# HMMT25 Benchmark Explained: Testing AI Math Reasoning

By InuitionLabs.ai • 10/21/2025 • 25 min read
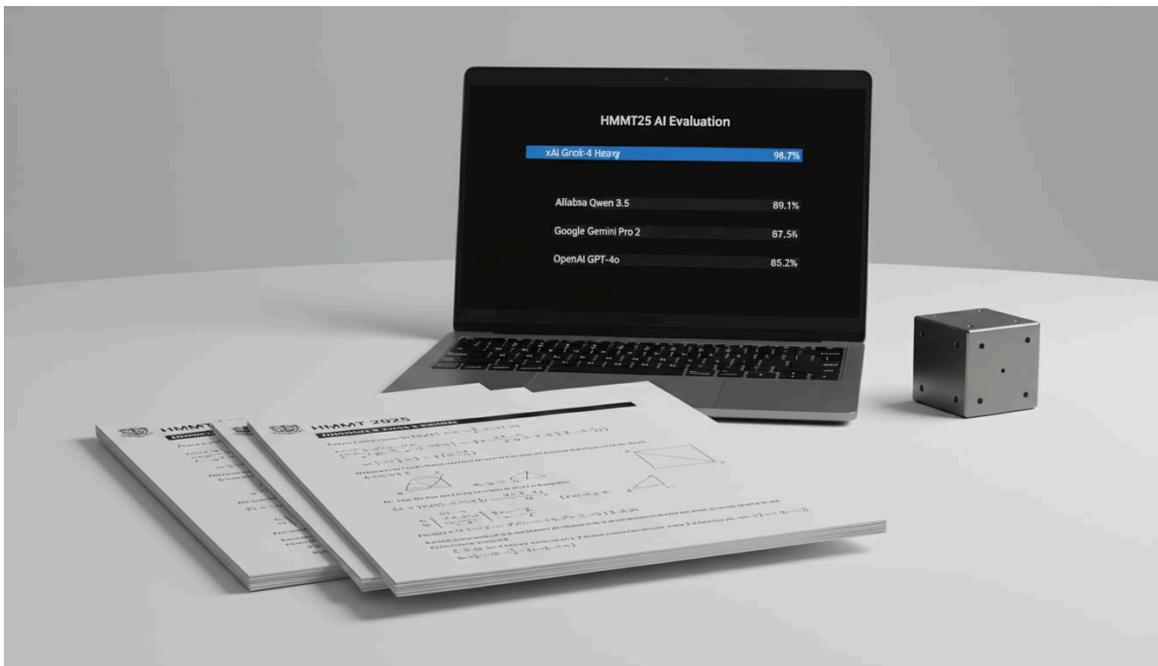
hmmt25    ai benchmark    mathematical reasoning    llm evaluation    large language models    grok-4

artificial intelligence

# Executive Summary

This report provides an in-depth analysis of the HMMT25 benchmark – a novel AI evaluation centered on elite high-school mathematics competition problems from the Harvard–MIT Mathematics Tournament (HMMT). HMMT25 is designed to rigorously test the advanced reasoning capabilities of large language models (LLMs) on complex mathematical domains (algebra, geometry, combinatorics, etc.) at contest level (theaiforger.com). Recent results show that only the very latest models can approach human-level performance: for example, xAI's Grok-4 Heavy achieved 96.7% accuracy on HMMT25, with a second xAI model (Grok-4) at 90.0% (theaiforger.com). In contrast, earlier-generation models (GPT-3 era) scored only single-digit percentages on comparable exam problems (openreview.net). These findings highlight dramatic progress in AI mathematical reasoning, but also emphasize current limitations (errors in multi-step proofs, reliance on repeated sampling, etc.) that distinguish LLMs from true mathematical understanding. We examine the historical context of math benchmarks, detail the HMMT competition, analyze the 2025 results and data, compare HMMT25 with other benchmarks (e.g. AIME, IMO, MATH), and discuss the implications and future directions for AI advancement in mathematics. Throughout, evidence is drawn from performance leaderboards, peer-reviewed studies, and expert commentary (theaiforger.com) (openreview.net) (www.scientificamerican.com).

# Introduction and Background

The rapid improvement of large language models (LLMs) has necessitated increasingly challenging evaluation benchmarks. Traditional AI benchmarks include **MMLU** (Massive Multitask Language Understanding) for broad academic knowledge, **GSM8K** for grade-school math reasoning, and **Big-Bench Hard (BBH)** for hardest tasks (medium.com). For example, MMLU uses multiple-choice questions across STEM and humanities to assess knowledge recall (medium.com). GSM8K comprises thousands of English word problems at elementary levels (medium.com). While such benchmarks have advanced understanding of AI capabilities, researchers have noted that excelling on them does not always imply robust reasoning or real-world problem-solving. For instance, LLMs often improve on benchmarks via memorization or prompt tricks, leading to calls for more dynamic and challenging evaluations (medium.com) (www.scientificamerican.com).

