

---

# RIEMANN-BENCH: A BENCHMARK FOR MOONSHOT MATHEMATICS

Surge AI Research

## ABSTRACT

1       Recent AI systems have achieved gold-medal-level performance on the In-  
2       ternational Mathematical Olympiad, demonstrating remarkable proficiency at  
3       competition-style problem solving. However, competition mathematics repre-  
4       sents only a narrow slice of mathematical reasoning: problems are drawn from  
5       limited domains, require minimal advanced machinery, and can often reward in-  
6       sightful tricks over deep theoretical knowledge. We introduce RIEMANN-BENCH,  
7       a private benchmark of 25 expert-curated problems designed to evaluate AI sys-  
8       tems on research-level mathematics that goes far beyond the olympiad frontier.  
9       Problems are authored by Ivy League mathematics professors, graduate students,  
10      and PhD-holding IMO medalists, and routinely took their authors weeks to solve  
11      independently. Each problem undergoes double-blind verification by two inde-  
12      pendent domain experts who must solve the problem from scratch, and yields a  
13      unique, closed-form solution assessed by programmatic verifiers. We evaluate  
14      frontier models as unconstrained research agents, with full access to coding tools,  
15      search, and open-ended reasoning, using an unbiased statistical estimator computed  
16      over 100 independent runs per problem. Our results reveal that all frontier models  
17      currently score below 10%, exposing a substantial gap between olympiad-level  
18      problem solving and genuine research-level mathematical reasoning. By keep-  
19      ing the benchmark fully private, we ensure that measured performance reflects  
20      authentic mathematical capability rather than memorization of training data.

## 21 1 INTRODUCTION

22      Five years ago, Surge helped create GSM8K (Cobbe et al., 2021), the first mathematical reasoning  
23      benchmark for large language models. At the time, the tasks focused on grade-school math and  
24      GSM8K became one of the most widely cited benchmarks because it exposed a fundamental gap  
25      between fluent language and actual reasoning.

26      The frontier has since moved dramatically. The year 2025 marked an important moment for AI and  
27      mathematics. Google DeepMind’s Gemini with Deep Think scored 35 out of 42 points on the 2025  
28      International Mathematical Olympiad (IMO), officially achieving gold-medal standard as certified  
29      by IMO coordinators (DeepMind, 2025). DeepSeekMath-V2 achieved gold-level performance on  
30      IMO 2025 and scored 118/120 on Putnam 2024 (DeepSeekMath-V2, 2025). These achievements  
31      followed a rapid progression: AlphaProof solved the hardest problem at IMO 2024 (AlphaProof,  
32      2025), and performance on the American Invitational Mathematics Examination (AIME) approached  
33      near-perfect accuracy, with o4-mini achieving 99.5% on AIME 2025 using tool access (OpenAI,  
34      2025). The emergence of reasoning-focused models such as OpenAI’s o1 (OpenAI, 2024) and  
35      DeepSeek-R1 (DeepSeek-AI, 2025), which apply reinforcement learning to develop extended chains  
36      of reasoning, has been a key driver of these gains.

37      The IMO is deliberately limited to four domains: Algebra, Combinatorics, Geometry, and Number  
38      Theory. These foundational areas were specifically chosen because they require minimal advanced  
39      machinery; calculus, for instance, is strictly excluded. Because the available tools are limited, IMO  
40      problems are designed to reward lateral thinking, often hinging on a single key insight that renders  
41      the solution tractable. While IMO problems are incredibly clever, they are fundamentally designed to  
42      be solved in a few hours using known tools. The distance between this style of problem solving and  
43      the sustained, multi-step theoretical reasoning characteristic of professional mathematical research is  
44      vast.

---

45 This gap motivates a new evaluation paradigm. While benchmarks such as GSM8K (Cobbe et al.,  
46 2021), MATH (Hendrycks et al., 2021a), and Omni-MATH (Gao et al., 2024) have progressively  
47 raised the difficulty bar, they remain largely confined to competition-style problems. The saturation  
48 of existing benchmarks, combined with growing evidence of data contamination in public evalu-  
49 ations (Mirzadeh et al., 2025; Srivastava et al., 2024; Oren et al., 2024), underscores the need for  
50 private, rigorously constructed benchmarks at the research frontier.

51 We introduce RIEMANN-BENCH, a benchmark of 25 extreme-tier mathematical problems designed  
52 to evaluate AI not on competition puzzles, but on PhD-level research mathematics. We collaborated  
53 with Ivy League mathematics professors, graduate students, and PhD-holding IMO medalists to  
54 gather problems they encountered in their own research. These problems routinely took their authors  
55 weeks to solve independently. Authors noted that their own graduate students and colleagues would  
56 struggle to solve these problems on their own.

