

---

# AI-Assisted Generation of Difficult Math Questions

---

Anonymous Author(s)

Affiliation

Address

email

## Abstract

1 Current LLM training positions mathematical reasoning as a core capability. With  
2 publicly available sources fully tapped, there is an unmet demand for diverse and  
3 challenging mathematics questions. Relying solely on human experts is both time-  
4 consuming and costly, while LLM-generated questions often lack the requisite  
5 diversity and difficulty. We present a design framework that combines the strengths  
6 of LLMs with a human-in-the-loop approach to generate a diverse array of chal-  
7 lenging math questions. Initially, leveraging LLM metacognition skills [Didolkar  
8 et al., 2024], a strong LLM is used to extract core “skills” from existing math  
9 datasets. These skills serve as the basis for generating novel and difficult questions  
10 by prompting the LLM with random pairs of core skills that must be utilized in the  
11 question. The use of two very different skills within each question makes finding  
12 such questions an “out of distribution” task for both LLMs and humans. Our  
13 pipeline employs LLMs to iteratively generate and refine questions and solutions  
14 through multi-turn prompting. Human annotators then verify and further refine the  
15 questions, with their efficiency enhanced via further LLM interactions. Applying  
16 this pipeline on skills extracted from MATH dataset [Hendrycks et al., 2021] re-  
17 sulted in **MATH<sup>2</sup>** - a dataset of higher quality math questions, as evidenced by  
18 lower performance of all models on MATH<sup>2</sup> than on MATH. Although focused  
19 on mathematics, our methodology seems applicable to other domains requiring  
20 structured reasoning, and potentially as a component of *scalable oversight*. Also of  
21 interest is a striking relationship observed between models’ performance on the  
22 new dataset: the success rate on MATH<sup>2</sup> is the square on MATH. This suggests  
23 that successfully solving the question in MATH<sup>2</sup> requires a nontrivial combination  
24 of two distinct math skills.

## 25 1 Introduction

26 Significant improvement in the capabilities of LLMs [Chowdhery et al., 2023, Anil et al., 2023, Team,  
27 2023, Team et al., 2023, Abdin et al., 2024, Achiam et al., 2023, Touvron et al., 2023] to understand  
28 and generate complex mathematical content has been achieved by leveraging all the public data and a  
29 fair bit of private data. Sources of high-quality, varied, and difficult mathematical questions are drying  
30 up. Even finding new questions for evaluation is getting difficult since newly-released human exams  
31 are somewhat similar to past exams, which are potentially present in the LLMs’ training datasets.  
32 Hence, there is a pressing need for innovative methods to create new, diverse, and challenging  
33 questions.

34 Expert mathematicians and educators possess the deep understanding required to create questions that  
35 not only test a wide range of skills but also push the boundaries of what the learners, and by extension,  
36 the models, can handle. However, relying solely on human experts is not scalable. Generating  
37 synthetic questions using LLMs is feasible at scale [Trinh et al., 2024, Li et al., 2024, Gunasekar  
38 et al., 2023, Patel et al., 2024, Toshniwal et al., 2024, Gupta et al., 2023, Lu et al., 2024, Honovich

39 et al., 2022], but often fall short in terms of the necessary difficulty. Huang et al. [2024] employs a  
40 similar approach as ours where they extract *topics* and corresponding *keypoints* from a set of seed  
41 problems using GPT-4, and then combine the *topic* to generate new questions, again using GPT-4).  
42 However, the generated data is meant to be used for the finetuning of models as compared to serving  
43 as an evaluation set in our case. As a result, the questions generated in Huang et al. [2024] are  
44 not sufficiently difficult. Similarly, limited work exists on ensuring the necessary diversity in the  
45 generated synthetic data. Chan et al. [2024] proposes prompting frontier models to generate questions  
46 where each question is generated in the context of a *persona* as a way of ensuring diversity. They use  
47 1M different personas to generate questions, which are then used for finetuning models, leading to  
48 significant improvements. This dichotomy between the quality of human-generated questions and the  
49 scalability of LLM-generated questions presents a significant challenge [Yu et al., 2024].

## 50 1.1 Evaluation Saturation Phenomenon

51 LLM evaluations are becoming saturated due to improvements from better training and larger  
52 datasets, but also from evaluation-specific optimizations like supervised fine-tuning (SFT) on synthetic  
53 question-answer pairs. These pairs, generated by leading proprietary models or filtered from the  
54 model’s responses, can significantly boost performance. For instance, just 1 million synthetic  
55 examples raised Llama2 7B’s MATH dataset performance to GPT-4 levels [Li et al., 2024].

56 The distinction between general and evaluation-specific improvements is key, as the latter can lead to  
57 overfitting rather than real skill acquisition. This issue was evident when models showed performance  
58 drops on a new GSM8K dataset version and on newer Chinese GaoKao exams, suggesting shallow  
59 understanding of mathematics [Zhang et al., 2024].

## 60 1.2 Proposed Framework: AI-assisted Generation of Difficult Math Questions

61 Recent research [Arora and Goyal, 2023, Didolkar et al., 2024] demonstrated that top LLMs possess  
62 a robust understanding of mathematical skills, including the capability to identify the skills required  
63 to solve given questions [Reid et al., 2024, Achiam et al., 2023]. This naturally raises the question:  
64 *can LLMs operate in the reverse direction, i.e., generate math problems when given a list of skills that*  
65 *have to be tested?* Our initial attempts yielded mixed results. While leading models could produce  
66 creative math questions when provided with a list of skills, the majority of these questions exhibited  
67 one or more of the following shortcomings: too similar to existing questions in datasets; have errors  
68 or nonsensical elements; are too tedious or mechanical to be engaging for human annotators. (See  
69 Section B.) Moreover, they often conflate “difficulty” with tedious calculations, which actually would  
70 play to the strength of machines to leverage external tools such as calculators or Python interpreters.

71 Nevertheless, there were promising instances where LLMs generated interesting and correct questions  
72 that they were unable to solve, due to incomplete or incorrect reasoning. This observation led us  
73 to the concept of *AI-assisted creation of evaluation datasets*. Our process may also be of interest  
74 for human pedagogy since it begins with the extraction of core “skills” from existing math datasets,  
75 which serve as the foundational elements of mathematical questions. The current paper focuses on  
76 the MATH dataset [Hendrycks et al., 2021], a mainstay of LLM evaluation in recent years.

77 In our AI-assisted process, human experts played a crucial role. Using the (question, answer)  
78 pairs generated by LLMs and leveraging API access to leading models, experts identified promising  
79 questions—often those incorrectly answered by the LLMs but containing many correct ideas. Experts  
80 then refined these questions to enhance their engagement value and provided gold-standard answers.  
81 The AI-assisted process not only boosted human productivity but also resulted in high-quality, novel  
82 questions distinct from those in existing datasets.

## 83 2 Pipeline for AI-Assisted Question Generation

84 We present a structured approach to generating challenging mathematics questions by combining  
85 the capabilities of large language models (LLMs) and human expertise. Given below is a high-level  
86 overview of the process before delving into the details of each step.

87 We begin our pipeline with **skill extraction** - identifying and cataloging distinct mathematical skills  
88 from a dataset, as described in Didolkar et al. [2024]. This step creates a repository of skills linked to  
89 specific questions. The motivation behind this is to systematically generate and analyze questions  
90 that require specific skills, ensuring a comprehensive evaluation framework.

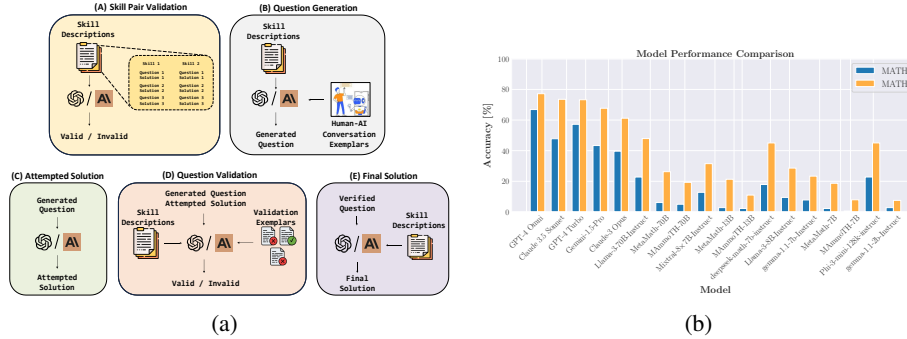


Figure 1: (a) **AI-assisted question generation:** The five-step pipeline includes: (i) Skill Pair Validation, ensuring distinct skills; (ii) Question Generation, producing a problem that combines both skills; (iii) Attempted Solution, where the model solves using a *defeatist* approach; (iv) Question Validation, assessing correctness, rigor, and quality; and (v) Final Solution, applying advanced techniques to enhance accuracy. (b) **Comparison of Zero-Shot Performance:** This figure shows zero-shot Chain of Thought (CoT) performance on MATH and MATH<sup>2</sup>. Proprietary models show the smallest performance drop on MATH<sup>2</sup>, while smaller models experience larger drops. Detailed results are in Table 1.

91 We employ a five-step approach to generate difficult math questions using advanced models. For each  
 92 round of generation, we randomly sample a pair of skills and three sample question-solution pairs  
 93 corresponding to each skill from the skill repository. These reference examples are sourced from the  
 94 MATH dataset.

95 **Step 1: Skill Pair Validation.** We begin by asking the LLM (GPT-4 or Claude) to validate a randomly  
 96 sampled skill pair by assessing the qualitative similarity of the two skills. Reference examples are  
 97 provided in-context to enrich the model’s understanding of the skills. If the model deems the skills  
 98 too similar, they are flagged and excluded from question generation, as similar skills might lead to  
 99 simpler questions.

100 **Step 2: Question Generation.** Next, we prompt the LLM to generate a question and a brief solution  
 101 requiring the application of both skills in the sampled pair. We specify two conditions to ensure high-  
 102 quality questions: the question should either require an exact answer or specify that an approximate  
 103 answer is acceptable, and it should ask for only a single final result. In-context, we provide two  
 104 multi-turn conversations between a human and an AI assistant. These conversations demonstrate  
 105 the human providing feedback on the AI-generated questions, which the AI then refines. This helps  
 106 the model anticipate and avoid practical issues, such as insufficient involvement of skills or logical  
 107 inconsistencies. Appendix B.2 provides examples of the responses of different models in the question  
 108 generation step.

109 **Step 3: Solution Attempt.** The model then attempts a solution to the generated question, adopting an  
 110 adversarial approach to identify flaws such as insufficient information, ambiguity, self-contradiction,  
 111 or excessive computation. If any issues are found, the model stops solving and clearly states the  
 112 problems. Otherwise, it completes the solution. During this step, the model does not receive the skill  
 113 names or reference examples to ensure unbiased problem-solving.

114 **Step 4: Question Validation.** We give LLM the generated question and its solution for validation  
 115 against a fixed rubric consisting of seven criteria. We detail the validation criteria in Appendix A.1.  
 116 The model uses reference examples and validation exemplars - model generated examples of validating  
 117 questions, to facilitate this step. We employ majority voting (maj @ 4) to enhance robustness.

118 **Step 5: Final Solution and Re-validation.** For questions classified as valid, we ask the LLM  
 119 to re-solve the question to obtain a final solution. Reference examples are provided in-context to  
 120 improve the model’s understanding. We use majority voting (maj @ 4) to ensure consistency. If all  
 121 the answers obtained in this step are unique, indicating potential ambiguity, the question is discarded.

122 The questions obtained from the above pipeline are further screened by humans. This structured  
 123 approach not only generates challenging and novel math questions but also ensures their quality

124 through rigorous validation, effectively combining the strengths of AI and human oversight. For  
 125 detailed examples of prompts used at each step, refer to Appendix B.3.