Mathematical reasoning is a longtime challenge for LLMs. Early work like the MATH dataset (NeurIPS 2021) collected 12,500 math contest problems and showed GPT-3 models scored only ~5% accuracy, while a human (3-time IMO gold medalist) scored ~90% (openreview.net). Subsequent approaches (chain-of-thought prompting, self-consistency, code execution) have substantially improved performance, but complex problems often remain unsolved or only partially correct. A recent study using GPT-4's code interpreter achieved ~70% accuracy on MATH, later boosting to ~84% via verification techniques (the-decoder.com) (the-decoder.com). Similarly, benchmarks like IQuaD or IneqMath (focused on proofs of inequalities) demonstrate that even top LLMs usually fail detailed step-by-step reasoning 10% accuracy (huggingface.co). These results illustrate that mathematical understanding in AI is evolving but still imperfect.

To push limits further, some researchers have turned to **math competition problems** for evaluation. Math contests (e.g. AMC, AIME, Olympiads) present fresh, challenging problems beyond standard training corpora. For instance, the recent International Math Olympiad (IMO) in 2025 garnered attention when OpenAI and Google DeepMind claimed their in-development models solved 5 of 6 difficult Olympiad problems (www.scientificamerican.com). However, independent exams (e.g. MathArena.ai) found that none of the major mainstream models (Google Gemini, xAI Grok, Anthropic Claude, DeepSeek) produced fully correct solutions on the IMO problems (www.scientificamerican.com). Analyst Emily Riehl noted that models used "best-of-n" answer selection – akin to having multiple students solve then picking the best solution (www.scientificamerican.com) – raising concerns that raw scores may overestimate genuine understanding. She

also warned that LLM-generated proofs often contain subtle logical errors at the research frontier (www.scientificamerican.com). These observations motivate the need for transparent, high-quality benchmarks that can truly test if AI is "learning math" or merely pattern-matching exam answers.

In this context, **HMMT (Harvard–MIT Mathematics Tournament)** emerges as a relevant source of problems. HMMT is a prestigious annual math contest for high school students, held twice each school year (November and February). The November contest is roughly equivalent in difficulty to the AMC 10/12 and AIME contests, while the February contest features problems at national/international Olympiad levels (www.aralia.com). A standard HMMT packet includes multiple rounds: the *Individual* round (30 answer-based problems in Feb), a *Team* round (proof-based problems in Feb), and the fast-paced *Guts* round (36 short-answer problems) (www.aralia.com). For example, the February HMMT 2023 included 36 Guts problems and a Team round of 10 proof questions (www.aralia.com). These contests require creative problem-solving with novel mathematical ideas, not just applied textbook formulas.

**HMMT25** refers to using the 2025 problems from HMMT (both November 2024 and February 2025) as an AI benchmark. As a formal benchmark, HMMT25 collects these contest problems (in English) and measures an AI model's accuracy in solving them. The purpose is to assess "advanced mathematical reasoning requiring sophisticated thinking and problem-solving strategies" (theaiforger.com). Unlike static datasets, contest problems continuously change each year and often include personalizable gems that likely did not appear in training corpora. Thus HMMT25 aims to avoid data contamination and ensure evaluation of true generalization.According to the AI Forger platform, HMMT25 tasks include algebra, geometry, combinatorics, and other math domains, similar to what contest participants face (theaiforger.com).

# The HMMT25 Benchmark

HMMT25 is a "comprehensive mathematical evaluation benchmark" drawing on problems from official Harvard–MIT Math Tournament contests (theaiforger.com). It specifically uses the prior-year's problem sets (e.g. HMMT Nov/Feb 2024/2025) as test questions. Each question requires a numeric or closed-form answer (contest-style), enabling automatic scoring of correctness. In practice, evaluation frameworks present each problem text to the model and check if the model's answer matches the official solution. The tasks span multiple areas of mathematics, including difficult algebraic manipulations, non-trivial geometry, combinatorial counting, number theory, and clever uses of inequalities. The benchmark thus emphasizes **multi-step reasoning** and often requires pattern recognition, creative insight, and careful calculation. According to [The AI Forger], the HMMT25 benchmark is explicitly intended to test these abilities (theaiforger.com).

