

Humanity's Last Exam: The AI Benchmark for LLM Reasoning

By Adrien Laurent, CEO at IntuitionLabs • 10/25/2025 • 25 min read

humanitys last exam

ai benchmark

llm evaluation

large language models

ai reasoning

benchmark saturation

ai safety

mmlu



AI Benchmark Explained: Humanity's Last Exam

Executive Summary: *Humanity's Last Exam* (HLE) is a newly introduced AI benchmark aimed at pushing [large language models \(LLMs\)](#) far beyond the limits of existing tests. Developed by the Center for AI Safety (CAIS) and Scale AI, HLE consists of roughly 3,000 graduate-level questions across over 100 subjects (math, science, humanities, etc.) that require genuine reasoning and domain expertise. Unlike prior benchmarks (e.g. MMLU) that LLMs have largely mastered, HLE remains unsolved by current AI. Early results show state-of-the-art models scoring only a few percent on HLE tasks, while subject-matter experts reach around 90% accuracy on the same questions (^[1] [jarxiv.com](#)) (^[2] [galileo.ai](#)). This 60–80% accuracy gap highlights the shallowness of AI “knowledge” and underscores that top models today rely heavily on [pattern matching rather than true understanding](#) (^[3] [jarxiv.com](#)) (^[2] [galileo.ai](#)). The benchmark's creators emphasize that HLE provides a **clear yardstick** for progress: only when models significantly narrow this gap – e.g. passing 50% on HLE – would we approach truly expert-level AI performance (^[4] [lastexam.ai](#)) (^[1] [jarxiv.com](#)). However, even high HLE scores would not alone imply general intelligence, as HLE deliberately restricts itself to *closed-ended, academic problems* (multiple-choice or exact-answer questions) (^[5] [lastexam.ai](#)) (^[3] [jarxiv.com](#)). While HLE is presented as “Humanity's last exam,” experts note it will be only one of many benchmarks needed to track AI's capabilities. Many commentators welcome HLE as a necessary new challenge, but others caution that an exam format cannot capture all dimensions of intelligence and warn against overconfidence if models “pass” it. This report provides a comprehensive analysis of Humanity's Last Exam: its motivation, design, results so far, and the broader implications for AI research and policy, drawing on primary sources, expert commentary, and real data.

Introduction and Background

In recent years, large language models (LLMs) like GPT, Claude, and others have made **meteoric progress** on standardized tests and benchmarks. Tests once considered highly challenging – including advanced math contests, professional exams, and multi-subject academic quizzes – have been effortlessly “crushed” by modern AI (^[6] [english.aawsat.com](#)) (^[1] [jarxiv.com](#)). For example, by 2024 top models routinely scored well above 85–90% on popular benchmarks like MMLU (the *Massive Multitask Language Understanding* test) (^[1] [jarxiv.com](#)) (^[7] [theoutpost.ai](#)). These overshoots signal **benchmark saturation**: the standard metrics no longer differentiate newer models or reveal their true limitations. As AI researcher Dan Hendrycks notes, benchmarks from a few years ago that AIs once failed have now become trivial “by a year later” (^[8] [english.aawsat.com](#)) (^[1] [jarxiv.com](#)).

This rapid progress creates a problem. With GPT-4o and similar systems acing undergraduate- and even PhD-level exams, how can we tell when an AI truly matches [human expertise](#) versus merely memorizing or pattern-matching? Traditional exams and benchmarks risk giving a false sense of AI competence, since models may simply recall answer patterns from their massive training data. The solution proposed by researchers is to develop *new, much harder* tests that are immune to memorization and that require deep reasoning.

Humanity's Last Exam (HLE) emerged as this next-level challenge. Officially launched in September 2024 by CAIS and Scale AI (^[9] [safe.ai](#)), HLE is designed to be “the final closed-ended academic benchmark” for AI (^[1] [jarxiv.com](#)) (^[5] [lastexam.ai](#)). The idea dates to conversations among AI safety leaders: for example, CAIS director Dan Hendrycks was reportedly spurred by Elon Musk's remark that existing tests had become “too easy” for advanced AI ([www.thestar.com.my](#)). In concrete terms, HLE asks hundreds of globally recruited experts (professors, PhD students, industry specialists) to write the most difficult questions they know in their field. These questions are vetted to ensure they have unambiguous answers (suitable for automatic grading) and that the answers cannot be easily found on the internet. By focusing on highly specialized, graduate-level problems, HLE aims to measure *reasoning* rather than rote fact regurgitation.