57 Our contributions are:

- 58 1. **Research-level mathematical benchmark.** We introduce RIEMANN-BENCH, comprising 25  
59 problems spanning multiple mathematical domains including areas that require understanding  
60 of variational principles, measure theory, stability analysis, manifolds, and advanced algebraic  
61 structures.
- 62 2. **Double-blind, from-scratch verification.** Every problem is independently verified by two  
63 domain experts who must solve the problem from scratch before confirming validity.
- 64 3. **Unconstrained agent evaluation.** Unlike existing benchmarks that force models into rigid,  
65 automated evaluation loops, RIEMANN-BENCH evaluates true, unconstrained AI research agents  
66 with full access to coding tools, search, and open-ended reasoning.
- 67 4. **Statistically rigorous evaluation.** We evaluate every problem against each frontier model 100  
68 times and compute pass rates using the unbiased estimator of Chen et al. (2021), providing  
69 stable and reproducible measurements of model capability.
- 70 5. **Comprehensive frontier evaluation.** We evaluate frontier models from major AI laboratories,  
71 finding that all currently score below 10%.
- 72 6. **Fully private and uncontaminated.** The dataset is kept strictly private to ensure a fully  
73 unbiased evaluation for all frontier labs.

## 74 2 RELATED WORK

### 75 2.1 MATHEMATICAL REASONING BENCHMARKS

76 The landscape of mathematical reasoning benchmarks has evolved rapidly, tracing a clear difficulty  
77 progression from elementary arithmetic to research-level problems.

78 **Elementary and competition-level benchmarks.** GSM8K (Cobbe et al., 2021) introduced 8,500  
79 grade-school math word problems requiring 2–8 reasoning steps, alongside the verifier-based evalua-  
80 tion paradigm. The MATH dataset (Hendrycks et al., 2021a) raised the bar significantly with 12,500  
81 competition-level problems across seven categories sourced from AMC, AIME, and other competi-  
82 tions. When introduced, the best models achieved roughly 7% accuracy; frontier models now exceed  
83 90%, rendering the benchmark effectively saturated. MMLU (Hendrycks et al., 2021b) includes  
84 mathematics-related subjects among its 57 domains, spanning elementary through college-level  
85 abstract algebra.

86 **Olympiad-level benchmarks.** Several recent benchmarks target olympiad-level difficulty. Omni-  
87 MATH (Gao et al., 2024) contains 4,428 problems across 33 sub-domains sourced from competitions  
88 including USAMO, APMO, and Putnam. OlympiadBench (He et al., 2024) provides 8,476 bilingual  
89 problems in mathematics and physics drawn from international olympiads. OlymMATH (Olym-  
90 MATH, 2025) targets olympiad-level reasoning across multiple difficulty tiers. JEEBench (Arora  
91 et al., 2023) uses 515 problems from India’s JEE-Advanced examination. MathOdyssey (Fang et  
92 al., 2024) contributes 387 expert-crafted problems spanning high school to university level. Math-  
93 Arena (Balunović et al., 2025) evaluates models on recently released competition problems with  
94 rigorous contamination controls. The AI Mathematical Olympiad (AIMO) Prize (AIMO Prize, 2023)

---

95 has further catalyzed progress by awarding prizes for publicly shared models that solve olympiad-level  
96 problems.

97 **Graduate and research-level benchmarks.** GPQA (Rein et al., 2024) provides 448 expert-crafted  
98 multiple-choice questions in physics, chemistry, and biology at a graduate level where domain  
99 experts achieve only 65% accuracy. Humanity’s Last Exam (Phan et al., 2026) crowdsourced  
100 3,000 expert-level questions across dozens of academic disciplines; frontier models scored below  
101 10% at launch. GHOSTS (Frieder et al., 2023) was among the first benchmarks curated by working  
102 mathematicians to target graduate-level mathematics. TheoremQA (Chen et al., 2023) tests knowledge  
103 of over 350 theorems across mathematics, physics, and finance. ARB (Sawada et al., 2023) targets  
104 graduate and expert-level reasoning across multiple domains. SciBench (Wang et al., 2024) evaluates  
105 college-level scientific problem solving with free-response questions drawn from standard textbooks.  
106 MathBench (Liu et al., 2024) spans 3,709 problems across five progressive stages from arithmetic to  
107 college mathematics.

108 **Formal mathematics benchmarks.** MiniF2F (Zheng et al., 2022) contains 488 problems formalized  
109 across Lean, Metamath, Isabelle, and HOL Light. ProofNet (Azerbaiyev et al., 2023) provides  
110 371 parallel examples of formal and natural-language theorem statements from undergraduate  
111 textbooks. PutnamBench (Tsoukalas et al., 2024) offers 1,692 hand-constructed formalizations of  
112 640 Putnam competition theorems. DeepSeek-Prover-V2 (DeepSeek-AI, 2025b) recently advanced  
113 formal theorem proving by combining reinforcement learning with subgoal decomposition in Lean 4.

114 **FrontierMath.** FrontierMath (Glazer et al., 2024) was developed by Epoch AI in collaboration  
115 with over 60 mathematicians including Fields Medalist Terence Tao. It contains approximately 350  
116 problems organized into four difficulty tiers covering most major branches of modern mathematics.  
117 When launched, all frontier models scored under 2%. EternalMath (EternalMath, 2026) has proposed  
118 an alternative paradigm in which the benchmark evolves alongside mathematical discovery to resist  
119 contamination over time. RIEMANN-BENCH complements FrontierMath in several ways: it is  
120 independently constructed, fully private, uses double-blind from-scratch expert verification, and  
121 evaluates models as unconstrained research agents rather than constraining them to rigid evaluation  
122 loops.