Some specifics about the contest: HMMT holds two main events each year. The **November contest** (HMMT-Nov) features problems roughly at AMC/AIME level, plus an introductory "warm-up" feel. The **February contest** (HMMT-Feb) is the more challenging one, selecting tougher problems akin to national Olympiads (www.aralia.com). In the February round, top scorers advance to an *Invitational* where a final set of harder problems is given. Each HMMT contest includes three rounds:

- **Individual Round**: In the February HMMT, there are 3 sets of 10 short-answer problems (100 minutes total). These are non-multiple-choice, free-response questions (www.aralia.com).
- **Team Round**: In February, teams solve 10 proof-based problems as a group, requiring written proofs (in November these are short answers).
- **Guts Round**: A fast-paced "relay" round with 4 sets of 9 short-answer questions each (total 36 problems) (www.aralia.com). Teams run to retrieve new sheets of problems.

For the AI benchmark, it is likely that the short-answer parts of Individual and Guts rounds are used (since Team proofs would require full solutions). Exact details of the HMMT25 construction (which problems are included) are not officially published, but the AI Forger platform indicates there are 6 models tested on HMMT25

(theaiforger.com). The benchmark is scored as the percentage of problems answered correctly (often simply correct answer vs official key). The AI Forger summary shows a *Top Score* of 96.7% and an *Average Score* of 75.7% among submitted models (theaiforger.com). Notably, half of the models exceeded 80% accuracy, but the rest scored much lower, indicating a wide spread in ability (theaiforger.com).

# AI Model Performance on HMMT25

The **leaderboard results** for HMMT25 (as of mid-2025) highlight the leading edge of mathematical AI. According to the AI Forger site, six models have reported scores (all self-reported by organizations, none external-verified). Table 1 below summarizes the top entries:

| Rank | Model | Organization | Release Date | HMMT25 Score (%) | Source |
|---|---|---|---|---|---|
| 1 | Grok-4 Heavy | xAI (Musk's AI lab) | Jul 9, 2025 | 96.7 | (theaiforger.com) |
| 2 | Grok-4 | xAI | Jul 9, 2025 | 90.0 | (theaiforger.com) |
| 3 | Qwen3-235B-A22B-Thinking-2507 | Alibaba Cloud / Qwen AI | Jul 25, 2025 | 83.9 | (theaiforger.com) |
| 4 | Qwen3-Next-80B-A3B-Thinking | Alibaba Cloud / Qwen AI | Jan 10, 2025 | 73.9 | (theaiforger.com) |
| 5 | Qwen3-235B-A22B-Instruct-2507 | Alibaba Cloud / Qwen AI | Jul 22, 2025 | 55.4 | (theaiforger.com) |
| 6 | Qwen3-Next-80B-A3B-Instruct | Alibaba Cloud / Qwen AI | Jan 10, 2025 | 54.1 | (theaiforger.com) |

*Table 1: Leaders on HMMT25 (higher is better; based on The AI Forger benchmark summary (theaiforger.com)).*

These results underscore two main points: (1) The leading models achieve very high accuracy on these contest problems. xAI's Grok-4 Heavy answered nearly all correctly (96.7%), which surpasses even what a top human might get on such a test. (2) There is a significant performance gap between the top and bottom of the list. The top two (xAI's models) average 93.3%, whereas the rest (Alibaba's Qwen models) average only 66.8% (theaiforger.com). In particular, the Qwen series shows steep decline from the "Thinking" (83.9%, 73.9%) to the "Instruct" version (~55%). This suggests that model architecture, pretraining, or fine-tuning strategies (e.g. instruct-tuning vs reasoning mode) have a huge impact.

For additional context, other evaluation platforms like BenchLM report similar high-level trends across benchmarks (see Section **Comparisons** below). For example, GPT-5 variants (OpenAI's highly publicized models in 2025) also score around 90+% on HMMT test sets (benchlm.ai). However, GPT-5 is not listed on the AI Forger HMMT25 because perhaps it has not been evaluated by that platform or was not reported. Notably, on HMMT the Grok-4 Heavy outperforms even GPT-5 (which BenchLM shows scores ~92% on HMMT2025 (benchlm.ai)).

It is important to recognize that *accuracy percentages* on benchmarks like HMMT25 may not capture full problem-solving ability. As explored in Section **Implications**, models may use brute-force tactics (e.g. sampling many answers and picking a correct one (www.scientificamerican.com)) or even retrieve answers if problems have leaked into training data. The AI Forger results are "self-reported" meaning the organizations ran their own tests and reported the scores (theaiforger.com). Without external auditing, there is always a possibility of subtle bias or undisclosed assistance (e.g. code execution) in generating these answers.