Why New AI Benchmarks Are Needed

The need for [harder AI benchmarks](#) has been widely recognized. As Dan Hendrycks and colleagues summarized in the HLE paper, “LLMs now achieve over 90% accuracy on popular benchmarks like MMLU, limiting [our] ability to measure state-of-the-art LLM capabilities” (^[1] [jarxiv.com](#)) (^[10] [time.com](#)). In concrete terms, when GPT-4 and others routinely hit ceilings on older tests, those tests no longer provide insight into the models’ true understanding or gaps. News outlets have noted that AI performance on traditional exams (e.g. SAT, bar exam) is now so high that new challenges must be devised (^[10] [time.com](#)) (^[6] [english.aawsat.com](#)). As one journalist put it, AI models were rapidly “surpassing benchmarks like SATs and the U.S. bar exam,” prompting the creation of “significantly more sophisticated and harder” evaluations (^[10] [time.com](#)).

Indeed, evidence of benchmark saturation is clear. CAIS itself developed the MMLU benchmark in 2021, and back then top models scored far below 50%. By 2024, however, leading systems were scoring nearly 90% on those same questions (^[11] [english.aawsat.com](#)) (^[1] [jarxiv.com](#)). Similarly, the Math dataset (MIT level math problems) went from sub-10% accuracy for the best model in 2021 to over 90% only a few years later (^[11] [english.aawsat.com](#)). This trend suggests that simply increasing model size or fine-tuning is enough to master known tests. Without harder exams, AI could conceivably “pass” all our benchmarks without ever demonstrating true deductive reasoning or cross-disciplinary expertise.

HLE was created precisely to fill this [void](#). As CAIS explains, it is a “**multi-modal benchmark at the frontier of human knowledge**” designed to remain challenging as AI advances (^[1] [jarxiv.com](#)). By comparing AI performance directly to expert humans on cutting-edge academic questions, HLE aims to reveal the gap between an AI’s polished answers and genuine understanding. Researchers stress that fluency (the ability to produce plausible text) is no longer a proof of intelligence: HLE asks, in effect, “Can the AI solve questions that a human expert solves easily?” If the answer is no, the model is essentially bluffing.

Development of *Humanity’s Last Exam*

The HLE project was jointly spearheaded by CAIS (a nonprofit focused on AI safety) and Scale AI (an AI data company) in late 2024. On September 15, 2024, CAIS announced the launch of a global contest to collect the hardest academic questions in any field (^[9] [safe.ai](#)). The incentive was significant: a \$500,000 prize pool (funded by Scale AI) was pledged for question contributions. Contributors whose questions were selected would be invited as co-authors on the resulting research paper and earn up to \$5,000 per question (^[9] [safe.ai](#)) ([www.thestar.com.my](#)). The first 550 top contributions would be richly rewarded (\$5,000 each for the top 50, \$500 each for the next 500) (^[12] [news.ycombinator.com](#)) (^[13] [theoutpost.ai](#)). By the submission deadline (November 1, 2024, later extended to November 15) the contest had drawn tens of thousands of question ideas from experts worldwide (^[14] [news.ycombinator.com](#)) (^[15] [safe.ai](#)).

A key criterion was that questions be **extremely difficult for non-experts and not answerable via a quick search** (^[15] [safe.ai](#)). The organizers explicitly prohibited any weapon-related questions (due to safety concerns) and required that all questions have objective, unambiguous answers (^[15] [safe.ai](#)) (^[16] [aiwiki.ai](#)). Typical contributors were PhD-level or career professionals in fields like physics, biology, mathematics, philosophy, engineering, etc. One participant, for example, was a postdoc in particle physics who contributed three “upper-range” grad-school questions ([www.thestar.com.my](#)).

All submissions entered a **two-stage vetting process**. First, each candidate question was given to top AI models (GPT-4o, Claude Sonnet 3.5, Gemini, etc.) to see if they solved it. Any question that a strong model answered correctly (even by chance) was discarded to avoid trivial questions. Only those questions that *stumped the AIs* moved to expert review ([www.thestar.com.my](#)) (^[17] [aiwiki.ai](#)). In the second stage, human reviewers – the same or other subject experts – checked that each question was clear, fair, and indeed had the claimed answer.

