
Beyond Memorization: Reasoning-Driven Synthesis as a Mitigation Strategy Against Benchmark Contamination

Anonymous Author(s)

Affiliation

Address

email

Abstract

1 Capability evaluation of large language models (LLMs) is increasingly shadowed
2 by rising concerns of data contamination that cast doubts on whether static bench-
3 marks measure genuine reasoning or mere memorization. We present an empirical
4 study using an infinitely scalable framework to synthesize research-level QA di-
5 rectly from arXiv papers, harnessing the natural temporal structure of research
6 publications where performance decay after knowledge cutoffs may indicate po-
7 tential contamination. We evaluated 4 frontier model represented by 2 models of
8 different knowledge cutoff dates per family on 1,643 multi-step reasoning questions
9 synthesized from 20,277 arXiv papers stratified over 26 months, covering at least 6
10 months before and after all cutoff dates. Our results consistently showed a lack of
11 significant performance decay near knowledge cutoff dates for models of various
12 sizes, developers, and release dates. We further performed a comparative analysis
13 with previous longitudinal studies that reported significant post-cutoff performance
14 decay using directly retrieved questions based on public data. we hypothesize
15 that the multi-step reasoning required by our synthesis pipeline offered additional
16 complexity that goes deeper than shallow memorization, which effectively serves a
17 mitigation strategy against benchmark contamination. We fully open source our
18 code and dataset to aid reproducibility and advocate for a paradigm shift that prior-
19 itize reasoning-driven synthesis to construct benchmarks over simply collecting
20 newly released questions periodically.

21 1 Introduction

22 *If you cannot measure it, you cannot improve it.*

23 —Lord Kelvin Thomson [1889]

24 The reliability of large language model (LLM) evaluations faces a critical threat from data con-
25 tamination, which could artificially inflate performance metrics while obscuring genuine reasoning
26 capabilities Dong et al. [2024]. Many previous work Ding et al. [2024], Li et al. [2025], Xu et al.
27 [2024] have demonstrated that contamination has reached concerning proportions across various
28 evaluation benchmarks, where models achieve inflated scores through memorized patterns rather than
29 authentic problem-solving abilities. This widespread contamination undermines the fundamental
30 premise of capability evaluation by rewarding benchmarking memorization over genuine reasoning
31 abilities.

32 The rapid saturation of established benchmarks further compounds this challenge, creating an
33 evaluation crisis where traditional metrics no longer differentiate between reasoning advances and
34 pattern matching. Legacy frameworks such as MATH Hendrycks et al. [2021] and GPQA Rein

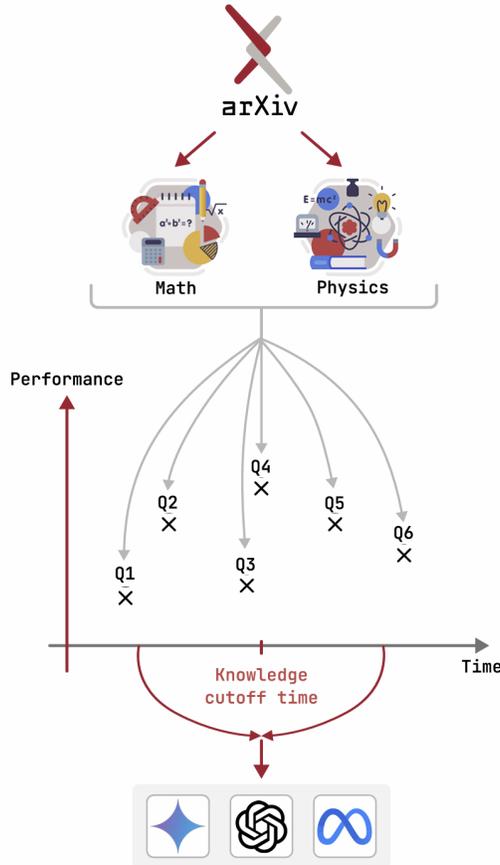


Figure 1: Overview of our framework for synthesizing research-level evaluation questions based on arXiv papers, which are temporally stratified to monthly window to assess frontier model reasoning capabilities before and after their knowledge cutoff dates. We hypothesize that potential contamination will lead to significant post-cutoff performance decay.