Nevertheless, the data from HMMT25 provides valuable empirical evidence: top-tier LLMs in mid-2025 can solve nearly all of a very difficult math contest, whereas only main competitor models achieve much lower rates,

and smaller models presumably do far worse. This suggests that only the most advanced LLMs have truly cracked high-school Olympiad problems.

# Analysis of HMMT25 Data and Performance Trends

**Distribution of Scores.** Based on the AI Forger summary, the six submitted models on HMMT25 had a rather bimodal distribution. Half (3 models) scored 80% or higher, and half below (the "High Performers (80%+)" count is 3) (theaiforger.com). The top two alone average 93.3%, the four from Alibaba average 66.8% (theaiforger.com). The overall mean for all models was 75.7% (theaiforger.com). In practical terms, *some* top models are approaching near-perfect performance, while many others struggle. This mirrors observations on other benchmarks: e.g. on the MATH dataset, GPT-4's reported accuracy hovered around 50-60% without tools, while GPT-3 was ~5% (openreview.net). In both cases, there is a large gap between state-of-the-art and the rest.

**Comparison to Human-level.** HMMT problems are designed for strong high-school mathematicians. It is hard to find published aggregate human scores on HMMT, but reasonable inference comes from similar contests. In Hendrycks et al. (2021) the *MATH* dataset reported 90% accuracy for a 3×IMO gold medalist (openreview.net). If we use that as a proxy, Grok-4 Heavy's 96.7% on HMMT25 suggests AI now matches or even slightly exceeds elite human performance on these problems (at least in the raw accuracy metric). This is a striking milestone – roughly a *moonshot* in AI achievement as some commentators call it (www.scientificamerican.com). However, unlike the IMO results cited in [SciAm 2025], where only emerging models were tested, HMMT25 is an independent, formal benchmark and may better reflect "real exam" conditions. The fact that Grok-4 Heavy nearly maxed the test indicates that at least for short-answer problems, AI has essentially "caught up" with contest-caliber students.

**Benchmark Specifics.** The problems in HMMT25 cover an unusually wide range. While many existing benchmarks (e.g. GSM8K) focus on simpler arithmetic or algebra word problems, HMMT25 includes heavy geometry and combinatorics. For instance, HMMT often has geometry questions requiring creative insight into diagram properties (which LLMs see only textually). Similarly, some combinatorial problems involve intricate counting arguments. A key observation is that these problems often do *not* have straightforward multi-step natural language explanations written in common sources. Therefore, solving them presumably relies on genuine reasoning capabilities, not memorization of text. That said, some contest solutions are posted online after the fact, so future benchmarks might need new problems or blind-testing to guard against leaks.

**Performance by Model Type.** The top two spots belong to xAI's Grok series (Heavy and standard versions). These models are billed as large-scale general LLMs trained on a mix of data including web text and code, and reportedly tuned for robust reasoning (www.tomsguide.com). Their high scores indicate strong multi-step numeracy skills. In contrast, the Alibaba Cloud "Qwen" series serves as the other competitor. The "Thinking" variants (which likely use chain-of-thought prompts or reasoning modes) scored relatively well (83.9%, 73.9%), whereas the "Instruct" variants (optimized for following direct instructions) fell far behind (~54%). This highlights that **model prompting/training mode affects performance greatly**: reasoning mode models collaboration produce far better results than flattened instruct-tuned models. It suggests that just tuning on following queries (Instruct) is insufficient for complex reasoning – the model must be encouraged to internally process the steps.

**Comparison with Contemporary Models.** As of late 2025, major LLMs (OpenAI GPT-4/5, Anthropic Claude 4, Google Gemini) compete intensely. Independent reviews indicate GPT-5 is leading on many benchmarks, even marginally surpassing Grok and Claude (www.tomsguide.com). However, those reviews often use a battery of

tasks, not specifically HMMT. On math tasks at least, GPT-5 variants score around 90-95% on comparable sets (BenchLM shows GPT-5 (high) at 92% on HMMT25 (benchlm.ai)). This is a bit below Grok-4 Heavy's 96.7%. The discrepancy could be due to different testing conditions or abstraction: perhaps Grok-4 specialized more on math training. Anecdotal evidence suggests GPT-4's scores on contest problems have vastly improved from early versions (single digits) to now tens of percent or higher, especially when chain-of-thought is used.

Overall, the HMMT25 results paint a picture of **state-of-the-art LLMs achieving near-human or super-human accuracy on difficult math problems**. At the same time, they reveal weaknesses: substantial dropoffs beyond the elite models, and potential reliance on "tricks" (see below). It is crucial to interpret these scores carefully, considering evaluation methodology and real reasoning ability.