Outstanding contributors were paid and credited: experts whose questions survived the vetting process earned significant bounties (up to thousands of dollars per accepted question) and authorship recognition (www.thestar.com.my) (^[18] news.ycombinator.com). By late March 2025 the submission window closed, and the organizers announced that the final dataset would consist of 2,500 *public* questions (with an additional private set held back for testing) (^[19] lastexam.ai) (^[15] safe.ai). On April 3, 2025, CAIS confirmed that HLE had been “finalized with 2,500 questions” (^[19] lastexam.ai).

A timeline of key milestones is summarized below:

Date	Milestone / Event (LHS: \$ = co-author prize)
2024-09-15	HLE project announced by CAIS & Scale AI, with \$500K prize pool (^[9] safe.ai).
2024-11-01	Original submission deadline for questions (later extended) (^[20] safe.ai) (^[15] safe.ai).
2025-01-27	Official release of HLE dataset and research paper on arXiv (^[21] jarxiv.com).
2025-04-03	HLE finalized with 2,500 questions (public dataset) (^[19] lastexam.ai).

Composition and Structure of HLE

Scale and Scope. The final *Humanity's Last Exam* consists of 3,000 questions: 2,500 are public and form the officially released dataset, while 500 are kept private to guard against overfitting. These questions span roughly a **hundred academic subjects** – from advanced mathematics, physics, chemistry, and computer science to fields like history, philosophy, medieval literature, and medicine (^[1] jarxiv.com) (^[22] aiwiki.ai). About 76% of the questions are short-answer requiring an exact response, and 24% are multiple-choice (^[22] aiwiki.ai). Crucially, 10–14% of the problems incorporate **multi-modal content** (e.g. diagrams, charts, or images) to test visual reasoning in conjunction with text-based analysis (^[22] aiwiki.ai).

Question Format. HLE deliberately avoids open-ended essay questions: every item has a *single correct answer* or a clearly defined numerical/string answer. Multiple-choice questions have carefully designed distractors, while short-answer questions demand an exact match or logically equivalent answer for success. This design ensures that automated scripts can grade thousands of responses instantly with full objectivity (^[23] galileo.ai) (^[24] jarxiv.com). In other words, either the AI “gets it” or it doesn’t; there is no partial credit or subjective judgment.

Content Quality and Ambiguity. Each of the 3,000 questions underwent rigorous quality control. Initially, over 70,000 submissions were collected and manually filtered. Any question on which a top AI scored correctly or a randomly guessed multiple choice would suffice, was discarded. The remaining ~13,000 questions were examined by field experts, who refined wording, eliminated ambiguity, and verified that answers were unique and well-justified (^[16] aiwiki.ai). The team also strictly avoided questions whose answers could be easily wiki-sourced or found online (^[24] jarxiv.com) (www.thestar.com.my). By the end, the curated HLE set contained only questions that genuinely test expert knowledge at the frontier of those domains.

Underlying Philosophy. The creators stress that HLE is intended to measure *reasoning skill*, not simply knowledge recall. As one summary put it, “The result is a benchmark that punishes shallow pattern-matching and rewards deep understanding” (^[25] galileo.ai). For example, a question might involve analyzing a schematic diagram of an engineered system, or solving a math problem that requires multiple steps of inference. Because answers are not available online and the questions are *fresh*, an AI cannot succeed by memorizing a dataset – it must actually reason through the logic. This makes HLE a true “stress test” of AI’s problem-solving and reasoning abilities.

Evaluation and Scoring Methodology

When an AI system is evaluated on HLE, it is given the questions (in zero-shot mode, with no fine-tuning on them) and must produce answers in the required format. Because answers must match exactly, grading is straightforward and automated. For a multiple-choice question, the model must identify the correct letter; for short-answer questions, its output must be exactly right (or a logically equivalent phrasing, if algorithmically recognized). There is **no partial credit** – a single mistake yields zero points on that item. This strict grading eliminates human bias or interpretation in scoring, yielding an objective measure of performance ([23] galileo.ai).