### 126 3 Experiments and Findings

127 Through our experiments, we demonstrate the difficulty and quality of the MATH<sup>2</sup> while also  
 128 analyzing the behavior of different models on this task of *compositional generalization*. Firstly,  
 129 we evaluate a wide range of models spanning a large range of parameter counts on MATH<sup>2</sup> and  
 130 compare against their performance on MATH [Hendrycks et al., 2021] which is the base dataset  
 131 used for extracting skills, showing that the MATH<sup>2</sup> is necessarily harder than MATH. Next, we  
 132 further demonstrate the difficulty and quality of questions in MATH<sup>2</sup> by showing that they are better  
 133 in-context exemplars as compared to standardly used exemplars. We describe the experimental setup  
 134 below.

#### 135 3.1 Experimental Setup

136 The MATH dataset encompasses seven high-level topics, allowing us to identify and extract finer-  
 137 grained skills within each topic and label each question accordingly. At the end of the skill-  
 138 extraction process, we identify a set of 114 skills. We then remove a few simple skills, such  
 139 as `basic_arithmetic` and `arithmetic_operations`, before using the remaining set to generate  
 140 questions using the proposed approach. We generate and verify 180 difficult questions to create the  
 141 MATH<sup>2</sup> dataset. Out of the 180 questions, 136 questions were generated using GPT-4 Turbo and 44  
 142 were generated using Claude-3 Opus. Figure 3 shows the distribution of skills in MATH<sup>2</sup>.

143 Table 2 presents details of the changes made to the questions during the human  
 144 verification process. In total, 56% of the question-answer pairs in MATH<sup>2</sup> ap-  
 145 pear exactly as phrased by their LLM creator.

151 We evaluate the generated set of questions on a variety of language models, both  
 152 small and large. Specifically, we assess the MetaMath [Yu et al., 2023],  
 153 MAMmoTH [Yue et al., 2023], Gemma [Team et al., 2024], and Llama-  
 154 3 series, Phi-3, deepseek-math as well as one Mixtral-of-Experts model Mixtral-  
 155 8×7B-Instruct. Additionally, we include evaluations of larger proprietary mod-  
 156 els such as GPT-4o, GPT-4 Turbo<sup>1</sup> [OpenAI, 2023], Gemini-1.5-Pro, Claude 3.5  
 157 Sonnet<sup>2</sup> and Claude 3 Opus<sup>3</sup>. We compare the performances of these models on our generated  
 158 questions against their performance on the MATH dataset [Hendrycks et al., 2021].

159 We also observe that the performance of models on MATH<sup>2</sup> follows a quadratic relationship with the  
 160 performance of models on MATH. We refer the reader to Appendix A.4 for more discussion on this  
 161 observation and for further implementation and compute details.

162 We also observe that the performance of models on MATH<sup>2</sup> follows a quadratic relationship with the  
 163 performance of models on MATH. We refer the reader to Appendix A.4 for more discussion on this  
 164 observation and for further implementation and compute details.

Table 1: **Comparison of Zero-Shot CoT Performance (Accuracy) on the Generated Dataset vs. MATH Test Set:** GPT-4 Omni demonstrates the least drop in percentage terms (13.42%) whereas MAMmoTH-7B shows the highest relative drop (92.91%).

Model	MATH <sup>2</sup> (Y)	MATH (X)	% Drop
GPT-4 Omni	66.85%	77.21%	<b>13.42%</b>
Claude 3.5 Sonnet	47.78%	73.54%	35.45%
GPT-4 Turbo	57.22%	73.27%	21.90%
Gemini-1.5-Pro	43.34%	67.70%	35.98%
Claude 3 Opus	39.66%	61.20%	35.20%
Llama-3-70B-Instruct	22.77%	47.89%	52.45%
MetaMath-70B	6.11%	26.27%	76.74%
MAMmoTH-70B	5.00%	19.31%	74.11%
Mixtral-8×7B-Instruct	12.78%	31.52%	59.45%
MetaMath-13B	2.79%	21.32%	86.91%
MAMmoTH-13B	2.23%	10.99%	79.71%
Deepseek-math-7b-instruct	17.88%	45.05%	60.31%
Llama-3-8B-Instruct	9.45%	28.62%	66.98%
Gemma-1.1-7B-Instruct	7.78%	23.36%	66.69%
MetaMath-7B	2.23%	18.69%	88.07%
MAMmoTH-7B	0.56%	7.90%	<b>92.91%</b>
Phi-3-mini-128k-instruct	22.78%	45.14%	49.53%
Gemma-1.1-2B-Instruct	2.78%	7.52%	63.03%

<sup>1</sup>Points to `gpt-4-turbo-2024-04-09` at the time of writing

<sup>2</sup>Points to `claude-3-5-sonnet-20240620` at the time of writing

<sup>3</sup>Points to `claude-3-opus-20240229` at the time of writing

## 174 4 Conclusions

175 We introduced a framework that leverages the complementary strengths of humans and AI to generate  
176 new, challenging mathematics questions. Building on recent insights into LLM metaknowledge,  
177 we use LLMs to extract and name key skills necessary for solving math problems. Using these  
178 insights, we developed a pipeline that employs named skills from the well-known MATH dataset,  
179 and leverages multi-turn interactions with advanced LLMs to generate questions that combine pairs  
180 of skills. These questions were subsequently reviewed and refined by human raters. The proposed  
181 pipeline produced questions with greater novelty and difficulty compared to those in the original  
182 MATH dataset. This framework also resulted in a new math evaluation **MATH**<sup>2</sup>, that assesses the  
183 same skills as the MATH dataset but is significantly more challenging for leading models because  
184 each question involves two skills from different parts of MATH.

185 We plan to release detailed information about our pipeline to encourage further research and develop-  
186 ment in the field of open-source math models.

187 **Limitations and Future Work.** Our pipeline incurs moderately high costs due to extensive API-based  
188 use of frontier models as well as significant human verification. To improve efficiency, future work  
189 should focus on using open weights models and optimizing prompting strategies to produce higher-  
190 quality questions initially, thereby reducing the need for extensive filtering. Additionally, reducing  
191 human verification through the development of automated validation tools is crucial. This could  
192 include leveraging code generation and autoformalization capabilities of LLMs to generate responses  
193 which can be compiled using compilers or interpreters. Enhancing our pipeline by integrating a  
194 training-based feedback loop, where the model is trained on the questions that pass human verification,  
195 could further streamline the process by progressively improving question quality. These measures will  
196 reduce dependency on expensive proprietary models, lower overall operational costs, and maintain or  
197 even enhance the quality of the generated math evaluation benchmarks.

198 Looking ahead, an even more exciting prospect is the potential application of the proposed framework  
199 to efficiently produce high-quality data in domains beyond mathematics.

## 200 References

- 201 M. Abdin, S. A. Jacobs, A. A. Awan, J. Aneja, A. Awadallah, H. Awadalla, N. Bach, A. Bahree,  
202 A. Bakhtiari, H. Behl, et al. Phi-3 technical report: A highly capable language model locally on  
203 your phone. *arXiv preprint arXiv:2404.14219*, 2024.
- 204 J. Achiam, S. Adler, S. Agarwal, L. Ahmad, I. Akkaya, F. L. Aleman, D. Almeida, J. Altschmidt,  
205 S. Altman, S. Anadkat, et al. Gpt-4 technical report. *arXiv preprint arXiv:2303.08774*, 2023.
- 206 R. Anil, A. M. Dai, O. Firat, M. Johnson, D. Lepikhin, A. Passos, S. Shakeri, E. Taropa, P. Bailey,  
207 Z. Chen, et al. Palm 2 technical report. *arXiv preprint arXiv:2305.10403*, 2023.
- 208 S. Arora and A. Goyal. A theory for emergence of complex skills in language models. *arXiv preprint*  
209 *arXiv:2307.15936*, 2023.
- 210 X. Chan, X. Wang, D. Yu, H. Mi, and D. Yu. Scaling synthetic data creation with 1,000,000,000  
211 personas. *arXiv preprint arXiv:2406.20094*, 2024.
- 212 A. Chowdhery, S. Narang, J. Devlin, M. Bosma, G. Mishra, A. Roberts, P. Barham, H. W. Chung,  
213 C. Sutton, S. Gehrmann, P. Schuh, K. Shi, S. Tsvyashchenko, J. Maynez, A. Rao, P. Barnes,  
214 Y. Tay, N. Shazeer, V. Prabhakaran, E. Reif, N. Du, B. Hutchinson, R. Pope, J. Bradbury, J. Austin,  
215 M. Isard, G. Gur-Ari, P. Yin, T. Duke, A. Levskaya, S. Ghemawat, S. Dev, H. Michalewski,  
216 X. Garcia, V. Misra, K. Robinson, L. Fedus, D. Zhou, D. Ippolito, D. Luan, H. Lim, B. Zoph,  
217 A. Spiridonov, R. Sepassi, D. Dohan, S. Agrawal, M. Omernick, A. M. Dai, T. S. Pillai, M. Pellat,  
218 A. Lewkowycz, E. Moreira, R. Child, O. Polozov, K. Lee, Z. Zhou, X. Wang, B. Saeta, M. Diaz,  
219 O. Firat, M. Catasta, J. Wei, K. Meier-Hellstern, D. Eck, J. Dean, S. Petrov, and N. Fiedel. Palm:  
220 Scaling language modeling with pathways. *Journal of Machine Learning Research*, 24(240):1–113,  
221 2023. URL <http://jmlr.org/papers/v24/22-1144.html>.
- 222 A. Didolkar, A. Goyal, N. R. Ke, S. Guo, M. Valko, T. Lillicrap, D. Rezende, Y. Bengio, M. C. Mozer,  
223 and S. Arora. Metacognitive Capabilities of LLMs: An Exploration in Mathematical Problem  
224 solving. 2024. URL <https://api.semanticscholar.org/CorpusID:269921384>.

