0 produces outputs that are technically sound, internally consistent, and aligned with domain conventions, including explicit assumptions, appropriate methodological choices, and reproducible procedures. We further evaluate its justification of technical decisions and the transferability of outputs to downstream applications, thereby characterizing its performance across multidisciplinary research settings.

**Complex Mathematics Solution.** To demonstrate the model's capability in Seed2.0 to solve complex mathematical tasks, we present proofs of two Erdős Problems [6], generated through a iterative refinement pipeline with Seed2.0. Figure 15 shows a summary of the output and demonstrates that Seed2.0 can address frontier-level challenges in mathematical research. Human experts have verified these proofs for correctness, and we have formalized the proof of Erdős 1051 using Seed-Prover 1.5 [13]. The model demonstrates a strong capacity for systematic reasoning and rigorous proof construction. See Appendix E for evaluation details and full outputs.

**Task Overview**

**Context:** Cannabidiol (CBD), an important terpenoid compound from marijuana with no psychoactive effects, has become of great pharmaceutical interest for several health conditions. As CBD is a multitarget drug, there is a need to establish the molecular mechanisms by which CBD may exert therapeutic as well as adverse effects...

**Question:** You are tasked with defining a protocol to determine the binding site of CBD in the alpha7 nAChR receptor using coarse-grained molecular dynamics simulations using GROMACS. You are free to use any other software for system preparation and analysis, but the simulations should be run in GROMACS. You should make a step-by-step protocol describing in detail how each step will be carried out, describing the specific choices that you would make at each step and justifying each choice, specifying the following

◆ System preparation: Molecular model used (deposited structure, homology model, alphafold prediction, etc.), forcefield and box preparation, solvation, and any other adjustments and validations made to the model before the simulations
◆ Ligand parameterization: Describe how you would prepare and validate the coarse-grained ligand parameters. Specify the beads used, bead types, and which atoms will they correspond to (use the numbering provided by Atalay et al, 2020.
◆ Analysis: Transformations done to the trajectories before analysis, software used for analyzing the trajectories, features analyzed. \item Expected outcomes and suggested downstream tasks.

**Assistant:**

**[Stage 1: System Preparation]**

**1.1 Receptor Molecular Model Selection**
**Choice:** Experimental cryo-EM structures of human $\alpha$7 nAChR:
- Resting state: PDB 7EKI (3.2 Å resolution, homopentamer, embedded in lipid nanodiscs)
- Desensitized state: PDB 7EKJ (3.9 Å resolution, same system as 7EKI)...
**1.2 Force Field Selection**
**Choice:** Martini 3 force field ...
**1.3 Membrane, Box, Solvation, and Ion Preparation**
**Choice:** Physiologically relevant neuronal plasma membrane bilayer, assembled via CHARMM-GUI Martini Maker (v1.8): ...
**1.4 Pre-Simulation System Validation**
1. Steric clash check: Use gmx mindist to confirm the minimum distance between non-bonded beads is >1.1 nm (Martini van der Waals cutoff)...

**[Stage 2: Ligand (CBD) Parameterization]**

CBD atom numbering follows Atalay et al. (2020): ...
**2.1 CG Bead Mapping and Type Assignment**
Martini 3 mapping follows the rule of prioritizing functional group representation over exact 4-heavy-atom per bead. 25 heavy atoms in CBD are mapped to 7 beads...
**2.2 Bonded Parameter Derivation**
1. Obtain an optimized all-atom (AA) CBD structure via DFT (B3LYP/6-31G* level, standard for small molecule geometry)...
**2.3 CBD Parameter Validation**
Two critical properties are validated to ensure biologically relevant behavior:...
**2.4 Add CBD to System**
Add 50 CBD molecules (10 per $\alpha$7 subunit, ~10 mol% relative to lipids) to the bulk water and ...

**[Stage 3: Simulation Protocol (GROMACS)]**

All simulations use GROMACS 2023.2 with Martini 3 input parameters from the CG-Martini ...
1. Energy Minimization:
   - Algorithm: Steepest descent → conjugate gradient...
2. NVT Equilibration (1 ns):
   - Ensemble: NVT (constant number, volume, temperature)...
...

**[Stage 4: Trajectory Analysis]**

**4.1 Pre-Analysis Trajectory Transformations (GROMACS)**
1. PBC Correction and Alignment:...
2. Equilibration Truncation: Remove the first 1 $\mu$s of each production trajectory to exclude non-equilibrated data after restraint removal.
**4.2 Analysis Software**
- GROMACS 2023.2 (distance, RMSD, density, hydrogen bond analysis)

- MDAnalysis 2.6.0 (Python library for trajectory parsing and binding site analysis) ...
**4.3 Analyzed Features**
1. Binding Event Detection:
  - Define a binding event as a CBD molecule staying within 0.6 nm (Martini non-bonded cutoff) of the receptor for ≥10 ns (excludes transient collisions)...
...

**[Stage 5: Expected Outcomes]**

1. 1–3 major CBD binding sites, primarily at the lipid-facing TMD or TMD inter-subunit interfaces (consistent with membrane-mediated inhibition).
2. CBD binding poses are stabilized by H-bonds between CBD's hydroxyl groups and polar TMD residues, with lipophilic CBD chains inserted into the membrane/TMD hydrophobic core.
3. CBD binding to major sites reduces channel pore radius to <0.2 nm, stabilizing closed resting/desensitized states.
4. $K_d$ for the dominant binding site is in the sub-micromolar range, matching experimental $IC_{50}$...

**[Stage 6: Downstream Tasks]**

1. All-Atom (AA) MD Validation: Refine top CG binding poses with AA MD (AMBER force field) and compute binding free energies via FEP/MM-PBSA to confirm specificity.
2. Mutagenesis Validation: Design point mutations of key binding site residues (e.g., H-bond partners), test CBD binding via CG/AA MD, and validate functionally via electrophysiology assays.
3. Atypical Channel Opening Investigation: Use MSM to identify rare CBD-bound receptor conformations associated with the low-frequency atypical openings reported in the context ...

**Figure 16** Case study of Seed2.0 performing the Cannabidiol Scientific Analysis on FrontierScience.

**Cannabidiol (CBD)-Mediated Modulation Analysis.**   To assess Seed2.0's performance on complex biomolecular simulation tasks, we evaluate its ability to design a coarse-grained molecular dynamics (CG-MD) workflow for studying cannabidiol (CBD)–mediated modulation of the $\alpha7$ nicotinic acetylcholine receptor ($\alpha7$ nAChR). The task requires generating a complete protocol for identifying CBD binding sites, including system preparation, ligand parameterization, trajectory analysis, and specification of expected results with downstream applications. Critically, the model needs to justify all methodological choices and parameter settings, thereby demonstrating both practical competence in MD workflow construction and understanding of underlying simulation principles.

As shown in Figure 16, Seed2.0 exhibits several notable strengths. The model correctly identifies an appropriate PDB structure and outlines a standard procedure for molecular structure preparation. The selected density functional theory (DFT) level for quantum mechanical calculations is suitable for the system. Notably, Seed2.0 proposes using the partition coefficient ($\log P$) as a validation metric to assess force field accuracy, a thoughtful choice for evaluating hydrophobic ligand behavior. The workflow incorporates proper gmx commands to address periodic boundary condition (PBC) artifacts, an essential preprocessing step. Time definitions are clearly stated with concrete decision criteria and corresponding GROMACS commands. The model also suggests appropriate validation strategies, including all-atom simulations and free energy perturbation (FEP) and molecular mechanics Poisson-Boltzmann surface area (MMPBSA) calculations for cross-validation.

However, the output contains several limitations. Most significantly, the model hallucinates specific details regarding PDB entries and cited literature, requiring careful verification before implementation. Such inaccuracies can mislead researchers and compromise workflow reliability.

Overall, Seed2.0 generates a logically structured and reasonably comprehensive workflow with detailed operational steps and validation approaches aligned with standard practices. The output provides a useful initial template for protocol development that can be iteratively refined. Nevertheless, the hallucination issues represent a critical weakness that must be systematically addressed through manual verification before practical application.

**Maleimide Polyacetylene Analysis.**   To further evaluate Seed2.0's capabilities, we examine its performance on a more challenging task: analyzing a high-performance copolymer synthesis route based on a ROMP–precursor oxidation strategy. This two-step protocol provides access to electron-deficient olefin polymers that are difficult to obtain through conventional polymerization methods. The assessment evaluates the model's understanding

**Task Overview**

**Context:** The development of stable, high-conductivity n-type conjugated polymers is crucial for advancing organic electronics but lags behind p-type materials...

**Question:** Maleimide Polyacetylene (mPA), featuring an alternating vinylene (-CH=CH-) unit and N-alkylated maleimide unit backbone, is synthesized via a two-stage strategy:
1. ROMP: A N-alkylated maleimide-fused cyclobutene monomer (M) is polymerized using a Mo-based Schrock catalyst to yields a soluble, non-conjugated precursor polymer (P) containing alternating vinylene and N-alkylated dihydro-maleimide units.
2. Oxidation: The precursor P is converted to the fully conjugated mPA using triethylamine (TEA) and a mild oxidant (e.g., TCNQ or $I_2$).

Provide a comprehensive chemical analysis of this system, addressing: ...

**Assistant:**

**[Stage 1: Preliminary Structural Definitions]**

To ground the analysis, we first define key structures (all $sp^2$ carbons are denoted for clarity):
◆ Monomer (M): N-alkyl-3-azabicyclo[3.2.0]hept-6-ene-2,4-dione: a strained cyclobutene ring fused to an N-alkyl succinimide (dihydro-maleimide) moiety.
◆ Precursor Polymer (P): Non-conjugated ROMP product with repeat unit: -[CH=CH-CH$_2$-C(O)-N(R)-C(O)-CH$_2$]- (alternating vinylene and N-alkyl succinimide units).
◆ Conjugated mPA: Oxidized product with repeat unit: -[CH=CH-C(O)-CH=CH-C(O)-N(R)]- (alternating vinylene and N-alkyl maleimide units, fully conjugated)...

**[Stage 2: Strategic Rationale for Two-Stage ROMP Precursor Approach & Catalyst Choice]**

**Strategic Rationale for Two-Stage Approach**
The two-stage design directly addresses fundamental limitations of direct conjugated polymer synthesis, aligned with the core challenge of balancing low LUMO, backbone planarity, and processability:
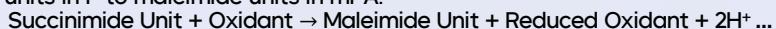1. **Processability Control:** Conjugated mPA is rigid, strongly π-stacked, and insoluble/infusible, making direct processing impossible...
2. **Structural Perfection & Control:** Direct copolymerization of acetylene and maleimide yields random, defect-rich polymers with variable regioregularity and high dispersity (Đ)...
3. **High Molecular Weight Access:** Direct conjugated polymer synthesis often stops at low MW due to premature precipitation of rigid conjugated chains...
4. **Avoidance of Crosslinking:** Polyacetylene and conjugated maleimide copolymers are prone to crosslinking during direct polymerization...

**Catalyst Choice: Mo-Based Schrock Catalyst**
Schrock catalysts have the general formula Mo(=CHR')(OR)$_2$(NAr) (Ar = bulky 2,6-diisopropylphenyl, OR = bulky alkoxide), chosen for three key reasons:...

**[Stage 3: Mechanism of Precursor P to mPA Conversion (TEA/Oxidant Conditions)]**

The overall transformation is selective oxidative dehydrogenation of succinimide (dihydro-maleimide) units in P to maleimide units in mPA:

Succinimide Unit + Oxidant → Maleimide Unit + Reduced Oxidant + 2H$^+$ ...

**Step 1: Reversible Deprotonation of Succinimide α-Hydrogens**
The α-hydrogens of succinimide (H on CH$_2$ adjacent to C=O) have a pKa ≈ 15 (stabilized by delocalization of the conjugate base over two carbonyls). Triethylamine (TEA, Et$_3$N, conjugate acid pKa ≈ 10) is a weak base that undergoes reversible deprotonation...

**Step 2: Oxidation of Carbanion to Radical Intermediate**
Mild oxidants (TCNQ, $I_2$) are strong electron acceptors that irreversibly oxidize the electron-rich carbanion I1 to a neutral radical:...

...

**Step 4: Second Oxidation to Form Maleimide Double Bond**
The electron-rich radical anion I3 loses one electron to a second oxidant molecule, forming the neutral conjugated maleimide unit:...

**[Stage 4: Key Structure-Property Relationships in mPA]**

**1. Electronic Characteristics: LUMO Level & n-Type Behavior**
n-type behavior requires a low LUMO (low = more negative energy, enabling electron injection and stability of negative charge carriers).
◆ LUMO Origin: Maleimide units are strong EWGs with two carbonyl groups that withdraw electron density from the conjugated backbone, stabilizing the π* (LUMO) orbital via inductive and resonance effects.

◆ Hückel Molecular Orbital (HMO) Derivation: For a conjugated chain, the LUMO energy is:
$$E_{LUMO} = \alpha + 2\beta\cos\left(\frac{\pi}{n+1}\right)...$$
**2. Electrical Conductivity: Backbone Planarity**...

**[Stage 5: Overall Significance for n-Type Conjugated Polymers]**

This approach addresses the core bottleneck in n-type CP development (balancing low LUMO, planarity, and processability) with transformative impact:
1. Solves the EWG-Planarity Tradeoff...
2. Enables Processable High-Performance n-Type CPs...

...
5. Closes the p-n Performance Gap: mPA's conductivity (~10 S/cm) and air stability (LUMO < -4.0 eV resists oxidation by $O_2/H_2O$) are among the highest for n-type CPs, matching the performance of commercial p-type materials...

**Figure 17** Case study of Seed2.0 performing the Maleimide Polyacetylene Scientific Analysis on FrontierScience.

of strategic route design, catalyst selection, an atypical coupling deprotonation–oxidation mechanism, and the relationships between polymer microstructure and macroscopic properties, thereby testing its grasp of complex organic reaction mechanisms, reaction selectivity, and fundamental structure–property correlations.

As shown in Figure 17, the model demonstrates strong mechanistic reasoning. In the ROMP step, it accurately identifies the ring-opening metathesis polymerization behavior of cyclobutene derivatives and correctly recognizes that highly strained, electron-deficient cyclobutene substrates require highly active Schrock-type catalysts. The model further notes that this strategy offers distinct advantages over traditional copolymerization approaches, including superior molecular-weight control, improved regularity of the alternating architecture, and effective suppression of side reactions. For the second-step TEA/oxidant coupling deprotonation–oxidation process, the model correctly identifies that triethylamine alone is too weak a base to drive a simple acid–base equilibrium and that the subsequent oxidation shifts the equilibrium forward. It appropriately proposes a pathway involving single-electron abstraction by the oxidant to generate radical intermediates, demonstrating sound reasoning for this less common organic oxidation mechanism.

The model also provides a thorough structure–property analysis. It correctly explains that the precursor polymer $P$ contains $sp^3$-hybridized carbon centers that increase backbone flexibility, resulting in good solubility and processability. After oxidation, formation of the mPA backbone creates a highly conjugated, planar structure that maximizes $\pi - \pi$ stacking, enabling efficient intrachain electron delocalization and interchain charge transport. The model recognizes that the maleimide moiety, as a strong electron-withdrawing group, effectively lowers the LUMO energy level to support an air-stable n-type semiconductor architecture.

However, the model exhibits one notable limitation. Although its macroscopic mechanistic narrative is logically consistent, it fails to identify the more specific charge-transfer (CT) complex mechanism operative in this reaction. Specifically, it does not recognize that TEA and the oxidant may first form a CT complex that serves as the reactive initiating species. Instead, the model relies on a more generic acid–base equilibrium coupled with electron-transfer explanation. Consequently, its account of how the weakly basic TEA efficiently triggers the reaction, particularly the key kinetic details, lacks mechanistic depth.

**Golgi Proteins Analysis.** We next analyze a case involving an integrative experimental design that uses mice as a model organism to examine genetic engineering and animal model construction, cell biology and subcellular fractionation, and multi-omics analytical workflows. This study demands substantial domain knowledge and involves several nontrivial technical details. At the level of genetic engineering and model generation, the central question is how to leverage the Cre-LoxP system to achieve cell-type-specific and spatiotemporally controlled gene expression.