## 123 2.2 AI PERFORMANCE ON MATHEMATICAL OLYMPIADS

124 AlphaGeometry (Trinh et al., 2024) solved 25 of 30 historical IMO geometry problems, matching the  
125 average gold medalist. AlphaProof (AlphaProof, 2025), combined with AlphaGeometry 2 (Chervonyi  
126 et al., 2025), scored 28/42 at IMO 2024 (silver medal; the gold cutoff was 29). By IMO 2025,  
127 Gemini with Deep Think became the first AI system officially certified at gold-medal standard  
128 by IMO coordinators (DeepMind, 2025). DeepSeekMath-V2 achieved gold-level scores on IMO  
129 2025 and CMO 2024, and scored 118/120 on Putnam 2024 (DeepSeekMath-V2, 2025). In formal  
130 mathematics, Axiom Math’s AxiomProver solved all 12 problems on the 2025 William Lowell  
131 Putnam Mathematical Competition with machine-verified proofs in Lean 4 (Axiom Math, 2025).

132 In parallel, language model reasoning capabilities have advanced rapidly. Chain-of-thought prompt-  
133 ing (Wei et al., 2022) demonstrated that eliciting intermediate reasoning steps substantially improves  
134 mathematical performance. Subsequent work on tree-structured reasoning (Yao et al., 2023) and tool-  
135 augmented approaches such as PAL (Gao et al., 2023) and ToRA (Gou et al., 2024) further expanded  
136 the problem-solving capabilities of language models. The introduction of dedicated reasoning models,  
137 beginning with OpenAI o1 (OpenAI, 2024) and followed by DeepSeek-R1 (DeepSeek-AI, 2025),  
138 marked a paradigm shift: these systems allocate substantial test-time compute to extended chains  
139 of reasoning, yielding dramatic gains on competition mathematics. Domain-specific mathematical  
140 training has also proven effective, with Minerva (Lewkowycz et al., 2022) demonstrating early  
141 gains through pretraining on mathematical corpora, and more recent systems such as DeepSeek-  
142 Math (DeepSeek-AI, 2024), Llemma (Azerbaiyev et al., 2024), InternLM-Math (Ying et al., 2024),  
143 and Qwen2.5-Math (Yang et al., 2024) advancing the state of the art for open-weight mathematical  
144 models.

---

## 145 3 BENCHMARK DESIGN

### 146 3.1 DESIGN PHILOSOPHY

147 RIEMANN-BENCH operates in a fundamentally different regime from competition-style benchmarks.  
148 While IMO problems are incredibly clever, they are designed to be solved in a few hours using known  
149 tools and mathematical machinery. RIEMANN-BENCH operates in the universe of PhD-level research:  
150 problems that demand deep domain knowledge, complex theory, and the synthesis of advanced  
151 mathematical machinery over long reasoning chains.

152 The problems are not open-ended conjectures; they have known, verifiable answers. But arriving at  
153 those answers demands the kind of multi-step reasoning and mastery of specialized mathematical  
154 frameworks that characterizes work at the research frontier. To succeed on RIEMANN-BENCH, an  
155 AI system cannot rely on pattern recognition or lateral thinking alone. It must navigate complex  
156 abstract definitions, apply advanced theorems, and sustain coherent reasoning across extended chains  
157 of increasing complexity.

### 158 3.2 PROBLEM CONSTRUCTION

159 RIEMANN-BENCH comprises 25 problems authored by Ivy League mathematics professors, graduate  
160 students, and PhD-holding IMO medalists with active research programs. Contributors were asked to  
161 draw on problems they encountered in their own research: problems that routinely took them weeks  
162 to solve independently. Authors noted that their own graduate students and colleagues would struggle  
163 to solve these problems on their own.

164 Each problem satisfies the following requirements:

- 165 • **Unambiguous answer.** Every problem yields a unique, closed-form solution. There is no partial  
166 credit and no subjective judgment: the answer is either correct or incorrect.
- 167 • **Programmatic verification.** When the solution admits multiple equivalent representations (for  
168 example, a rational number that may be expressed in different forms), programmatic verifiers  
169 assess correctness automatically.
- 170 • **Research-level difficulty.** Problems require deep domain knowledge and multi-step theoretical  
171 reasoning that goes substantially beyond what is testable in competition settings.
- 172 • **Self-contained statement.** Problems are presented with all necessary definitions and context, so  
173 that a mathematician in the relevant field can solve them without external references.

### 174 3.3 VERIFICATION PROTOCOL

175 A benchmark is only as trustworthy as its verification process. Every problem in RIEMANN-BENCH  
176 was subjected to a strict double-blind, from-scratch verification protocol:

- 177 1. Two independent domain experts, who were not shown the author’s solution in advance, are  
178 assigned to verify each problem.
- 179 2. Each verifier must solve the problem from scratch and arrive at the correct answer through their  
180 own reasoning before confirming a problem’s validity.
- 181 3. The verifiers also assess problem quality, checking for ambiguity, underspecification, and  
182 appropriate difficulty calibration.