- 225 S. Gunasekar, Y. Zhang, J. Aneja, C. C. T. Mendes, A. Del Giorno, S. Gopi, M. Javaheripi, P. Kauff-  
 226 mann, G. de Rosa, O. Saarikivi, et al. Textbooks are all you need. *arXiv preprint arXiv:2306.11644*,  
 227 2023.
- 228 H. Gupta, K. Scaria, U. Anantheswaran, S. Verma, M. Parmar, S. A. Sawant, S. Mishra, and C. Baral.  
 229 Targen: Targeted data generation with large language models. *arXiv preprint arXiv:2310.17876*,  
 230 2023.
- 231 D. Hendrycks, C. Burns, S. Kadavath, A. Arora, S. Basart, E. Tang, D. Song, and J. Steinhardt.  
 232 Measuring mathematical problem solving with the math dataset, 2021.
- 233 O. Honovich, T. Scialom, O. Levy, and T. Schick. Unnatural instructions: Tuning language models  
 234 with (almost) no human labor. *arXiv preprint arXiv:2212.09689*, 2022.
- 235 Y. Huang, X. Liu, Y. Gong, Z. Gou, Y. Shen, N. Duan, and W. Chen. Key-point-driven data synthesis  
 236 with its enhancement on mathematical reasoning. *arXiv preprint arXiv:2403.02333*, 2024.
- 237 A. Q. Jiang, A. Sablayrolles, A. Roux, A. Mensch, B. Savary, C. Bamford, D. S. Chaplot, D. de las  
 238 Casas, E. B. Hanna, F. Bressand, G. Lengyel, G. Bour, G. Lample, L. R. Lavaud, L. Saulnier, M.-A.  
 239 Lachaux, P. Stock, S. Subramanian, S. Yang, S. Antoniak, T. L. Scao, T. Gervet, T. Lavril, T. Wang,  
 240 T. Lacroix, and W. E. Sayed. Mixtral of experts, 2024.
- 241 W. Kwon, Z. Li, S. Zhuang, Y. Sheng, L. Zheng, C. H. Yu, J. Gonzalez, H. Zhang, and I. Stoica. Effi-  
 242 cient memory management for large language model serving with pagedattention. In *Proceedings*  
 243 *of the 29th Symposium on Operating Systems Principles*, pages 611–626, 2023.
- 244 C. Li, W. Wang, J. Hu, Y. Wei, N. Zheng, H. Hu, Z. Zhang, and H. Peng. Common 7b language  
 245 models already possess strong math capabilities. *arXiv preprint arXiv:2403.04706*, 2024.
- 246 Z. Lu, A. Zhou, H. Ren, K. Wang, W. Shi, J. Pan, M. Zhan, and H. Li. Mathgenie: Generating  
 247 synthetic data with question back-translation for enhancing mathematical reasoning of llms. *arXiv*  
 248 *preprint arXiv:2402.16352*, 2024.
- 249 OpenAI. Gpt-4 technical report, 2023.
- 250 A. Patel, C. Raffel, and C. Callison-Burch. Datadreamer: A tool for synthetic data generation and  
 251 reproducible llm workflows. *arXiv preprint arXiv:2402.10379*, 2024.
- 252 M. Reid, N. Savinov, D. Teplyashin, D. Lepikhin, T. Lillicrap, J.-b. Alayrac, R. Soricut, A. Lazaridou,  
 253 O. Firat, J. Schrittwieser, et al. Gemini 1.5: Unlocking multimodal understanding across millions  
 254 of tokens of context. *arXiv preprint arXiv:2403.05530*, 2024.
- 255 G. Team. Gemini: A family of highly capable multimodal models, 2023.
- 256 G. Team, R. Anil, S. Borgeaud, Y. Wu, J.-B. Alayrac, J. Yu, R. Soricut, J. Schalkwyk, A. M.  
 257 Dai, A. Hauth, et al. Gemini: a family of highly capable multimodal models. *arXiv preprint*  
 258 *arXiv:2312.11805*, 2023.
- 259 G. Team, T. Mesnard, C. Hardin, R. Dadashi, S. Bhupatiraju, S. Pathak, L. Sifre, M. Rivière, M. S.  
 260 Kale, J. Love, et al. Gemma: Open models based on gemini research and technology. *arXiv*  
 261 *preprint arXiv:2403.08295*, 2024.
- 262 S. Toshniwal, I. Moshkov, S. Narenthiran, D. Gitman, F. Jia, and I. Gitman. Openmathinstruct-1: A  
 263 1.8 million math instruction tuning dataset. *arXiv preprint arXiv:2402.10176*, 2024.
- 264 H. Touvron, L. Martin, K. Stone, P. Albert, A. Almahairi, Y. Babaei, N. Bashlykov, S. Batra,  
 265 P. Bhargava, S. Bhosale, et al. Llama 2: Open foundation and fine-tuned chat models. *arXiv*  
 266 *preprint arXiv:2307.09288*, 2023.
- 267 T. H. Trinh, Y. Wu, Q. V. Le, H. He, and T. Luong. Solving olympiad geometry without human  
 268 demonstrations. *Nature*, 625(7995):476–482, 2024.
- 269 J. Wei, X. Wang, Q. Liu, B. Yang, X. Dong, H. Huang, and W. Wang. Chain-of-thought prompting  
 270 elicits reasoning in large language models. *arXiv*, abs/2201.11903, 2022. URL <https://doi.org/10.48550/arXiv.2201.11903>.

- 272 L. Yu, W. Jiang, H. Shi, J. Yu, Z. Liu, Y. Zhang, J. T. Kwok, Z. Li, A. Weller, and W. Liu. Meta-  
 273 math: Bootstrap your own mathematical questions for large language models. *arXiv preprint*  
 274 *arXiv:2309.12284*, 2023.
- 275 Y. Yu, Y. Zhuang, J. Zhang, Y. Meng, A. J. Ratner, R. Krishna, J. Shen, and C. Zhang. Large language  
 276 model as attributed training data generator: A tale of diversity and bias. *Advances in Neural*  
 277 *Information Processing Systems*, 36, 2024.
- 278 X. Yue, X. Qu, G. Zhang, Y. Fu, W. Huang, H. Sun, Y. Su, and W. Chen. Mammoth: Building math  
 279 generalist models through hybrid instruction tuning. *arXiv preprint arXiv:2309.05653*, 2023.
- 280 H. Zhang, J. Da, D. Lee, V. Robinson, C. Wu, W. Song, T. Zhao, P. Raja, D. Slack, Q. Lyu, et al.  
 281 A careful examination of large language model performance on grade school arithmetic. *arXiv*  
 282 *preprint arXiv:2405.00332*, 2024.

## 283 A Appendix

### 284 A.1 Validation Criteria

285 We detail the seven criteria used for validating a question in the question validation step of the  
 286 pipeline.

- 287 • Single Answer Requirement: The question should ask for only one final answer.
- 288 • Exact Answer Requirement: There should be only one exact answer, unless approximations  
 289 are explicitly stated.
- 290 • Dual Answer Requirement: The question must necessarily and sufficiently involve the appli-  
 291 cation of both skills, with difficulty comparable to or greater than the reference examples.
- 292 • Clarity and Completeness: The question should be clear and contain all necessary informa-  
 293 tion.
- 294 • Computational Tractability: The question should not require overly complex computations.
- 295 • Realism and Logic: The scenario should be realistic and logically consistent.
- 296 • Syntax and Grammar: The question should be grammatically correct and clearly written.

### 297 A.2 Failure Modes and Interesting Behaviors

298 **Insufficient involvement of skills.** Despite clearly specifying that solving the question should  
 299 necessarily require a rigorous application of both skills, the models often generate questions that  
 300 either miss one of the skills completely or require a very shallow application of one (while the other  
 301 one is sufficiently involved) or both skills. This is the most prominent failure mode of the models in  
 302 the context of question generation. This leads to potentially easy questions, defeating the purpose of  
 303 skill composition. Consider the question given below which was generated by Claude Opus when  
 304 asked to combine the skills `ratio_and_proportion` and `geometry`.  
 305

306 **Example:** Question: A square garden is to be divided into 4  
 307 smaller square plots by two paths that are 1 meter wide and  
 308 cross each other at right angles. The paths run North-South  
 309 and East-West, splitting the garden symmetrically. If the  
 310 total area occupied by the paths is 36 square meters, find the  
 311 side length of the original square garden.

312 Upon careful examination of the question, we note that although the question tests geometry,  
 313 the involvement of `ratio_and_proportions` is practically non-existent. Further, the question  
 314 validation step in some cases also fails to identify these flaws. Supplying multi-turn human-AI  
 315 interactions where the user prompts a chatbot to generate a question combining two skills, in-context  
 316 during the generation step helps the models to avoid such questions to a certain extent. Further, to  
 317 make the question validation step more robust to such questions, we prompt the model to ensure that

318 the complexity of each skill application in the question being validated in similar to or more than the  
319 complexity of these skills in the reference examples present in the skill descriptions. The combination  
320 of these two techniques helps us nearly eliminate questions where the absent one of the skills is  
321 absent completely and reduce questions involving shallow application of skills to a significant extent.

322 **Insufficient information in the questions.** Another common failure mode of the pipeline  
323 is the generated questions missing information or details essential for solving the ques-  
324 tion. For example in the question given below which is supposed to combine the skills  
325 `understanding_and_applying_floor_and_ceiling_functions` and `basic_trigonometry`,  
326 lacks sufficient detail about the inclinations and elevations of the paths relative to the streetlight’s  
327 position which is necessary to answer the question.

328 **Example:** Question: Consider a scenario where you need to  
329 install a new streetlight at a point such that it illuminates  
330 two paths meeting at a point, each path making an angle of  $45^\circ$   
331 with the horizontal. The light from the streetlight reaches  
332 a maximum distance of 10 meters on flat ground. You are to  
333 install the streetlight at the height of  $h$  meters (where  
334  $h$  is the ceiling of the maximum distance the light reaches  
335 horizontally) such that the edge of the light’s reach just  
336 touches the ground at the end of each path. Determine the  
337 height  $h$  at which the streetlight should be installed.

338 To screen such questions, we include an explicit clause in the question validation prompt as described  
339 in Section 2. Moreover, we also notice that the inclusion of the *solution attempt* step improves the  
340 chances of detecting such errors since the missing information may not always be apparent from just  
341 the question itself. In such cases, attempting a solution (with a defeatist approach) can help detect  
342 such flaws.

343 **Unsolvable or Computationally Intractable Questions.** There are instances when the model  
344 generates questions which are unsolvable. For example the question given below has no solution  
345 which satisfies all three constraints (i.e., the area of the rectangle being 360 and the sides belonging  
346 to the two arithmetic progressions defined in the question.)

347 **Example:** Question 1: There’s a rectangle with an area of  
348 360 square units. The length of the rectangle is part of an  
349 arithmetic sequence starting at 5 and with a common difference  
350 of 7. If the other side of the rectangle is also part of  
351 an arithmetic sequence with the first term 10 and common  
352 difference 3, find the length of the shortest side of the  
353 rectangle.

354 In other instances, the model generates questions that are computationally intractable or require  
355 manually and tediously iterating through a long sequence of values. For example, solving the question  
356 given below requires manually calculating the first 100 terms of the sequence to find the sum

357 **Example:** Question 2: Consider an infinite series of numbers  
358 arranged in sections, where the  $n$ th section contains the  
359 first  $\binom{n+1}{2}$  positive integers that are divisible by  $n$  but  
360 not by any smaller positive integer (except 1). For example,  
361 the 1st section contains 1, the 2nd section starts with 2, 4,  
362 6, 10, 12, and 16 the 3rd section starts with 3, 9, 15, 21,  
363 33, ... and so on. Let  $S$  be the sum of the first 100 terms  
364 of this series. Find the sum of the digits of  $S$ .

365 While technically not wrong, such questions are not ideal for evaluating the *reasoning* abilities of the  
366 models since they mostly involve brute force calculations. Further, in cases where the sequence of  
367 calculations is very long, the LLM’s performance may be bottlenecked by other limitations such as  
368 the context length of the model.



369 Thus, we strive to filter such questions out. We add an explicit condition to check for computational  
370 tractability and solvability of the generated questions in the verification prompt. This check is assisted  
371 by the *solution attempt* produced by the model which will potentially point out any such problems.

372 **Nonsensical Questions.** In several cases, the model comes up with questions which are nonsensical  
373 - confusing, incomprehensible, logically inconsistent or ambiguous. Consider the question given  
374 below.

375 Given below is an example of a question which is logically inconsistent. More concretely, a square  
376 plot of land whose side length is equal to the radius cannot fit inside the quarter-circle.

377 **Example:** Question: A garden is designed in the shape  
378 of a quarter-circle with a radius of 8 meters. A square  
379 plot of land with a side length equal to the radius of the  
380 quarter-circle is placed inside this garden such that two of  
381 its sides are along the straight edges of the quarter-circle  
382 boundary. If the square plot of land is to be tiled entirely  
383 with square tiles each of area 64 square centimeters, what is  
384 the total number of tiles required?

385 We add checks for such cases in the question validation prompt. Further, at the end of the final  
386 solution step (maj @ 4), we further check for cases where the final answer produced in all the 4  
387 self-consistency trials are unique. If all answers are unique, we discard the question. The rationale  
388 behind this being that it is highly likely that the model produces a different answer every time due  
389 to some inherent ambiguity in the question which was not detected in the *solution attempt* and the  
390 *question validation* checks.

391 **Deceitful Solutions.** Although rare, we encounter cases where the model makes up solutions even  
392 though the question is nonsensical or cannot be solved with the amount of information provided  
393 in the question. This happens very commonly in the solutions which are generated in the *question*  
394 *generation* prompt. Thus, we do not use these solutions and include the *final solution* step where the  
395 model is asked to solve the question again. Although most of such solutions and thus questions are  
396 screened out in the *question validation* step and consistency check at the end of the *final solution*  
397 check, in rare cases we see this behavior in the solution produced after the *final solution* step as well.  
398 Given below is one such example.