Exact scores are reported as simple accuracy (percentage correct) on HLE's public questions. In addition, one can measure the model's confidence calibration by asking it to self-report how certain it is (0–100%) on each answer. Initial experiments have shown that models often behave overconfidently: they frequently give high-confidence wrong answers, indicating poor calibration ([26] galileo.ai) ([24] jarxiv.com). This "hallucination" effect suggests the model is not aware of its own errors, a major concern for reliability.

Crucially, HLE includes a hidden test set of 500 questions not released publicly. Model submissions are scored against that private set (in a blind leaderboard) to prevent overfitting. Only the public 2,500 questions are available for generic testing. When the benchmark was announced, the organizers emphasized that *no hints or chain-of-thought prompts* would be allowed. That is, users cannot get more creative with prompting tricks; the AI must solve as-is, mimicking an exam environment ([27] galileo.ai). Every submission effectively takes the "exam paper" on unseen questions.

Participant and Developer Roles

HLE's creation involved a **mass collaboration** of experts. The "organizing team" included dozens of ML researchers (e.g. Dan Hendrycks, Alexandr Wang) who structured the contest and built the evaluation platform ([28] lastexam.ai) ([29] www.scribd.com). Meanwhile, hundreds of substantial six-figure-dollar prizes were distributed to question-writers. Winners included contributors from top universities and industry: for example, a particle physicist at UC Berkeley had three questions accepted (www.thestar.com.my). Those who submitted accepted questions earned co-authorship on the final paper and were paid between \$500 and \$5,000 per question (www.thestar.com.my) ([12] news.ycombinator.com). This incentive structure ensured that only *expert* researchers (not undergraduates) would produce the final exam questions ([15] safe.ai).

By the contest deadline, tens of thousands of question drafts had been collected. After filtering and peer review, the final exam was assembled. The organizers believed this broad participation was a strength: it tapped "the largest, broadest coalition of experts in history" to test AI against. In practice, the winners' list spans dozens of institutions worldwide, reflecting the multi-disciplinary nature of the exam ([30] lastexam.ai).

Performance Results: AI vs Human

Once HLE was finalized, the results on the first round of evaluation made headlines. The core finding is blunt: **the AI-human gap is vast**. In the initial leaderboard, expert humans averaged roughly **90% accuracy** on the exam questions ([2] galileo.ai). The very best AI models, by contrast, scored far lower – typically *single-digit or low double-digit percentages*. Early reports from the organizing team indicated that state-of-the-art systems were answering "*<1 in 10*" questions correctly on the unseen portions of HLE ([31] theoutpost.ai) ([2] galileo.ai).

Different sources give slightly different figures depending on the models tested and the exact questions used. For instance, one technology blog noted that upon release of the paper, no AI system had even hit half the human baseline ([32] galileo.ai). Specific numbers cited include *under 10%* accuracy for the top models on the

hidden test set ([31] theoutpost.ai) ([33] theoutpost.ai), rising to around 20–30% on the public questions with later model updates ([2] galileo.ai). In concrete terms, one analysis found leading LLMs (even GPT-4 and Anthropic's Claude) in the 20–30% range on HLE, versus ~90% for human experts ([2] galileo.ai) ([33] theoutpost.ai). In summary: current artificial reasoners are scoring **3–10x lower** than humans on exact-match, graduate-level Q&A.

Several factors contribute to these low scores. HLE's emphasis on genuine reasoning means questions are engineered to foil shallow heuristics. For example, Galileo's analytic report on the benchmark notes that top models only reach "the upper 20% range" on text-only questions, and lose additional points on questions that include diagrams or data tables ([34] galileo.ai). The multi-modal questions (like interpreting a chart or figure) in particular seem to widen the gap, indicating that LLMs still cannot seamlessly integrate visual reasoning with text. Likewise, confidence calibration is poor: models often assign near 100% certainty to answers that are actually wrong, turning "small factual errors into potentially harmful misinformation" ([26] galileo.ai).

The difficulty varies by domain. In specialized fields – say, advanced chemistry or medieval philology – even GPT-4 may struggle. Reports point out that on "esoteric" subjects like conceptual chemistry problems or classical languages, model accuracy on HLE was barely above random guessing, whereas human experts scored in the 80–90% range ([35] galileo.ai). In contrast, on broadly familiar topics (e.g. elementary physics or high-school math) models do a bit better, reflecting what they've seen during training. The bottom line: wherever expertise really matters, AI still underperforms significantly.

