# Delphi Model Evaluation Report: General Reasoning in Commercial Fast LLMs

**Gensyn AI Team**

We track the performance of the fast variants of commercial LLMs by popular frontier labs. Over three weeks, we will run evaluations across 11 general AI reasoning benchmarks curated to capture broad real-world reasoning ability. The score of each benchmark is normalized to a 0–100 scale, and a model's final score is the average across all of them, giving a single, clean metric of practical capability.

## Entrant models

We evaluate the following models, using the default API accesses provided and versions current as of February 3, 2026.
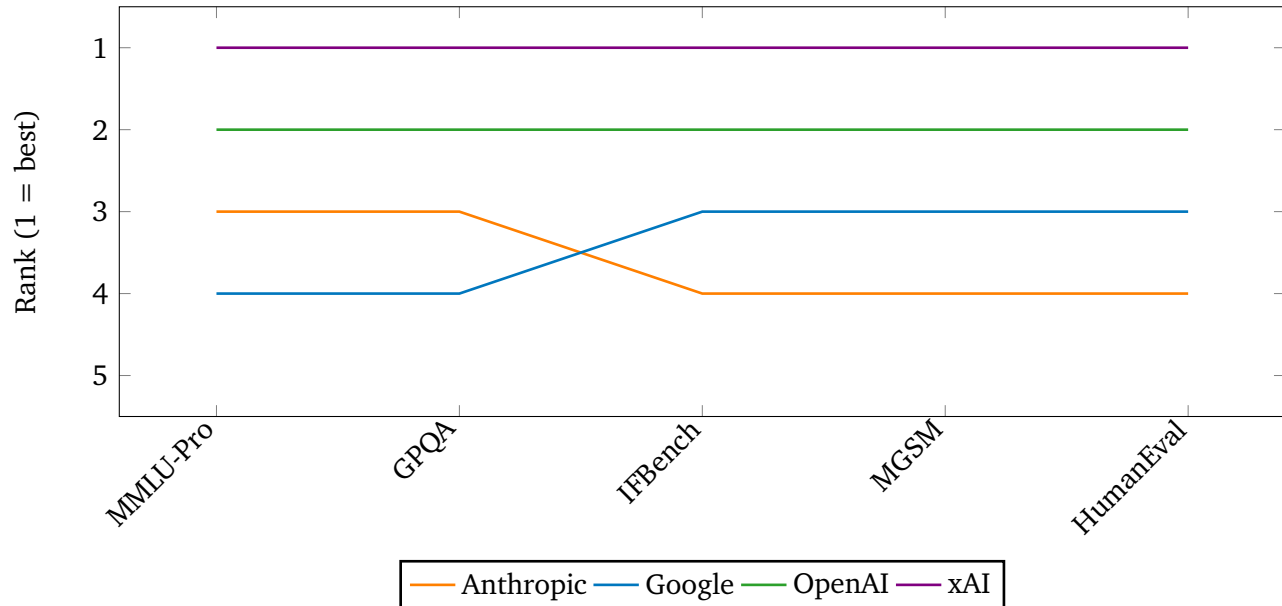
1. Claude Haiku 4.5 from Anthropic.
   - Cost: $1/input MTok, $5/output MTok.
2. Gemini 3 Flash Preview from Google.
   - Cost: $0.50/input MTok, $3/output MTok.
3. GPT 5 Mini from OpenAI.
   - Cost: $0.25/input MTok, $2/output MTok.
4. Grok 4.1 Fast Reasoning from xAI.
   - Cost: $0.20/input MTok, $0.50/output MTok.

**Model settings:** We use the default temperature setting (which is typically 1.0). Unless otherwise noted for a given benchmark, we also used default settings for thinking level (also referred to as thinking budget or effort, depending on the provider).

## Average scores

### Running ranks after each benchmark



| Rank | Family | Model | Avg. score | Δ rank | Notes |
|------|--------|-------|------------|--------|-------|
| 1 | xAI | Grok 4.1 Fast Reasoning | 81.23 | - | - |
| 2 | OpenAI | GPT 5 Mini | 80.51 | - | - |
| 3 | Google | Gemini 3 Flash Preview | 73.30 | - | - |
| 4 | Anthropic | Claude Haiku 4.5 | 70.41 | - | - |

*Avg. score* is the mean of normalized benchmark scores. Δ *rank* is the change in rank compared to the last evaluation.