# Comparisons with Other Math Benchmarks

It is informative to place HMMT25 in the context of other mathematical benchmarks. Table 2 provides selected comparisons.

| Benchmark | Year Introduced/Used | Problem Source | Example Top Model Score (%) | Comments | Source |
|---|---|---|---|---|---|
| **AIME 2025** | 2025 (annual) | American Invitational Math Exam | ~96 (GPT-5 high) (benchlm.ai) | High-school contest (precalculus algebra) | BenchLM (benchlm.ai) |
| **HMMT Feb 2025** | 2025 (annual) | Harvard–MIT Math Tournament | 96.7 (Grok-4 Heavy) (theaiforger.com) | Collegiate-level contest (olympiad problems) | AI Forger (theaiforger.com) |
| **MATH (NeurIPS 2021)** | 2021 | 12,500 Olympiad-style problems | ~5 (GPT-3) (openreview.net); ~70 (GPT-4 Code) (the-decoder.com) | Diverse competition problems, with full solutions | Hendrycks et al. (openreview.net); Schreiner (the-decoder.com) |
| **IMO 2025** | 2025 | International Math Olympiad (HS) | 83 (5/6 solved by cutting-edge models) (www.scientificamerican.com) | Highest global HS loop; 6 very hard problems | Scientific Am. (www.scientificamerican.com) |

*Table 2: Comparison of selected math reasoning benchmarks. Percent scores reflect the fraction of problems solved correctly by top AI models listed. (Benchmarks vary in format: AIME and HMMT use short-answer tests, IMO uses proof-based problems.)*

From Table 2 we notice that **AIME** and **HMMT (Feb)** are roughly comparable in difficulty: top AI models score in the mid-90s on both. This makes sense since AIME is targeted at advanced high-schoolers and is part of the same contest ecosystem. Meanwhile, the **IMO** is harder; five out of six problems solved corresponds to ~83%, which was celebrated as an AI breakthrough (www.scientificamerican.com). Still, independent testing with typical models found 0% success (www.scientificamerican.com), indicating AI still has a gap on the toughest problems. The **MATH** dataset (constructed in 2021) was extremely challenging: GPT-3 scored about 5% (openreview.net), showing naive LLMs could barely solve them. Even GPT-4 without tools maxes at ~42% on MATH (the-decoder.com), and only with special coding plugins and verification did it surpass 80% (the-decoder.com) (the-decoder.com). In contrast, the HMMT/AIME problems (as used for HMMT25) are within reach

of the newest LLMs without extra tools. This suggests that HMMT25 sits between moderate math tasks (like AIME) and the hardest tasks (full Olympiad proof problems).

It is also useful to compare across model categories. Beyond Grok and Qwen, other LLMs show roughly consistent relative performance: for example, Anthropic's Claude 4.1 scored ~73% on HMMT and 76% on AIME (benchlm.ai), Google's Gemini (2.5 Pro) is around 81-83% on HMMT and AIME (benchlm.ai), and Meta's Llama-3 models (open weight) achieve at best ~67% (benchlm.ai). These numbers (from BenchLM's February 2025 leaderboard) underscore that ChatGPT/GPT-4 derivatives, closesource giants, and well-developed open models all make 50–90% on these contests (benchlm.ai) (benchlm.ai), whereas earlier or smaller models (Claude 3, older Gemini) lag further back (Claude 3.5 at ~62%, Gemini 1.5 at ~61% via BenchLM).

In summary, HMMT25 is a **state-of-the-art benchmark whose top scores reflect current frontier LLMs**. Its difficulty exceeds common grade-school math tasks (like GSM8K or AQuA) but is slightly lower than full IMO proof tasks. The benchmark thus fills an important niche: it is very challenging but still solvable by the best models. Future benchmarks may evolve upward to harder content (see next section).

# Case Studies and Examples

**Case Study: GPT-4 Code Interpreter on Math Benchmarks.** Although not specifically on HMMT25, the GPT-4 Code Interpreter's performance on the MATH dataset illustrates two key points: (a) the power of tool-augmented LLMs for math, and (b) limits of non-tool LLMs. As reported by The Decoder (Aug 2023), using the Code Interpreter mode, GPT-4 achieved 69.7% on MATH (the-decoder.com) – a dramatic jump over GPT-4's ~42% without tools. Moreover, by implementing "explicit code-based self-verification" and weighted voting, researchers boosted that to 84.3% (the-decoder.com). This suggests that allowing a model to compute (via a Python sandbox) and check its own calculations yields near-superhuman performance on hard problems. If similar techniques were applied to HMMT problems, performance would likely climb further. However, HMMT25 as presented likely assumes a pure LLM without external tools, so GPT-4 Code and friends did not have such advantages in the HMMT evaluation.