183 This double-blind protocol goes substantially beyond the standard practice of relying on prob-  
184 lem authors to self-certify their solutions. It provides a strong guarantee that each problem has a  
185 unique correct answer that can be independently derived by multiple experts. Problems that failed  
186 verification—due to ambiguity, errors, or insufficient difficulty—were revised or excluded.

### 187 3.4 PRIVACY

188 To ensure a fully unbiased evaluation for all frontier labs, the dataset is kept strictly private. Public  
189 benchmarks, however well-intentioned, are vulnerable to leakage and contamination (Xu et al.,

---

190 2024; Zhou et al., 2023). A benchmark that has been seen, even indirectly, is a benchmark that has  
191 been compromised. Labs wishing to evaluate their models may submit them through a controlled  
192 evaluation service.

### 193 3.5 UNCONSTRAINED AGENT EVALUATION

194 A distinctive feature of RIEMANN-BENCH is its evaluation protocol for AI systems. Existing  
195 benchmarks can force models into rigid, automated evaluation loops. RIEMANN-BENCH evaluates  
196 unconstrained AI research agents. Models are given full access to coding tools (Python interpreter),  
197 search capabilities, and open-ended reasoning with no artificial constraints on interaction format or  
198 token budget. This design reflects our belief that measuring research-level mathematical capability  
199 requires allowing models to operate as they would in a genuine research setting.

## 200 4 ILLUSTRATIVE PROBLEM

201 To convey the character and difficulty of RIEMANN-BENCH, we present one sample problem. This  
202 problem was selected because it illustrates several key properties of the benchmark: it involves  
203 advanced mathematical objects, requires deep familiarity with specialized theory, and demands  
204 sustained multi-step reasoning that a domain expert estimated would take 40–50 hours to complete  
205 from scratch.

206 **Problem overview.** The problem concerns the classification of multibasic  $A$ -modules over the  
207 ring of Hahn series with real-valued valuation and residue field  $\mathbb{F}_2$ . The field  $K$  of Hahn series  
208 in indeterminate  $t$  with value group  $\mathbb{R}$  is considered as a module over its subring  $A$  of elements  
209 with non-negative valuation. Special  $A$ -modules, termed *basic* (quotients of submodules of  $K$ ) and  
210 *multibasic* (finite direct sums of basic modules), are defined, with the property that every multibasic  
211  $A$ -module has a unique decomposition into a direct sum of basic submodules. The problem asks for  
212 the number of distinct isomorphism classes of multibasic  $A$ -modules  $M$  satisfying three structural  
213 conditions involving the endomorphism ring and a dimension function on associated  $\mathbb{F}_2$ -vector spaces.

214 **Discussion.** The setup is highly technical and involves advanced mathematical objects. Hahn series  
215 with real-valued valuation are formal infinite series in which the exponents may be any real numbers.  
216 A key insight in the solution is that submodules of the Hahn series field behave very simply with  
217 respect to the valuation: any submodule is determined by which powers  $t^q$  it contains, forcing the  
218 submodule to correspond to a cut in  $\mathbb{R}$ . As a result, every basic module  $L/N$  must come from a small  
219 list of canonical possibilities. Because multibasic modules decompose uniquely into basic pieces,  
220 the classification problem reduces to determining which combinations of these building blocks are  
221 allowed.

222 The three conditions in the problem then act as filters, each eliminating a different family of candi-  
223 dates through a qualitatively distinct algebraic mechanism. The solution draws on a diversity of  
224 mathematical ideas: classifying  $A$ -submodules of  $K$  via the valuation, applying the tensor-hom  
225 adjunction to determine how tensor products and hom functors interact with multibasic modules,  
226 using ring-theoretic properties to constrain which basic summands can appear, and finally reducing to  
227 a finite case analysis with a combinatorial count.

## 228 5 EXPERIMENTAL SETUP

### 229 5.1 MODELS EVALUATED

230 We evaluated major frontier AI models.

231 All models were evaluated as unconstrained research agents through their respective APIs, with full  
232 access to coding tools (Python interpreter), search capabilities, and open-ended reasoning. No artificial  
233 constraints were imposed on interaction format or token budget. This setup ensures that measured  
234 performance gaps reflect genuine mathematical reasoning limitations rather than implementation  
235 bottlenecks.

Table 1: Frontier model performance on RIEMANN-BENCH. Pass rates are computed from 100 independent runs per problem using the unbiased estimator of Chen et al. (2021). All models were evaluated as unconstrained agents with access to coding tools and search.

Model	Lab	pass@1 (%)
Gemini 3.1 Pro	Google	6
Claude Opus 4.6	Anthropic	6
Gemini 3 Pro	Google	4
Kimi K2.5	Moonshot AI	4
DeepSeek V3.2	DeepSeek	3
GPT 5.2	OpenAI	2
Claude Opus 4.5	Anthropic	2

## 236 5.2 EVALUATION PROTOCOL

237 For each problem, we ran every model 100 times independently and computed pass rates using the  
 238 unbiased statistical estimator introduced by Chen et al. (2021). Given  $n$  total samples and  $c$  correct  
 239 samples, the pass@ $k$  estimator is:

$$\text{pass}@k = 1 - \frac{\binom{n-c}{k}}{\binom{n}{k}}$$

240 This estimator provides unbiased measurements of model capability at any sampling budget  $k$ ,  
 241 computed from a fixed pool of  $n = 100$  independent attempts. The resulting difficulty assessments  
 242 are not based on a small number of attempts or selected runs; they reflect stable, reproducible  
 243 measurements.