## Per-Benchmark Breakdown

| Benchmark | Claude Haiku 4.5 | Gemini 3 Flash Preview | GPT 5 Mini | Grok 4.1 Fast Reasoning |
|-----------|------------------|------------------------|------------|-------------------------|
| MMLU-Pro (11) | 78.86 | 52.50 | 81.65 | 85.08 |
| GPQA-Diamond (9) | 62.63 | 66.41 | 77.27 | 84.34 |
| IFBench (8) | 28.03 | 61.09 | 59.05 | 52.18 |
| MGSM (10) | 91.09 | 92.00 | 90.11 | 89.45 |
| HumanEval (6) | 91.46 | 94.51 | 94.51 | 95.12 |
| ? | ▬ | ▬ | ▬ | ▬ |
| ? | ▬ | ▬ | ▬ | ▬ |
| ? | ▬ | ▬ | ▬ | ▬ |
| ? | ▬ | ▬ | ▬ | ▬ |
| ? | ▬ | ▬ | ▬ | ▬ |
| ? | ▬ | ▬ | ▬ | ▬ |

## Benchmark 5: HumanEval

The HumanEval benchmark (6) assesses code generation capabilities by making models complete Python function implementations given a function signature and docstring. It consists of 164 handwritten

programming problems covering fundamental algorithms, data structures, string manipulation, and mathematical operations. We use the EvalPlus framework(7), which extends the original benchmark with comprehensive test suites containing approximately 80 times more test cases per problem to rigorously evaluate both functional correctness against edge cases. HumanEval is a valuable benchmark as its compact, specification-based challenges effectively measure how well a model can convert natural language descriptions into working code. It offers clear insights into its reasoning capabilities and ability to generalize when generating programs.

**Experimental setup.** We evaluate all 164 problems from HumanEval, using the EvalPlus library (2), which extends the original benchmark with comprehensive test suites. For pass@k evaluation, we generate k=10 samples per problem with temperature=0.2, yielding 1,640 total samples. Results are presented as pass@10 metrics. Note that the models are given a maximum output length of 2048 tokens.

We make the following observation:

• All models performed very well, scoring between 91% and 96%. As with Benchmark 4: MGSM, the HumanEval dataset is likely a part of the models' training data.

## Benchmark 4: Multilingual Grade School Math Benchmark (MGSM)

The Multilingual Grade School Math (MGSM) benchmark (10) assesses mathematical reasoning abilities across eleven languages. It is built from 250 elementary-level math word problems drawn from the English GSM8K dataset (5), which were professionally translated into Bengali, Chinese, French, German, Japanese, Russian, Spanish, Swahili, Telugu, and Thai. The tasks involve multi-step arithmetic reasoning that a middle-schooler should be able to solve.

**Experimental setup.** We evaluate all 250 problems in each of the eleven languages, yielding 2,750 total instances, and score performance using exact match on the final extracted numerical answer after normalizing for commas and trailing decimal zeros. We use the evaluation scripts from the `openai/simple-evals` repository (3), adapted for local model inference using vLLM. Prompts are formatted in the target language with language-specific instructions that request reasoning followed by a numeric answer in a specified format. Models are sampled with greedy decoding (temperature=0.0) and a maximum generation length of 2,048 tokens. Reported results correspond to the mean accuracy across the 11 languages.

We make the following observation:

• All models performed very well, scoring between 89% and 92%. The MGSM dataset is likely a part of the models' training data.

## Benchmark 3: IF Bench

The Instruction Following benchmark (IFBench) from AllenAI assesses instruction-following across diverse tasks like counting, formatting, and text manipulation (8). The models are evaluated on how well they follow constraints in their responses, such as "Include keyword beneath in the 23rd sentence, as the 34th word of that sentence," "Use at least 3 different coordinating conjunctions in the response," or "The second word in your response and the second to last word in your response should be the word vibrant."

**Experimental setup.** We use the 294-question single-turn dataset with 5 repetitions per question, evaluated via pass@1 scoring (1). Responses are scored using the repository's official evaluation code in loose mode, which tolerates formatting variations and extraneous text by testing multiple output formats (removing leading/trailing lines, asterisks, etc.). Results reflect average prompt-level accuracy

across all questions and runs. In our experiments, models are sampled with temperature 0.7, top-p 0.95, and a 4096 token limit.

We make the following observations:

• Gemini 3 Flash Preview bounced back with the highest score. This is the first benchmark where Grok 4.1 Fast Reasoning is not at the top.

• Claude Haiku 4.5 performed significantly worse than the others.