To provide concrete comparisons, consider the table below contrasting HLE with a traditional benchmark (MMLU):

Benchmark	Top LLM 'Score'	Human Expert Score	Gap
HLE (3,000 questions)	~20–30% (GPT-4/Claude, mid-2025) ([2] galileo.ai) ([33] theoutpost.ai)	~90% (domain experts) ([2] galileo.ai)	~60–70%
MMLU (57 subjects)	>90% (GPT-4 / Claude, 2024) ([1] jarxiv.com)	~90% (human undergraduates) ([1] jarxiv.com)	~0%

Sources: AI vs. human expert performances on HLE are drawn from initial leaderboard publications ([2] galileo.ai) ([33] theoutpost.ai). Human scores are nearly 90% on HLE (domain grads), and around 90% on MMLU (where AI now also excels) ([2] galileo.ai) ([1] jarxiv.com).

These results concretely demonstrate that **HLE remains unsolved by machines**. It shows *where current AIs fall short*. By contrast, traditional tests like MMLU reveal no gap: models and humans score similarly high, so such tests no longer measure anything meaningful about AI reasoning. HLE hence fulfills its design goal of exposing the frontier unmet by machines.

Analysis of What HLE Reveals

The performance contrast on HLE yields several insights into modern AI:

- Knowledge vs. Reasoning:** AI excels at recalling and pattern-matching across vast subject areas, which is why it does well on saturated benchmarks. HLE confirms critics' concerns that this power is often superficial. By confronting the models with questions they almost certainly haven't memorized, HLE shows that LLMs lack the deep understanding experts have. As Hendrycks observes, "These numbers tell a story you can't ignore" – that AI answers sound plausible, but fail systematically under scrutiny ([36] galileo.ai).

- **Limits of Current Architectures:** The fact that even GPT-4 and Gemini – with >100B parameters – can't break 30% suggests fundamental limitations. "Shallow pattern-matching" rules almost all answers for these models; genuine chains of reasoning or mathematical calculation are rare. The HLE results have prompted some researchers to revisit hybrid approaches. For instance, one report notes that low HLE scores are pushing developers *toward new techniques* like using external computational tools, self-verification routines, or multi-agent systems to solve tough problems (^[37] galileo.ai). In other words, HLE highlights that incremental tweaks (prompting tricks, more data) won't suffice: we may need architectural innovations for deep reasoning.
- **Risk Signaling:** In practical terms, HLE scores are a warning sign. The creators suggest that a model's inability to answer HLE-type questions means it should not be trusted unsupervised in domains requiring expert judgment (^[38] galileo.ai) (^[39] lastexam.ai). For high-stakes applications (medical diagnosis, scientific research, legal advice), the consistent failure on HLE signals that current LLM outputs "fluent-sounding" statements might contain critical errors. The benchmark provides quantifiable evidence of these weaknesses, helping stakeholders gauge *how much improvement is needed* before deploying AI in sensitive roles.
- **Progress Over Time:** Interestingly, while the gap remains large, it is **shrinking** as expected. The HLE website notes that benchmarks tend to saturate quickly once released (^[4] lastexam.ai). Indeed, shortly after the benchmark's publication, model developers reported incremental gains. For example, OpenAI's updates to GPT and Anthropic's Claude narrowed the gap by a few percentage points by mid-2025. One analysis found the top AI score rising from ~10% to nearly 30% within months (still far from 90%) (^[2] galileo.ai) (^[33] theoutpost.ai). The HLE organizers acknowledge this trend, predicting that models "could exceed 50% accuracy on HLE by the end of 2025" if current advances continue (^[40] lastexam.ai). Even so, they caution that passing 50% would only show expert-level closed-check questions performance, not true general intelligence (^[5] lastexam.ai).
- **No Shortcut Immunity:** HLE confirms that as questions become more complex, knowledge cutoff strategies (like fine-tuning on existing questions or chain-of-thought prompts) give diminishing returns. The benchmark's kept-private questions ensure models can't simply "overfit" to a public test suite. Every new HLE release thus genuinely challenges models to learn new reasoning skills, not memorize answers.