## 244 6 RESULTS

### 245 6.1 OVERALL PERFORMANCE

246 Table 1 and Figure 1 present our primary results across all 25 problems.

247 The central finding is stark: **all frontier models currently score below 10% on RIEMANN-BENCH,**  
 248 **even when operating as unconstrained research agents with full access to coding tools and**  
 249 **search.** This is in sharp contrast to olympiad-level benchmarks, where the same generation of  
 250 models approaches or exceeds human gold-medal performance. The gap confirms that research-level  
 251 mathematics, the kind of sustained, multi-step theoretical reasoning that characterizes PhD-level  
 252 work, remains beyond current model capabilities.

### 253 6.2 COMPARISON WITH COMPETITION-LEVEL PERFORMANCE

254 To contextualize these results, we note the performance of the same model generation on competition-  
 255 level mathematics. The models evaluated here, or their close variants, achieve near-perfect scores  
 256 on AIME problems and gold-medal-level performance on IMO problems. The dramatic drop from  
 257 near-100% on AIME to below 10% on RIEMANN-BENCH illustrates the qualitative difference  
 258 between competition mathematics and research-level problems. Competition problems, however  
 259 difficult, can often be resolved through a single key insight applied with elementary tools. RIEMANN-  
 260 BENCH problems require sustained theoretical reasoning over weeks of effort, drawing on specialized  
 261 knowledge that extends well beyond the competition canon.

### 262 6.3 ANALYSIS OF A REPRESENTATIVE FAILURE MODE

263 Beyond aggregate pass rates, qualitative analysis of model failures reveals important patterns in how  
 264 current systems break down on research-level mathematics. We present a representative failure on  
 265 the illustrative problem from Section 4 to demonstrate an important class of errors.

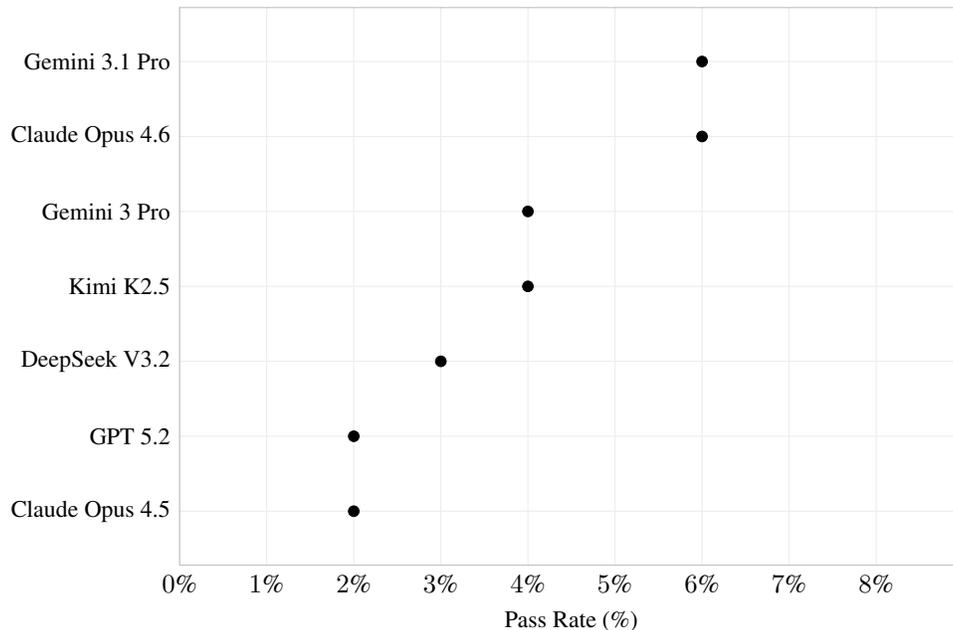


Figure 1: pass@1 across frontier models on RIEMANN-BENCH. All models score below 10%, confirming that research-level mathematics remains substantially beyond current capabilities.

266 **Model failure.** Rather than working within the  $A$ -module framework specified by the problem, the  
 267 model reinterpreted the entire problem in terms of an inapplicable theory of “generalized scales.”  
 268 Specifically, it incorrectly treated conditions (i) and (ii)—constraints on the endomorphism ring of  
 269  $M$ —as the definition of a “basic scale,” and misinterpreted condition (iii) as specifying the “support”  
 270 of  $M$ , when in fact the relevant notion of support is intrinsic to the Hahn series construction. To  
 271 justify its reasoning, the model fabricated a nonexistent classification theorem, attributing it to a  
 272 fictitious reference (“M. Getz, Theorem 4.14 on Generalized Scales”). Applying this fabricated  
 273 theorem, the model arrived at an answer of  $2^{299}$ , which is off by orders of magnitude from the correct  
 274 answer.