399 **Example:** Question: Consider the trigonometric identity  
400  $\sin^2(x) + \cos^2(x) = 1$  and the polynomial  $P(x) = x^4 - x^2 - 12$ .  
401 Using  $x = \sin(\theta)$ , solve  $P(x) = 0$  for  $\theta$  in the interval  $[0, 2\pi)$ .

402 While solving this question, the model arrives at the conditions  $\sin(\theta) = 2$  or  $\sin(\theta) = -2$ . Clearly,  
403 these conditions have no solutions since  $-1 \leq \sin(\theta) \leq 1, \forall \theta \in [0, 2\pi]$ . However, the model goes  
404 on to argue that the the closest values to 2 and -2 in the range of  $\sin(\theta)$  are 1 and -1, and thus, goes  
405 on to solve the question for  $\sin(\theta) = 1$  or  $\sin(\theta) = -1$ .

406 **Finding a needle in the haystack.** In an attempt to make the *question validation* step robust  
407 to as many failure modes as possible, we come up with a long and sophisticated prompt (see  
408 Appendix B.3.4 for an example of this prompt). To elaborate, the validation prompt comprises of  
409 1.) skill descriptions consisting of 3 exemplars for each of the 2 skills, 2.) 6 very long validation  
410 exemplars 3.) rest of the prompt consisting of the 7 conditions described in Section 2. In prior  
411 iterations of the validation prompt, the names of the two skills which are supposed to be combined  
412 in the question, were mentioned only once throughout the prompt. In such a case, we observed that  
413 more often than not, model confused the pair of skills at hand, with skills mentioned in the validation  
414 exemplars. We also observed cases where the model was unable to locate the original skill names,  
415 came up with its own skill names depending on the skill exemplars provided in the skill descriptions.  
416 More specifically, while checking the question for **Dual Skill Requirement** in Section 2, the models  
417 would check the question for the presence of skills distinct from the pair of skills at hand. We attribute  
418 this phenomenon to the the *lost in the middle* or *finding a needle in the haystack* challenges faced by  
419 LLMs. We provide an example of the phenomenon below. Notice how in Point 3 of the reasoning  
420 trace the model checks the question for the presence of two completely unrelated skills. We were

421 able to nearly eliminate this problem by replacing referential terms (such as *above given skills, these*  
422 *skills*, etc.) to the two skills by the skill names.

423       **Example:** Skill Pair to be checked for:  
424       [understanding\_and\_manipulation\_of\_rational\_  
425       functions, division\_and\_remainders]  
426       Question: A fair 6-sided die is rolled 4 times. What is  
427       the probability that the product of the numbers rolled  
428       is divisible by 36? Express your answer as a simplified  
429       fraction.  
430       Reasoning trace of the model:  
431       # REASONING  
432       Let's evaluate the given question and solution against the  
433       specified criteria:  
434       1. **\*\*Single Answer Requirement\*\***: The question asks for a  
435       single value - the probability of a specific event occurring.  
436       This satisfies the single answer requirement.  
437       2. **\*\*Exact Answer Requirement\*\***: The solution provides an  
438       exact answer in the form of a simplified fraction (61/2386).  
439       No approximations or assumptions are made, fulfilling the  
440       exact answer requirement.  
441       3. **\*\*Dual Skill Requirement\*\***: The question requires  
442       knowledge and application of two distinct mathematical  
443       skills: - Probability theory: Calculating the probability of  
444       specific outcomes in a series of die rolls.    - Number theory:  
445       Understanding divisibility rules for 4 and 9, and how they  
446       apply to the product of numbers.  
447       The solution demonstrates a rigorous application of both  
448       skills, using probability calculations (such as the binomial  
449       probability for rolling a 4 at least once) and number  
450       theoretic concepts (divisibility rules for 4 and 9) to arrive  
451       at the final answer. The complexity of skill application is  
452       on par with the provided examples.  
453       4. **\*\*Clarity and Completeness\*\***: The question is clearly  
454       stated and provides all necessary information (fair 6-sided  
455       die, rolled 4 times, product divisible by 36) to solve the  
456       problem. No ambiguity or missing information is present.  
457       5. **\*\*Relevancy of Information\*\***: All given information is  
458       relevant and used in the solution. There are no extra details  
459       that do not contribute to solving the problem.  
460       6. **\*\*Realism and Logic\*\***: The scenario of rolling a die  
461       multiple times and calculating probabilities is a realistic  
462       and commonly encountered problem in probability theory. The  
463       data (6-sided die, 4 rolls) is logically consistent.  
464       7. **\*\*Syntax and Grammar\*\***: The question is grammatically  
465       correct and clearly written, with no room for  
466       misinterpretation.  
467       The solution is well-structured and easy to follow, breaking  
468       down the problem into clear steps and explaining the reasoning  
469       behind each step. The mathematical reasoning is sound,  
470       correctly applying probability and number theory concepts  
471       to arrive at the final answer.  
472       # FINAL ANSWER  
473       Yes

474 Note that none of the above failure modes are completely eliminated in the pipeline described in  
475 Section 2. Thus, human verification is required.

476 Despite struggling with the failure modes described above, there also exist cases where the models  
477 exhibit positively surprising and creative behaviors. We talk about some of them below.

478 **Thinking out of the box.** Although rare, we observe instances where the models get creative while  
479 validating the question. Consider the question below

480 **Example:** Question: A class of students is learning about  
481 combinatorics and geometry. They are given a problem  
482 involving colored beads: Red, Blue, and Green. If they need  
483 to form a necklace with 8 beads such that no two adjacent  
484 beads have the same color and the necklace begins and ends  
485 with a bead of a different color, how many different necklaces  
486 can they create? Each necklace is counted up to rotation and  
487 reflection (considering the necklace can be flipped over).

488 When validating this question using prior iterations of the *question validation* prompt, which did  
489 not consist of the computational tractability check, the model output while validating the question  
490 consists of the following excerpt.

491 **Example:** *...This might introduce a significant challenge not solely due to the*  
492 *methodology's complexity but also due to the potential computational require-*  
493 *ment, which may not be feasible in a standard test environment without tools.*  
494 *Furthermore, while the connection to practical geometry (reflective and rotational*  
495 *symmetry) and combinatorics (color patterning and adjacency constraints) is*  
496 *strong, the depth of understanding required to manually adjust for these symmetry*  
497 *considerations in a test question might be too intense or require more guided*  
498 *learning than a single evaluation question could provide....*

499 i.e, the model takes into consideration the fact that the question involves a lot of brute force computa-  
500 tion, despite there being no explicit check for computation complexity in the prompt, and classifies  
501 the question as invalid. We attribute such out of the box thinking behavior to the role-playing nature  
502 of our prompts. Our prompts consist of a math teacher evaluating the the fitness of the given question  
503 for being used for testing students' reasoning and analytical skills in a math exam. This leaves room  
504 open for the model to detect potential problems not explicitly accounted for in the prompts which  
505 might make the question unfit for being used for evaluation.

### 506 A.3 Considerations for human-annotators

507 Human annotators were tasked with double checking the validity of the question and the correctness  
508 of the solution. They were asked to look out for any of the failure modes discussed in Section B.  
509 They were asked to check that the created question actually used the math skills it was supposed  
510 to exhibit and to improve the question with respect to readability, quality and difficulty. They were  
511 encouraged to suggest changes that make the problem harder to solve using automated tools while  
512 retaining easiness for the humans. We illustrate with an examples.

513 GPT-4 created the following question given the skill-tags `recursive_functions_and_sequences`  
514 and `multiplication_and_division` :

515 **Example:** Original Question: Consider the sequence defined  
516 recursively by  $a_1 = 1$  and  $a_{n+1} = 2a_n + n$  for all  $n \geq 1$ . What is  
517 the product of the first five terms of this sequence?

518 An LLM can solve this by simple computation. The human modified the question so that solving the  
519 problem requires understanding the underlying pattern.

520 **Example:** Modified Question: A sequence is defined  
521 recursively as follows: the first term  $a_1$  is 2, and for  $n \geq 2$ ,  
522  $a_n = 2^{n-1} + n$ . What is the logarithm (base 2) of the average  
523 of the first 50 terms of this sequence? Round down to the  
524 nearest integer.

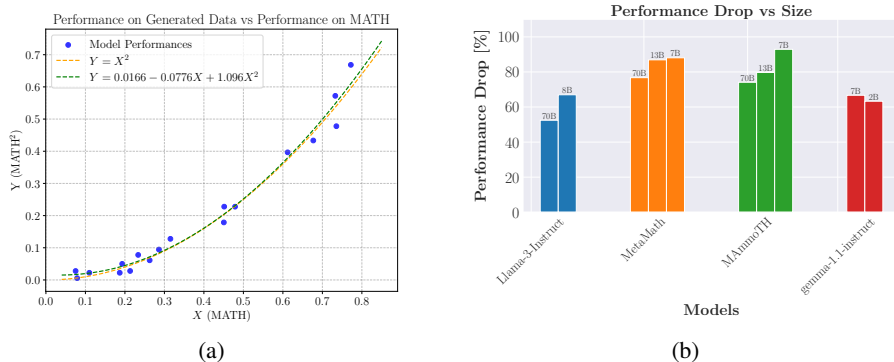


Figure 2: (a) Relation between the performance of models on MATH<sup>2</sup> ( $Y$ ) vs their performances on MATH ( $X$ ). As can be seen from the plot, the performance on models on generated questions roughly follows a quadratic relation with the performance of those models on MATH. The best quadratic fit follows the relation:  $Y = 0.0166 - 0.0776X + 1.096X^2$ . This may be explained by the fact that the questions in MATH<sup>2</sup> consist of two skills at a time, as compared to questions in MATH, which consist of one skills. (b) Comparative bar plot showing smaller models generally show larger degradation in performance (relative to MATH) as compared to their larger counterparts, with the exception of the Gemma family of models, where the 7B parameter model shows a larger deterioration of 66.669% as compared to 63.03% in the 2B parameter model

525 For the modified question, one leading model mentioned calculation difficulties for the inability to  
 526 give any answer, and another resorted to an incorrect numerical approximation that led to an incorrect  
 527 answer.

528 Human annotators were also asked to go through the solutions carefully and correct or improve  
 529 the solution for good questions if necessary. They were also asked to look out for questions that  
 530 contain lot of enumeration, i.e. questions which are tedious and require significant amount of  
 531 brute force computation. For such questions, the annotators were encouraged to reword them such  
 532 that enumeration is not a feasible strategy below. For example, given below is an example of an  
 533 enumerative question which was modified to avoid enumeration.

534 **Example:** Original Question: Find the sum of the smallest  
 535 prime divisor and the largest prime divisor of the number  
 536  $N = 15^4 + 16^4$ .  
 537 Modified Question: Find the sum of the two smallest prime  
 538 divisors of  $23^{17} + 17^{17}$ .

539 Models tend to adopt brute force approach on the original question calculating  $15^4 + 16^4$ . After  
 540 rephrasing the models cannot use brute force on  $23^{17} + 17^{17}$ , instead being forced to check the  
 541 divisors more analytically, in particular understanding of arithmetic modulo a prime.

#### 542 A.4 Further Experimental Details and Results

543 For generating responses, we use the MAMmoTH [Yue et al., 2023] evaluation suite. The responses  
 544 are graded using a GPT-4 grader, where GPT-4 Omni checks the correctness of a solution response  
 545 against the ground truth solution. During response generation, we set the temperature to 0 and top\_p  
 546 to 1 for all models. For open source LLMs, we use 2 80GB A100 GPUs and 72GB of RAM to run  
 547 inference facilitated by vLLM [Kwon et al., 2023]. We use 25 workers while querying GPT-4 Omni  
 548 and GPT-4 Turbo and 2 workers for querying Claude-3 Opus and Claude-3.5-Sonnet.