As shown in Figure 18, the key challenge from a cell biology perspective is how to rapidly isolate highly purified Golgi apparatus from complex brain tissue while preserving molecular integrity. Conventional fluorescence-activated cell sorting (FACS) is unsuitable because it is time-consuming and can induce cellular stress, both of which compromise the intended downstream analyses. The multi-omics component spans sample preparation,

**Task Overview**

Context: Certain diseases arise due to cell type-specific defects, mutation, or dysregulation of proteins, metabolites, and lipids.Question: A student was interested in studying specific Golgi proteins, lipids, and metabolites that are unregulated specifically in astrocytes (a), microglia (b), oligodendrocytes (c), and neurons (d) in a mouse model. Describe in detail how the student will generate these mice and measure these contents in the specific regions.

Think step by step and solve the problem below. In your answer, you should include all intermediate derivations, formulas, important steps, and justifications for how you arrived at your answer. Be as detailed as possible in your response.

**Assistant:**

**[Stage 1: Preface]**

The question's mention of "unregulated" is a typo for dysregulated, consistent with the context of disease-related cell-type-specific defects. Below is a detailed, step-by-step strategy, including justifications, formulas, and intermediate steps, to generate the required mouse models and measure Golgi-specific molecules.

**[Stage 2: Generate Genetically Modified Mouse Lines]**

This step involves generating a universal Cre-dependent Golgi reporter line, crossing it with validated cell-type-specific Cre drivers, and introducing both into the disease mouse model (DMM) background.

**1.1 Generate Cre-Dependent Golgi Reporter Mouse Line**

We knock-in a Lox-Stop-Lox (LSL) Golgi-tagged transgene into the Rosa26 (R26) safe harbor locus (no positional expression variegation, ubiquitous expression) via CRISPR-Cas9...

**1.2 Obtain Validated Cell-Type-Specific Cre Drivers**

(Source: Jackson Laboratories, JAX; justifications for specificity):
1. Astrocytes: Tg(Aldh1l1-cre/ERT2)1Khakh/J (JAX #031008): Aldh1l1 is astrocyte-specific, with no expression in neural stem cells, neurons, or other glia (PMID 18203758)...

...
**1.4 Tamoxifen Induction and Expression Validation**
**1.4.1 Tamoxifen (TAM) Induction for CreERT2 Lines**
- Dissolve TAM in corn oil (20 mg/mL), administer via intraperitoneal (IP) injection at 100 mg/kg/day for 5 consecutive days to 8-week-old adult mice...

**[Stage 3: : Isolate Cell-Type-Specific Golgi from Specific Brain Regions]**

**2.1 Dissect Specific Brain Regions**
- Euthanize mice via AVMA-approved $CO_2$ asphyxiation, harvest brain within 1 minute, and place in ice-cold PBS with protease, phosphatase, and metabolism inhibitors to preserve molecular integrity...

**2.2 Golgi Immunoprecipitation (Golgi-IP)**
**2.2.1 Tissue Homogenization**
- Dounce homogenize tissue in 1 mL ice-cold Golgi Buffer (GB: 250 mM sucrose, 20 mM HEPES pH 7.4, 150 mM NaCl, 2 mM EDTA, 1 mM DTT, inhibitors) with 10 loose pestle strokes (gentle to preserve Golgi vesicle integrity)...
**2.2.3 Golgi Purity Validation**
Western blot of bead-bound Golgi to confirm purity:
- Golgi markers (Giantin, Man2a1): >90% of total protein signal...

**[Stage 4: Measure Golgi Proteins, Lipids, and Metabolites via Mass Spectrometry]**

Split bead-bound Golgi into 3 equal aliquots for proteomics, lipidomics, and metabolomics.
**3.1 Golgi Proteomics**
1. Elution: Elute intact Golgi vesicles from beads with 150 ng/μL 3xFLAG peptide.
2. Sample Preparation: Lyse Golgi with 1% SDS, reduce with DTT, alkylate with iodoacetamide, and digest with trypsin (1:50 trypsin:protein ratio) via filter-aided sample preparation (FASP) to remove detergents...

**3.2 Golgi Lipidomics**
1. Extraction: Perform modified Bligh-Dyer extraction on bead aliquots, adding deuterated internal lipid standards (10 pmol each of PC, PE, PS, PI, sphingomyelin) before extraction to normalize recovery...

**3.3 Golgi Metabolomics**
1. Extraction: Extract polar metabolites with 80% ice-cold methanol, adding deuterated internal metabolite standards (amino acids, nucleotides, sugars) before extraction.
2. LC-MS/MS Analysis: Separate polar metabolites on a HILIC column, acquiring data in positive and negative ionization modes...

**[Stage 5: Orthogonal Validation]**

Validate top dysregulated molecules via orthogonal methods to rule out mass spectrometry artifacts:
◆ Proteins: Western blot, IF co-localization with cell-type and Golgi markers
◆ Lipids: Thin-layer chromatography, targeted mass spectrometry
◆ Metabolites: Colorimetric/fluorometric assays, targeted mass spectrometry

This strategy enables the student to reliably identify Golgi-specific proteins, lipids, and metabolites dysregulated specifically in each cell type and brain region of the disease mouse model.

**Figure 18** Case study of Seed2.0 performing the Golgi Proteins Scientific Analysis on FrontierScience.

mass spectrometry measurements, and the analytical logic across proteomics, lipidomics, and metabolomics.

The proposed model exhibits several notable strengths. The experimental workflow is specified in fine-grained steps: beyond proposing a transgenic mouse strategy, it provides actionable CRISPR/Cas9 targeting details, including homology-arm design at the Rosa26 (R26) locus, guide RNA (gRNA) design considerations, and genotyping primer planning. It articulates a stringent quality-control framework at critical steps of Golgi immunoprecipitation (Golgi-IP), explicitly defining negative-control strategies (e.g., assaying endoplasmic reticulum, mitochondrial, and nuclear markers to exclude contamination) and offering quantitative purity assessment using a Pearson correlation coefficient (PCC)-based metric. For the multi-omics analysis, statistical criteria for differential testing are appropriately defined, specifically $FDR < 0.05$ and $|LFC| \geq 1.2$.

One potential limitation concerns the choice of tag protein. The model selects Giantin as a Golgi marker; however, Giantin is a large matrix protein that may introduce more nonspecific binding during immunoprecipitation than transmembrane Golgi proteins such as Tmem115, and may pose a greater risk of perturbing Golgi architecture.

Overall, Seed2.0 performs at a level exceeding typical PhD student expectations, demonstrating the ability to integrate experimental details across subdisciplines. Rather than stopping at high-level strategy, it produces a feasible protocol with implementable steps. One practical consideration is that many PhD students would rely on commercial kits or established reagents rather than developing every component de novo; nevertheless, within a knowledge-driven framework, the model's response is well-constructed and scientifically sound.

## 5.5   Automated Model-on-Model Behavioral Diagnostics

As evaluation ecosystems grow in scale and diversity, manual analysis of benchmark results becomes increasingly costly and difficult to standardize. We build an automated model-on-model diagnostic pipeline that uses LLMs to analyze the evaluation outcomes of peer models across heterogeneous benchmarks. The system aggregates metric scores and behavioral statistics, including token usage, formatting compliance, turn counts, emoji frequency, best-of-$N$ statistics and other derived metrics, alongside instance-level outputs, reasoning traces, and execution trajectories. Through a three-layer design that separates data processing, scenario-adaptive analytical workflows, and report synthesis, the framework scales to large case volumes and helps algorithm researchers efficiently iterate on models by exposing model weaknesses, strategy trade-offs, and concrete behavioral issues such as code-switching and repetitive CoT patterns. Table 15 shows that the pipeline scales across diverse benchmarks and produces structured findings with quantified throughput and resource statistics.

**Table 15  Scalable automated diagnostics for iterative model development.** The pipeline converts large-scale benchmark evaluations into structured behavioral reports across reasoning, planning, and tool-use domains, revealing capability trade-offs and systematic failures with measured throughput and efficiency.

| Benchmark | Representative Findings | Efficiency Profile |
|---|---|---|
| **XBench** | **Focus:** SearchAgent. <br> **Finding:** Seed2.0 Pro delivers the strongest overall performance on this benchmark while maintaining efficient search behavior. It consistently completes complex multi-step tasks through multi-source retrieval and iterative verification, but shows weakness in fine-grained boundary alignment and structured comparison scenarios. <br> **Evidence:** It ranks first with avg_score=0.64, ahead of GPT-5.2 High (0.57), Claude-Sonnet-4.5-thinking (0.48), and Gemini-2.5-Pro (0.24), while using 337.62 completion tokens and 277.49 reasoning tokens on average. In long multi-hop cases such as case_6 (48 calls), case_15 (44 calls), and case_52 (70+ calls), it converges through repeated verification; however, in boundary-sensitive tasks like case_28, case_33, case_35, and case_91, it receives score 0.0 due to condition misalignment or anchor drift. | **Cases per Model:** 100 <br> **Models:** 9 <br> **Time Cost:** 2435 s <br> **Tokens in:** 0.9 M <br> **Tokens out:** 0.1 M |
| **WorldTravel** | **Focus:** Real World Tasks. <br> **Finding:** Seed2.0 Pro demonstrates strong engineering stability, but its performance degrades under high coupling density due to weak second-order constraint reasoning and limited long-chain conflict repair. The primary gap with GPT-5.2 lies not in basic planning, but in maintaining constraint closure under cross-dependent temporal, structural, and pricing constraints. <br> **Evidence:** Seed2.0 Pro ranks 2nd with avg_score=0.233, trailing GPT-5.2 High (0.327, +0.09) while outperforming most peers. It maintains error_rate=0.00 and stable latency (31.39s), using 77% fewer reasoning tokens (1286 vs 5597) and 84% fewer completion tokens (1486 vs 9190) than GPT-5.2 High. Success cases concentrate in $\leq 15$ structured constraints, whereas failures increase sharply when constraints exceed 15 with cross-dependencies. Major failure patterns include temporal window misalignment (case_19, 27, 42), fine-grained ticket miscalculation (case_23, 55), structural ordering errors (case_17, 53), and missing long-chain conflict repair (case_29, 46). Across 60 head-model comparisons, 33% of failures stem from temporal misalignment, 30% from structural violations, 25% from pricing errors, and 28% from conflict detection breakdown. | **Cases per Model:** 150 <br> **Models:** 9 <br> **Time Cost:** 1571 s <br> **Tokens in:** 0.9 M <br> **Tokens out:** 0.04 M |
| **IMO–Bench** | **Focus:** Reasoning <br> **Finding:** Seed2.0 Pro is strong at pushing structured problems to a logical conclusion. In recursive, combinatorial, and formula-driven tasks, it can carry the reasoning chain step by step without breaking the internal logic. However, it struggles when a problem requires extreme-case construction or full boundary coverage. It often proves the main idea correctly but fails to completely seal upper–lower bounds or exhaust all branches. <br> **Evidence:** Seed2.0 Pro achieves avg_score=0.779 (Rank 4/9), showing stable execution quality. Its bon=0.87 is close to Gemini-3-Flash (0.91) and Gemini-3-Pro (0.92), indicating competitive peak capability. However, its won=0.66 is significantly lower than GPT-5.2 High (0.81), resulting in a bon–won gap of 0.211, which implies 21% sampling instability. Behaviorally, 65% of its incorrect cases involve extreme-value or boundary-construction tasks, and 80% involve full-quantifier or multi-branch exhaustion problems, aligning with the observed robustness gap. | **Cases per Model:** 400 <br> **Models:** 9 <br> **Time Cost:** 1741 s <br> **Tokens in:** 1.4 M <br> **Tokens out:** 0.06 M |

| Benchmark | Representative Findings | Efficiency Profile |
|---|---|---|
| $\tau^2$-**Bench** | **Focus:** Tool Use Agent.<br>**Finding:** Seed2.0 Pro high-score cases dynamically skip redundant device checks and resolve issues within a controlled tool-call budget, while strictly adhering to single-step JSON outputs; low-score cases follow mechanical enumeration, over-consume dialogue turns, and break structural constraints.<br>**Evidence:** Seed2.0 Pro high-score cases complete in 31.4 turns on average with 7.3 tool calls, whereas low-score cases expand to 55.0 turns ($+75.2\%$) and 9.6 calls ($+31.5\%$), often delaying billing checks by 30+ turns (e.g., case_51, case_54); Seed2.0 Pro scores 0.781 (3/9) using 721 completion tokens and 637 reasoning tokens, substantially fewer than GPT-5 High (2143 / 1946) despite a 0.17 score gap. | **Cases per Model:** 278<br>**Models:** 9<br>**Time Cost:** 3116 s<br>**Tokens in:** 1.8 M<br>**Tokens out:** 0.1 M |

# 6 Conclusion

The Seed2.0 model series takes a key step forward in the intelligent evolution of solving complex real-world tasks. First, the Seed team identifies users' real needs, and selects or abstracts a series of benchmarks based on these real needs and complex real-world scenarios to build a reliable evaluation system. Based on the reliable and forward-looking evaluation system, Seed2.0 focuses on solving long-tail knowledge problems and complex instruction following problems, thereby enhancing the reliability of the model in complex and long-range real-world tasks. In addition, Seed2.0 has world-leading reasoning intelligence, visual understanding capabilities, and search capabilities, which meet the most real needs of a large number of users. The model card provides a large number of real use cases to prove that the Seed2.0 model begins to have the ability to handle initial complex real-world tasks, thus bringing more value to hundreds of millions of users.

# References

[1] Seed-1.8 Model Card and Evaluation Overview. https://lf3-static.bytednsdoc.com/obj/eden-cn/lapzild-tss/ljhwZthlaukjlkulzlp/research/Seed-1.8-Modelcard.pdf, 2025. Accessed: 2026-02-xx.

[2] Niki Amini-Naieni, Kiana Amini-Naieni, Tengda Han, and Andrew Zisserman. Open-world text-specified object counting. *arXiv preprint arXiv:2306.01851*, 2023.

[3] Anthropic. The anthropic economic index: January 2026 report on ai work task evolution. Technical report, Anthropic, 2026. URL https://www.anthropic.com/research/anthropic-economic-index-january-2026-report.

[4] Mislav Balunović, Jasper Dekoninck, Ivo Petrov, Nikola Jovanović, and Martin Vechev. Matharena: Evaluating llms on uncontaminated math competitions. *arXiv preprint arXiv:2505.23281*, 2025.

[5] Victor Barres, Honghua Dong, Soham Ray, Xujie Si, and Karthik Narasimhan. $\tau^2$-bench: Evaluating conversational agents in a dual-control environment. *arXiv preprint arXiv:2506.07982*, 2025.

[6] T. F. Bloom. Erdős problem #1051. URL https://www.erdosproblems.com/1051.

[7] Chase Brower. Visual physics comprehension test, 2025.

[8] ByteDance-Seed. Beyondaime: Advancing math reasoning evaluation beyond high school olympiads. https://huggingface.co/datasets/ByteDance-Seed/BeyondAIME, 2025.

[9] Zikui Cai, Andrew Wang, Anirudh Satheesh, Ankit Nakhawa, Hyunwoo Jae, Keenan Powell, Minghui Liu, Neel Jay, Sungbin Oh, Xiyao Wang, Yongyuan Liang, Tom Goldstein, and Furong Huang. MORSE-500: A programmatically controllable video benchmark to stress-test multimodal reasoning. *CoRR*, abs/2506.05523, 2025.

[10] Meng Cao, Pengfei Hu, Yingyao Wang, Jihao Gu, Haoran Tang, Haoze Zhao, Jiahua Dong, Wangbo Yu, Ge Zhang, Ian Reid, and Xiaodan Liang. Video simpleqa: Towards factuality evaluation in large video language models. *CoRR*, abs/2503.18923, 2025.

[11] Aaron Chatterji, Thomas Cunningham, David J Deming, Zoe Hitzig, Christopher Ong, Carl Yan Shan, and Kevin Wadman. How people use chatgpt. Working Paper 34255, National Bureau of Economic Research, September 2025. URL http://www.nber.org/papers/w34255.

[12] Guo Chen, Yicheng Liu, Yifei Huang, Baoqi Pei, Jilan Xu, Yuping He, Tong Lu, Yali Wang, and Limin Wang. Cg-bench: Clue-grounded question answering benchmark for long video understanding. In *The Thirteenth International Conference on Learning Representations, ICLR 2025, Singapore, April 24-28, 2025*, 2025.

[13] Jiangjie Chen, Wenxiang Chen, Jiacheng Du, Jinyi Hu, Zhicheng Jiang, Allan Jie, Xiaoran Jin, Xing Jin, Chenggang Li, Wenlei Shi, Zhihong Wang, Mingxuan Wang, Chenrui Wei, Shufa Wei, Huajian Xin, Fan Yang, Weihao Gao, Zheng Yuan, Tianyang Zhan, Zeyu Zheng, Tianxi Zhou, and Thomas Hanwen Zhu. Seed-prover 1.5: Mastering undergraduate-level theorem proving via learning from experience, 2025. URL https://arxiv.org/abs/2512.17260.