Example Questions

To illustrate HLE's difficulty, consider a couple of sample questions publicly shared by the creators. These examples – chosen by journalists – show the depth expected:

- **Biology (Hummingbird Anatomy):** *"Hummingbirds within Apodiformes uniquely have a bilaterally paired oval bone, a sesamoid embedded in the caudolateral portion of the expanded, cruciate aponeurosis of insertion of m. depressor caudae. How many paired tendons are supported by this sesamoid bone? Answer with a number."* (This question demands specialized anatomical knowledge of avian musculature (www.thestar.com.my). An AI would need to synthesize obscure facts and count anatomical features.)
- **Physics (Rod and Tension Problem):** *"A block is placed on a horizontal frictionless rail and attached to a rigid, massless rod of length R . A mass W is attached at the rod's other end. Initially the mass is directly above the block; it is given an infinitesimal push so that the rod can rotate 360° . When the rod is horizontal, it carries tension T_1 ; when vertical below the block, tension T_2 . (Tension could be negative for compression.) What is the value of $(T_1 - T_2)/W$?"* (This multi-step mechanics problem requires understanding energy, motion constraints, and equilibrium phases (www.thestar.com.my). Even physics PhD students would need to work it out carefully.)

These types of questions typify HLE's level: highly technical, graduate-level material from academic exams. They are well beyond the casual knowledge an AI picks up from text. Models tested on these examples reportedly failed to produce correct answers (indeed, even humans decoding the official answer might struggle without domain expertise). Such examples underscore that passing natural language interviews or writing essays is very different from solving rigorous subject-matter puzzles.

Criticisms and Perspectives

While HLE has been generally well-received as a needed stress test, several experts offer cautionary viewpoints:

- **Not a Measure of General Intelligence:** A recurring criticism is that exam performance measures only one narrow aspect of intelligence. As one commentator noted, a human from 2,000 years ago would fail HLE (not knowing modern science or language), yet that doesn't mean their intelligence was lower than an ape's (^[41] [news.ycombinator.com](#)). HLE explicitly tests *knowledge-loaded reasoning* in academic disciplines, so a model might "pass" HLE without true common-sense, emotional, or creative intelligence. In short, HLE is valuable but far from a complete picture of AGI. As Thestar columnist Kevin Roose emphasized, HLE creators themselves admit it covers *closed-ended* problems and likely *won't capture* open-ended creativity or genuine research ability (^[5] [lastexam.ai](#)). Even if a model eventually aced HLE, it wouldn't prove it could independently conduct scientific research.
- **Overconfidence Risk:** Some worry that calling HLE the "last exam" is rhetorically problematic. Big Think columnist Ethan Siegel argues that naming it as such might "fool humans into believing that an AI possesses capabilities it, in fact, does not" (^[42] [bigthink.com](#)). In his view, HLE's multiple-choice and short-answer format may be too "lax" – an AI might game the format (for example by using process-of-elimination) in ways that do not reflect real understanding, yet still technically get the right answer. Siegel cautions that this could engender overconfidence: observers might think an AI is smarter than it is simply because it manages to hit the correct options by manipulation rather than insight. He stresses that intelligence is multifaceted and an exam of this sort can miss key dimensions (^[42] [bigthink.com](#)) (^[43] [news.ycombinator.com](#)).
- **"Benchmark Quality" Debate:** Within the AI research community, some argue that no single test can define intelligence. As one Hacker News commenter put it, focusing on exam-style questions may be "arrogant" – once a model masters one set of questions, new ones can always be devised (^[44] [news.ycombinator.com](#)). Others point out that HLE is still a static dataset; once known, future models could train on it. The designers mitigate this via hidden questions and leaderboard protocols, but skeptics note that over time models inevitably adapt to any public test. In that sense, HLE may simply become *another* obsolete benchmark eventually, rather than a truly "final" exam.
- **Contest Management Critique:** Some question contributors publicly complained that the contest administration was unclear or unfair (^[18] [news.ycombinator.com](#)). These issues (extended deadlines, prize distribution) relate more to logistics than to the benchmark itself, but they did spark debate. The organizers responded by thanking participants and noting that around **50 winners** (of the intended 550 prizes) came from those who met the original deadline, and that the goal was inclusivity. Importantly, such operational disputes don't affect the technical validity of HLE, but they do illustrate the challenges of crowd-sourcing benchmarks at scale.