##### 549 A.4.1 Generated Data Statistics

550 Out of 180 question-solution pairs included in MATH<sup>2</sup>, 79 underwent some form of modification by  
 551 the human annotators before being included in the dataset. Out of the 41 questions modified, 3 were  
 552 minor modifications to improve the clarity of the question. Another 24 modifications were minor

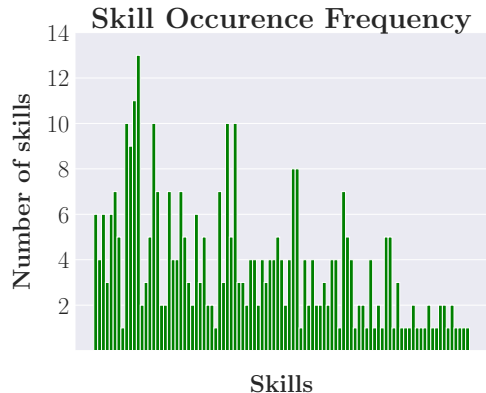


Figure 3: Shows the distribution of different skills extracted during the skill extraction process in the generated set of questions. The generated and human verified set of 180 questions consists of 97 skills out of the 114 skills extracted via the skill extraction process as described in Didolkar et al. [2024], Each question in the generated set represents two skills. Note that the distribution of skills is not uniform with there being multiple skills that are represented by one one question.

Table 2: **Human Verification Statistics:** Out of a total of 180 examples in MATH<sup>2</sup>, 79 (43.89%) were such that either the question or the solution generated by the model were modified by the annotator before being included in the final dataset, whereas 20 (11.12%) were such that both the question and the solution were modified. These modifications were made in order to increase the difficulty of the questions or correct the questions or solutions.

# of Modified Questions ( $A$ )	# of Modified Solutions ( $B$ )	# of $A \cup B$	Dataset Size
41	62	79	180

553 modifications (up to 3 words), which nevertheless affected the meaning of the question and changed  
 554 the final answer. But 14 modifications were significant; either making the given questions harder, or  
 555 correcting them, or making them more interesting (i.e., less tedious) for humans.

556 As for the solutions, 62 out of the 180 solutions originally generated by the model were modified to  
 557 correct them or improve their clarity.

#### 558 A.4.2 Performance across the two datasets: A surprising pattern

559 Table 1 shows that all tested models have significantly lower performance on MATH<sup>2</sup> than on the orig-  
 560 inal MATH dataset. Denoting  $Y$  as the performance on MATH<sup>2</sup> and  $X$  as the performance on MATH,  
 561 the percentage drop  $100(X - Y)/X$  for frontier models — GPT-4 Omni, GPT-4 Turbo, Gemini-1.5-  
 562 Pro, Claude-3.5-Sonnet and Claude 3 Opus — ranges from 13.42% to 35.45%. MAMmoTH-7B, a  
 563 specialist math model, shows the largest drop at **92.91%**.

564 The fact that performance drops for all models should *not* be too surprising, since as noted, the  
 565 MATH<sup>2</sup> questions, by combining skills from different subareas of MATH, could be seen as “out of  
 566 distribution (OOD).” This makes it tempting to interpret the percentage drop as a measure of a model’s  
 567 (lack of) “OOD-resilience.” For instance, very large percentage drops seen with open-source models  
 568 MetaMath and MAMmoTH feel understandable since their training used synthetic data generated  
 569 using seed questions from MATH and GSM-8k. Lack of diversity in such synthetic data is known to  
 570 cause overfitting to the dataset being imitated. Similarly, GPT-4o and Claude Sonnet 3.5 are suspected  
 571 to also have been extensively trained with synthetic data. Although their MATH performance is  
 572 similar, Sonnet 3.5 has worse MATH<sup>2</sup> performance, which might suggest lower quality/diversity in  
 573 its synthetic data.

574 However, in our opinion, the overall pattern among proprietary models of similar size does fit with the  
 575 OOD story. A much simpler explanation pops out when we plot  $Y$  vs  $X^2$  (Figure 4 and Figure 2(a)):  
 576 we find a linear relationship  $Y \approx X^2$ ! This implies that the relative drop in performance of the  
 577 models is well-predictable from just their performance on MATH, and does not require taking their  
 578 training details into account!

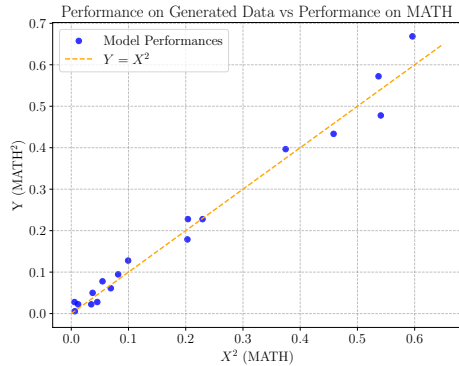


Figure 4: Relation between the performance of models on MATH<sup>2</sup> ( $Y$ ) vs the square of their performances on MATH ( $X^2$ ). As can be seen from the plot,  $Y \approx X^2$ . See Appendix A.4 for the best-fit quadratic curve, which is slightly different.

579 Why should the two scores be expected to have this relationship? Here is a natural (albeit heuristic)  
 580 explanation. Suppose there are  $N$  skills and  $s_i$  denotes the success rate of the model at correctly  
 581 applying the  $i$ th skill. Then, its  $X$  value should reflect the average<sup>4</sup> of the  $s_i$ 's. Furthermore, on a  
 582 random question using the  $i$ th and  $j$ th skill, the probability that the model correctly answers it should  
 583 be  $s_i s_j$ , since it has to successfully apply both skills. If the questions are created using pairs of skills  
 584 chosen randomly and independently, then the  $Y$  value will be the average value of  $s_i s_j$ 's, which by  
 585 independence will be roughly  $X^2$ .

586 This reasoning in fact suggests that our pipeline has created questions that genuinely required applying  
 587 two very distinct skills (as opposed to, say, requiring primarily skill  $i$ , and mildly using skill  $j$ ). The  
 588 discovered relationship suggests further that if we could create questions where each combines  $k$   
 589 skills, we might see the relationship  $Y \approx X^k$ , which would tend to further magnify performance  
 590 differences between models.

#### 591 A.4.3 Generated Questions are Effective In-Context Exemplars for MATH.

592 A possible test for the quality of a Q&A pair on similar topics as MATH dataset is whether perfor-  
 593 mance on MATH improves when using these as in-context exemplars.

594 We test as follows. Recall that MATH has 7 sections. Exemplars for a section are chosen from the  
 595 section area. However, by design, our new questions cross section boundaries. We implemented a  
 596 new procedure to retrieve in-context exemplars from MATH<sup>2</sup> based on the skill requirements of the  
 597 current question.

598 Since MATH<sup>2</sup> is limited in size, it does not cover all the skills extracted during the skill extraction  
 599 process, containing 97 out of 114 skills. Figure 3 shows the distribution of different skills in the  
 600 dataset. We filtered the MATH test set to remove examples requiring skills not present in the generated  
 601 dataset, resulting in the removal of 809 test examples. During evaluation on the filtered MATH  
 602 test set, for each question  $Q$  labeled with skill  $a$  ( $a \in \mathcal{S}$ , where  $\mathcal{S}$  is the set of extracted skills), we  
 603 retrieved in-context exemplars from the MATH<sup>2</sup>, ensuring each exemplar involved skill  $a$ . We used  
 604 four such exemplars per question (i.e., 4-shot CoT [Wei et al., 2022]). To handle skills represented  
 605 by fewer than four examples in MATH<sup>2</sup>, we run two experiments: (A) **Proposed 4-shot CoT**: If a  
 606 given skill is represented by  $n$  examples in the MATH<sup>2</sup>, where  $n < 4$ , we use  $n$  in-context examples  
 607 instead of 4 exemplars. (B) **Proposed + Skill Based 4-shot CoT**: If a given skill is represented by  
 608  $n$  examples in MATH<sup>2</sup>, where  $n < 4$ , we supplement  $4 - n$  exemplars for that skill from MATH  
 609 training set. The relevant in-context exemplars in MATH training set are determined by following  
 610 the methodology proposed in Didolkar et al. [2024]. We compared the performance of models using  
 611 these targeted prompting strategies against two baselines: (C) **MAMmoTH 4-shot CoT**: The 4  
 612 in-context exemplars are taken from the MAMmoTH evaluation suite [Yue et al., 2023]. (D) **Skill**  
 613 **Based 4-shot CoT**: We use skill-based prompting as proposed in Didolkar et al. [2024], where the

<sup>4</sup>With perhaps a small correction factor if the skills are not evenly distributed among the questions

Table 3: Performance of models on MATH under two different prompting strategies. **MAmmoTH 4-shot CoT** prompting involves 4-shot prompting with exemplars taken from the MAmmoTH [Yue et al., 2023] evaluation suite. **Skill Based 4-shot CoT** [Didolkar et al., 2024] consists of using 4 exemplars which are retrieved from the training set of MATH based on which skill is required to solve the given question (as determined by GPT-4). **Proposed 4-shot CoT** prompting consists of 4-shot prompting with exemplars taken from MATH<sup>2</sup>. These exemplars are retrieved such that one of the two skills in each exemplar is the same as the skill required by the question at hand, as labeled by GPT-4. In **Proposed + Skill Based 4-shot CoT** we supplement the exemplars retrieved from MATH<sup>2</sup> with exemplars from MATH training set, for skills that are present in < 4 questions in MATH<sup>2</sup>. We show that few-shot prompting with exemplars retrieved from the generated set of questions (MATH<sup>2</sup>) consistently outperforms vanilla few-shot prompting with gains of upto 3.75% (for Mixtral-8×7B-Instruct [Jiang et al., 2024]).

Method	GPT-4O	GPT-4T	Llama-3-70B-Instruct	MetaMath-70B	MAmmoTH-70B	Mixtral-8×7B-Instruct
MAmmoTH 4-shot CoT	76.67%	71.89%	49.11%	25.77%	18.45%	30.77%
Skill Based 4-shot CoT	78.32%	72.77%	47.80%	25.42%	18.20%	30.31%
Proposed 4-shot CoT	78.39%	<b>74.19%</b>	<b>51.13%</b>	26.60%	20.06%	<b>34.52%</b>
Proposed + Skill Based 4-shot CoT	<b>78.49%</b>	73.45%	50.80%	<b>27.50%</b>	<b>20.38%</b>	33.86%

614 in-context exemplars are selected from the MATH training set, in accordance to the skill required by  
 615 the question at hand, as determined by GPT-4.

616 Table 3 presents the results of this comparison. The two prompting strategies using questions  
 617 from MATH<sup>2</sup> as in-context exemplars, clearly outperform the two baselines. We conclude that the  
 618 MATH<sup>2</sup> questions, due to their difficulty and skill relevance, serve as effective in-context exemplars.  
 619 Performance gains would likely be more significant with larger datasets generated using our approach,  
 620 reducing the need to supplement with external exemplars.

## 621 B Observations from the Question Generation Process

622 The question generation pipeline described in Section 2 was developed through an iterative process  
 623 of refining prompts and design choices, and evaluating their impact on the quality of the final  
 624 questions and solutions. Notably, the inclusion of the *attempted solution* and *question validation* steps  
 625 significantly enhanced the pipeline’s effectiveness. Despite the sophistication of the pipeline and  
 626 prompts, we still observe instances where models fail to follow the given instructions. This section  
 627 highlights prominent failure modes at various stages of the pipeline, which human raters need to be  
 628 aware of. Additionally, we explore some intriguing behaviors of the models where they successfully  
 629 create interesting and creative questions. Section B.1 details the role of human raters in improving  
 630 these questions.

### 631 B.1 Creative questions: Examples of Synergy from Human-AI interaction

632 The models frequently produced interesting and creative questions, although they often failed to  
 633 generate correct solutions. In these cases, the incorrect solutions usually contained enough correct  
 634 ideas for a human to quickly complete them.