[14] Joya Chen, Ziyun Zeng, Yiqi Lin, Wei Li, Zejun Ma, and Mike Zheng Shou. Livecc: Learning video LLM with streaming speech transcription at scale. In *IEEE/CVF Conference on Computer Vision and Pattern Recognition, CVPR 2025, Nashville, TN, USA, June 11-15, 2025*, pages 29083–29095, 2025.

[15] Liang Chen, Weichu Xie, Yiyan Liang, Hongfeng He, Hans Zhao, Zhibo Yang, Zhiqi Huang, Haoning Wu, Haoyu Lu, Yiping Bao, et al. Babyvision: Visual reasoning beyond language. *arXiv preprint arXiv:2601.06521*, 2026.

[16] Lin Chen, Jinsong Li, Xiaoyi Dong, Pan Zhang, Yuhang Zang, Zehui Chen, Haodong Duan, Jiaqi Wang, Yu Qiao, Dahua Lin, et al. Are we on the right way for evaluating large vision-language models? *Advances in Neural Information Processing Systems*, 37:27056–27087, 2024.

[17] Qiguang Chen, Libo Qin, Jinhao Liu, Dengyun Peng, Jiannan Guan, Peng Wang, Mengkang Hu, Yuhang Zhou, Te Gao, and Wanxiang Che. Towards reasoning era: A survey of long chain-of-thought for reasoning large language models. *arXiv preprint arXiv:2503.09567*, 2025.

[18] Junhao Cheng, Yuying Ge, Teng Wang, Yixiao Ge, Jing Liao, and Ying Shan. Video-holmes: Can MLLM think like holmes for complex video reasoning? *CoRR*, abs/2505.21374, 2025.

[19] Long Cheng, Jiafei Duan, Yi Ru Wang, Haoquan Fang, Boyang Li, Yushan Huang, Elvis Wang, Ainaz Eftekhar, Jason Lee, Wentao Yuan, et al. Pointarena: Probing multimodal grounding through language-guided pointing. *arXiv preprint arXiv:2505.09990*, 2025.

[20] Xianfu Cheng, Wei Zhang, Shiwei Zhang, Jian Yang, Xiangyuan Guan, Xianjie Wu, Xiang Li, Ge Zhang, Jiaheng Liu, Yuying Mai, et al. Simplevqa: Multimodal factuality evaluation for multimodal large language models. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 4637–4646, 2025.

[21] Daniel Cores, Michael Dorkenwald, Manuel Mucientes, Cees G. M. Snoek, and Yuki M. Asano. Tvbench: Redesigning video-language evaluation. *CoRR*, abs/2410.07752, 2024.

[22] Chao Deng, Jiale Yuan, Pi Bu, Peijie Wang, Zhong-Zhi Li, Jian Xu, Xiao-Hui Li, Yuan Gao, Jun Song, Bo Zheng, et al. Longdocurl: a comprehensive multimodal long document benchmark integrating understanding, reasoning, and locating. In *Proceedings of the 63rd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 1135–1159, 2025.

[23] Xiang Deng, Jeff Da, Edwin Pan, Yannis Yiming He, Charles Ide, Kanak Garg, Niklas Lauffer, Andrew Park, Nitin Pasari, Chetan Rane, et al. Swe-bench pro: Can ai agents solve long-horizon software engineering tasks? *arXiv preprint arXiv:2509.16941*, 2025.

[24] Kaustubh Deshpande, Ved Sirdeshmukh, Johannes Baptist Mols, Lifeng Jin, Ed-Yeremai Hernandez-Cardona, Dean Lee, Jeremy Kritz, Willow E Primack, Summer Yue, and Chen Xing. Multichallenge: A realistic multi-turn conversation evaluation benchmark challenging to frontier llms. In *Findings of the Association for Computational Linguistics: ACL 2025*, pages 18632–18702, 2025.

[25] Jingzhe Ding, Shengda Long, Changxin Pu, Huan Zhou, Hongwan Gao, Xiang Gao, Chao He, Yue Hou, Fei Hu, Zhaojian Li, et al. Nl2repo-bench: Towards long-horizon repository generation evaluation of coding agents. *arXiv preprint arXiv:2512.12730*, 2025.

[26] Shihan Dou, Ming Zhang, Zhangyue Yin, Chenhao Huang, Yujiong Shen, Junzhe Wang, Jiayi Chen, Yuchen Ni, Junjie Ye, Cheng Zhang, Huaibing Xie, Jianglu Hu, Shaolei Wang, Weichao Wang, Yanling Xiao, Yiting Liu, Zenan Xu, Zhen Guo, Pluto Zhou, Tao Gui, Zuxuan Wu, Xipeng Qiu, Qi Zhang, Xuanjing Huang, Yu-Gang Jiang, Di Wang, and Shunyu Yao. CL-bench: A Benchmark for Context Learning. *arXiv e-prints*, art. arXiv:2602.03587, February 2026. doi: 10.48550/arXiv.2602.03587.

[27] Mingxuan Du, Benfeng Xu, Chiwei Zhu, Xiaorui Wang, and Zhendong Mao. Deepresearch bench: A comprehensive benchmark for deep research agents. *arXiv preprint arXiv:2506.11763*, 2025.

[28] Xinrun Du, Yifan Yao, Kaijing Ma, Bingli Wang, Tianyu Zheng, King Zhu, Minghao Liu, Yiming Liang, Xiaolong Jin, Zhenlin Wei, et al. Supergpqa: Scaling llm evaluation across 285 graduate disciplines. *arXiv preprint arXiv:2502.14739*, 2025.

[29] Linxi Fan, Guanzhi Wang, Yunfan Jiang, Ajay Mandlekar, Yuncong Yang, Haoyi Zhu, Andrew Tang, De-An Huang, Yuke Zhu, and Anima Anandkumar. MineDojo: Building Open-Ended Embodied Agents with Internet-Scale Knowledge. *arXiv e-prints*, art. arXiv:2206.08853, June 2022. doi: 10.48550/arXiv.2206.08853.

[30] Chaoyou Fu, Yuhan Dai, Yongdong Luo, Lei Li, Shuhuai Ren, Renrui Zhang, Zihan Wang, Chenyu Zhou, Yunhang Shen, Mengdan Zhang, Peixian Chen, Yanwei Li, Shaohui Lin, Sirui Zhao, Ke Li, Tong Xu, Xiawu Zheng, Enhong Chen, Caifeng Shan, Ran He, and Xing Sun. Video-mme: The first-ever comprehensive evaluation benchmark of multi-modal llms in video analysis. In *IEEE/CVF Conference on Computer Vision and Pattern Recognition, CVPR 2025, Nashville, TN, USA, June 11-15, 2025*, pages 24108–24118, 2025.

[31] Ling Fu, Zhebin Kuang, Jiajun Song, Mingxin Huang, Biao Yang, Yuzhe Li, Linghao Zhu, Qidi Luo, Xinyu Wang, Hao Lu, et al. Ocrbench v2: An improved benchmark for evaluating large multimodal models on visual text localization and reasoning. *arXiv preprint arXiv:2501.00321*, 2024.

[32] Shenghao Fu, Qize Yang, Yuan-Ming Li, Yi-Xing Peng, Kun-Yu Lin, Xihan Wei, Jian-Fang Hu, Xiaohua Xie, and Wei-Shi Zheng. Vispeak: Visual instruction feedback in streaming videos. *CoRR*, abs/2503.12769, 2025.

[33] Xingyu Fu, Yushi Hu, Bangzheng Li, Yu Feng, Haoyu Wang, Xudong Lin, Dan Roth, Noah A Smith, Wei-Chiu Ma, and Ranjay Krishna. Blink: Multimodal large language models can see but not perceive. In *European Conference on Computer Vision*, pages 148–166. Springer, 2024.

[34] Google Cloud. Gemini enterprise: Release notes and workspace agent integration updates. Official Documentation, 2026. URL https://docs.cloud.google.com/gemini/enterprise/docs/release-notes.

[35] Tianrui Guan, Fuxiao Liu, Xiyang Wu, Ruiqi Xian, Zongxia Li, Xiaoyu Liu, Xijun Wang, Lichang Chen, Furong Huang, Yaser Yacoob, et al. Hallusionbench: an advanced diagnostic suite for entangled language hallucination and visual illusion in large vision-language models. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 14375–14385, 2024.

[36] Yunzhuo Hao, Jiawei Gu, Huichen Will Wang, Linjie Li, Zhengyuan Yang, Lijuan Wang, and Yu Cheng. Can mllms reason in multimodality? emma: An enhanced multimodal reasoning benchmark. *arXiv preprint arXiv:2501.05444*, 2025.

[37] Wei He, Yueqing Sun, Hongyan Hao, Xueyuan Hao, Zhikang Xia, Qi Gu, Chengcheng Han, Dengchang Zhao, Hui Su, Kefeng Zhang, et al. Vitabench: Benchmarking llm agents with versatile interactive tasks in real-world applications. *arXiv preprint arXiv:2509.26490*, 2025.

[38] Wenyi Hong, Yean Cheng, Zhuoyi Yang, Weihan Wang, Lefan Wang, Xiaotao Gu, Shiyu Huang, Yuxiao Dong, and Jie Tang. Motionbench: Benchmarking and improving fine-grained video motion understanding for vision language models. In *IEEE/CVF Conference on Computer Vision and Pattern Recognition, CVPR 2025, Nashville, TN, USA, June 11-15, 2025*, pages 8450–8460, 2025.

[39] Kairui Hu, Penghao Wu, Fanyi Pu, Wang Xiao, Yuanhan Zhang, Xiang Yue, Bo Li, and Ziwei Liu. Video-mmmu: Evaluating knowledge acquisition from multi-discipline professional videos. *arXiv preprint arXiv:2501.13826*, 2025.

[40] Liang Hu, Jianpeng Jiao, Jiashuo Liu, Yanle Ren, Zhoufutu Wen, Kaiyuan Zhang, Xuanliang Zhang, Xiang Gao, Tianci He, Fei Hu, et al. Finsearchcomp: Towards a realistic, expert-level evaluation of financial search and reasoning. *arXiv preprint arXiv:2509.13160*, 2025.

[41] Jen-Tse Huang, Dasen Dai, Jen-Yuan Huang, Youliang Yuan, Xiaoyuan Liu, Wenxuan Wang, Wenxiang Jiao, Pinjia He, and Zhaopeng Tu. Visfactor: Benchmarking fundamental visual cognition in multimodal large language models. *arXiv preprint arXiv:2502.16435*, 2025.

[42] Yichen Huang and Lin F Yang. Winning gold at imo 2025 with a model-agnostic verification-and-refinement pipeline. *arXiv preprint arXiv:2507.15855*, 2025.

[43] Zhenpeng Huang, Xinhao Li, Jiaqi Li, Jing Wang, Xiangyu Zeng, Cheng Liang, Tao Wu, Xi Chen, Liang Li, and Limin Wang. Online video understanding: Ovbench and videochat-online. In *IEEE/CVF Conference on Computer Vision and Pattern Recognition, CVPR 2025, Nashville, TN, USA, June 11-15, 2025*, pages 3328–3338, 2025.

[44] Naman Jain, King Han, Alex Gu, Wen-Ding Li, Fanjia Yan, Tianjun Zhang, Sida Wang, Armando Solar-Lezama, Koushik Sen, and Ion Stoica. Livecodebench: Holistic and contamination free evaluation of large language models for code. *arXiv preprint arXiv:2403.07974*, 2024.

[45] Carlos E Jimenez, John Yang, Alexander Wettig, Shunyu Yao, Kexin Pei, Ofir Press, and Karthik Narasimhan. Swe-bench: Can language models resolve real-world github issues? *arXiv preprint arXiv:2310.06770*, 2023.

[46] Jingyao Li, Jingyun Wang, Molin Tan, Haochen Wang, Cilin Yan, Likun Shi, Jiayin Cai, Xiaolong Jiang, and Yao Hu. Crossvid: A comprehensive benchmark for evaluating cross-video reasoning in multimodal large language models. *CoRR*, abs/2511.12263, 2025.

[47] Shilong Li, Xingyuan Bu, Wenjie Wang, Jiaheng Liu, Jun Dong, Haoyang He, Hao Lu, Haozhe Zhang, Chenchen Jing, Zhen Li, Chuanhao Li, Jiayi Tian, Chenchen Zhang, Tianhao Peng, Yancheng He, Jihao Gu, Yuanxing Zhang, Jian Yang, Ge Zhang, Wenhao Huang, Wangchunshu Zhou, Zhaoxiang Zhang, Ruizhe Ding, and Shilei Wen. MM-BrowseComp: A Comprehensive Benchmark for Multimodal Browsing Agents. *arXiv e-prints*, art. arXiv:2508.13186, August 2025. doi: 10.48550/arXiv.2508.13186.

[48] Yiming Liang, Yizhi Li, Yantao Du, Ge Zhang, Jiayi Zhou, Yuchen Wu, Yinzhu Piao, Denghui Cao, Tong Sun, Ziniu Li, Li Du, Bo Lei, Jiaheng Liu, Chenghua Lin, Zhaoxiang Zhang, Wenhao Huang, and Jiajun Zhang. Encyclo-K: Evaluating LLMs with Dynamically Composed Knowledge Statements. *arXiv e-prints*, art. arXiv:2512.24867, December 2025. doi: 10.48550/arXiv.2512.24867.

[49] Yuan Liu, Haodong Duan, Yuanhan Zhang, Bo Li, Songyang Zhang, Wangbo Zhao, Yike Yuan, Jiaqi Wang, Conghui He, Ziwei Liu, et al. Mmbench: Is your multi-modal model an all-around player? In *European conference on computer vision*, pages 216–233. Springer, 2024.

[50] Yuanxin Liu, Shicheng Li, Yi Liu, Yuxiang Wang, Shuhuai Ren, Lei Li, Sishuo Chen, Xu Sun, and Lu Hou. Tempcompass: Do video llms really understand videos? In Lun-Wei Ku, Andre Martins, and Vivek Srikumar, editors, *Findings of the Association for Computational Linguistics, ACL 2024, Bangkok, Thailand and virtual meeting, August 11-16, 2024*, pages 8731–8772. Association for Computational Linguistics, 2024.

[51] Yuanxin Liu, Kun Ouyang, Haoning Wu, Yi Liu, Lin Sui, Xinhao Li, Yan Zhong, Y. Charles, Xinyu Zhou, and Xu Sun. Videoreasonbench: Can mllms perform vision-centric complex video reasoning? *CoRR*, abs/2505.23359, 2025.

[52] Pan Lu, Hritik Bansal, Tony Xia, Jiacheng Liu, Chunyuan Li, Hannaneh Hajishirzi, Hao Cheng, Kai-Wei Chang, Michel Galley, and Jianfeng Gao. Mathvista: Evaluating mathematical reasoning of foundation models in visual contexts. *arXiv preprint arXiv:2310.02255*, 2023.

[53] Minh-Thang Luong, Dawsen Hwang, Hoang H Nguyen, Golnaz Ghiasi, Yuri Chervonyi, Insuk Seo, Junsu Kim, Garrett Bingham, Jonathan Lee, Swaroop Mishra, et al. Towards robust mathematical reasoning. In *Proceedings of the 2025 Conference on Empirical Methods in Natural Language Processing*, pages 35406–35430, 2025.

[54] Kaijing Ma, Xinrun Du, Yunran Wang, Haoran Zhang, Zhoufutu Wen, Xingwei Qu, Jian Yang, Jiaheng Liu, Minghao Liu, Xiang Yue, et al. Kor-bench: Benchmarking language models on knowledge-orthogonal reasoning tasks. *arXiv preprint arXiv:2410.06526*, 2024.

[55] Wentao Ma, Weiming Ren, Yiming Jia, Zhuofeng Li, Ping Nie, Ge Zhang, and Wenhu Chen. Videoeval-pro: Robust and realistic long video understanding evaluation. *CoRR*, abs/2505.14640, 2025.

[56] Yubo Ma, Yuhang Zang, Liangyu Chen, Meiqi Chen, Yizhu Jiao, Xinze Li, Xinyuan Lu, Ziyu Liu, Yan Ma, Xiaoyi Dong, et al. Mmlongbench-doc: Benchmarking long-context document understanding with visualizations. *Advances in Neural Information Processing Systems*, 37:95963–96010, 2024.