35 et al. [2023] have experienced dramatic performance saturation across model generations. This
 36 phenomenon reflects the broader inadequacy of static evaluation paradigms in measuring rapidly
 37 evolving model capabilities.

38 Existing mitigation strategies bear fundamental limitations in scalability and effectiveness. Manual
 39 curation approaches, exemplified by FrontierMath Glazer et al. [2024], demand extensive expert
 40 effort for problem construction and validation, creating practical barriers to frequent updates and
 41 systematic coverage. Periodic collection methods such as LiveBench White et al. [2025] promised to
 42 complete regular content updates but suffer from short shelf life and resource-intensive maintenance
 43 requirements on long time horizon. These approaches fail to address the core challenge of creating
 44 evaluation frameworks that scale automatically with advancing model capabilities while maintaining
 45 contamination resistance.

46 To address these challenges, we leverage the natural temporal structure of research publication
 47 to create a contamination-resistant evaluation framework that scales automatically with scientific
 48 progress. Academic publications follow strict chronological ordering, making it impossible for
 49 models to encounter evaluation content synthesized from papers published after their training cutoff
 50 dates. This temporal design enables systematic contamination detection through performance decay
 51 analysis, where post-cutoff performance decay would indicate more reliance on memorization rather
 52 than genuine reasoning capabilities.

53 Contributions

- 54 • We leverage a fully automated, infinitely scalable framework to synthesize 1,643 multi-step
55 reasoning questions from 20,277 arXiv math and physics papers covering 26 months with
56 comprehensive temporal coverage spanning at least 6 months before and after all model
57 knowledge cutoffs.
 - 58 • We comprehensively evaluated all 8 frontier models from 4 leading developers with fine-
59 grained monthly temporal granularity across math and physics domain to improve robustness
60 of our findings across various models and domains.
 - 61 • We systematically compare longitudinal analysis based on retrieval vs. synthesis-based
62 benchmarks to empirically show that reasoning-driven synthesis offers stronger contamina-
63 tion resistance compared to questions directly harvested from public sources.
- 64 Our approach addresses the fundamental trade-off between evaluation scalability and integrity that
65 shadows current benchmarking paradigm. We further advocate for multi-step reasoning-driven
66 transformation that creates increased cognitive distance against simple pattern matching, thereby
67 resist contamination by design.

68 2 Related Work

69 **Research-Level Reasoning Benchmarks** GPQA Rein et al. [2023] present PhD-level challenges
70 across biology, physics, and chemistry, but its static nature makes it vulnerable to memorization
71 as training corpora expand. CURIE Anonymous [2025] advances multitask scientific evaluation
72 across specialized domains requiring expert knowledge synthesis, yet lacks temporal stratification
73 mechanisms for contamination detection. While HLE Jain et al. [2025] and FrontierMath Glazer
74 et al. [2024] achieve research-level difficulty through expert curation, they suffer from scalability
75 bottlenecks and lack systematic contamination detection that our automated synthesis pipeline
76 addresses through temporal stratification.

77 **Data Synthesis from Research Papers** Scientific paper-based evaluation represents promising
78 directions for scalable benchmark construction. LiveXiv Shabtay et al. [2025] pioneered temporal
79 evaluation through chronologically ordered arXiv papers, establishing contamination detection princi-
80 ples, but focuses on factual extraction rather than complex reasoning chains. SPIQA Anonymous
81 [2024] demonstrates comprehensive multimodal evaluation from scientific papers, yet emphasizes
82 comprehension over multi-step mathematical reasoning that creates cognitive barriers against memo-
83 rization. OpenScholar Ajith et al. [2024], DIGESTables Skarlinski et al. [2024], and SciQAG Wang
84 et al. [2024] advance cross-domain synthesis with sophisticated metrics but lack systematic temporal
85 stratification for contamination detection. PaperQA Lala et al. [2023] combines question answering
86 with document retrieval effectively, yet introduces contamination vectors through retrieval-based
87 approaches that our synthesis methodology avoids.