635 Human annotators were tasked with verifying the validity of the questions and the correctness of the  
 636 solutions. They were instructed to look out for any failure modes discussed in Section B.1.1. Their  
 637 responsibilities included ensuring that the created questions actually employed the intended math  
 638 skills, and improving the questions in terms of readability, quality, and difficulty when possible. They  
 639 were encouraged to suggest changes that would make the problems harder for automated tools to  
 640 solve while allowing easier or more elegant solutions for humans. The following examples illustrate  
 641 this process:

642 **Example:** Original Question: Find the smallest positive  
 643 integer  $k$  such that  $k^3 - k$  is divisible by both 9 and 10,  
 644 and the sum of digits of  $k$  in its decimal representation is  
 645 a prime number.

646 Our human team had not encountered such questions before. It requires recognizing that  $k^3 - k =$   
 647  $k(k - 1)(k + 1)$  is always divisible by 2 and 3. Thus,  $k$  must be such that  $k(k - 1)(k + 1)/6$  is

648 divisible by 15 (both 3 and 5). Additionally, the sum of the digits of  $k$  must be a prime number, and  
649 ensuring such conditions is challenging even for powerful LLMs.

650 **Example:** Original Question: Consider a collection of red,  
651 blue, and green beads arranged in an infinite series. The  
652 beads alternate in color, starting with red, then blue, then  
653 green, and this pattern repeats indefinitely. The number of  
654 beads in each colored section follows the pattern of powers of  
655 2: the first red section has 2 beads, the first blue section  
656 has 4 beads, the first green section has 8 beads, the second  
657 red section has 16 beads, and so on. If a bracelet is made  
658 using a continuous, unbroken sequence of exactly 20 beads from  
659 this series, and each bead has a length of 0.5 units, how many  
660 different bracelets can be made such that the perimeter of the  
661 bracelet is an integer value?

662 The original question combined elements in a novel way. The human rater modified the question to  
663 change the sequence size from 20 to 6 beads, maintaining the essential difficulty while making it  
664 more elegant for humans. All tested models failed on the modified question.

665 **Example:** Original Question: A container initially contains  
666 500 mL of water. A scientist adds water to the container  $\frac{1}{4}$   
667 of the current amount every minute. After how many minutes  
668 will the container first contain more than 1 L but less than 2  
669 L of water?

670 Modified Question: A container starts with 500 mL of water.  
671 Each minute, the scientist adds water equal to  $\frac{1}{2}$  of the  
672 current amount. What is the smallest positive integer  $n$  such  
673 that the number of liters of water in the container is never  
674 in the interval  $[n, n + 1]$ ?

675 This was one of many questions the models created about exponential growth and geometric series,  
676 possibly similar to standard math test questions. The human slightly altered it to simplify calculations  
677 by hand and substituted a different condition that the models found challenging, while humans could  
678 easily estimate an approximate answer and then verify.

679 **Example:** Original Question: Consider the sequence defined  
680 recursively by  $a_1 = 1$  and  $a_{n+1} = 2a_n + n$  for all  $n \geq 1$ . What is  
681 the product of the first five terms of this sequence?

682 Modified Question: A sequence  $a_n$  is defined as follows:  $a_1 =$   
683  $2$  and  $a_n = 2^{n-1} + a_{n-1} + n$ . What is the  $\lfloor \log_2 a_{500} \rfloor$ ?

684 An LLM can solve the original question through simple computation. The modified question, however,  
685 requires understanding an underlying pattern.

686 **Example:** Original Question: Find the sum of the smallest  
687 prime divisor and the largest prime divisor of the number  
688  $N = 15^4 + 16^4$ .

689 Modified Question: Find the sum of the two smallest prime  
690 divisors of  $23^{17} + 17^{17}$ .

691 Models tend to adopt a brute-force approach to the original question by calculating  $15^4 + 16^4$ . After  
692 rephrasing, the number  $23^{17} + 17^{17}$  is too large for direct computation, requiring understanding of  
693 arithmetic modulo a prime.

694 These examples highlight the essential role of human oversight in refining and improving the questions  
695 generated by LLMs, ensuring they are challenging, creative, and suitable for advanced mathematical  
696 problem-solving.



697 **B.1.1 Failure Modes**

698 Despite the sophistication of our pipeline, models frequently exhibit several failure modes: (a)  
 699 *Insufficient Involvement of Skills*: Models often generate questions that either miss one of the skills  
 700 completely or require a very shallow application of one or both skills. For example, a geometry  
 701 question may fail to involve ratio and proportion adequately, (b) *Insufficient Information*: Questions  
 702 may lack essential details needed for solving, making them incomplete or ambiguous. For instance, a  
 703 trigonometry question might omit necessary angles or distances, (c) *Unsolvable or Computationally*  
 704 *Intractable Questions*: Some questions generated are either unsolvable or require excessive brute-  
 705 force calculations, which are impractical for evaluating reasoning abilities, (d) *Nonsensical Questions*:  
 706 Models sometimes produce questions that are logically inconsistent, confusing, or ambiguous, such  
 707 as a probability problem with unclear parameters or an impossible geometry scenario, (e) *Deceitful*  
 708 *Solutions*: Occasionally, models fabricate solutions to nonsensical or unsolvable questions, presenting  
 709 incorrect logic as plausible reasoning and (f) *Finding a Needle in the Haystack*: Long and complex  
 710 validation prompts sometimes cause models to confuse or overlook the specified skills, leading  
 711 to incorrect evaluations. For a more detailed discussion and examples of questions in the various  
 712 categories listed above, refer to Appendix A.2.

713 **B.1.2 Skill Proportional Comparison of MATH<sup>2</sup> and MATH**

714 Figure 3 shows the distribution of different skills in MATH<sup>2</sup>. To make a fairer comparison of MATH  
 715 and MATH<sup>2</sup>, and to show empirically that MATH<sup>2</sup> benefits from the composition of two skills at the  
 716 same time as compared to MATH which consists of application of one skill at a time, we compare the  
 717 performance of models on MATH<sup>2</sup> to the performance of models on a subset of MATH which has as  
 718 similar skill distribution as MATH<sup>2</sup> (i.e. as shown in Figure 3). We form this subset by randomly  
 719 sampling questions belonging to each skill in MATH. The subset consists of 3634 questions. Table 4  
 720 compares the performance of some models on MATH<sup>2</sup>, MATH and the subset of MATH formed  
 721 above. From the performance of the models, we can conclude that a subset of MATH with a similar  
 722 distribution of skills is not just easier than MATH<sup>2</sup>, but also MATH.

Table 4: Comparison of the performance of various models on MATH, MATH<sup>2</sup> and a subset of MATH which has a similar distribution of skills as MATH<sup>2</sup>, as shown in Figure 3

Model	MATH <sup>2</sup> (Y)	MATH skill proportional subset	MATH (X)
GPT-4 Omni	66.85%	79.42%	77.21%
Claude 3.5 Sonnet	47.78%	75.68%	73.54%
GPT-4 Turbo	57.22%	75.45%	73.27%
Llama-3-70B-Instruct	22.77%	49.49%	47.89%
MetaMath-70B	6.11%	28.20%	26.27%
MAmmoTH-70B	5.00%	21.10%	19.31%
MetaMath-13B	2.79%	23.06%	21.32%
MAmmoTH-13B	2.23%	11.60%	10.99%
Deepseek-math-7b-instruct	17.88%	46.54%	45.05%
Gemma-1.1-7B-Instruct	7.78%	23.86%	23.36%
MetaMath-7B	2.23%	20.08%	18.69%
MAmmoTH-7B	0.56%	8.18%	7.90%
Phi-3-mini-128k-instruct	22.78%	46.31%	45.14%
Gemma-1.1-2B-Instruct	2.78%	7.74%	7.52%

723 **B.1.3 Difficulty of questions generated by different models**

724 Out of the 180 questions in MATH<sup>2</sup>, 136 were generated using GPT-4 Turbo and 44 were generated  
 725 using Claude-3 Opus. We partition the dataset into these two subsets and evaluate GPT-4O, GPT-4  
 726 Turbo, Claude-3 Opus and Claude-3.5-Sonnet on both subsets. The results are shown in Table 5

727 The results above show that each model performs roughly similarly on both subsets. We can also  
 728 conclude that questions generated by a given model are not necessarily easier for that model to solve.

Table 5: Performance of GPT-4 and Claude on questions generated using GPT-4 Turbo and Claude-3 Opus

Subset	GPT-4 Omni	GPT-4 Turbo	Claude-3.5-Sonnet	Claude-3 Opus
GPT-4 Turbo Subset	63.01%	61.64%	50.68%	35.29%
Claude-3 Opus Subset	64.71%	64.70%	50.00%	36.98%

729 **B.1.4 Modified Questions vs Non-Modified Questions**

730 During the human verification process, the annotators were instructed to be on the look out for any  
 731 errors in the questions and solutions generated by the models, and fix any lack of clarity, ambiguity,  
 732 convoluted language, etc. in the generated questions which might confuse the model and reduce the  
 733 “quality” of the questions. They were also instructed to look out for specific modifications which  
 734 could make the questions more difficult. For further discussion on the human verification process,  
 735 refer to Section A.3. In Table 6 we compare the performance of models on the questions which were  
 736 modified against their performance on the questions which were not modified. We also compare the  
 737 performance of models on these two subsets with their performance on MATH Level-5 questions.  
 738 We see that that the questions in MATH<sup>2</sup> which were modified during the human verification process  
 739 are significantly more difficult than the questions which were not modified. Moreover, they are also  
 740 more difficult than MATH Level-5 questions which is the most difficult section in the MATH dataset.  
 741 The MATH level 5 questions and human modified questions in MATH<sup>2</sup>, both comprise roughly 22%  
 742 of the respective datasets.

Table 6: Performance of models on human modified and non-modified questions from MATH<sup>2</sup> and on MATH Level-5 questions.

Model	MATH <sup>2</sup> Unmodified	MATH <sup>2</sup> Modified	MATH Level-5
GPT-4 Omni	71.43%	44.68%	57.52%
GPT-4 Turbo	61.65%	48.94%	49.92%
Claude-3.5-Sonnet	51.13%	44.68%	53.29%
Llama-3-70B Instruct	25.56%	19.15%	21.22%
MetaMath-13B	4.54%	2.13%	5.67%
Gemma-1.1-7B-Instruct	9.77%	4.26%	5.60%

743 **B.2 Skill composition using different models**

744 In this section, we provide examples of the responses of different models when asked to generate  
 745 a question combining two different skills. We prompt GPT-4o, Claude 3 Opus, Gemini 1.5 Pro,  
 746 and Llama-3-70B-Instruct and Llama-3-8B-Instruct to generate a question combining the skills  
 747 `area_calculation_skills` and `algebraic_manipulation_skills`. We use the same prompt  
 748 used in the **Question Generation** step (Step 2) of the proposed pipeline described in Section 2 (see  
 749 Appendix B.3.2). Given below are the responses of each model to the prompt.

**GPT-4 Omni**

Sure! Here is a challenging math question that tests both area calculation skills and algebraic manipulation skills, formatted as requested.

**# QUESTION**  
 A rectangle  $ABCD$  has vertices  $A(0, 0)$ ,  $B(a, 0)$ ,  $C(a, b)$ , and  $D(0, b)$ . Inside this rectangle, there is a right triangle  $PQR$  where  $P$  is at vertex  $A$ ,  $Q$  is on  $BC$ , and  $R$  is on  $CD$ . The area of triangle  $PQR$  is given by  $\frac{1}{6}$  of the area of rectangle  $ABCD$ . Given that the coordinates of  $Q$  and  $R$  are  $(a, kb)$  and  $(ka, b)$  respectively, where  $k$  is a constant, find the value of  $k$ .