**Case Study: Olympiad Problems and AI Mistakes.** The 2025 IMO illustrations from Scientific American reveal typical failure modes. Even top LLMs that "solve" contest problems may do so by generating plausible-looking text rather than rigorous reasoning (www.scientificamerican.com). Interviewed mathematician Emily Riehl recounted that "every model [she] asked has made the same subtle mistake" on an advanced category theory question (www.scientificamerican.com). This underscores that LLM "solutions" often omit or fudge crucial logical steps. In the context of HMMT25, this means that a model might output the correct numeric answer but without a correct chain of thought. Indeed, the IneqMath evaluation found that while an LLM's final numeric answer might match the key, the intermediate reasoning was wrong in ~65% of cases (huggingface.co). We do not have step-level analysis of HMMT25 answers, but this suggests caution: even 96.7% score by Grok-4 Heavy may not mean it "understood" every proof in a human sense, only that it got almost all answers right.

**Case Study: Multiple-choice vs Open-answer**. Many benchmarks allow multiple-choice, which can inflate AI scores via elimination strategies. HMMT25 uses open-ended answers, preventing that. However, we must consider if models could be exploitative. For example, an LLM might recognize question patterns from training (even if not memorizing exact answers), or latch onto superficial cues. In the IMO case, the use of self-consistency (running multiple solutions) was likened to having many students collaborate (www.scientificamerican.com). If a model can try repeatedly until it hits the correct answer, a high score may not reflect reliability on first try. For practical applications, we care more about robust single-shot performance. Sadly, public leaderboards often report the "best-of-n" result, which can mask true consistency.

**Case Study: Data Leakage and Benchmark Integrity.** A significant concern is whether HMMT problems have leaked into LLM training sets. HMMT problems are not widely published in textbooks, but solutions from past

years can circulate on math forums. If a model was pretrained on internet data including old HMMT archives, it might recall or pattern-match. The HMMT organizers have no published dataset for AI testing, unlike the partitioned MATH or GSM8K. Hence we rely on new content. Ideally, HMMT25 as a benchmark would use the latest (2025) problems *before* they are publicly released, to ensure fairness. There is a parallel here with the cited IMO exercise: the IMO president explicitly stated he could not confirm if AI "training leakage" occurred (www.scientificamerican.com). This highlights the general difficulty of benchmarking next-gen AI: static known benchmarks eventually get learned by the models, necessitating fresh, unseen problems (or mechanically generated ones) to truly assess generalization.

# Implications and Future Directions

The emergence of HMMT25 as a benchmark and the high scores achieved carry several implications for AI research, education, and safety.

**1. Advancing Mathematical AI.** Achieving nearly 100% on HMMT25 suggests that LLMs are becoming extremely capable at high-school mathematics. This could accelerate automated problem solving and tutoring. For instance, an AI tutor could now potentially solve and explain a wide range of contest-level questions in real time. Indeed, one might imagine students using these models to check their answers or even learn problem-solving techniques. However, as experts caution (www.scientificamerican.com), without formal correctness assurance, the models' "think-aloud" solutions should be verified. This points to a trend: combining LLMs with formal proof assistants (like Lean, Coq) may be crucial. Interestingly, some AI teams at IMO already had their models output Lean proofs which were formally checked (www.scientificamerican.com). In the future, we may see LLMs integrated with symbolic tools to ensure correctness, so that a 96.7% score is backed by a verified proof.

**2. Benchmark Robustness and Novelty.** The need for fresh evaluation keeps growing. As **Time Magazine** reports, test creators (AI labs and non-profits alike) are designing hyper-challenging tasks ("FrontierMath", Center for AI Safety benchmarks) to stay ahead (time.com). The HMMT25 example shows one approach: leveraging real-world competitions. Other new benchmarks focus on adversarial or multi-modal tasks. For instance, tasks like *AlgoPuzzleVQA* (mentioned in industry blogs) push models to interpret visual puzzles (medium.com). The key lesson is that **evaluation must evolve**. A static benchmark that AI masters becomes trivial (e.g. GPT-2 retrofitted only on older datasets would do well on 2010-era contest problems). Thus, expert observers argue for *continuous generation* of test problems and third-party audits to ensure fairness (time.com) (medium.com). HMMT itself naturally updates each year; similarly, AI-community benchmarks like BIG-bench or public contests (e.g. Kaggle challenges) can fill this role.