275 **Broader pattern.** This failure exemplifies a recurring pattern observed across RIEMANN-BENCH  
 276 evaluations: when confronted with problems requiring specialized theoretical frameworks, models  
 277 may substitute a superficially related but inapplicable framework and fabricate supporting results to  
 278 complete the reasoning chain. The model’s output reads as structurally coherent—it identifies the  
 279 problem as a classification task, proposes a theoretical framework, invokes a theorem, and computes  
 280 a numerical answer, but the entire reasoning chain is built on a misidentified foundation.

## 281 7 DISCUSSION

### 282 7.1 WHY COMPETITION MATH IS NOT ENOUGH

283 The contrast between olympiad-level and research-level performance reveals a fundamental distinction  
 284 in mathematical reasoning. Problems in RIEMANN-BENCH demand deep familiarity with specialized  
 285 theory, the ability to chain together multiple advanced results, and sustained computation through  
 286 complex algebraic or analytic manipulations. A model that can identify the key insight in an IMO  
 287 problem may nonetheless be unable to navigate the Eynard–Orantin topological recursion or classify  
 288 multibasic modules over Hahn series rings.

289 This distinction carries important implications for how we interpret benchmark saturation. The  
 290 saturation of competition-level benchmarks does not imply that mathematical reasoning is solved.  
 291 It indicates that a particular style of mathematical reasoning, one that rewards lateral thinking with  
 292 more elementary tools, is within reach of current systems. The broader and deeper landscape

---

293 of mathematics, encompassing the specialized theory and extended reasoning chains required for  
294 research-level work, remains largely beyond current capabilities.

## 295 7.2 IMPLICATIONS FOR AI-ASSISTED MATHEMATICAL RESEARCH

296 The results carry both encouraging and cautionary implications for the prospect of AI systems  
297 contributing to mathematical research. The rapid generational improvement observed on competition  
298 mathematics suggests that continued scaling and targeted training can yield meaningful progress. On  
299 the other hand, performance below 10% on RIEMANN-BENCH means that even the best models fail  
300 on the vast majority of research-level problems, making current systems unreliable as autonomous  
301 mathematical reasoning agents.

302 A more realistic near-term application may be AI-assisted research, in which human mathematicians  
303 use AI systems as computational assistants for specific subtasks while verifying outputs against  
304 their own expertise. This mirrors the trajectory observed in other domains of AI deployment, where  
305 practitioners deliberately constrain agent autonomy to maintain reliability (Pan et al., 2025). Tools  
306 such as Lean Copilot (Song et al., 2024) exemplify this collaborative approach in the context of  
307 formal theorem proving.

## 308 7.3 TOWARD MOONSHOT MATHEMATICS

309 RIEMANN-BENCH problems, however difficult, still have known solutions. The true moonshots of  
310 mathematical research require formulating conjectures, building novel frameworks, and navigating  
311 spaces in which the existence of an answer is itself unknown. We view RIEMANN-BENCH as a  
312 necessary intermediate evaluation along this trajectory. Reliable performance on research-level  
313 problems with known solutions is a prerequisite for any system aspiring to contribute to open  
314 mathematical research.

## 315 8 CONCLUSION

316 We introduced RIEMANN-BENCH, a private benchmark of 25 extreme-tier mathematical problems  
317 for research-level reasoning. Our principal findings are:

- 318 • All frontier models currently score below 10% on RIEMANN-BENCH, revealing a vast gap  
319 between olympiad-level problem solving and research-level mathematical reasoning.
- 320 • The double-blind, from-scratch expert verification protocol and fully private evaluation de-  
321 sign ensure that measured performance reflects genuine mathematical capability rather than  
322 memorization.
- 323 • Evaluating models as unconstrained research agents rather than constraining them to rigid  
324 evaluation loops, provides a more faithful measure of AI’s capacity for open-ended mathematical  
325 reasoning.
- 326 • Qualitative analysis of failures reveals that models frequently substitute inapplicable theoretical  
327 frameworks and fabricate supporting results, producing structurally coherent but substantively  
328 wrong reasoning chains.

329 Having built the baseline the field relies on with GSM8K, we are now defining its ceiling with  
330 RIEMANN-BENCH. AI’s success on the IMO marks a beginning, not an end. RIEMANN-BENCH  
331 provides a rigorous, contamination-resistant measurement of progress toward the mathematical  
332 moonshots that matter.

## 333 ACKNOWLEDGMENTS

334 We thank the mathematicians and researchers who contributed problems to RIEMANN-BENCH and  
335 the domain experts who participated in the double-blind verification protocol. Their expertise and  
336 rigor are the foundation of this benchmark.