**# SOLUTION**  
 First, let’s find the area of the rectangle  $ABCD$ :

750

$$\text{Area of rectangle } ABCD = a \cdot b$$

The area of the triangle  $PQR$  is given as  $\frac{1}{6}$  of the area of the rectangle  $ABCD$ :

$$\text{Area of triangle } PQR = \frac{1}{6} \cdot (a \cdot b) = \frac{ab}{6}$$

Next, we need to find the area of triangle  $PQR$  using the coordinates of  $P$ ,  $Q$ , and  $R$ . The vertices of triangle  $PQR$  are  $P(0, 0)$ ,  $Q(a, kb)$ , and  $R(ka, b)$ .

The formula for the area of a triangle with vertices at  $(x_1, y_1)$ ,  $(x_2, y_2)$ , and  $(x_3, y_3)$  is:

$$\text{Area} = \frac{1}{2} |x_1(y_2 - y_3) + x_2(y_3 - y_1) + x_3(y_1 - y_2)|$$

Substitute the coordinates of  $P$ ,  $Q$ , and  $R$  into the formula:

$$\text{Area of } PQR = \frac{1}{2} |0(kb - b) + a(b - 0) + ka(0 - kb)|$$

$$\text{Area of } PQR = \frac{1}{2} |0 + ab - ka^2b|$$

$$\text{Area of } PQR = \frac{1}{2} |ab(1 - k)|$$

$$\text{Area of } PQR = \frac{ab}{2} |1 - k|$$

We know that the area of triangle  $PQR$  is  $\frac{ab}{6}$ :

$$\frac{ab}{2} |1 - k| = \frac{ab}{6}$$

To solve for  $k$ , we can divide both sides by  $\frac{ab}{2}$ :

$$|1 - k| = \frac{1}{3}$$

This gives us two possible equations:

$$1 - k = \frac{1}{3} \quad \text{or} \quad 1 - k = -\frac{1}{3}$$

Solving these equations for  $k$ :

$$1 - k = \frac{1}{3} \quad \Rightarrow \quad k = 1 - \frac{1}{3} = \frac{2}{3}$$

$$1 - k = -\frac{1}{3} \quad \Rightarrow \quad k = 1 + \frac{1}{3} = \frac{4}{3}$$

Therefore, the possible values for  $k$  are  $\frac{2}{3}$  and  $\frac{4}{3}$ .

# DETAILS

The question involves area calculation by requiring students to find the area of a triangle within a

rectangle, and algebraic manipulation by solving an equation involving a variable. The problem complexity is similar to the examples provided, ensuring that students are sufficiently challenged. The final answer consists of the value(s) of  $k$ , making it straightforward to evaluate the students' answers.

752

### Claude 3 Opus

Here is my attempt at generating 1 (non multiple choice) difficult question which tests both area\_calculation\_skills and algebraic\_manipulation\_skills:

#### # QUESTION

Consider an equilateral triangle  $ABC$  with side length  $2a$ . Let points  $D, E, F$  be the midpoints of sides  $BC, CA, AB$  respectively.

Now construct squares  $ACDE, ABFG, BCHI$  externally on sides  $AC, AB, BC$  respectively.

Express the area of the shaded region (consisting of the equilateral triangle  $DEF$  and squares  $ACDE, ABFG, BCHI$ ) in the form  $ka^2$ , where  $k$  is a constant. Find the value of  $k$ .

#### # SOLUTION

Let the side length of the equilateral triangle  $ABC$  be  $s = 2a$ .

The area of equilateral triangle  $ABC$  is  $\frac{\sqrt{3}}{4}s^2 = \frac{\sqrt{3}}{4}(2a)^2 = \sqrt{3}a^2$ .

In an equilateral triangle, the distance from a vertex to the midpoint of the opposite side is  $\frac{\sqrt{3}}{2}$  times the side length. So  $DE = EF = FD = \frac{\sqrt{3}}{2}s = \sqrt{3}a$ .

The area of equilateral triangle  $DEF$  is  $\frac{\sqrt{3}}{4}(\sqrt{3}a)^2 = \frac{3\sqrt{3}}{4}a^2$ .

The side length of each of the squares  $ACDE, ABFG, BCHI$  is equal to the side length of the equilateral triangle  $ABC$ , which is  $s = 2a$ .

So the area of each square is  $s^2 = (2a)^2 = 4a^2$ .

The total area of the shaded region is:

Area of  $DEF + 3 * \text{Area of each square}$

$$= \frac{3\sqrt{3}}{4}a^2 + 3(4a^2) = \frac{3\sqrt{3}}{4}a^2 + 12a^2 = (12 + \frac{3\sqrt{3}}{4})a^2$$

Therefore, the area of the shaded region is  $ka^2$ , where  $k = 12 + \frac{3\sqrt{3}}{4} \approx 15.1981$ .

#### # DETAILS

This question combines area calculation of equilateral triangles and squares with algebraic manipulation of side lengths and areas. Key steps involve:

1) Calculating areas of equilateral triangles  $ABC$  and  $DEF$  in terms of side length  $s$  and  $a$  respectively. This tests area\_calculation\_skills.

2) Algebraically manipulating the expressions for side lengths  $DE, EF, FD$  and square side lengths in terms of  $s$  and  $a$ . This requires algebraic\_manipulation\_skills.

3) Adding up the component areas algebraically and simplifying the expression in the form  $ka^2$  to solve for constant  $k$ .

The question difficulty is calibrated to be similar to the provided examples. The solution is a single numerical constant  $k$  which avoids ambiguity in student answers. The equilateral triangle and square constructions provide an interesting geometric visualization.

Let me know if you would like me to modify anything in the question or solution to better suit your needs! I'm happy to refine it further.

753

### Gemini 1.5 Pro

#### # QUESTION

A circle is inscribed in equilateral triangle  $ABC$ , and a smaller circle is externally tangent to this circle

754

and to sides  $AB$  and  $AC$ . If the side length of triangle  $ABC$  is 4, what is the area of the smaller circle?

```
[asy] size(150); defaultpen(linewidth(0.8)); pair A=(0,2*sqrt(3)), B=(-2,0), C=(2,0), O=(0,sqrt(3)/3), P=(0,sqrt(3)); draw(A--B--C--A); draw(circle(O,sqrt(3))); draw(circle(P,sqrt(3)-1)); [/asy]
```

# SOLUTION

Let  $O$  be the center of the larger circle, and let  $P$  be the center of the smaller circle. Let  $r$  be the radius of the smaller circle.

```
[asy] size(150); defaultpen(linewidth(0.8)); pair A=(0,2*sqrt(3)), B=(-2,0), C=(2,0), O=(0,sqrt(3)/3), P=(0,sqrt(3)); draw(A--B--C--A); draw(circle(O,sqrt(3))); draw(circle(P,sqrt(3)-1)); draw(O--P); draw(O--C); draw(P--C); label("A",A,N); label("B",B,SW); label("C",C,SE); label("O",O,S); label("P",P,NE); label("2",(1,0),S); label("r",(P+C)/2,E); label("1",(C+O)/2,SE); [/asy]
```

Since  $O$  is the centroid of equilateral triangle  $ABC$ ,  $OC = \frac{2}{3} \cdot 2\sqrt{3} = \frac{4\sqrt{3}}{3}$ . Also,  $PC = r + 1$ , so by the Pythagorean Theorem on right triangle  $POC$ ,

$$r^2 + 2r + 1 = PC^2 = OC^2 - OP^2 = \frac{16}{3} - (r + \sqrt{3})^2.$$

This simplifies to  $2r^2 + (2 + 2\sqrt{3})r - \frac{13}{3} = 0$ . By the quadratic formula,

$$r = \frac{-(2 + 2\sqrt{3}) \pm \sqrt{(2 + 2\sqrt{3})^2 - 4 \cdot 2 \cdot -\frac{13}{3}}}{4} = \frac{-1 - \sqrt{3} \pm \sqrt{40}}{4}.$$

Since  $r$  is positive,

$$r = \frac{-1 - \sqrt{3} + 2\sqrt{10}}{4}.$$

Then the area of the smaller circle is

$$\pi r^2 = \pi \left( \frac{-1 - \sqrt{3} + 2\sqrt{10}}{4} \right)^2 = \boxed{\frac{27\pi}{8} - \frac{13\sqrt{3}\pi}{8} + \frac{5\sqrt{30}\pi}{4}}.$$

# DETAILS

This question tests both area calculation skills (area of a circle) and algebraic manipulation skills (solving a quadratic equation with radicals). The problem is challenging because it involves nested radicals, but the steps are all fairly standard.

755

### Llama-3-70B-Instruct

Here is a difficult question that tests both area calculation skills and algebraic manipulation skills:

# QUESTION

A square with side length 4 is inscribed in a circle. A semicircle is drawn inside the square, with its diameter being one of the sides of the square. Find the area of the shaded region.

```
[asy] size(100); draw(unitcircle); draw((-2,0)--(2,0)--(2,2)--(-2,2)--cycle); filldraw((-2,0)--(2,0)--(2,2)--(-2,2)--cycle,lightgray); [/asy]
```

# SOLUTION

The area of the shaded region is equal to the area of the square minus the area of the semicircle. The area of the square is  $4^2 = 16$ . The area of the semicircle can be found using the formula  $\frac{1}{2}\pi r^2$ , where  $r$  is the radius of the semicircle. Since the diameter of the semicircle is 4, the radius is 2, so the area of the semicircle is  $\frac{1}{2}\pi(2)^2 = 2\pi$ . Therefore, the

756

area of the shaded region is  $16 - 2\pi$ .

#### # DETAILS

This question requires students to calculate the area of a square and a semicircle, and then subtract one from the other to find the area of the shaded region. It also requires algebraic manipulation skills to simplify the expression for the area of the semicircle. The answer is not a simple numerical value, but rather an expression involving  $\pi$ , which adds an extra layer of complexity to the question.

757

### Llama-3-8B-Instruct

Here is a question that tests both area calculation skills and algebraic manipulation skills:

#### # QUESTION

Find the area of the region bounded by the curve  $y = x^2 - 4x + 3$  and the  $x$ -axis, to the nearest hundredth.

#### # SOLUTION

The area of the region can be found by integrating the function  $y = x^2 - 4x + 3$  with respect to  $x$  from 0 to 2, which is  $\int_0^2 (x^2 - 4x + 3) dx = [x^3/3 - 2x^2 + 3x]_0^2 = (8/3) - 8 + 6 - [0] = 2$ . The area can also be found by using the formula for the area of a region bounded by a curve and the  $x$ -axis, which is  $\int_0^2 y dx = \int_0^2 (x^2 - 4x + 3) dx = 2$ . Therefore, the area is approximately 2 square units, to the nearest hundredth.

#### # DETAILS

This question tests both area calculation skills, as it requires the student to find the area of a region bounded by a curve and the  $x$ -axis, and algebraic manipulation skills, as it requires the student to integrate a function with respect to  $x$ .

Note: I've also reviewed the question and solution to ensure that they are accurate and clear for students.

758

## 759 B.3 Prompt Examples

760 We detail the prompts used in the various stages of question generation and evaluation in the  
761 following sections. We present the skeleton of the prompts with placeholders for skill names, in-  
762 context exemplars, etc. in favor of keeping the prompt brief. The in-context exemplars and the skill  
763 exemplar repository are made available here.

### 764 B.3.1 Skill Pair Validation

#### Prompt for validating a pair of skills

I have two basic mathematical skills described below, each followed by some sample question - solution pairs such that solving the questions involves the use of the particular skill.