## Benchmark 2: GPQA Diamond

The GPQA (Graduate-Level Google-Proof Q&A) Diamond dataset consists of highly challenging multiple-choice problems in biology, physics, and chemistry, requiring scientific knowledge and reasoning abilities to answer. Each question is authored by subject-matter experts and is intentionally formulated to be difficult for non-experts to solve, even with access to online resources. The most difficult "Diamond" split includes 198 questions for which both expert annotators selected the correct answer, while most non-experts did not.

**Experimental setup.** We evaluate the models following the guidelines from the benchmark authors (9) and using the OpenAI simple-eval library (3). Every question is presented 10 times, each time with the answer choices in a different randomized order. The models are allowed a maximum output length of 2028 tokens.

We make the following observations:

• Grok 4.1 Fast Reasoning achieves the highest score again.

• Claude Haiku 4.5 has a good score in the first benchmark, but scored the lowest in this one.

## Benchmark 1: MMLU-Pro

MMLU-Pro (11) is an enhanced benchmark for evaluating language models, building on the original MMLU dataset with significantly harder, reasoning-intensive questions and 10 answer options instead of 4. It contains over 12,000 curated questions from academic sources spanning 14 fields, including Biology, Computer Science, Mathematics, Physics, and Law. Experiments demonstrate MMLU-Pro substantially increases difficult with model accuracies dropping 16–33% compared to standard MMLU. Notably, Chain-of-Thought prompting yields greater improvements on MMLU-Pro than on the original benchmark, indicating the dataset demands deeper, more structured reasoning.

We evaluate the models following the guidelines from the benchmark authors (on `TIGER-Lab/MMLU-Pro repository` (4)). The models are allowed a maximum output length of 2048 tokens.

We observed the following from the performance of the Delphi entrant models.

• Gemini 3 Flash Preview performed significantly worse than the other three, scoring around 52%, whereas the remaining models scored between 78% and 85%.

## References

[1] allenai/IFBench: A new, challenging benchmark for precise instruction following. https://github.com/allenai/IFBench.

[2] evalplus/evalplus: Rigourous evaluation of LLM-synthesized code. https://github.com/evalplus/evalplus/tree/master.

[3] OpenAI simple-eval: a lightweight library for evaluating language models. `https://github.com/openai/simple-evals/tree/main`.

[4] TIGER-AI-Lab/MMLU-Pro: A More Robust and Challenging Multi-Task Language Understanding Benchmark. `https://github.com/TIGER-AI-Lab/MMLU-Pro`, 2024.

[5] Karl Cobbe, Vineet Kosaraju, Mohammad Bavarian, Mark Chen, Heewoo Jun, Lukasz Kaiser, Matthias Plappert, Jerry Tworek, Jacob Hilton, Reiichiro Nakano, Christopher Hesse, and John Schulman. Training verifiers to solve math word problems. *arXiv preprint arXiv:2110.14168*, 2021.

[6] Xueying Du, Mingwei Liu, Kaixin Wang, Hanlin Wang, Junwei Liu, Yixuan Chen, Jiayi Feng, Chaofeng Sha, Xin Peng, and Yiling Lou. Evaluating large language models in class-level code generation. In *Proceedings of the IEEE/ACM 46th International Conference on Software Engineering*, ICSE '24, New York, NY, USA, 2024. Association for Computing Machinery.

[7] Jiawei Liu, Chunqiu Steven Xia, Yuyao Wang, and Lingming Zhang. Is your code generated by chatgpt really correct? rigorous evaluation of large language models for code generation. NIPS '23, Red Hook, NY, USA, 2023. Curran Associates Inc.

[8] Valentina Pyatkin, Saumya Malik, Victoria Graf, Hamish Ivison, Shengyi Huang, Pradeep Dasigi, Nathan Lambert, and Hannaneh Hajishirzi. Generalizing verifiable instruction following, 2025.

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

[10] Freda Shi, Mirac Suzgun, Markus Freitag, Xuezhi Wang, Suraj Srivats, Soroush Vosoughi, Hyung Won Chung, Yi Tay, Sebastian Ruder, Denny Zhou, Dipanjan Das, and Jason Wei. Language models are multilingual chain-of-thought reasoners, 2022.

[11] Y. Wang, X. Ma, G. Zhang, Y. Ni, A. Chandra, S. Guo, W. Ren, A. Arulraj, X. He, Z. Jiang, T. Li, M. Ku, K. Wang, A. Zhuang, R. Fan, X. Yue, and W. Chen. Mmlu-pro: A more robust and challenging multi-task language understanding benchmark. *arXiv preprint arXiv:2406.01574*, 2024.