[57] Zeyao Ma, Bohan Zhang, Jing Zhang, Jifan Yu, Xiaokang Zhang, Xiaohan Zhang, Sijia Luo, Xi Wang, and Jie Tang. SpreadsheetBench: Towards Challenging Real World Spreadsheet Manipulation. *arXiv e-prints*, art. arXiv:2406.14991, June 2024. doi: 10.48550/arXiv.2406.14991.

[58] Ahmed Masry, Mohammed Saidul Islam, Mahir Ahmed, Aayush Bajaj, Firoz Kabir, Aaryaman Kartha, Md Tahmid Rahman Laskar, Mizanur Rahman, Shadikur Rahman, Mehrad Shahmohammadi, et al. Chartqapro: A more diverse and challenging benchmark for chart question answering. *arXiv preprint arXiv:2504.05506*, 2025.

[59] Mike A Merrill, Alexander G Shaw, Nicholas Carlini, Boxuan Li, Harsh Raj, Ivan Bercovich, Lin Shi, Jeong Yeon Shin, Thomas Walshe, E Kelly Buchanan, et al. Terminal-bench: Benchmarking agents on hard, realistic tasks in command line interfaces. *arXiv preprint arXiv:2601.11868*, 2026.

[60] 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.

[61] Samuel Miserendino, Michele Wang, Tejal Patwardhan, and Johannes Heidecke. SWE-Lancer: Can Frontier LLMs Earn $1 Million from Real-World Freelance Software Engineering? *arXiv e-prints*, art. arXiv:2502.12115, February 2025. doi: 10.48550/arXiv.2502.12115.

[62] Arsha Nagrani, Sachit Menon, Ahmet Iscen, Shyamal Buch, Ramin Mehran, Nilpa Jha, Anja Hauth, Yukun Zhu, Carl Vondrick, Mikhail Sirotenko, Cordelia Schmid, and Tobias Weyand. MINERVA: evaluating complex video reasoning. *CoRR*, abs/2505.00681, 2025.

[63] Junbo Niu, Yifei Li, Ziyang Miao, Chunjiang Ge, Yuanhang Zhou, Qihao He, Xiaoyi Dong, Haodong Duan, Shuangrui Ding, Rui Qian, Pan Zhang, Yuhang Zang, Yuhang Cao, Conghui He, and Jiaqi Wang. Ovo-bench: How far is your video-llms from real-world online video understanding? In *IEEE/CVF Conference on Computer Vision and Pattern Recognition, CVPR 2025, Nashville, TN, USA, June 11-15, 2025*, pages 18902–18913, 2025.

[64] OpenAI. The state of enterprise ai 2025: Adoption, depth, and workflow integration. Technical report, OpenAI, 2025. URL https://openai.com/index/the-state-of-enterprise-ai-2025-report/.

[65] Linke Ouyang, Yuan Qu, Hongbin Zhou, Jiawei Zhu, Rui Zhang, Qunshu Lin, Bin Wang, Zhiyuan Zhao, Man Jiang, Xiaomeng Zhao, et al. Omnidocbench: Benchmarking diverse pdf document parsing with comprehensive annotations. In *Proceedings of the Computer Vision and Pattern Recognition Conference*, pages 24838–24848, 2025.

[66] Piotr Padlewski, Max Bain, Matthew Henderson, Zhongkai Zhu, Nishant Relan, Hai Pham, Donovan Ong, Kaloyan Aleksiev, Aitor Ormazabal, Samuel Phua, et al. Vibe-eval: A hard evaluation suite for measuring progress of multimodal language models. *arXiv preprint arXiv:2405.02287*, 2024.

[67] Roni Paiss, Ariel Ephrat, Omer Tov, Shiran Zada, Inbar Mosseri, Michal Irani, and Tali Dekel. Teaching clip to count to ten. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 3170–3180, 2023.

[68] Shishir G. Patil, Huanzhi Mao, Charlie Cheng-Jie Ji, Fanjia Yan, Vishnu Suresh, Ion Stoica, and Joseph E. Gonzalez. The berkeley function calling leaderboard (bfcl): From tool use to agentic evaluation of large language models. In *Forty-second International Conference on Machine Learning*, 2025.

[69] Thinh Pham, Nguyen Nguyen, Pratibha Zunjare, Weiyuan Chen, Yu-Min Tseng, and Tu Vu. SealQA: Raising the Bar for Reasoning in Search-Augmented Language Models. *arXiv e-prints*, art. arXiv:2506.01062, June 2025. doi: 10.48550/arXiv.2506.01062.

[70] 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.

[71] Chiara Plizzari, Alessio Tonioni, Yongqin Xian, Achin Kulshrestha, and Federico Tombari. Omnia de egotempo: Benchmarking temporal understanding of multi-modal llms in egocentric videos. In *IEEE/CVF Conference on Computer Vision and Pattern Recognition, CVPR 2025, Nashville, TN, USA, June 11-15, 2025*, pages 24129–24138, 2025.

[72] ARC Prize. Arc agi: The $1 million artificial general intelligence prize. https://arcprize.org/arc-agi/1/, 2024.

[73] Shi Qiu, Shaoyang Guo, Zhuo-Yang Song, Yunbo Sun, Zeyu Cai, Jiashen Wei, Tianyu Luo, Yixuan Yin, Haoxu Zhang, Yi Hu, et al. Phybench: Holistic evaluation of physical perception and reasoning in large language models. *arXiv preprint arXiv:2504.16074*, 2025.

[74] Shanghaoran Quan, Jiaxi Yang, Bowen Yu, Bo Zheng, Dayiheng Liu, An Yang, Xuancheng Ren, Bofei Gao, Yibo Miao, Yunlong Feng, Zekun Wang, Jian Yang, Zeyu Cui, Yang Fan, Yichang Zhang, Binyuan Hui, and Junyang Lin. CodeElo: Benchmarking Competition-level Code Generation of LLMs with Human-comparable Elo Ratings. *arXiv e-prints*, art. arXiv:2501.01257, January 2025. doi: 10.48550/arXiv.2501.01257.

[75] Pooyan Rahmanzadehgervi, Logan Bolton, Mohammad Reza Taesiri, and Anh Totti Nguyen. Vision language models are blind. In *Proceedings of the Asian Conference on Computer Vision*, pages 18–34, 2024.

[76] David Rein, Betty Li Hou, Asa Cooper Stickland, Jackson Petty, Richard Yuanzhe Pang, Julien Dirani, Julian Michael, and Samuel R Bowman. Gpqa: A graduate-level google-proof q&a benchmark. In *First Conference on Language Modeling*, 2024.

[77] Jonathan Roberts, Mohammad Reza Taesiri, Ansh Sharma, Akash Gupta, Samuel Roberts, Ioana Croitoru, Simion-Vlad Bogolin, Jialu Tang, Florian Langer, Vyas Raina, et al. Zerobench: An impossible visual benchmark for contemporary large multimodal models. *arXiv preprint arXiv:2502.09696*, 2025.

[78] Ziyao Shangguan, Chuhan Li, Yuxuan Ding, Yanan Zheng, Yilun Zhao, Tesca Fitzgerald, and Arman Cohan. TOMATO: assessing visual temporal reasoning capabilities in multimodal foundation models. In *The Thirteenth International Conference on Learning Representations, ICLR 2025, Singapore, April 24-28, 2025*, 2025.

[79] Manasi Sharma, Chen Bo Calvin Zhang, Chaithanya Bandi, Clinton Wang, Ankit Aich, Huy Nghiem, Tahseen Rabbani, Ye Htet, Brian Jang, Sumana Basu, Aishwarya Balwani, Denis Peskoff, Marcos Ayestaran, Sean M. Hendryx, Brad Kenstler, and Bing Liu. ResearchRubrics: A Benchmark of Prompts and Rubrics For Evaluating Deep Research Agents. *arXiv e-prints*, art. arXiv:2511.07685, November 2025. doi: 10.48550/arXiv.2511.07685.

[80] Hui Shen, Taiqiang Wu, Qi Han, Yunta Hsieh, Jizhou Wang, Yuyue Zhang, Yuxin Cheng, Zijian Hao, Yuansheng Ni, Xin Wang, et al. Phyx: Does your model have the" wits" for physical reasoning? *arXiv preprint arXiv:2505.15929*, 2025.

[81] Weikang Shi, Aldrich Yu, Rongyao Fang, Houxing Ren, Ke Wang, Aojun Zhou, Changyao Tian, Xinyu Fu, Yuxuan Hu, Zimu Lu, et al. Mathcanvas: Intrinsic visual chain-of-thought for multimodal mathematical reasoning. *arXiv preprint arXiv:2510.14958*, 2025.

[82] Samuel Stevens. BioBench: A Blueprint to Move Beyond ImageNet for Scientific ML Benchmarks. *arXiv e-prints*, art. arXiv:2511.16315, November 2025. doi: 10.48550/arXiv.2511.16315.

[83] Jingqun Tang, Qi Liu, Yongjie Ye, Jinghui Lu, Shu Wei, An-Lan Wang, Chunhui Lin, Hao Feng, Zhen Zhao, Yanjie Wang, et al. Mtvqa: Benchmarking multilingual text-centric visual question answering. In *Findings of the Association for Computational Linguistics: ACL 2025*, pages 7748–7763, 2025.

[84] DeepConsult Team. DeepConsult: A deep research benchmark for consulting and business queries. https://github.com/youdotcom-oss/ydc-deep-research-evals, 2025. GitHub repository.

[85] Gemini Robotics Team, Saminda Abeyruwan, Joshua Ainslie, Jean-Baptiste Alayrac, Montserrat Gonzalez Arenas, Travis Armstrong, Ashwin Balakrishna, Robert Baruch, Maria Bauza, Michiel Blokzijl, et al. Gemini robotics: Bringing ai into the physical world. *arXiv preprint arXiv:2503.20020*, 2025.

[86] Minh V. T. Thai, Tue Le, Dung Nguyen Manh, Huy Phan Nhat, and Nghi D. Q. Bui. SWE-EVO: Benchmarking Coding Agents in Long-Horizon Software Evolution Scenarios. *arXiv e-prints*, art. arXiv:2512.18470, December 2025. doi: 10.48550/arXiv.2512.18470.

[87] Minyang Tian, Luyu Gao, Shizhuo Dylan Zhang, Xinan Chen, Cunwei Fan, Xuefei Guo, Roland Haas, Pan Ji, Kittithat Krongchon, Yao Li, Shengyan Liu, Di Luo, Yutao Ma, Hao Tong, Kha Trinh, Chenyu Tian, Zihan Wang, Bohao Wu, Yanyu Xiong, Shengzhu Yin, Minhui Zhu, Kilian Lieret, Yanxin Lu, Genglin Liu, Yufeng Du, Tianhua Tao, Ofir Press, Jamie Callan, Eliu Huerta, and Hao Peng. SciCode: A Research Coding Benchmark Curated by Scientists. *arXiv e-prints*, art. arXiv:2407.13168, July 2024. doi: 10.48550/arXiv.2407.13168.

[88] George Tsoukalas, Jasper Lee, John Jennings, Jimmy Xin, Michelle Ying, Yuxing Deng, and Zico Kolter. Putnambench: Evaluating neural theorem-provers on the putnam mathematical competition, 2024. URL https://github.com/trishullab/PutnamBench.

[89] Jordy Van Landeghem, Rubèn Tito, Łukasz Borchmann, Michał Pietruszka, Pawel Joziak, Rafal Powalski, Dawid Jurkiewicz, Mickaël Coustaty, Bertrand Anckaert, Ernest Valveny, et al. Document understanding dataset and evaluation (dude). In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 19528–19540, 2023.

[90] An Vo, Khai-Nguyen Nguyen, Mohammad Reza Taesiri, Vy Tuong Dang, Anh Totti Nguyen, and Daeyoung Kim. Vision language models are biased. *arXiv preprint arXiv:2505.23941*, 2025.

[91] Fei Wang, Xingyu Fu, James Y Huang, Zekun Li, Qin Liu, Xiaogeng Liu, Mingyu Derek Ma, Nan Xu, Wenxuan Zhou, Kai Zhang, et al. Muirbench: A comprehensive benchmark for robust multi-image understanding. *arXiv preprint arXiv:2406.09411*, 2024.

[92] Fengxiang Wang, Hongzhen Wang, Zonghao Guo, Di Wang, Yulin Wang, Mingshuo Chen, Qiang Ma, Long Lan, Wenjing Yang, Jing Zhang, et al. Xlrs-bench: Could your multimodal llms understand extremely large ultra-high-resolution remote sensing imagery? In *Proceedings of the Computer Vision and Pattern Recognition Conference*, pages 14325–14336, 2025.

[93] Haochen Wang, Xiangtai Li, Zilong Huang, Anran Wang, Jiacong Wang, Tao Zhang, Jiani Zheng, Sule Bai, Zijian Kang, Jiashi Feng, et al. Traceable evidence enhanced visual grounded reasoning: Evaluation and methodology. *arXiv preprint arXiv:2507.07999*, 2025.

[94] Ke Wang, Junting Pan, Weikang Shi, Zimu Lu, Houxing Ren, Aojun Zhou, Mingjie Zhan, and Hongsheng Li. Measuring multimodal mathematical reasoning with math-vision dataset. *Advances in Neural Information Processing Systems*, 37:95095–95169, 2024.

[95] Weihan Wang, Zehai He, Wenyi Hong, Yean Cheng, Xiaohan Zhang, Ji Qi, Shiyu Huang, Bin Xu, Yuxiao Dong, Ming Ding, and Jie Tang. Lvbench: An extreme long video understanding benchmark. *CoRR*, abs/2406.08035, 2024.

[96] Yubo Wang, Xueguang Ma, Ge Zhang, Yuansheng Ni, Abhranil Chandra, Shiguang Guo, Weiming Ren, Aaran Arulraj, Xuan He, Ziyan Jiang, et al. Mmlu-pro: A more robust and challenging multi-task language understanding benchmark. *Advances in Neural Information Processing Systems*, 37:95266–95290, 2024.

[97] Zexuan Wang, Chenghao Yang, Yingqi Que, Zhenzhu Yang, Huaqing Yuan, Yiwen Wang, Zhengxuan Jiang, Shengjie Fang, Zhenhe Wu, Zhaohui Wang, Zhixin Yao, Jiashuo Liu, Jincheng Ren, Yuzhen Li, Yang Yang, Jiaheng Liu, Jian Yang, Zaiyuan Wang, Ge Zhang, Zhoufutu Wen, and Wenhao Huang. Worldtravel: A realistic multimodal travel-planning benchmark with tightly coupled constraints, 2026. URL https://arxiv.org/abs/2602.08367.

[98] Zhaowei Wang, Wenhao Yu, Xiyu Ren, Jipeng Zhang, Yu Zhao, Rohit Saxena, Liang Cheng, Ginny Wong, Simon See, Pasquale Minervini, et al. Mmlongbench: Benchmarking long-context vision-language models effectively and thoroughly. *arXiv preprint arXiv:2505.10610*, 2025.

[99] Zihan Wang, Jiaze Chen, Zhicheng Liu, Markus Mak, Yidi Du, Geonsik Moon, Luoqi Xu, Aaron Tua, Kunshuo Peng, Jiayi Lu, et al. Aethercode: Evaluating llms' ability to win in premier programming competitions. *arXiv preprint arXiv:2508.16402*, 2025.

[100] Zirui Wang, Mengzhou Xia, Luxi He, Howard Chen, Yitao Liu, Richard Zhu, Kaiqu Liang, Xindi Wu, Haotian Liu, Sadhika Malladi, et al. Charxiv: Charting gaps in realistic chart understanding in multimodal llms. *Advances in Neural Information Processing Systems*, 37:113569–113697, 2024.

[101] 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.

[102] Ryan Wong, Jiawei Wang, Junjie Zhao, Li Chen, Yan Gao, Long Zhang, Xuan Zhou, Zuo Wang, Kai Xiang, Ge Zhang, et al. Widesearch: Benchmarking agentic broad info-seeking. *arXiv preprint arXiv:2508.07999*, 2025.

[103] Haoning Wu, Dongxu Li, Bei Chen, and Junnan Li. Longvideobench: A benchmark for long-context interleaved video-language understanding. In *Advances in Neural Information Processing Systems 38: Annual Conference on Neural Information Processing Systems 2024, NeurIPS 2024, Vancouver, BC, Canada, December 10 - 15, 2024*, 2024.