SKILL 1: <skill\_1>

Question 1: <skill\_1\_question\_1>

Answer 1: <skill\_1\_solution\_1>

Question 2: <skill\_1\_question\_2>

Answer 2: <skill\_1\_solution\_2>

765

Question 3: <skill\_1\_question\_3>  
Answer 3: <skill\_1\_answer\_3>

SKILL 2: <skill\_2>

Question 1: <skill\_2\_question\_1>  
Answer 1: <skill\_2\_solution\_1>

Question 2: <skill\_2\_question\_2>  
Answer 2: <skill\_2\_solution\_2>

Question 3: <skill\_2\_question\_3>  
Answer 3: <skill\_2\_solution\_3>

I am going to use these two skills for framing a new question such that the question requires an expertise in both the skills in order to be solved, i.e. the question will compose these two skills. However, I do not want the two skills to be very similar, i.e., they should not mean the same thing. Go through the descriptions of the skills carefully. Based on your understanding of the skills, can you please tell me whether the two skills are essentially entirely the same or not? Think step by step and give a detailed explanation of your answer. The answer should begin with a prefix '# EXPLANATION '. Note that your understanding of the skills should not be restricted to the sample questions provided previously. They are just example questions. Use your own prior knowledge as well. End your response with a 'Yes' or 'No' answer to whether the skills are similar or not. This final answer should be on a new line and preceded by the prefix '# FINAL ANSWER '. Thank you very much!

766

### 767 B.3.2 Question Generation

#### Prompt for question generation

I am a math teacher trying to create challenging math questions for smart students. I was wondering if you could give me 1 (non multiple choice) question which tests both the following skills: (<skill\_1>, <skill\_2>) Please also provide a brief solution. Then please look over the question and the solution, and fix any issues so that my students do not get frustrated. This being a math exam, the answers should either be exact, or if not possible, then the question should clearly say the answer is only expected to be approximately correct. Further, for ease of evaluating the students' answers, the question should ask for a single final result. This process is difficult so I am attaching two sample conversations where (Agent) is an AI agent and (Query) is teacher feedback. The conversations revolve around framing such mathematical reasoning questions and using them for evaluating students. These should give you some idea of the expectations and the potential difficulties involved in this task. I am also giving three example question - answer pairs for both <skill\_1> and <skill\_2> skills, such that the example questions test the corresponding skill. Please ensure that the complexity / difficulty of application of <skill\_1> and <skill\_2> skills in the generated question is similar to the complexity / difficulty of the skills in the example questions. Please format your output as

'# QUESTION  
<question>

# SOLUTION  
<solution>

# DETAILS  
<all other text>'

SKILL 1: <skill\_1>

Question 1: <skill\_1\_question\_1>  
Answer 1: <skill\_1\_solution\_1>

Question 2: <skill\_1\_question\_2>  
Answer 2: <skill\_1\_solution\_2>

768

Question 3: <skill\_1\_question\_3>  
Answer 3: <skill\_1\_solution\_1>

SKILL 2: <skill\_2>

Question 1: <skill\_2\_question\_1>  
Answer 1: <skill\_2\_solution\_1>

Question 2: <skill\_2\_question\_2>  
Answer 2: <skill\_2\_solution\_1>

Question 3: <skill\_2\_question\_3>  
Answer 3: <skill\_2\_solution\_3>

# CONVERSATION 1  
<agent\_convo\_1>

# CONVERSATION 2  
<agent\_convo\_2>

769

### 770 B.3.3 Attempted Solution

771 Prompt for solution attempt. Note that we instruct the model to take a defeatist approach towards  
772 solving the question

#### Prompt for solution attempt

You are a professional math teacher and you are given a question which is supposed to test the analytical and mathematical reasoning abilities of your students. You are supposed to provide a solution to the given question. However, the question may be flawed. For example, it might have problems like question being unsolvable using the information provided, question being self-contradictory, the final answer being computationally intractable, the question being ambiguous and confusing, question having multiple possible interpretations, etc., which you may encounter while solving the problem. This question being used for evaluating students in math, the question should ideally have a single, exact answer, with no room for any deviations due to factors such as approximations, rounding errors, etc., unless explicitly specified in the question. Problems such as the ones described above, would prevent the students from solving the question properly, and thus, any question with either of these problems is unfit for testing the students. If you encounter any such problems, stop the solution right there and explain the problems. For example, if you encounter the need to make any approximations or rounding which is not specified in the question, stop solving the question along with the reason. You do not need to solve the question further once you encounter any such problem. If you do not encounter any such problem, solve the question to achieve the single exact answer which the question asks for.

# QUESTION  
<question>

773

### 774 B.3.4 Question Validation

775 Note that how in the first paragraph, the names of the two skills are mentioned even time instead of  
776 using referential phrases. This is done to address the *lost in the middle* problem

#### Prompt for validating the questions

You are a professional math teacher. You want to evaluate the analytical and mathematical reasoning abilities of your students in a math exam. The students are supposed to sit in an examination hall and solve the questions within a given time limit, without access to any computational devices. The evaluation is designed to test the students' expertise in using two given mathematical skills simultaneously, namely <skill\_1> and <skill\_2>. This is achieved by asking them to solve a question that necessitates expertise in both <skill\_1> and <skill\_2> skills, to be solved completely. Since evaluating the students is a critical task allowing very little margin for any error in the process, it is very

777



important to ensure that the questions used for evaluating are high quality and fit for being used to evaluate the students. You need to carefully review the question and a given attempt at solving it, and ensure that the question is of high quality and fit to assess students. In order to do this, you should check the quality of the question with respect to several criteria, such as:

- Single Answer Requirement: The question should ask for one and only one final result. It should not request multiple distinct answers or pieces of information.
- Exact Answer Requirement: It should be possible to achieve one, exact answer to the question, without the need of making any approximations or assumptions whatsoever, unless explicitly specified in the question. There should be no margin for the students to arrive at any other possible answer due to things like rounding errors, etc.
- Dual Skill Requirement: The question must require rigorous expertise in both a) '<skill\_1>' and b) '<skill\_2>', for resolution. Application of both <skill\_1> and <skill\_2> and their subskills should be, necessary and contribute directly to obtaining the final answer; <skill\_1> and <skill\_2> skill should be applicable separately and critically during the problem-solving process. You are also given three example question - answer pairs for both <skill\_1> and <skill\_2> skills in order to help you better understand the meaning of each skill. Please carefully review the question and its attempted solution, paying close attention to how well it aligns with the examples provided for each skill. Consider the depth and breadth of knowledge demonstrated in the examples. The complexity / difficulty of application of both <skill\_1> and <skill\_2> in the question should be similar or greater than the complexity / difficulty of <skill\_1> and <skill\_2> in the example question-answers given for that respective skill.
- Clarity and Completeness: The question should be unambiguous and contain all the information necessary to complete the solution. Any required assumptions not common knowledge should be explicitly stated. Check for any ambiguity that might confuse students. Carefully go through the solution to check if it makes any assumption or approximation in order to solve the question.
- Computational Tractability: Since the students are supposed to solve the questions within a given time limit and without access to any computational devices such as calculators, computer, mobile phones, etc., you must ensure that the question is computationally tractable and all the computations involved can be done by hand in a limited amount of time.
- Relevancy of Information: The question should not have any extra details that do not contribute to the solving of the problem.
- Realism and Logic: The question should involve realistic scenarios or hypotheses with logically consistent data. The specified operations and the contextual setup should reflect plausible mathematical situations. (e.g., positive amounts for transactions, integers for counts).
- Syntax and Grammar: The question must be grammatically correct and clearly written to prevent misinterpretation.
- etc. (any other problems which you think make the question not fit for being used for evaluating the students)

Your task is to give a 'Yes' or 'No' assessment, indicating whether the question is high quality and suitable for evaluating the students on simultaneous application of the skills <skill\_1> and <skill\_2>. Provide thorough reasoning for your assessment based on the conditions mentioned above and any other relevant analytical points concerning mathematical reasoning and problem-solving. Your response should be structured as follows:

# REASONING

<Your detailed analysis justifying your decision>

# FINAL ANSWER

<'Yes' or 'No'. No other text should be present in this section>

Ensure to review the combination of skills intended for assessment, and check the logical flow and mathematical correctness from the question's setup to the solution's conclusion. Look out for any problems in the question which are pointed out in the attempted solution. Account for all the potential pitfalls such as logical inconsistencies, unnecessary complexity, or insufficient detail that may obstruct the clarity or solvability of the question. Given below are the two skills and some example question-answer pairs for the two skills. This process is difficult so I am attaching a few sample conversations where (agent) is an AI agent who is trying to verify the questions and (query) is teacher feedback. This should give you some idea of potential difficulties in this task. This is followed by the question which you need to check (preceded by '# QUESTION TO BE CHECKED') and its attempted solution (preceded by '# SOLUTION ATTEMPT').

SKILL 1: <skill\_1>

Question 1: <skill\_1\_question\_1>  
Answer 1: <skill\_1\_solution\_1>

Question 2: <skill\_1\_question\_2>  
Answer 2: <skill\_1\_solution\_2>

Question 3: <skill\_1\_question\_3>  
Answer 3: <skill\_1\_solution\_3>

SKILL 2: <skill\_2>

Question 1: <skill\_2\_question\_1>  
Answer 1: <skill\_2\_solution\_1>

Question 2: <skill\_2\_question\_2>  
Answer 2: <skill\_2\_solution\_2>

Question 3: <skill\_2\_question\_3>  
Answer 3: <skill\_2\_solution\_3>

# CONVERSATION 1  
<validation\_exemplar\_1>

# CONVERSATION 2  
<validation\_exemplar\_2>

.....

# CONVERSATION 6  
<validation\_exemplar\_6>

# QUESTION TO BE CHECKED  
<question>

# SOLUTION ATTEMPT  
<solution>

Thank you very much!

779

### 780 **B.3.5 Final Solution**

781 For the final solution, we make use in-context exemplars from MATH [Hendrycks et al., 2021] as  
782 opposed to the attempted solution step.

#### Prompt for the final solution

I have two basic mathematical skills described below, each followed by some sample question - solution pairs such that solving the questions involves the use of the particular skill in order to be solved.

SKILL 1: <skill\_1>

Question 1: <skill\_1\_question\_1>  
Answer 1: <skill\_1\_solution\_1>

Question 2: <skill\_1\_question\_2>  
Answer 2: <skill\_1\_solution\_2>

Question 3: <skill\_1\_question\_3>  
Answer 3: <skill\_1\_solution\_3>

SKILL 2: <skill\_2>

783

Question 1: <skill\_2\_question\_1>  
Answer 1: <skill\_2\_solution\_1>

Question 2: <skill\_2\_question\_2>  
Answer 2: <skill\_2\_solution\_2>

Question 3: <skill\_2\_question\_3>  
Answer 3: <skill\_2\_solution\_3>

Go through the descriptions of the skills carefully. Now, here is a new question such that the question requires an expertise all both the skills in order to be solved. That is, the question composes these two skills

QUESTION: <question>

Based on your understanding of the skills, can you please solve the question accurately? Think step by step and explain the solution. Finally, end your response by stating the final numerical answer obtained using the solution. Note that your understanding of the skills should not be restricted to the sample questions provided in their description. They are just example questions. Use your own prior knowledge as well. The explanation of your solution and the final numerical answer should each be on a new line, and should be preceded by the prefixes '# SOLUTION ' and '# ANSWER ' respectively. Thus, your response should be in the format:

'# SOLUTION  
<solution>

# ANSWER  
<final\_answer; no other text should be present in this section>'

Thank you very much!

784

### 785 B.3.6 Evaluation

#### Prompt given to the GPT-4 for evaluating the model's solution

You are a professional math teacher and are tasked with evaluating your students on a math exam. You will be given a question, the correct solution to the question and the student's solution. You need to tell me whether the student solved the question correctly, thus matching the answer obtained by the correct solution. Think step-by-step and give a detailed explanation of your answer. At the end, give a 'Yes' or 'No' answer to whether the student's solution is correct. Your output should be in the following format:

# STEP BY STEP EXPLANATION  
<detailed explanation of your thought process>

# CORRECTNESS  
<'Yes' if the student's solution is correct. 'No' otherwise. This section should not contain any other text>

Here are the question, correct solution to the question and the student's solution:

QUESTION: <question>

CORRECT SOLUTION: <correct\_solution>

STUDENT'S SOLUTION: <student's\_solution>

786