



Delphi Model Evaluation Report: General Reasoning in Middleweight LLMs

Gensyn AI Team

We track the performance of “middleweight” variants of open source LLMs published by popular frontier labs. Over three weeks we will run evaluations across 11 general AI reasoning benchmarks that have been curated to capture broad real-world reasoning ability.

Entrant models

We evaluate the following models, as provided on HuggingFace, using the default datatype provided and versions current as of December 7, 2025.

- Qwen3-30B-A3B-Instruct-2507 from Qwen.
- gpt-oss-20b from OpenAI.
- gemma-3-27b-it from Google.
- GLM-4-32B-0414 from Zai-Org.
- Falcon-H1-34B-Instruct from TII-UAE.

Evaluation setup: We perform inference with temperature 0.0 on machines with NVIDIA H100 GPUs using vLLM, unless specified.

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 observations:

- OpenAI’s gpt-oss-20b scored the highest, with a pass@10 score of 90%. Note that even frontier models, such as GPT-5 and Claude 3.5 Sonnet, score under 95%.
- The remaining models all scored between 75% and 85%.



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 ($\text{temperature}=0.0$) and a maximum generation length of 2,048 tokens. Reported results correspond to the mean accuracy across the 11 languages.

We make the following observations:

- Google’s gemma-3-27b-it significantly outperformed the others, scoring around 88% while OpenAI’s gpt-oss-20b got second place with around 77%.
- For the first time, Qwen3-30B-A3B-Instruct-2507 performed the worst in a benchmark, scoring around 70%.

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 $\text{pass}@1$ 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:

- The model that scored the lowest in the prior benchmarks, gpt-oss-20b, significantly outperformed the others in this evaluation, scoring about 55%.
- The other four models were close in performance, scoring between 32% and 40%.

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:

- The top four models all achieved accuracies within 4.5 percentage points of one another.
- The lowest-performing model, gpt-oss-20b, again scored noticeably lower than the other models.

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

- Qwen3-30B-A3B-Instruct-2507 performed the best, scoring nearly 75%.
- The first place model scored less than 10% ahead of the fourth place model.
- The last place model, gpt-oss-20b, scored noticeably lower than the other models.

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.

Overall scores

The running average after 5 evaluations.

Rank	Family	Model	Avg. score	Δ rank	Notes
1	Google	gemma-3-27b-it	63.04	-	-
2	Qwen	Qwen3-30B-A3B-Instruct-2507	62.46	↑ 1	-
3	OpenAI	gpt-oss-20b	61.54	↑ 2	First time not at the bottom.
4	TIIuae	Falcon-H1-34B-Instruct	60.96	↓ 2	-
5	Zai-Org	GLM-4-32B-0414	59.81	↓ 1	-

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

Per-Benchmark Breakdown

Thus far we have revealed 5 out of 11 evaluations



Benchmark	gpt-oss-20b	Qwen3-30B-A3B-Instruct-2507	Falcon-H1-34B-Instruct	GLM-4-32B-0414	gemma-3-27b-it
MMLU-Pro (11)	47.88	74.76	72.27	67.52	65.19
GPQA-Diamond (9)	37.22	46.92	49.60	46.97	45.66
IFBench (8)	54.56	34.83	32.59	31.16	40.48
MGSM (10)	77.16	70.40	73.56	73.49	88.25
HumanEval (10)	90.90	85.40	76.80	79.90	75.60
[?]	■■■	■■■	■■■	■■■	■■■
[?]	■■■	■■■	■■■	■■■	■■■
[?]	■■■	■■■	■■■	■■■	■■■
[?]	■■■	■■■	■■■	■■■	■■■
[?]	■■■	■■■	■■■	■■■	■■■
[?]	■■■	■■■	■■■	■■■	■■■

Rows with “[?]” as the benchmark will be revealed in the near future as we run those evaluations. Stay tuned for more.⁴



References

- [1] allenai/IFBench: A new, challenging benchmark for precise instruction following. <https://github.com/allenai/IFBench>.
- [2] evalplus/evalplus: Rigorous 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.