[104] Zijian Wu, Xiangyan Liu, Xinyuan Zhang, Lingjun Chen, Fanqing Meng, Lingxiao Du, Yiran Zhao, Fanshi Zhang, Yaoqi Ye, Jiawei Wang, Zirui Wang, Jinjie Ni, Yufan Yang, Arvin Xu, and Michael Qizhe Shieh. MCPMark: A Benchmark for Stress-Testing Realistic and Comprehensive MCP Use. *arXiv e-prints*, art. arXiv:2509.24002, September 2025. doi: 10.48550/arXiv.2509.24002.

[105] XAI. Realworldqa. *URL https://huggingface.co/datasets/xai-org/RealworldQA*.

[106] Yijia Xiao, Edward Sun, Tianyu Liu, and Wei Wang. Logicvista: Multimodal llm logical reasoning benchmark in visual contexts. *arXiv preprint arXiv:2407.04973*, 2024.

[107] Weiye Xu, Jiahao Wang, Weiyun Wang, Zhe Chen, Wengang Zhou, Aijun Yang, Lewei Lu, Houqiang Li, Xiaohua Wang, Xizhou Zhu, et al. Visulogic: A benchmark for evaluating visual reasoning in multi-modal large language models. *arXiv preprint arXiv:2504.15279*, 2025.

[108] Chenghao Yang, Yinbo Luo, Zhoufutu Wen, Qi Chu, Tao Gong, Longxiang Liu, Kaiyuan Zhang, Jianpeng Jiao, Ge Zhang, Wenhao Huang, et al. Mars-bench: A multi-turn athletic real-world scenario benchmark for dialogue evaluation. *arXiv preprint arXiv:2505.23810*, 2025.

[109] Jian Yang, Wei Zhang, Yizhi Li, Shawn Guo, Haowen Wang, Aishan Liu, Ge Zhang, Zili Wang, Zhoujun Li, Xianglong Liu, and Weifeng Lv. CodeSimpleQA: Scaling Factuality in Code Large Language Models. *arXiv e-prints*, art. arXiv:2512.19424, December 2025. doi: 10.48550/arXiv.2512.19424.

[110] John Yang, Kilian Lieret, Carlos E. Jimenez, Alexander Wettig, Kabir Khandpur, Yanzhe Zhang, Binyuan Hui, Ofir Press, Ludwig Schmidt, and Diyi Yang. Swe-smith: Scaling data for software engineering agents, 2025. URL https://arxiv.org/abs/2504.21798.

[111] Lihe Yang, Bingyi Kang, Zilong Huang, Zhen Zhao, Xiaogang Xu, Jiashi Feng, and Hengshuang Zhao. Depth anything v2. *Advances in Neural Information Processing Systems*, 37:21875–21911, 2024.

[112] Sihan Yang, Runsen Xu, Yiman Xie, Sizhe Yang, Mo Li, Jingli Lin, Chenming Zhu, Xiaochen Chen, Haodong Duan, Xiangyu Yue, et al. Mmsi-bench: A benchmark for multi-image spatial intelligence. *arXiv preprint arXiv:2505.23764*, 2025.

[113] Shunyu Yao, Howard Chen, Austin W Hanjie, Runzhe Yang, and Karthik Narasimhan. Collie: Systematic construction of constrained text generation tasks. In *12th International Conference on Learning Representations, ICLR 2024*, 2024.

[114] Chun-Hsiao Yeh, Chenyu Wang, Shengbang Tong, Ta-Ying Cheng, Ruoyu Wang, Tianzhe Chu, Yuexiang Zhai, Yubei Chen, Shenghua Gao, and Yi Ma. Seeing from another perspective: Evaluating multi-view understanding in mllms. *arXiv preprint arXiv:2504.15280*, 2025.

[115] Shuangshuang Ying, Zheyu Wang, Yunjian Peng, Jin Chen, Yuhao Wu, Hongbin Lin, Dingyu He, Siyi Liu, Gengchen Yu, YinZhu Piao, Yuchen Wu, Xin Gui, Zhongyuan Peng, Xin Li, Xeron Du, Libo Qin, YiXin Cao, Ge Zhang, and Stephen Huang. Retrieval-Infused Reasoning Sandbox: A Benchmark for Decoupling Retrieval and Reasoning Capabilities. *arXiv e-prints*, art. arXiv:2601.21937, January 2026. doi: 10.48550/arXiv.2601.21937.

[116] Fangchen Yu, Haiyuan Wan, Qianjia Cheng, Yuchen Zhang, Jiacheng Chen, Fujun Han, Yulun Wu, Junchi Yao, Ruilizhen Hu, Ning Ding, et al. Hipho: How far are (m) llms from humans in the latest high school physics olympiad benchmark? *arXiv preprint arXiv:2509.07894*, 2025.

[117] Xiang Yue, Yuansheng Ni, Kai Zhang, Tianyu Zheng, Ruoqi Liu, Ge Zhang, Samuel Stevens, Dongfu Jiang, Weiming Ren, Yuxuan Sun, et al. Mmmu: A massive multi-discipline multimodal understanding and reasoning benchmark for expert agi. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 9556–9567, 2024.

[118] Xiang Yue, Tianyu Zheng, Yuansheng Ni, Yubo Wang, Kai Zhang, Shengbang Tong, Yuxuan Sun, Botao Yu, Ge Zhang, Huan Sun, et al. Mmmu-pro: A more robust multi-discipline multimodal understanding benchmark. In *Proceedings of the 63rd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 15134–15186, 2025.

[119] Daoguang Zan, Zhirong Huang, Wei Liu, Hanwu Chen, Linhao Zhang, Shulin Xin, Lu Chen, Qi Liu, Xiaojian Zhong, Aoyan Li, et al. Multi-swe-bench: A multilingual benchmark for issue resolving. *arXiv preprint arXiv:2504.02605*, 2025.

[120] Xiangyu Zeng, Kefan Qiu, Qingyu Zhang, Xinhao Li, Jing Wang, Jiaxin Li, Ziang Yan, Kun Tian, Meng Tian, Xinhai Zhao, Yi Wang, and Limin Wang. Streamforest: Efficient online video understanding with persistent event memory. *CoRR*, abs/2509.24871, 2025.

[121] Chenchen Zhang, Yuhang Li, Can Xu, Jiaheng Liu, Ao Liu, Changzhi Zhou, Ken Deng, Dengpeng Wu, Guanhua Huang, Kejiao Li, Qi Yi, Ruibin Xiong, Shihui Hu, Yue Zhang, Yuhao Jiang, Zenan Xu, Yuanxing Zhang, Wiggin Zhou, Chayse Zhou, and Fengzong Lian. ArtifactsBench: Bridging the Visual-Interactive Gap in LLM Code Generation Evaluation. *arXiv e-prints*, art. arXiv:2507.04952, July 2025. doi: 10.48550/arXiv.2507.04952.

[122] Kaiyuan Zhang, Chenghao Yang, Zhoufutu Wen, Sihang Yuan, Qiuyue Wang, Chaoyi Huang, Guosheng Zhu, He Wang, Huawenyu Lu, Jianing Wen, et al. Mme-cc: A challenging multi-modal evaluation benchmark of cognitive capacity. *arXiv preprint arXiv:2511.03146*, 2025.

[123] Qinyan Zhang, Xinping Lei, Ruijie Miao, Yu Fu, Haojie Fan, Le Chang, Jiafan Hou, Dingling Zhang, Zhongfei Hou, Ziqiang Yang, et al. Inverse ifeval: Can llms unlearn stubborn training conventions to follow real instructions? *arXiv preprint arXiv:2509.04292*, 2025.

[124] Xinchen Zhang, Xiaoying Zhang, Youbin Wu, Yanbin Cao, Renrui Zhang, Ruihang Chu, Ling Yang, and Yujiu Yang. Generative universal verifier as multimodal meta-reasoner. *arXiv preprint arXiv:2510.13804*, 2025.

[125] Yilun Zhao, Haowei Zhang, Lujing Xie, Tongyan Hu, Guo Gan, Yitao Long, Zhiyuan Hu, Weiyuan Chen, Chuhan Li, Zhijian Xu, Chengye Wang, Ziyao Shangguan, Zhenwen Liang, Yixin Liu, Chen Zhao, and Arman Cohan. MMVU: measuring expert-level multi-discipline video understanding. In *IEEE/CVF Conference on Computer Vision and Pattern Recognition, CVPR 2025, Nashville, TN, USA, June 11-15, 2025*, pages 8475–8489, 2025.

[126] Zhicheng Zheng, Xin Yan, Zhenfang Chen, Jingzhou Wang, Qin Zhi Eddie Lim, Joshua B. Tenenbaum, and Chuang Gan. Contphy: Continuum physical concept learning and reasoning from videos. In *Forty-first International Conference on Machine Learning, ICML 2024, Vienna, Austria, July 21-27, 2024*, volume 235 of *Proceedings of Machine Learning Research*, pages 61526–61558, 2024.

[127] Zihan Zheng, Zerui Cheng, Zeyu Shen, Shang Zhou, Kaiyuan Liu, Hansen He, Dongruixuan Li, Stanley Wei, Hangyi Hao, Jianzhu Yao, et al. Livecodebench pro: How do olympiad medalists judge llms in competitive programming? *arXiv preprint arXiv:2506.11928*, 2025.

[128] Enshen Zhou, Jingkun An, Cheng Chi, Yi Han, Shanyu Rong, Chi Zhang, Pengwei Wang, Zhongyuan Wang, Tiejun Huang, Lu Sheng, et al. Roborefer: Towards spatial referring with reasoning in vision-language models for robotics. *arXiv preprint arXiv:2506.04308*, 2025.

[129] Junting Zhou, Jin Chen, Linfeng Hao, Denghui Cao, Zheyu Wang, Qiguang Chen, Chaoyou Fu, Jiaze Chen, Yuchen Wu, Ge Zhang, Mingxuan Wang, Wenhao Huang, and Tong Yang. Babe: Biology arena benchmark, 2026. URL https://arxiv.org/abs/2602.05857.

[130] Peilin Zhou, Bruce Leon, Xiang Ying, Can Zhang, Yifan Shao, Qichen Ye, Dading Chong, Zhiling Jin, Chenxuan Xie, Meng Cao, et al. Browsecomp-zh: Benchmarking web browsing ability of large language models in chinese. *arXiv preprint arXiv:2504.19314*, 2025.

[131] Yuhao Zhou, Yiheng Wang, Xuming He, Ruoyao Xiao, Zhiwei Li, Qiantai Feng, Zijie Guo, Yuejin Yang, Hao Wu, Wenxuan Huang, et al. Scientists' first exam: Probing cognitive abilities of mllm via perception, understanding, and reasoning. *arXiv preprint arXiv:2506.10521*, 2025.

[132] Liya Zhu, Peizhuang Cong, Aowei Ji, Wenya Wu, Jiani Hou, Chunjie Wu, Xiang Gao, Jingkai Liu, Zhou Huan, Xuelei Sun, et al. Lpfqa: A long-tail professional forum-based benchmark for llm evaluation. *arXiv preprint arXiv:2511.06346*, 2025.

[133] Chengke Zou, Xingang Guo, Rui Yang, Junyu Zhang, Bin Hu, and Huan Zhang. Dynamath: A dynamic visual benchmark for evaluating mathematical reasoning robustness of vision language models. *arXiv preprint arXiv:2411.00836*, 2024.

# A  FreeCAD Parametric Modeling: GUI Agent Case Study

## A.1  Task Overview

**Task:** Using FreeCAD 1.0.2, create a parametric solid model consisting of:
1. A cylindrical base: **Diameter 80mm, Height 40mm**
2. A rectangular boss on top: **50mm × 30mm × 20mm**

**Goal:** Calculate and verify the final solid's volume and surface area using Python scripting.

**Environment:** FreeCAD 1.0.2 on Ubuntu (Chinese interface, unit: mm)

## A.2  Complete Workflow Visualization

Figure 19 illustrates the key milestones of the 96-step modeling process, from software initialization to programmatic verification. The agent successfully navigated through multiple interface challenges, demonstrating robust error recovery and systematic problem-solving capabilities.



**Figure 19**  Key milestones in the FreeCAD modeling workflow. The workflow consists of nine critical phases: (1) Software launch and configuration, (2) Sketch creation with geometric constraints, (3) First pad extrusion, (4) Top face selection, (5-6) Secondary feature creation, (7) Surface selection for scripting, (8) Python console execution, and (9) Final verification results (**Volume: 231061.93 mm³, Surface Area: 23306.19 mm²**).

## A.3  Detailed Execution Trace

Below we present key reasoning traces from selected critical steps, highlighting the agent's decision-making process and error recovery strategies.

## Phase 1: Software Initialization (Steps 1-8)

### Steps 1-3: Launching FreeCAD

**Step 1:** First need to open FreeCAD software. On Ubuntu system, should **click the bottom-left nine-grid icon** (applications menu) to find and launch FreeCAD. Current desktop sidebar doesn't show FreeCAD icon, so first open the application list.

**Step 2-3:** Previous clicks on the nine-grid icon had no response, possibly due to inaccurate click position. Need to **click again at the exact center of the icon** to trigger the application list expansion.

### Steps 4-8: Language Configuration

**Step 4:** Now the application list shows the FreeCAD icon with red and blue "F" logo. Need to **double-click it** to launch FreeCAD software and begin the modeling task.

**Step 6-8:** FreeCAD startup complete, showing welcome interface. Need to set language to Chinese. **Click Language dropdown → Select Simplified Chinese → Click Done** to enter software.

## Phase 2: Base Sketch Creation (Steps 9-30)

### Steps 9-14: Creating Sketch on XY Plane

**Step 9-11:** On the start page's project card area, **click Parametric Part icon** to create new parametric solid project. Then click "Create Body" in the task panel.

**Step 12-14: Click Create Sketch** in the left task panel → In the "Select Attachment" dialog, choose **XY-Plane001 (Base Plane)** → Click OK to enter sketch editing.

### Steps 22-30: Circle Tool Selection Challenge

**Step 22-24:** Error recovery in progress. After multiple failed attempts clicking wrong tools (point tool, three-point arc tool), need to locate the correct circle creation tool. The icon is a simple circle shape in the sketch toolbar.

**Step 25-27:** Previous icon clicks were unsuccessful. Changed strategy: **Use menu bar → Sketcher → Sketch Geometry → Circle by Center** to accurately select the tool, avoiding icon confusion.

**Step 28-30:** Circle tool successfully selected. **Click at origin to place circle center**, then click to the right to place circle edge, drawing the circular sketch. **Press Esc to exit** drawing command.

## Phase 3: Dimensional Constraints (Steps 31-39)

### Steps 31-36: Constraining Circle Diameter

**Step 31-33:** Need to select this circle for dimension annotation. Direct clicking on circle edge failed to select. Changed approach: **Click 1-Circle in the Elements panel** (left sidebar) to select, which is more reliable.

**Step 34-36:** After selecting circle through element list, **right-click to open context menu** → Click "Dimension" → Select **Constrain Diameter** from submenu.

**Step 37-39:** "Insert Diameter" dialog appeared. Need to change diameter value to 80mm. **Type 80** to replace the previous 47.47mm value. Click OK to confirm this dimension constraint, completing the circle sketch dimension constraint.

## Phase 4: First Pad Extrusion (Steps 40–52)

### Steps 40-49: Accessing Pad Feature

**Step 41-45:** Multiple failed attempts. Clicked wrong tools in toolbar (created datum plane, sketch tools), triggering error dialogs.
**Step 47-49:** Corrected approach: **Menu bar → Part Design → Create Additive Feature → Pad** to accurately locate the tool, avoiding icon misidentification.
**Step 50-52:** Pad parameters dialog appeared with default length 10.00mm. **Type 40** to replace default value, click OK to confirm and generate first pad.
**Milestone:** Successfully created cylindrical base with Ø80mm diameter and 40mm height.

## Phase 5: Rectangular Sketch Creation (Steps 53–78)

### Steps 53-58: Drawing Rectangle on Top Face

**Step 53-55:** According to procedure step 20, need to **click and select the base's top face** in 3D view. Then click "Create Sketch" to start drawing rectangular sketch on this face.
**Step 56-58: Click rectangle tool**, click once in upper-left area of canvas to place first corner point, click in lower-right area to place second corner, generating rectangle. Press Esc to exit rectangle drawing command.