Overall, the consensus appears to be: HLE is a **valuable tool**, but not a panacea. Most experts agree that even as HLE sets a higher bar, AI progress will continue. The benchmark will likely be updated and supplemented. Indeed, HLE's own maintainers acknowledge that "it is the last academic exam we need to give to models, but it is far from the last benchmark for AI" (^[45] [lastexam.ai](#)). In other words, passing HLE (whenever that happens) won't be the end of testing LLMs – new tasks (perhaps entirely different in nature) will be needed to probe other dimensions of AI.

Implications and Future Directions

The advent of HLE has several important implications for AI research, deployment, and policy:

- **Benchmarking as Standard Practice:** For developers, HLE establishes a **new standard test**. Teams working on advanced models (OpenAI, Anthropic, Google, etc.) will likely use HLE regularly to gauge progress. Because it is open and research groups share scores on leaderboards, HLE creates a **common yardstick**. Policymakers and safety auditors can also refer to HLE scores when assessing how "smart" a system is. As one CAIS blog pointed out, having a standard measure helps inform regulatory discussion and trust: "This enables more informed discussions about development trajectories, potential risks, and necessary governance measures" (^[39] [lastexam.ai](#)).
- **Guiding Architecture and Strategy:** The HLE results indicate research priorities. For AI practitioners, the low scores on HLE send a signal: continue improving reasoning and knowledge retrieval. Approaches like chain-of-thought prompting, external tools, or specialized solvers might be needed to crack these questions. Indeed, some models already incorporate external calculators and Python interpreters to tackle reasoning tasks. HLE will likely influence how LLMs are engineered, pushing for better role of symbolic reasoning, modularity, or structured memory for factual data.

- **Educational and Societal Perspective:** HLE also serves as a reality check on hype. Public discussion of AI can swing between utopian and dystopian extremes. Quantitative benchmarks like HLE ground the conversation: they show concretely that **despite rapid advances, humans still outperform machines on core intellectual tasks**. For education and workforce planning, this means expert jobs (math research, specialized sciences, graduate-level teaching, etc.) are not imminent casualties of AI – at least until these gaps close significantly.
- **Future Benchmarks and Exams:** The phrasing “Last Exam” is provocative, implying futurist meaning. Pragmatically, if AI continues to improve on HLE, researchers will create still harder versions or alternative evaluations. Indeed, other groups are developing similar ideas. For example, the AI community has proposed benchmarks like “FrontierMath” (for advanced math) or domain-specific tests (biology UMass exam, etc.) ([10] time.com). In summary, we can expect an *evolving suite* of tests, of which HLE is the latest.
- **Limitations and Complementary Tests:** Experts emphasize that HLE should be used in **combination** with other safety evaluations. Passing an HLE-like test doesn't guarantee safety or general competence. The creators themselves note HLE doesn't address issues like bias, hallucinations in open dialogue, or creativity ([5] lastexam.ai). Thus, developers are advised to also run practical simulations, red-teaming, and domain-specific sign-off processes before deployment. The “last exam” is just a piece of the overall verification puzzle.

Finally, the establishment of HLE implies a **socio-technical shift** in AI development. Instead of secret internal benchmarks, major AI labs increasingly engage in open, community benchmarks (MMLU, HLE, etc.) with public leaderboards. This transparency fosters collaboration and comparability across models. It also channels public attention and funding onto the most pressing challenges. The very fact that HLE's existence made headlines (Reuters, TIME, Wired, etc.) shows that AI evaluation is now a central topic of technical and public discourse ([6] english.aawsat.com) ([10] time.com).

Conclusion

Humanity's Last Exam represents a milestone in AI benchmarking. It is by design the most challenging closed-ended test assembled for machine intelligence to date, requiring genuine graduate-level reasoning across hundreds of subjects. The initial outcomes – expert humans near 90% vs. AI near 10–30% – vividly demonstrate that current large models, despite their fluency, still fall far short of human experts on deep understanding tasks.

This benchmark matters not merely as a number, but as a **diagnostic tool**. It reveals specific weaknesses (multi-step logic, visual reasoning, niche knowledge) and highlights where model improvements must focus. Moreover, by publicly quantifying the AI-human gap, HLE can help calibrate expectations: it reminds us that, for now, AI is often only “as smart as a strong undergraduate” on familiar topics, not a genuine polymath in the wild.