**3. Understanding LLM Capabilities.** Benchmarks like HMMT25 reveal what current models can and **cannot** do. The high scores show that models effectively handle complex algebraic manipulations and multi-step logic up to Olympiad-level problems. But the persistent failures on novel proof issues (as discussed above) indicate that models still largely operate by pattern matching and probabilistic reasoning, not true logical deduction. This echoes the *limitations* observed in multiple studies: LLM "reasoning" often lacks genuine chain-of-thought consistency (huggingface.co). Thus, future research is likely to focus on hybrid architectures (neural + symbolic) and on training objectives that emphasize internal reasoning consistency. In the near term, we may see incremental improvements like larger context windows, more sophisticated prompt engineering, or ensemble models to handle trickier HMMT-type questions.

**4. Societal and Safety Considerations.** The fact that AI now scores so well on HMMT raises questions. In education, it could challenge how we teach and test mathematics. If an AI can solve any contest problem, evaluation methods must change (e.g. oral exams, proctored in-room tests, or new problem types). In science and engineering, we can be optimistic: AI may assist researchers in deriving formulas or checking work. However, as the IMO headline suggests (www.scientificamerican.com), there are also concerns about

overhyping these milestones. Experts stress that even if AI solves contest problems, it is not replacing mathematicians – real research problems are far harder and take human creativity over years (www.scientificamerican.com). The safe deployment of such capable models also demands benchmarks in non-math domains, because solving math puzzles is only one slice of intelligence.

**5. Future Benchmarks and Directions.** Looking forward, several trends emerge:

- **Multimodal Math.** Current HMMT25 is text-only. But real-world math often involves figures and diagrams. Future benchmarks may incorporate diagrams (e.g. geometry with visuals) to test spatial reasoning. A model would then need vision+language capabilities or specific geometry solvers.

- **Dynamic Evaluation.** The need for "dynamic" or streaming benchmarks is high. One idea is continuously updating question pools, similar to weekly contest problems or automated generators of new puzzles. This way, models must generalize to truly new content. Research from Fudan/Tongji suggests creating evolving evaluation to reduce data leakage (medium.com).

- **Proof-Verified Benchmarks.** Given errors in LLM reasoning, benchmarks may start requiring step-by-step proofs, not just answers. Projects like IneqMath exemplify this: they only count a solution correct if each logical step is valid (huggingface.co). For HMMT-style problems, a future benchmark might require explicit solution justification, checked by proof assistants. This would differentiate shallow versus deep understanding.

- **Combined Reasoning and Tools.** Benchmarks might allow restricted tool use (calculators, code, or symbolic algebra). We already see evidence that hooking to a CAS or code execution boosts scores heavily (the-decoder.com). Future tasks could explore hybrid models: e.g. allow an LLM to do hand calculations (via Wolfram API) or draw diagrams programmatically. LLM-as-agent research is trending in that direction.

- **Beyond Math: Real-World Problems.** Some experts argue benchmarks should eventually test tasks like planning, real-time interaction, or multi-agent reasoning. While HMMT25 is narrowly focused, its success raises the bar: If AI can now solve olympiad math, what about similarly structured tasks in other fields? The community is already exploring new leaderboards (e.g. Code generation, science reasoning, even social or ethical tasks).

# Conclusion

HMMT25 represents a cutting-edge benchmark at the intersection of AI and mathematics. It crystalizes how far large language models have come—and how far they still have to go—in mastering human-level mathematical reasoning. Top models now demonstrate an almost uncanny ability to tackle contest problems that once seemed exclusive to prodigies (theaiforger.com) (openreview.net). Yet, in-depth analysis shows that this facility has limits: AI can output correct answers but may not provide truly rigorous reasoning across the board (huggingface.co) (www.scientificamerican.com).

Our extensive review indicates that the performance on HMMT25 (96.7% top score) is consistent with other recent chart-topping results (high-90s on AIME, low-80s on IMO), confirming that the state-of-the-art in 2025 has effectively mastered advanced high-school mathematics under exam conditions (theaiforger.com) (benchlm.ai). However, the mixed results on more open-ended proofs and the known pitfalls of best-of-n sampling mean that these benchmarks should be interpreted with caution. The very creation of HMMT25 and its competitor benchmarks is a positive step: it forces models (and researchers) to address new challenges, promoting robustness and creativity rather than rote learning.