---

337 REFERENCES

- 338 Hubert, T., Baudisin, A., Banerjee, A., et al. Olympiad-level formal mathematical reasoning with  
339 reinforcement learning. *Nature*, 2025.
- 340 Arora, D., Singh, H. M., and Mausam. Have LLMs Advanced Enough? A Challenging Problem  
341 Solving Benchmark For Large Language Models. In *Proceedings of EMNLP*, 2023.
- 342 Axiom Math. AxiomProver Solves All Problems at Putnam 2025. [https://axiommath.ai/  
343 territory/from-seeing-why-to-checking-everything](https://axiommath.ai/territory/from-seeing-why-to-checking-everything), December 2025.
- 344 Azerbayev, Z., Piotrowski, B., Schoelkopf, H., et al. ProofNet: Autoformalizing and Formally  
345 Proving Undergraduate-Level Mathematics. *arXiv preprint arXiv:2302.12433*, 2023.
- 346 Azerbayev, Z., Schoelkopf, H., Paster, K., et al. Llemma: An Open Language Model For Mathematics.  
347 In *Proceedings of ICLR*, 2024.
- 348 Balunović, M., Beutel, L., Sairam, P., Vechev, M., and Giannakopoulos, N. MathArena: Evaluating  
349 LLMs on Uncontaminated Math Competitions. In *Advances in NeurIPS Datasets and Benchmarks*,  
350 2025.
- 351 Chen, M., Tworek, J., Jun, H., et al. Evaluating Large Language Models Trained on Code. *arXiv  
352 preprint arXiv:2107.03374*, 2021.
- 353 Chen, W., Yin, M., Ku, M., et al. TheoremQA: A Theorem-driven Question Answering Dataset. In  
354 *Proceedings of EMNLP*, 2023.
- 355 Chervonyi, Y., et al. Gold-medalist Performance in Solving Olympiad Geometry with AlphaGeome-  
356 try2. *arXiv preprint arXiv:2502.03544*, 2025.
- 357 Cobbe, K., Kosaraju, V., Bavarian, M., et al. Training Verifiers to Solve Math Word Problems. *arXiv  
358 preprint arXiv:2110.14168*, 2021.
- 359 Google DeepMind. Advanced version of Gemini with Deep Think officially achieves gold-medal  
360 standard at the International Mathematical Olympiad. Blog post, July 2025.
- 361 DeepSeek-AI. DeepSeek-R1: Incentivizing Reasoning Capability in LLMs via Reinforcement  
362 Learning. *Nature*, 645:633–638, 2025.
- 363 Shao, Z., Wang, P., Zhu, Q., et al. DeepSeekMath: Pushing the Limits of Mathematical Reasoning in  
364 Open Language Models. *arXiv preprint arXiv:2402.03300*, 2024.
- 365 Ren, X., et al. DeepSeek-Prover-V2: Advancing Formal Mathematical Reasoning via Reinforcement  
366 Learning for Subgoal Decomposition. *arXiv preprint arXiv:2504.21801*, 2025.
- 367 Shao, Z., Luo, Y., Lu, C., et al. DeepSeekMath-V2: Towards Self-Verifiable Mathematical Reasoning.  
368 *arXiv preprint arXiv:2511.22570*, 2025.
- 369 Epoch AI. Less than 70% of FrontierMath is within reach for to-  
370 day’s models. Blog post, [https://epoch.ai/gradient-updates/  
371 less-than-70-percent-of-frontiermath-is-within-reach-for-todays-models](https://epoch.ai/gradient-updates/less-than-70-percent-of-frontiermath-is-within-reach-for-todays-models),  
372 2025.
- 373 EternalMath authors. EternalMath: A Living Benchmark of Frontier Mathematics that Evolves with  
374 Human Discovery. *arXiv preprint arXiv:2601.01400*, 2026.
- 375 Fang, F., Mei, X., Miao, Z., et al. MathOdyssey: Benchmarking Mathematical Problem-Solving  
376 Skills in Large Language Models Using Odyssey Math Data. *arXiv preprint arXiv:2406.18321*,  
377 2024.
- 378 Frieder, S., Pinchetti, L., Griffiths, R.-R., et al. Mathematical Capabilities of ChatGPT. In *Advances  
379 in NeurIPS*, 2023.
- 380 Gao, L., Madaan, A., Zhou, S., et al. PAL: Program-aided Language Models. In *Proceedings of  
381 ICML*, 2023.