At the same time, Humanity's Last Exam is only one chapter in a larger story. Its creators explicitly note that even perfect performance on HLE would not mean machines have become truly autonomous researchers or “AGI” ([5] lastexam.ai). Intelligence is multi-dimensional, and many aspects (creativity, consciousness, empathy) elude any standardized test. Critics rightly argue that we should not equate high exam scores with general intelligence ([42] bigthink.com) ([43] news.ycombinator.com). The title “Last Exam” is hence partly hyperbole – things will always evolve and new challenges will appear.

Looking ahead, HLE sets a precedent. We anticipate that as AI grows stronger, we will see **continually more advanced benchmarks**. These might test open-ended problem solving, collaborative planning, or moral reasoning. For now, HLE frames the frontier of academic reasoning: it is “last” in the sense that it is currently the hardest hurdle, but surely not final.

In summary, *Humanity's Last Exam* illuminates both our progress and our limits. It tells us that despite astonishing gains in AI, machines today can't “think like a human expert” across diverse domains. Only when future models significantly close the remaining gap will we have truly transformed our relationship with machine

intelligence. Until then, Humanity's Last Exam stands as a rigorous check on AI advances – and a cautious reminder that true understanding remains uniquely human (at least for now).

Word Count: *Extremely Comprehensive (Draft)*

References: Key information and quotes in this report are drawn from expert blogs, news outlets, and the official HLE documentation ([6] english.aawsat.com) ([9] safe.ai) ([1] jarxiv.com) ([2] galileo.ai) ([31] theoutpost.ai) ([5] lastexam.ai) ([42] bigthink.com) (www.thestar.com.my) ([10] time.com) ([20] safe.ai) (details in inline citations). Each factual claim is supported by at least one cited source.

External Sources

- [1] <https://jarxiv.com/2025/01/27/humanitys-last-exam/#:~:Bench...>
- [2] <https://galileo.ai/blog/humanitys-last-exam-ai-benchmark#:~:Think...>
- [3] <https://jarxiv.com/2025/01/27/humanitys-last-exam/#:~:bench...>
- [4] <https://lastexam.ai/?ueid=ecd0511fe8770f943568c96e37600be1#:~:While...>
- [5] <https://lastexam.ai/?ueid=ecd0511fe8770f943568c96e37600be1#:~:progr...>
- [6] <https://english.aawsat.com/technology/5061676-ai-experts-ready-%E2%80%98humanity%E2%80%99s-last-exam%E2%80%99-stump-powerful-tech#:~:Dubbe...>
- [7] <https://theoutpost.ai/news-story/new-ai-benchmark-humanity-s-last-exam-stumps-top-models-revealing-limits-of-current-ai-11040/#:~:3%20,...>
- [8] <https://english.aawsat.com/technology/5061676-ai-experts-ready-%E2%80%98humanity%E2%80%99s-last-exam%E2%80%99-stump-powerful-tech#:~:At%20...>
- [9] <https://safe.ai/blog/humanitys-last-exam#:~:%E2%8...>
- [10] <https://time.com/7203729/ai-evaluations-safety/#:~:AI%20...>
- [11] <https://english.aawsat.com/technology/5061676-ai-experts-ready-%E2%80%98humanity%E2%80%99s-last-exam%E2%80%99-stump-powerful-tech#:~:Hendr...>
- [12] <https://news.ycombinator.com/item?id=42806105#:~:exam%...>
- [13] <https://theoutpost.ai/news-story/new-ai-benchmark-humanity-s-last-exam-stumps-top-models-revealing-limits-of-current-ai-11040/#:~:Also%...>
- [14] <https://news.ycombinator.com/item?id=42806105#:~:compe...>
- [15] <https://safe.ai/blog/humanitys-last-exam#:~:,LLMs...>
- [16] https://aiwiki.ai/wiki/Humanity%27s_Last_Exam#:~:Quest...
- [17] https://aiwiki.ai/wiki/Humanity%27s_Last_Exam#:~:resis...
- [18] <https://news.ycombinator.com/item?id=42806105#:~:gets%...>
- [19] <https://lastexam.ai/?ueid=ecd0511fe8770f943568c96e37600be1#:~:,work...>
- [20] <https://safe.ai/blog/humanitys-last-exam#:~:match...>
- [21] <https://jarxiv.com/2025/01/27/humanitys-last-exam/#:~:%E8%A...>

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.

© 2025 IntuitionLabs.ai. All rights reserved.