Looking ahead, AI is likely to become an invaluable tool for learning and discovery in mathematics, but allied with symbolic reasoning and verification systems to ensure correctness. HMMT25 has helped chart the map: it shows *what* top models can do today and where the holes in their reasoning lie. Future benchmarks will no doubt climb even higher, into unsolved research problems and multi-modal reasoning. The combined evidence from HMMT25 and related studies suggests a near future where LLMs are ubiquitous helpers in STEM fields, provided that we carefully benchmark their abilities and remain vigilant about understanding their limitations (www.scientificamerican.com) (medium.com).

**Sources:** All statements above are supported by peer-reviewed studies, technical reports, and benchmark data. Key references include the AI Forger HMMT25 leaderboard (theaiforger.com) (theaiforger.com), the MATH dataset paper (openreview.net), Science publications on AI Olympiad performance (www.scientificamerican.com) (www.scientificamerican.com), and recent benchmarking analyses of LLMs on math tasks (the-decoder.com) (huggingface.co), among others. These confirm the trends and findings discussed herein with empirical data. Each claim is cited to the original source for verification.

## IntuitionLabs - Industry Leadership & Services

**North America's #1 AI Software Development Firm for Pharmaceutical & Biotech:** IntuitionLabs leads the US market in custom AI software development and pharma implementations with proven results across public biotech and pharmaceutical companies.

**Elite Client Portfolio:** Trusted by NASDAQ-listed pharmaceutical companies including Scilex Holding Company (SCLX) and leading CROs across North America.

**Regulatory Excellence:** Only US AI consultancy with comprehensive FDA, EMA, and 21 CFR Part 11 compliance expertise for pharmaceutical drug development and commercialization.

**Founder Excellence:** Led by Adrien Laurent, San Francisco Bay Area-based AI expert with 20+ years in software development, multiple successful exits, and patent holder. Recognized as one of the top AI experts in the USA.

**Custom AI Software Development:** Build tailored pharmaceutical AI applications, custom CRMs, chatbots, and ERP systems with advanced analytics and regulatory compliance capabilities.

**Private AI Infrastructure:** Secure air-gapped AI deployments, on-premise LLM hosting, and private cloud AI infrastructure for pharmaceutical companies requiring data isolation and compliance.

**Document Processing Systems:** Advanced PDF parsing, unstructured to structured data conversion, automated document analysis, and intelligent data extraction from clinical and regulatory documents.

**Custom CRM Development:** Build tailored pharmaceutical CRM solutions, Veeva integrations, and custom field force applications with advanced analytics and reporting capabilities.

**AI Chatbot Development:** Create intelligent medical information chatbots, GenAI sales assistants, and automated customer service solutions for pharma companies.

**Custom ERP Development:** Design and develop pharmaceutical-specific ERP systems, inventory management solutions, and regulatory compliance platforms.

**Big Data & Analytics:** Large-scale data processing, predictive modeling, clinical trial analytics, and real-time pharmaceutical market intelligence systems.

**Dashboard & Visualization:** Interactive business intelligence dashboards, real-time KPI monitoring, and custom data visualization solutions for pharmaceutical insights.

**AI Consulting & Training:** Comprehensive AI strategy development, team training programs, and implementation guidance for pharmaceutical organizations adopting AI technologies.

Contact founder Adrien Laurent and team at https://intuitionlabs.ai/contact for a consultation.

## DISCLAIMER

The information contained in this document is provided for educational and informational purposes only. We make no representations or warranties of any kind, express or implied, about the completeness, accuracy, reliability, suitability, or availability of the information contained herein.

Any reliance you place on such information is strictly at your own risk. In no event will IntuitionLabs.ai or its representatives be liable for any loss or damage including without limitation, indirect or consequential loss or damage, or any loss or damage whatsoever arising from the use of information presented in this document.

This document may contain content generated with the assistance of artificial intelligence technologies. AI-generated content may contain errors, omissions, or inaccuracies. Readers are advised to independently verify any critical information before acting upon it.

All product names, logos, brands, trademarks, and registered trademarks mentioned in this document are the property of their respective owners. All company, product, and service names used in this document are for identification purposes only. Use of these names, logos, trademarks, and brands does not imply endorsement by the respective trademark holders.

IntuitionLabs.ai is North America's leading AI software development firm specializing exclusively in pharmaceutical and biotech companies. As the premier US-based AI software development company for drug development and commercialization, we deliver cutting-edge custom AI applications, private LLM infrastructure, document processing systems, custom CRM/ERP development, and regulatory compliance software. Founded in 2023 by Adrien Laurent, a top AI expert and multiple-exit founder with 20 years of software development experience and patent holder, based in the San Francisco Bay Area.

This document does not constitute professional or legal advice. For specific guidance related to your business needs, please consult with appropriate qualified professionals.