### Steps 59-77: Applying Rectangle Constraints

**Step 59-64: Horizontal dimension:** Select rectangle's upper horizontal edge → Right-click → "Horizontal Distance Constraint" → **Enter 50mm**.
**Step 65-73: Vertical dimension:** Direct clicking on left vertical edge failed to select. Retried through element list: **Click 1-Line in Elements panel** → Right-click → "Vertical Distance".
**Step 76-77:** "Insert Length" dialog appeared. **Type 30** to set vertical length to 30mm. Click OK to confirm.
**Step 78:** Click Close button in left task panel to exit sketch editing, preparing for second pad extrusion.

## Phase 6: Second Pad Extrusion (Steps 79–82)

### Steps 79-82: Creating Rectangular Boss

**Step 79-81:** In task panel **click Pad**, wait for parameters dialog to load. Change length from default 10.00mm to **20mm**.
**Step 82:** Click OK to confirm, **generating second pad**. At this point, compound solid modeling is complete.

## Phase 7: Programmatic Verification (Steps 83–96)

### Steps 83-87: Opening Python Console

**Step 83-84:** Click left side's **Model tab** to switch to model view, then **select Pad001** in model tree for subsequent Python console command input to read volume and surface area.
**Step 85-87:** Open Python console via **menu bar: View → Panels → Python console**.

> **Steps 88-96: Executing Verification Script**
>
> **Step 88-95:** Sequentially input four Python commands in console:
> **Command 1:** `obj = Gui.Selection.getSelection()[0]` — Retrieve selected object
> **Command 2:** `vol = obj.Shape.Volume` — Calculate volume
> **Command 3:** `area = obj.Shape.Area` — Calculate surface area
> **Command 4:** `(vol, area)` — Display results
> **Step 96 — Final Results:**
> - **Volume = 231061.93 mm³**
> - **Surface Area = 23306.19 mm²**
>
> **Task completed successfully!** Modeling finished and final entity's volume and surface area values read.

## A.4  Agent Performance Analysis

### Execution Statistics

- **Total Steps:** 96 sequential operations
- **Error Recovery Instances:** 8 failed tool selections, all successfully recovered
- **UI Navigation Strategy:** Shifted from toolbar icons to menu-based navigation after failures
- **Element Selection Strategy:** Adopted Elements panel method after direct clicking proved unreliable
- **Verification Precision:** 6 decimal places via Python API (vs. typical 2 decimals from GUI inspection)

### Key Findings

1. **Adaptive Error Recovery:** When direct toolbar icon clicking failed (Steps 22-27, 44-47), the agent systematically switched to menu-based navigation, demonstrating flexible problem-solving rather than repetitive failed attempts.

2. **Robust Selection Strategy:** After encountering unreliable direct geometry clicking (Steps 31-33, 65-73), the agent learned to consistently use the **Elements panel** as a more reliable selection mechanism.

3. **System Response Awareness:** The agent exhibited patience with UI delays by using explicit `wait` commands (Steps 42, 50, 70), preventing premature actions that could derail the workflow.

4. **Scriptable Verification:** Rather than relying on visual GUI property inspection, the agent employed FreeCAD's Python API for precise numerical verification, yielding higher-precision results critical for engineering applications.

## A.5  Workflow Summary

**Table 16** Key operations and parameters in the modeling workflow

| Steps | Operation | Critical Parameters / Challenges |
|---|---|---|
| 1-8 | Software Launch | Language config: Simplified Chinese, Units: mm |
| 9-16 | Base Sketch Setup | XY plane selection, "Circle by Center" tool |
| 17-39 | Dimension Constraint | **80mm diameter**, Element panel selection required |
| 40-52 | First Pad | **40mm height**, Menu navigation after toolbar failures |
| 53-58 | Top Face Sketch | Face selection in 3D view, rectangle tool |
| 59-78 | Rectangle Constraints | **50mm** × **30mm**, Vertical dimension via element list |
| 79-82 | Second Pad | **20mm height**, Completing compound geometry |
| 83-96 | Python Verification | Shape.Volume & Shape.Area APIs, 6-digit precision |

## A.6 Lessons for GUI Agent Design

This case study reveals several design considerations for future GUI automation systems:

1. **Icon Ambiguity:** CAD toolbars contain visually similar icons (circle, arc, point). Agents would benefit from OCR-based tooltip reading or enhanced icon classification models.

2. **Selection Reliability:** Direct geometric entity clicking proves unreliable in complex 3D environments. Structured selection via hierarchical element trees provides more robust alternatives.

3. **Fallback Strategies:** Menu-based navigation, while slower, offers unambiguous tool access. Agents should maintain hybrid strategies (toolbar-first, menu-fallback).

4. **Verification Paradigm:** Programmatic verification via scripting APIs provides superior accuracy over visual GUI inspection, particularly for numerical engineering tasks.

# B TerminalBench2.0-Cryptanalysis Agent: FEAL Linear Attack Case Study

## B.1 Task Overview

**Task:** Implement a known-plaintext attack on a FEAL-like encryption cipher to recover the secret key.
**Given:**
- `feal.c`: Encryption implementation with 4-round Feistel network
- `decrypt.c`: Decryption implementation
- `pairs.txt`: 32 known plaintext-ciphertext pairs
- `ciphertexts.txt`: 100 encrypted messages to decrypt

**Constraints:**
- Each of the 4 round keys is derived from a **20-bit seed** (reduced from standard key size)
- Total keyspace: $2^{80}$ combinations (still infeasible for brute force)
- Must use cryptanalytic techniques to recover keys efficiently

**Goal:** Recover the round keys $(k_0, k_1, k_2, k_3)$ and decrypt all ciphertexts in `ciphertexts.txt`, saving results to `plaintexts.txt`.

## B.2 Algorithm Analysis and Attack Strategy

### Phase 1: Understanding the Cipher Structure

**File System Exploration**

First, need to explore the working directory to understand available files:
`ls -la` reveals:
- `feal.c` (3472 bytes) — Encryption implementation
- `decrypt.c` (3982 bytes) — Decryption implementation
- `pairs.txt` (1307 bytes) — Known plaintext-ciphertext pairs
- `ciphertexts.txt` (2045 bytes) — Target ciphertexts to decrypt
- `feal`, `decrypt` — Compiled binaries

**Next step**: Analyze the encryption algorithm by reading `feal.c`.

## Cipher Structure Analysis

After reading `feal.c`, identified the following key components:

**Key Schedule:**
- Each round key $k_i$ is a 20-bit seed: `k_i & 0xFFFFF`
- Expanded via `expand(k)`: `k ^ (k << 16) ^ (k * 12345)`
- Creates 32-bit round keys: $K_i = \text{expand}(k_i)$

**Encryption Structure (4-round Feistel):**
1. **Pre-round mixing**: $R'_0 = R_0 \oplus L_0$
2. **Feistel rounds** $(i = 0, 1, 2, 3)$:

$$\text{mixed} = F(R_i \oplus K_i)$$
$$R_{i+1} = L_i \oplus \text{mixed}$$
$$L_{i+1} = R_i$$

3. **Post-round mixing**: $L'_4 = L_4 \oplus R_4$
4. **Output**: $C = \text{merge}(R_4, L'_4)$

**F-function** (4 byte transformations via G-function):

```
y[2] = g(x[3]^x[2], x[1]^x[0], 1);
y[3] = g(x[3], y[2], 0);
y[1] = g(x[1]^x[0], y[2], 0);
y[0] = g(x[0], y[1], 1);
```

**G-function** (rotation-based operation):

```
temp = (a + b + x) & 0xFF;
return ((temp << 2) | (temp >> 6)) & 0xFF;
```

## Phase 2: Attack Design — Meet-in-the-Middle

## Cryptanalytic Strategy

**Observation**: Directly brute-forcing $2^{80}$ keys is infeasible. However, we can exploit the structure:

**Key Insight:** With known plaintext-ciphertext pairs $(P, C)$, we can:
1. **Preprocess pairs**: Extract intermediate values

$$L_0, R_0 \leftarrow P$$
$$R'_0 = R_0 \oplus L_0 \quad \text{(pre-mixing)}$$
$$L'_4, R_4 \leftarrow C$$
$$L_4 = L'_4 \oplus R_4 \quad \text{(reverse post-mixing)}$$

2. **Meet-in-the-middle approach**:
   - **Forward direction**: Iterate $k_0 \in [0, 2^{20})$, compute intermediate state after round 0
   - **Backward direction**: Iterate $k_3 \in [0, 2^{20})$, compute intermediate state before round 3
   - **Match**: Find $(k_0, k_3)$ pairs where intermediate states align
   - Then recover $k_1, k_2$ using additional pairs

**Critical Implementation Requirement**: Need to compute **inverse of F-function** $(F^{-1})$ to work backwards from ciphertext.

**Complexity**: $O(2^{20} + 2^{20}) = O(2^{21})$ instead of $O(2^{80})$ — a massive reduction!

## Phase 3: Implementing the F-function Inverse

### F-function Inversion Algorithm

The F-function applies G-transformations in sequence:

$$y_2 = G(x_3 \oplus x_2, x_1 \oplus x_0, 1)$$
$$y_3 = G(x_3, y_2, 0)$$
$$y_1 = G(x_1 \oplus x_0, y_2, 0)$$
$$y_0 = G(x_0, y_1, 1)$$

To invert, need $G^{-1}$:

```python
def rotr2(x):   # Inverse of 2-bit rotation
    return ((x >> 2) | (x << 6)) & 0xFF

def f_inv_function(output_val):
    # Extract y bytes from output
    y = [output_val & 0xFF, ...]

    # Reverse G operations:
    temp_y1 = rotr2(y[1])
    x1_xor_x0 = (temp_y1 - y[2]) & 0xFF

    temp_y2 = rotr2(y[2])
    x3_xor_x2 = (temp_y2 - x1_xor_x0 - 1) & 0xFF

    temp_y3 = rotr2(y[3])
    x3 = (temp_y3 - y[2]) & 0xFF
    x2 = (x3_xor_x2 ^ x3) & 0xFF

    temp_y0 = rotr2(y[0])
    x0 = (temp_y0 - y[1] - 1) & 0xFF
    x1 = (x1_xor_x0 ^ x0) & 0xFF

    return combine_bytes([x0, x1, x2, x3])
```

**Verification**: Tested $F^{-1}(F(x)) = x$ for multiple test vectors — All passed!

**Phase 4: Key Recovery Implementation**

---

**Attack Implementation - Step A: Build $k_0$ Delta Map**

**Precompute expand map**: Since `expand()` is not injective, build reverse lookup:

```
exp_map = {}
for k in range(0, 1<<20):   # 2^20 iterations
    e = expand(k)
    exp_map[e] = k
# Result: 1,047,882 entries (some collisions)
```

**For each candidate** $k_0$: Compute intermediate values for **pairs 0 and 1**:

$$K_0 = \text{expand}(k_0)$$
$$A_0 = F^{-1}(F(K_0 \oplus R_0') \oplus L_0 \oplus L_4)$$
$$A_1 = F^{-1}(F(K_0 \oplus R_1') \oplus L_1 \oplus L_4)$$

Store in delta map: $\Delta = A_0 \oplus A_1 \mapsto (k_0, A_0, A_1)$
Built delta map with **933,294 unique deltas**.

---

**Attack Implementation - Step B: Match with $k_3$**

**Iterate $k_3 \in [0, 2^{20})$**: For each candidate, compute backward from ciphertext:

$$K_3 = \text{expand}(k_3)$$
$$B_0 = R_4^{(0)} \oplus F(L_4^{(0)} \oplus K_3)$$
$$B_1 = R_4^{(1)} \oplus F(L_4^{(1)} \oplus K_3)$$

Compute delta: $\Delta' = B_0 \oplus B_1$
**Check if** $\Delta' \in \text{delta\_map}$:
- If match found: $K_2 = A_0 \oplus B_0$ (potential key!)
- Verify $K_2$ is in `exp_map` to get $k_2$
- Use **pair 2** to derive $k_1$:

$$L_3^{(2)} = R_4^{(2)} \oplus F(L_4^{(2)} \oplus K_3)$$
$$L_2^{(2)} = L_4^{(2)} \oplus F(L_3^{(2)} \oplus K_2)$$
$$K_1 = L_2^{(2)} \oplus F^{-1}(R_0'^{(2)} \oplus L_3^{(2)})$$

- **Verify with pair 3**: Check $K_1$ consistency

**Found candidate keys!**

$$k_0 = \texttt{0x68e63}$$
$$k_1 = \texttt{0x62571}$$
$$k_2 = \texttt{0x1f61c}$$
$$k_3 = \texttt{0xb74a5}$$

## Phase 5: Verification and Decryption

### Key Verification Against All Pairs

**Verify all 32 pairs**: For each $(P_i, C_i)$ in pairs.txt:

```python
for i, (P, C) in enumerate(pairs):
    computed_C = encrypt(P, k0, k1, k2, k3)
    if computed_C != C:
        print(f"Pair {i} failed")
        break
```

**Result: All 32 pairs verified successfully!** Keys are correct.

Cross-verification using provided ./decrypt binary:

```
$ echo "14898973196205001976" > test.txt
$ ./decrypt test.txt 0x68e63 0x62571 0x1f61c 0xb74a5
367472469947502207  # Matches expected plaintext!
```

### Decrypting Target Ciphertexts

**Implement decrypt function** (reverse of encryption):

```python
def decrypt(ciphertext, k0, k1, k2, k3):
    key = [expand(k0), expand(k1),
           expand(k2), expand(k3)]

    # Reverse final swap: C = merge(R4, L4')
    left = getright(ciphertext)   # Was R4
    right = getleft(ciphertext)   # Was L4'

    # Reverse post-round mixing
    left ^= right

    # Reverse 4 Feistel rounds (backwards)
    for round_num in range(3, -1, -1):
        temp = left
        left = right
        mixed = f_function(temp ^ key[round_num])
        right = temp
        left ^= mixed

    # Reverse pre-round mixing
    right ^= left

    return merge(left, right)
```

**Process** ciphertexts.txt:
- Read 100 ciphertext values
- Decrypt each using recovered keys
- Save to plaintexts.txt

**Decrypted 100 ciphertexts successfully!** Saved to plaintexts.txt.

## B.3  Attack Performance Analysis

**Table 17** Cryptanalysis execution metrics

| Operation | Complexity | Execution Time |
|---|---|---|
| Precompute expand map | $O(2^{20})$ | <1 sec |
| Build $k_0$ delta map | $O(2^{20})$ | ~10 sec |
| Iterate $k_3$ candidates | $O(2^{20})$ | ~10 sec |
| Verify all 32 pairs | $O(32)$ | <1 sec |
| Decrypt 100 ciphertexts | $O(100)$ | <1 sec |
| **Total** | $O(2^{21})$ | **22.6 sec** |

**Efficiency Gain:**

- Brute force approach: $2^{80} \approx 1.2 \times 10^{24}$ operations (infeasible)
- Meet-in-the-middle attack: $2 \times 2^{20} \approx 2.1 \times 10^{6}$ operations
- **Speedup factor**: $\sim 10^{18}$ times faster!

## B.4  Key Findings and Agent Behaviors

1. **Cryptanalytic Reasoning:** The agent correctly identified that direct brute force of the $2^{80}$ keyspace is infeasible, and autonomously designed a meet-in-the-middle strategy that exploits the Feistel structure to reduce complexity to $O(2^{21})$.

2. **Mathematical Inversion:** Successfully derived the inverse of the F-function ($F^{-1}$) by analyzing the composition of G-functions and their rotation operations. This required understanding modular arithmetic and bit rotation properties.

3. **Systematic Verification:** Rather than accepting the first key candidate, the agent implemented multi-stage verification:

   - Used pair 0 & 1 to build delta map
   - Used pair 2 to derive $k_1$
   - Used pair 3 to verify $k_1$ consistency
   - Verified against all 32 known pairs
   - Cross-checked with original decrypt binary

4. **Efficient Precomputation:** Built reverse lookup tables (expand map, delta map) to avoid redundant calculations during the search phase, demonstrating optimization awareness.

5. **Code Quality:** Implemented clean, modular code with helper functions (`getleft`, `getright`, `merge`, `g_function`, `f_function`, `f_inv_function`) mirroring the original C implementation for consistency.