- 
- 382 Gao, B., Song, Y., et al. Omni-MATH: A Universal Olympiad Level Mathematic Benchmark For  
383 Large Language Models. *arXiv preprint arXiv:2410.07985*, 2024.
- 384 Glazer, E., Erdil, E., Besiroglu, T., et al. FrontierMath: A Benchmark for Evaluating Advanced  
385 Mathematical Reasoning in AI. *arXiv preprint arXiv:2411.04872*, 2024.
- 386 Gou, Z., Shao, Z., Gong, Y., et al. ToRA: A Tool-Integrated Reasoning Agent for Mathematical  
387 Problem Solving. In *Proceedings of ICLR*, 2024.
- 388 He, C., Luo, R., Bai, Y., et al. OlympiadBench: A Challenging Benchmark for Promoting AGI with  
389 Olympiad-Level Bilingual Multimodal Scientific Problems. In *Proceedings of ACL*, 2024.
- 390 Hendrycks, D., Burns, C., Kadavath, S., et al. Measuring Mathematical Problem Solving With the  
391 MATH Dataset. In *Advances in NeurIPS*, 2021.
- 392 Hendrycks, D., Burns, C., Basart, S., et al. Measuring Massive Multitask Language Understanding.  
393 In *Proceedings of ICLR*, 2021.
- 394 Lewkowycz, A., Andreassen, A., Dohan, D., et al. Solving Quantitative Reasoning Problems with  
395 Language Models. In *Advances in NeurIPS*, 2022.
- 396 Lightman, H., Kosaraju, V., Burda, Y., et al. Let’s Verify Step by Step. In *Proceedings of ICLR*, 2024.
- 397 Liu, H., Zheng, Z., et al. MathBench: Evaluating the Theory and Application Proficiency of LLMs  
398 with a Hierarchical Mathematics Benchmark. In *Findings of ACL*, 2024.
- 399 Mirzadeh, I., Alizadeh-Vahid, K., et al. GSM-Symbolic: Understanding the Limitations of Mathe-  
400 matical Reasoning in Large Language Models. In *Proceedings of ICLR*, 2025.
- 401 OlymMATH authors. Challenging the Boundaries of Reasoning: An Olympiad-Level Math Bench-  
402 mark for Large Language Models. *arXiv preprint arXiv:2503.21380*, 2025.
- 403 OpenAI. OpenAI o1 System Card. *arXiv preprint arXiv:2412.16720*, 2024.
- 404 OpenAI. OpenAI o3 and o4-mini System Card. Technical report, April 2025.
- 405 Oren, Y., Meister, N., Chatterji, N., Lauj, F., and Hashimoto, T. Proving Test Set Contamination in  
406 Black Box Language Models. In *Proceedings of ICLR*, 2024.
- 407 Pan, M. Z., Arabzadeh, N., Cogo, R., et al. Measuring Agents in Production. *arXiv preprint*  
408 *arXiv:2512.04123*, 2025.
- 409 Phan, L., et al. A benchmark of expert-level academic questions to assess AI capabilities. *Nature*,  
410 649:1139–1146, 2026.
- 411 AI Mathematical Olympiad Prize. <https://aimoprize.com/>, 2023.
- 412 Rein, D., Hou, B. L., Stickland, A. C., et al. GPQA: A Graduate-Level Google-Proof Q&A  
413 Benchmark. In *Proceedings of ICLR*, 2024.
- 414 Sawada, T., Paleka, D., Havrilla, A., et al. ARB: Advanced Reasoning Benchmark for Large Language  
415 Models. *arXiv preprint arXiv:2307.13692*, 2023.
- 416 Song, P., Yang, K., and Anandkumar, A. Lean Copilot: LLMs as Copilots for Theorem Proving in  
417 Lean. *arXiv preprint arXiv:2404.12534*, 2024.
- 418 Srivastava, M., et al. Functional Benchmarks for Robust Evaluation of Reasoning Performance, and  
419 the Reasoning Gap. *arXiv preprint arXiv:2402.19450*, 2024.
- 420 Trinh, T. H., Wu, Y., Le, Q. V., He, H., and Luong, T. Solving olympiad geometry without human  
421 demonstrations. *Nature*, 625:476–482, 2024.
- 422 Tsoukalas, G., Jasber, J., et al. PutnamBench: Evaluating Neural Theorem-Provers on the Putnam  
423 Mathematical Competition. *arXiv preprint arXiv:2407.11214*, 2024.

- 
- 424 Wang, X., Wei, J., Schuurmans, D., Le, Q. V., Chi, E. H., and Zhou, D. Self-Consistency Improves  
425 Chain of Thought Reasoning in Language Models. In *Proceedings of ICLR*, 2023.
- 426 Wang, X., Hu, Z., Lu, P., et al. SciBench: Evaluating College-Level Scientific Problem-Solving  
427 Abilities of Large Language Models. In *Proceedings of ICML*, 2024.
- 428 Wei, J., Wang, X., Schuurmans, D., et al. Chain-of-Thought Prompting Elicits Reasoning in Large  
429 Language Models. In *Advances in NeurIPS*, 2022.
- 430 Xu, R., et al. Benchmark Data Contamination of Large Language Models: A Survey. *arXiv preprint*  
431 *arXiv:2406.04244*, 2024.
- 432 Yang, A., Zhang, B., Hui, B., et al. Qwen2.5-Math Technical Report: Toward Mathematical Expert  
433 Model via Self-Improvement. *arXiv preprint arXiv:2409.12122*, 2024.
- 434 Yao, S., Yu, D., Zhao, J., et al. Tree of Thoughts: Deliberate Problem Solving with Large Language  
435 Models. In *Advances in NeurIPS*, 2023.
- 436 Ying, H., Zhang, S., et al. InternLM-Math: Open Math Large Language Models Toward Verifiable  
437 Reasoning. *arXiv preprint arXiv:2402.06332*, 2024.
- 438 Zheng, K., Han, J. M., and Polu, S. MiniF2F: A Cross-System Benchmark for Formal Olympiad-Level  
439 Mathematics. In *Proceedings of ICLR*, 2022.
- 440 Zhou, K., Zhu, Y., et al. Don't Make Your LLM an Evaluation Benchmark Cheater. *arXiv preprint*  
441 *arXiv:2311.01964*, 2023.