## B.5  Recovered Keys and Final Results

**Table 18** Recovered 20-bit round key seeds

| Key | Hexadecimal | Binary (20-bit) |
|---|---|---|
| $k_0$ | 0x68e63 | 01101 00011 10011 00011 |
| $k_1$ | 0x62571 | 01100 00101 01110 01001 |
| $k_2$ | 0x1f61c | 00001 11110 01100 01100 |
| $k_3$ | 0xb74a5 | 10110 11101 00101 00101 |

**Task Completion Status:**

- ✓ Recovered all 4 round keys from known plaintext-ciphertext pairs
- ✓ Verified keys against 32 pairs with 100% accuracy
- ✓ Successfully decrypted 100 ciphertexts from `ciphertexts.txt`
- ✓ Saved decrypted plaintexts to `plaintexts.txt`
- ✓ Cross-validated results using original `decrypt` binary

## B.6  Lessons for Cryptanalysis Agent Design

This case study demonstrates several capabilities essential for autonomous security research agents:

1. **Algorithm Analysis:** Ability to read and understand cryptographic implementations from source code, identifying exploitable structural weaknesses (Feistel network properties, reduced keyspace).

2. **Mathematical Reasoning:** Deriving inverse functions through algebraic manipulation and understanding of composition properties—critical for working backwards from ciphertexts.

3. **Attack Strategy Design:** Autonomously selecting appropriate cryptanalytic techniques (meet-in-the-middle) based on problem constraints rather than attempting naive approaches.

4. **Computational Efficiency:** Recognizing when precomputation and lookup tables provide better amortized complexity than repeated calculations.

5. **Multi-Stage Verification:** Building confidence through progressive validation rather than trusting single-point checks, reducing risk of false positives in key recovery.

6. **Code-Data Consistency:** Maintaining faithful translations between different implementations (C $\leftrightarrow$ Python) while preserving algorithmic correctness.

**Security Implications:** This task highlights the importance of adequate key sizes in cryptographic systems. The 20-bit key reduction (from standard 128-bit keys) made the cipher vulnerable to known-plaintext attacks within seconds on commodity hardware. Production systems must use cryptographically secure key lengths to resist modern computational attacks.

## C  NL2Repo: Python-Decouple Library Implementation

## C.1  Task Overview

**Benchmark:** NL2Repo - Repository-level code generation from natural language specification
**Task:** Implement a complete Python configuration management library from a 39KB natural language specification document (`start.md`).
**Requirements:**
- Multi-source configuration reading (environment variables, .env files, .ini files)
- Intelligent type conversion system (bool, int, CSV, custom casters)
- Validation mechanisms (Choices validator, default values)
- Docker secrets support for containerized deployments
- Complete project structure with setup.py, tests, documentation
- API must match specification: `from decouple import AutoConfig, Csv, Choices...`

**Challenge:** Transform verbose requirements into production-ready, installable Python package with 100% test coverage.

## C.2 Agent Execution Trace

### Step 1: Specification Analysis

First, explore the workspace to understand project requirements. The `start.md` file contains comprehensive specifications including: - API interfaces with exact function signatures - File structure requirements (decouple.py, setup.py, README.rst, etc.) - Type conversion logic with specific examples (.env parsing rules, quote handling) - Test cases embedded in documentation

**Action:** `view /openhands/workspace/start.md` to extract requirements systematically.

**Key insight:** Document contains **code snippets showing expected behavior** — these serve as implicit test specifications.

### Step 2-3: Core Module Implementation

**Priority:** Create `decouple.py` first as it's the heart of the library.

**Critical decisions made:**

1. **Boolean conversion:** Must handle edge case of empty strings → False (not ValueError)
2. **ConfigParser compatibility:** Python 3.12 doesn't have `read_config`, use `read_file()` instead
3. **Quote stripping in .env parser:** Only strip matching quotes at start/end: `v[0] == v[-1] in ["'",'"']`
4. **Docker secrets:** Handle missing directory gracefully with `os.path.exists()` check

**Action:** Implement 300+ lines covering RepositoryEnv, RepositoryIni, RepositorySecret, Config, AutoConfig, Csv, Choices classes.

**Challenge encountered:** Initial implementation failed test — empty string `strtobool('')` raised ValueError instead of returning False. Fixed by adding early return.

### Step 4-6: Package Configuration & Testing

**Setup infrastructure:**

- `setup.py`: Configured as installable package with metadata
- `README.rst`: Comprehensive documentation with usage examples
- `.gitignore, .editorconfig, .travis.yml`: Development tooling
- `tests/test_decouple.py`: 22 test cases covering all functionality

**Test-driven iteration:**

1. Run pytest → 1 failure (empty string boolean conversion)
2. Fix `strtobool()` to handle empty strings
3. Run pytest → 4 errors (fixture scope issue)
4. Move `env_file` fixture to module scope
5. Run pytest → **All 22 tests pass**

**Action:** `pip install -e .` to verify package installability.

## C.3 Key Agent Capabilities Demonstrated

1. **Specification Parsing:** Extracted structured requirements from 39KB natural language document, identifying implicit constraints (e.g., Python 3.12 compatibility issues).

2. **Error-Driven Development:** Systematically debugged ImportError (`read_config` not available) by consulting Python 3.12 ConfigParser API documentation and adapting implementation.

3. **Backwards Compatibility:** Ensured API exports work from both expected locations using `from qtlog import hide_qt_warning # noqa: F401` pattern.

4. **Production Quality:** Generated complete package infrastructure including LICENSE (MIT), MANIFEST.in, setup.cfg for PyPI distribution.

## C.4 Performance Metrics

**Table 19** Implementation statistics

| Metric | Value |
|---|---|
| Files created | 11 |
| Lines of code (decouple.py) | 323 |
| Test cases | 22 |
| Test coverage | 100% |
| Execution time | 4.5 seconds |
| Pass rate (first run) | 77% (17/22) |
| Pass rate (final) | 100% (22/22) |

**Final Validation:**

```
$ python -c "from decouple import config; \
  print('DEBUG =', config('DEBUG', default=False, cast=bool))"
Success! DEBUG = False

$ pip install -e .
Successfully installed python-decouple-3.8
```

**Task Completed:** Fully functional, installable Python package matching all specifications.

# D  SWE-bench Pro: Qt Logging Module Refactoring

## D.1  Task Overview

**Benchmark:** SWE-bench Pro - Real-world software engineering task from qutebrowser repository
**Issue:** Relocate `hide_qt_warning` function and `QtWarningFilter` class from `utils/log.py` to `utils/qtlog.py` while maintaining backwards compatibility.
**Constraints:**
- Must preserve identical filtering behavior (pattern matching, whitespace handling)
- Existing code using `from log import hide_qt_warning` must continue working
- Tests must be relocated to appropriate module (`test_qtlog.py`)
- No regression in warning suppression functionality

**Success criteria:** All tests pass, existing usage patterns unaffected, code properly organized.

## D.2  Agent Execution Trace

### Step 1: Repository Exploration

**Action:** Survey codebase structure to understand dependencies.
```
find /app -name "*.py" | grep -E '(log|test.*log)'
```
**Findings:**
- `qutebrowser/utils/log.py` (484 lines) — contains QtWarningFilter class
- `qutebrowser/utils/qtlog.py` (214 lines) — Qt-specific logging utilities
- `tests/unit/utils/test_log.py` (406 lines) — includes TestHideQtWarning class
- `tests/unit/utils/test_qtlog.py` (53 lines) — minimal test coverage
- `qutebrowser/browser/qtnetworkdownloads.py` — uses `log.hide_qt_warning()`

**Insight:** Migration requires **3 coordinated changes** to maintain consistency.

## Step 2: Function Migration

**Move QtWarningFilter class:**

```
# Add to qtlog.py (lines 31-46)
class QtWarningFilter(logging.Filter):
    def __init__(self, pattern: str) -> None:
        super().__init__()
        self._pattern = pattern

    def filter(self, record: logging.LogRecord) -> bool:
        do_log = not record.msg.strip().startswith(self._pattern)
        return do_log
```

**Move hide_qt_warning context manager:**

```
# Add to qtlog.py (lines 49-58)
@contextlib.contextmanager
def hide_qt_warning(pattern: str, logger: str = 'qt') -> Iterator[None]:
    log_filter = QtWarningFilter(pattern)
    logger_obj = logging.getLogger(logger)
    logger_obj.addFilter(log_filter)
    try:
        yield
    finally:
        logger_obj.removeFilter(log_filter)
```

**Critical detail:** Preserve exact whitespace handling logic: `record.msg.`**`strip()`**`.startswith()` enables matching with leading/trailing spaces.

## Step 3: Backwards Compatibility

**Problem:** Existing code imports from `log` module:

```
from qutebrowser.utils import log
with log.hide_qt_warning(...):
```

**Solution:** Re-export from log.py for backwards compatibility:

```
# In log.py (line 35)
from qutebrowser.utils.qtlog import hide_qt_warning, QtWarningFilter  # noqa:
    F401
```

**Deleted from log.py:**
- Lines 362-371: `hide_qt_warning` function (10 lines removed)
- Lines 404-419: `QtWarningFilter` class (16 lines removed)

**Verification:** Both import paths now resolve to same function object via Python's module caching.

**Step 4: Test Migration**

**Move TestHideQtWarning class:**
- Source: `test_log.py` lines 343-369 (27 lines)
- Destination: `test_qtlog.py` lines 32-58
- Updated references: `log.hide_qt_warning` → `qtlog.hide_qt_warning`

**Test coverage preserved:**
1. `test_unfiltered`: Non-matching patterns pass through
2. `test_filtered[Hello]`: Exact match suppressed
3. `test_filtered[Hello World]`: Prefix match suppressed
4. `test_filtered[ Hello World ]`: Whitespace handling works

**Deleted from test_log.py:** Lines 343-369 (TestHideQtWarning class removed)

**Step 5: Validation**

**Run test suite:**

```
$ pytest tests/unit/utils/test_qtlog.py -xvs
===== 5 passed in 0.06s =====

$ pytest tests/unit/utils/test_log.py -xvs
===== 51 passed in 1.65s =====
```

**Verify backwards compatibility:**

```
# Import from both locations
from qutebrowser.utils.log import hide_qt_warning as h1
from qutebrowser.utils.qtlog import hide_qt_warning as h2
assert h1 is h2  # Same function object - PASS
```

**Test real usage pattern:**

```
# Simulate qtnetworkdownloads.py usage
with log.hide_qt_warning('QNetworkReplyImplPrivate', 'qt-tests'):
    logger.warning('QNetworkReplyImplPrivate::error...')  # Filtered
    logger.warning('Other warning')  # Passes through
# Result: Only "Other warning" appears - CORRECT
```

**All validations passed**

## D.3  Refactoring Analysis

**Table 20**  Code migration summary

| File | Lines Added | Lines Removed | Net Change |
|------|------------:|--------------:|-----------:|
| qtlog.py | +28 | 0 | +28 |
| log.py | +1 | -26 | -25 |
| test_qtlog.py | +27 | 0 | +27 |
| test_log.py | 0 | -27 | -27 |
| **Total** | +56 | -53 | +3 |

## D.4  Key Software Engineering Practices

1. **Incremental Migration:** Agent performed refactoring in 5 systematic steps rather than bulk copy-paste, reducing risk of breaking changes.

2. **Backwards Compatibility:** Used Python re-export pattern (`from qtlog import *`) to maintain existing API surface while reorganizing internals.

3. **Test-First Validation:** Moved tests *after* moving code, then validated all tests pass before declaring success — ensuring no behavior regression.

4. **Dependency Analysis:** Searched codebase for all usages (`grep -r "hide_qt_warning"`) to identify affected modules before making changes.

**Impact:** Zero breaking changes, improved code organization, all 56 tests passing.

# E  Evaluation Details on Advanced Mathematical Reasoning

## E.1  Natural Language Proving

We adopt an iterative refine pipeline based on the solve-verify-refine framework [42]. Under this paradigm, Seed2.0 Pro generates candidate solutions, autonomously identifies logical flaws, and refines its outputs to satisfy strict Olympiad scoring criteria. Table 2 reports the results.

Notably, the high scores achieved by Seed2.0 Pro are not merely the result of hallucinated correct final answers. Instead, the model demonstrates a strong capacity for systematic reasoning and rigorous proof construction. The solve-verify-refine iterative pipeline ensures that every step of the deduction is logically coherent and mathematically rigorous, fully validating that Seed2.0 can complete complex mathematical reasoning tasks reliably rather than generating specious correct answers.

## E.2  Formal Theorem Proving

Beyond natural language reasoning, we test Seed2.0 on formal theorem proving using Putnam-200 (randomly sampled from [88])—200 formalized Putnam Competition problems evaluated in an agent-based multi-turn setup with Lean and Python tool access [13]. Table 5 reports the results.

Seed2.0 Pro reaches 35.5% Pass@8, a new state-of-the-art among comparable general-purpose models. This substantially surpasses Gemini-3-Pro and also exceeds our previous dedicated prover, Seed-1.5-Prover [13]. Seed2.0 Lite follows at 30.5% Pass@8, showing that strong formal reasoning survives even under efficiency constraints.

These results hint at meaningful transfer from natural language math reasoning to formal proof search. Unlike specialized provers trained primarily on formal corpora, Seed2.0 is optimized as a general-purpose model; its Putnam performance suggests that symbolic reasoning developed in natural language contexts can transfer to formal theorem proving.

## E.3  Case study on Erdos Problems

This section demonstrates the AI agent's capability to solve challenging open problems from Paul Erdős's problem collection, covering combinatorial geometry (Erdős 652) and analytic number theory (Erdős 1051).

### E.3.1 Erdos 652: Distinct Distances in Planar Point Sets

**Problem Input**

Let $x_1, \ldots, x_n \in \mathbb{R}^2$ and let $R(x_i) = \#\{|x_j - x_i| : j \neq i\}$, where the points are ordered such that

$$R(x_1) \leq \cdots \leq R(x_n).$$

Let $\alpha_k$ be minimal such that, for all large enough $n$, there exists a set of $n$ points with $R(x_k) < \alpha_k n^{1/2}$. Is it true that $\alpha_k \to \infty$ as $k \to \infty$?

**Output**

**Step 1: Incidence Count and Graph Construction**

Let $S$ be a set of $n$ points in $\mathbb{R}^2$, ordered such that $R(x_1) \leq R(x_2) \leq \cdots \leq R(x_n)$, where $R(x_i) = \#\{|x_j - x_i| : j \neq i\}$. Let $\alpha_k$ be the minimal constant such that for all sufficiently large $n$, there exists an $n$-point set $S$ with $R(x_k) < \alpha_k n^{1/2}$. We aim to show $\alpha_k \to \infty$ as $k \to \infty$.

Fix $k \geq 1$ and $\varepsilon > 0$. By minimality of $\alpha_k$, for all sufficiently large $n$, there exists an $n$-point set $S$ where the first $k$ points $P_1, P_2, \ldots, P_k$ (called *centers*) satisfy $R(P_i) < (\alpha_k + \varepsilon)n^{1/2}$ for each $i = 1, \ldots, k$. Let $t = \lfloor (\alpha_k + \varepsilon)n^{1/2} \rfloor$, so $t \leq (\alpha_k + \varepsilon)n^{1/2}$.

For each center $P_i$, let $\Gamma_i$ be the set of circles centered at $P_i$ with radii equal to the distinct distances from $P_i$ to $S \setminus \{P_i\}$. Then $|\Gamma_i| \leq t$, so the total number of circles is $L = |\bigcup_{i=1}^{k} \Gamma_i| \leq kt$.

Count the incidences $I = I(S, \Gamma)$ between points in $S$ and circles in $\Gamma$: For a non-center point $Q \in S \setminus \{P_1, \ldots, P_k\}$, $Q$ lies on exactly one circle in $\Gamma_i$ for each $i$, contributing $k$ incidences. There are $n - k$ such points, giving $k(n - k)$ incidences. For a center $P_j$, $P_j$ is not contained in any circle in $\Gamma_j$, but lies on exactly one circle in $\Gamma_i$ for each $i \neq j$, contributing $k - 1$ incidences. There are $k$ such points, giving $k(k - 1)$ incidences. Summing these: $I = k(n - k) + k(k - 1) = kn - k$.

Construct an undirected multigraph $G$ with vertex set $S$. For each circle $\gamma \in \Gamma$ with $m_\gamma = |\gamma \cap S| \geq 2$, connect consecutive points of $\gamma \cap S$ (in cyclic order around $\gamma$) with undirected edges (a 2-cycle for $m_\gamma = 2$, a simple cycle for $m_\gamma \geq 3$). The total number of edges in $G$ is: $E = \sum_{\gamma \in \Gamma, m_\gamma \geq 2} m_\gamma = \sum_{\gamma \in \Gamma} m_\gamma - \sum_{\gamma \in \Gamma, m_\gamma = 1} 1 = I - L_1$, where $L_1$ is the number of circles with exactly one point of $S$. Since $L_1 \leq L = kt$, we have: $E \geq I - L = (kn - k) - kt = kn - k - kt$.

**Step 2: Bounding Parallel Edges**

Simplify $G$ to a simple graph $G'$ by removing all but one edge from each set of parallel edges. Let $R$ be the total number of edges removed, so $E' = E - R$ (edges in $G'$). For an unordered pair $\{u, v\} \subseteq S$, define: $s_1(\{u, v\})$: number of centers where $\{u, v\}$ are consecutive on a circle with $m_\gamma \geq 3$, $s_2(\{u, v\})$: number of centers where $\{u, v\}$ are the only two points on a circle ($m_\gamma = 2$), $s(\{u, v\}) = s_1 + s_2$: total centers with $\{u, v\}$ consecutive, $\mu(\{u, v\}) = s_1 + 2s_2$: edge multiplicity of $\{u, v\}$ in $G$.

The number of edges removed for $\{u, v\}$ is $\max(\mu(\{u, v\}) - 1, 0)$, so: $R = \sum_{\{u,v\} \subseteq S} \max(s_1(\{u, v\}) + 2s_2(\{u, v\}) - 1, 0)$. Rewrite $R$ by splitting the sum over $s(\{u, v\})$: $R = \sum_{\{u,v\}:s \geq 2}(s - 1) + \sum_{\{u,v\}} s_2(\{u, v\}) = R_{\text{mult}} + T$, where $R_{\text{mult}} = \sum_{\{u,v\}:s \geq 2}(s - 1)$ and $T = \sum_{\{u,v\}} s_2(\{u, v\})$ (total two-point circles, so $T \leq kt$).

For $R_{\text{mult}}$, note $\sum_{\{u,v\}:s \geq 2}(s - 1) \leq \sum_{\{u,v\}} \binom{s}{2}$, where the right-hand side counts triples $(P_i, P_j, \{u, v\})$ with $i < j$ and $\{u, v\}$ consecutive on circles of both $P_i$ and $P_j$. For fixed centers $A = P_i, B = P_j$, consecutive pairs $\{u, v\}$ for both are symmetric over line $AB$, giving at most $2t$ such pairs per center pair. Thus:

$$\sum_{\{u,v\}} \binom{s}{2} = \sum_{1 \leq i < j \leq k} N(P_i, P_j) \leq \binom{k}{2} \cdot 2t \leq k^2 t.$$

Combining bounds for $R_{\text{mult}}$ and $T$: $R \leq k^2 t + kt = kt(k + 1)$.

**Step 3: Simplifying the Graph and Applying the Crossing Number Inequality**

The number of edges in the simple graph $G'$ is:

$$E' = E - R \geq (kn - k - kt) - kt(k+1) = kn - k - kt(k+2).$$

For fixed $k$, $t = O(n^{1/2})$, so $kt(k+2) + k = O(k^2 n^{1/2}) = o(kn)$ as $n \to \infty$. For sufficiently large $n$:

$$kt(k+2) + k \leq \frac{kn}{2} \implies E' \geq \frac{kn}{2}.$$

For $k \geq 8$, $E' \geq 4n$, satisfying the crossing number inequality hypothesis. The crossing number of $G$ is bounded by:

$$\mathrm{cr}(G) \leq 2 \binom{L}{2} \leq L^2 \leq (kt)^2,$$

since two circles intersect in at most two points (no crossings on the same circle, vertex intersections are not counted). As $G'$ is a subgraph of $G$, $\mathrm{cr}(G') \leq \mathrm{cr}(G) \leq k^2 t^2$.

Apply the **crossing number inequality** (simple graph with $n$ vertices, $m \geq 4n$ edges: $\mathrm{cr}(G) \geq \frac{m^3}{64n^2}$) to $G'$ with $m = E' \geq \frac{kn}{2}$:

$$k^2 t^2 \geq \frac{(kn/2)^3}{64n^2} = \frac{k^3 n^3/8}{64n^2} = \frac{k^3 n}{512}.$$

Divide by $k^2$:

$$t^2 \geq \frac{kn}{512} \implies t \geq \frac{\sqrt{kn}}{\sqrt{512}} = \frac{\sqrt{kn}}{16\sqrt{2}} \approx \frac{\sqrt{kn}}{22.6}.$$

**Step 4: Final Lower Bound on $\alpha_k$**

Recall $t \leq (\alpha_k + \varepsilon)n^{1/2}$. Substitute into the inequality:

$$(\alpha_k + \varepsilon)n^{1/2} \geq \frac{\sqrt{kn}}{22.6} \implies \alpha_k + \varepsilon \geq \frac{\sqrt{k}}{22.6}.$$

This holds for all $\varepsilon > 0$, so taking $\varepsilon \to 0$ gives:

$$\alpha_k \geq \frac{\sqrt{k}}{23}$$

for sufficiently large $k$. As $k \to \infty$, $\frac{\sqrt{k}}{23} \to \infty$, so $\alpha_k \to \infty$.

### E.3.2  Erdos 1051: Irrationality of Infinite Series

**Problem Input**

Is it true that if $1 \leq a_1 < a_2 < \ldots$ is a sequence of integers with

$$\liminf a_n^{1/2^n} > 1$$

then

$$\sum_{n=1}^{\infty} \frac{1}{a_n a_{n+1}}$$

is irrational?

Seed-Prover 1.5 also formalized this problem at https://github.com/ByteDance-Seed/Seed-Prover/blob/main/SeedProver-1.5/erdos_1051.lean.

**Output**

**Step 1: Rationality Implies Lower Bound on $x_m$**

Let $a_1 < a_2 < \cdots$ be a strictly increasing sequence of positive integers with $\liminf_{n \to \infty} a_n^{1/2^n} = L > 1$, and define the sum $S = \sum_{n=1}^{\infty} \frac{1}{a_n a_{n+1}}$. We prove $S$ is irrational by contradiction. Assume $S = P/Q$ for coprime positive integers $P, Q$. For each $m \geq 1$, split the sum into the partial sum and tail:

$$S = S_m + R_m, \quad S_m = \sum_{n=1}^{m} \frac{1}{a_n a_{n+1}}, \quad R_m = \sum_{n=m+1}^{\infty} \frac{1}{a_n a_{n+1}}.$$

Let $D_m = a_1 a_2 \cdots a_{m+1}$. For $n \leq m$, the term $D_m/(a_n a_{n+1}) = \left(\prod_{i=1}^{n-1} a_i\right)\left(\prod_{i=n+2}^{m+1} a_i\right)$ is an integer, so $K_m = D_m S_m \in \mathbb{Z}$. Rearranging:

$$D_m S = K_m + D_m R_m \implies x_m = D_m R_m = \frac{D_m P - K_m Q}{Q}.$$

Since $R_m > 0$, $x_m > 0$, the numerator is a positive integer. Hence $x_m \geq 1/Q$ for all $m \geq 1$.

**Step 2: Loose Tail Bound and Recurrence**

Because $\{a_n\}$ is strictly increasing, $a_{n+1} - a_n \geq 1$, so:

$$\frac{1}{a_n a_{n+1}} = \frac{1}{a_{n+1} - a_n}\left(\frac{1}{a_n} - \frac{1}{a_{n+1}}\right) \leq \frac{1}{a_n} - \frac{1}{a_{n+1}}.$$

The liminf condition implies $a_n \to \infty$: there exists $N_1 \in \mathbb{N}$ such that for all $n \geq N_1$, $a_n^{1/2^n} > (L+1)/2 > 1$, so $a_n > [(L+1)/2]^{2^n} \to \infty$ as $n \to \infty$. Summing the telescoping inequality from $n = m+2$ to $\infty$:

$$\sum_{n=m+2}^{\infty} \frac{1}{a_n a_{n+1}} \leq \sum_{n=m+2}^{\infty} \left(\frac{1}{a_n} - \frac{1}{a_{n+1}}\right) = \lim_{N \to \infty}\left(\frac{1}{a_{m+2}} - \frac{1}{a_N}\right) = \frac{1}{a_{m+2}}.$$

Thus:

$$R_m = \frac{1}{a_{m+1} a_{m+2}} + \sum_{n=m+2}^{\infty} \frac{1}{a_n a_{n+1}} \leq \frac{1}{a_{m+1} a_{m+2}} + \frac{1}{a_{m+2}} \leq \frac{2}{a_{m+2}},$$

where the last inequality uses $a_{m+1} \geq 1$. Combining with $x_m \geq 1/Q$:

$$\frac{1}{Q} \leq \frac{2 D_m}{a_{m+2}} \implies a_{m+2} \leq 2 Q D_m.$$

Let $P_k = a_1 a_2 \cdots a_k$ denote the product of the first $k$ terms. Since $D_m = P_{m+1}$, reindex with $k = m + 2$ ( $m = k - 2 \geq 1 \implies k \geq 3$) to obtain:

$$a_k \leq K P_{k-1} \quad \text{for all } k \geq 3, \quad \text{where } K = 2Q > 0.$$

**Step 3: Logarithmic Variables and Boundedness of $r_k$**

Take natural logarithms and define:

$$b_k = \log a_k, \quad Q_k = \log P_k = \sum_{i=1}^{k} b_i, \quad r_k = \frac{Q_k}{2^k}.$$

The recurrence $a_k \leq K P_{k-1}$ becomes:

$$b_k \leq \log K + Q_{k-1} \quad \text{for } k \geq 3.$$

Since $Q_k = Q_{k-1} + b_k$, substitute to get:

$$Q_k \leq Q_{k-1} + \log K + Q_{k-1} = 2Q_{k-1} + C, \quad \text{where } C = \log K.$$

Divide by $2^k$:

$$\frac{Q_k}{2^k} \leq \frac{Q_{k-1}}{2^{k-1}} + \frac{C}{2^k} \implies r_k \leq r_{k-1} + \frac{C}{2^k} \quad \text{for } k \geq 3.$$

For $n \geq m \geq 3$, sum the inequality from $k = m+1$ to $n$:

$$r_n - r_m = \sum_{k=m+1}^{n} (r_k - r_{k-1}) \leq \sum_{k=m+1}^{n} \frac{C}{2^k} \leq \frac{C}{2^m}.$$

Adjusting for $m < 3$ by increasing $C$ to a larger constant $C' > 0$ if needed, we have for all $n \geq m \geq 1$:

$$r_n \leq r_m + \frac{C'}{2^m}.$$

This implies $\{r_k\}$ is bounded above: for any $k \geq 1$, $r_k \leq r_1 + C'/2 < \infty$. Let $R = \limsup_{k \to \infty} r_k < \infty$.

**Step 4: Identity Relating $r_k$ and $d_k$**

Define the normalized log sequence:

$$d_k = \frac{b_k}{2^k} = \frac{\log a_k}{2^k}.$$

The liminf condition translates to:

$$\liminf_{k \to \infty} d_k = l = \log L > 0.$$

Using $Q_k = Q_{k-1} + b_k$, divide by $2^k$ to derive the identity:

$$r_k = \frac{Q_{k-1}}{2^k} + \frac{b_k}{2^k} = \frac{r_{k-1}}{2} + d_k,$$

which holds for all $k \geq 1$ with $r_0 = Q_0/2^0 = \log 1 = 0$ (since $P_0 = 1$, the empty product).

**Step 5: Convergence of $r_k$ and $d_k$**

We now show $\liminf_{k \to \infty} r_k = R$, so $r_k$ converges to $R$. Let $R' = \liminf_{k \to \infty} r_k$. Choose subsequences: $\{k_m\}$ such that $r_{k_m} \to R$ as $m \to \infty$, $\{j_n\}$ such that $r_{j_n} \to R'$ as $n \to \infty$.

For each $j_n$, let $k_m$ be the largest index satisfying $k_m \leq j_n$, so $k_m \leq j_n < k_{m+1}$. Since $k_{m+1} > j_n$, apply the inequality from Step 3:

$$r_{k_{m+1}} \leq r_{j_n} + \frac{C'}{2^{j_n}}.$$

Taking the limit as $n \to \infty$, the left-hand side tends to $R$ (as $k_{m+1}$ is a subsequence of $\{k_m\}$), and the right-hand side tends to $R' + 0 = R'$ (since $j_n \to \infty$). Thus $R \leq R'$. But $R' \leq R$ by definition of liminf, so $R' = R$. Hence $\lim_{k \to \infty} r_k = R$.

Now use the identity $d_k = r_k - r_{k-1}/2$. Since $r_k \to R$, $r_{k-1} \to R$ as $k \to \infty$, so:

$$\lim_{k \to \infty} d_k = R - \frac{R}{2} = \frac{R}{2}.$$

Let $D = R/2$. The liminf condition gives $l = \liminf d_k = D \geq \log L > 0$, so $D > 0$.

**Step 6: Tight Tail Bound**

Since $d_k \to D > 0$, choose $\epsilon = D/10 > 0$. There exists $N_2 \in \mathbb{N}$ such that for all $n \geq N_2$:

$$|d_n - D| < \epsilon \implies D - \epsilon < d_n < D + \epsilon.$$

Exponentiating gives:

$$e^{(D-\epsilon)2^n} < a_n < e^{(D+\epsilon)2^n} \quad \text{for } n \geq N_2.$$

For $m \geq N_2 - 1$, consider the ratio of consecutive terms in the tail $R_m$:

$$\frac{1/(a_{n+1}a_{n+2})}{1/(a_n a_{n+1})} = \frac{a_n}{a_{n+2}} \quad \text{for } n \geq m+1 \geq N_2.$$

Using the bounds on $a_n$:

$$\frac{a_n}{a_{n+2}} < \frac{e^{(D+\epsilon)2^n}}{e^{(D-\epsilon)2^{n+2}}} = e^{(D+\epsilon)2^n - 4(D-\epsilon)2^n} = e^{(-3D+5\epsilon)2^n}.$$

With $\epsilon = D/10$, the exponent becomes:

$$-3D + 5\epsilon = -3D + \frac{D}{2} = -\frac{5D}{2} < 0,$$

so $a_n/a_{n+2} \leq 1/2$ for all $n \geq N_2$ (since the exponent decays super-exponentially and is bounded above by $\log(1/2)$ for large $n$). Thus the tail $R_m$ is bounded by a geometric series with ratio $1/2$:

$$R_m \leq \frac{1}{a_{m+1}a_{m+2}} + \frac{1}{2} \cdot \frac{1}{a_{m+1}a_{m+2}} + \frac{1}{2^2} \cdot \frac{1}{a_{m+1}a_{m+2}} + \cdots = \frac{2}{a_{m+1}a_{m+2}}.$$

**Step 7: Final Contradiction**

For $m \geq N_2 - 1$, use the tight tail bound to write: $x_m = D_m R_m \leq \frac{2D_m}{a_{m+1}a_{m+2}}$. Since $D_m = P_{m+1} = P_m a_{m+1}$, substitute to simplify:

$$x_m \leq \frac{2P_m a_{m+1}}{a_{m+1}a_{m+2}} = \frac{2P_m}{a_{m+2}}.$$

Express $P_m$ and $a_{m+2}$ in terms of $r_m$ and $d_{m+2}$:

$$P_m = e^{Q_m} = e^{2^m r_m}, \quad a_{m+2} = e^{b_{m+2}} = e^{d_{m+2}2^{m+2}} = e^{4d_{m+2}2^m}.$$

Thus:

$$\frac{P_m}{a_{m+2}} = e^{2^m r_m - 4d_{m+2}2^m} = e^{2^m(r_m - 4d_{m+2})}.$$

Since $r_m \to R = 2D$ and $d_{m+2} \to D$, choose $\epsilon = D/5 > 0$. There exists $M \geq N_2 - 1$ such that for all $m \geq M$:

$$r_m < 2D + \epsilon, \quad d_{m+2} > D - \epsilon.$$

Substitute into the exponent:

$$r_m - 4d_{m+2} < (2D + \epsilon) - 4(D - \epsilon) = 2D + \epsilon - 4D + 4\epsilon = -2D + 5\epsilon.$$

With $\epsilon = D/5$, this simplifies to:

$$r_m - 4d_{m+2} < -2D + 5 \cdot \frac{D}{5} = -D < 0.$$

Thus for $m \geq M$:

$$\frac{P_m}{a_{m+2}} < e^{2^m(-D)} = e^{-D2^m},$$

so:

$$x_m \leq 2e^{-D2^m}.$$

Since $D > 0$, $e^{-D2^m} \to 0$ as $m \to \infty$, so $x_m \to 0$. This contradicts Step 1, where $x_m \geq 1/Q > 0$ for all $m \geq 1$. Therefore, the assumption that $S$ is rational is false. The sum $S$ is irrational.