88 **Longitudinal Analysis on Data Contamination** Systematic contamination detection reveals
89 widespread evaluation compromise, establishing temporal analysis as essential for benchmark in-
90 tegrity, yet existing methodologies focus on retrieval-based approaches vulnerable to memorization.
91 Comprehensive surveys Ding et al. [2024], Li et al. [2025], Xu et al. [2024] document pervasive
92 benchmark compromise while ContaminationCDD Fan et al. [2024] provides computational detection
93 through output distribution analysis. Temporal boundary research Roberts et al. [2024] demonstrates
94 effectiveness in Mathematical Olympiad contexts with performance decay on post-cutoff problems,
95 yet relies on publicly available competition problems vulnerable to training data inclusion. Pro-
96 gramming applications through LiveCodeBench Jain et al. [2024] and DyCodeEval Chen et al.
97 [2025] demonstrate performance drops on post-training data but focus on code generation rather than
98 mathematical reasoning synthesis. Our contribution uniquely combines monthly temporal resolution
99 with reasoning-driven synthesis that generates multi-step problems requiring genuine mathematical
100 understanding rather than memorized solution patterns.

Mathematics
arXiv:2406.19979v2 - On quantitative convergence for stochastic processes: Crossings, fluctuations and martingales
Main Field: math.PR (Probability)
<p>Theorem: Given $\varepsilon, K > 0$ and $g : \mathbb{N} \rightarrow \mathbb{N}$ there exists some N dependent only on these parameters such that whenever $\{x_n\}$ is a monotone sequence in $[-K, K]$, there exists some $n \leq N$ such that $x_i - x_j < \varepsilon$ for all $n \leq i \leq j \leq n + g(n)$. Moreover, we can assign to N the following concrete value: $N := \tilde{g}^{(\lceil 2K/\varepsilon \rceil)}(0)$ for $\tilde{g}(n) := n + g(n)$, where $\tilde{g}^{(i)}$ denotes the ith iteration of \tilde{g}.</p>
<p>Q: Given parameters $\varepsilon > 0, K > 0$, and a function $g : \mathbb{N} \rightarrow \mathbb{N}$, what explicit bound N in terms of ε, K, and g guarantees that for any monotone sequence (x_n) in $[-K, K]$ there exists some $n \leq N$ such that $x_i - x_j < \varepsilon$ for all $n \leq i \leq j \leq n + g(n)$?</p>
Physics
arXiv:2407.02415v2 - The Symplectic Schur Process
Main Field: math-ph (Mathematical Physics)
<p>Theorem: Consider the rescaling</p> $i(\tau) = \left\lfloor \frac{9n}{8} + \sqrt{\frac{27n}{64}} \tau \right\rfloor, \quad u(\alpha) = -n + \left\lfloor \left(\frac{n}{12}\right)^{\frac{1}{4}} \alpha \right\rfloor.$ <p>For any fixed $k \in \mathbb{Z}_{\geq 1}$ and $\tau_1, \dots, \tau_k \in \mathbb{R}$ and $\alpha_1, \dots, \alpha_k \in \mathbb{R}_+$, let $i_\ell = i(\tau_\ell), u_\ell = u(\alpha_\ell)$, then</p> $\lim_{n \rightarrow \infty} \det_{1 \leq \ell, \ell' \leq k} \left[\left(\frac{n}{12}\right)^{\frac{1}{4}} \left(\delta_{i_\ell, i_{\ell'}} \delta_{u_\ell, u_{\ell'}} - K^{\text{SSP}}(i_\ell, u_\ell; i_{\ell'}, u_{\ell'}) \right) \right] = \det_{1 \leq \ell, \ell' \leq k} [\mathcal{K}(\tau_\ell, \alpha_\ell; \tau_{\ell'}, \alpha_{\ell'})],$ <p>where \mathcal{K} is the Pearcey-like kernel.</p>
<p>Q: Let $i(\tau) = \left\lfloor \frac{9n}{8} + \sqrt{\frac{27n}{64}} \tau \right\rfloor$ and $u(\alpha) = -n + \left\lfloor \left(\frac{n}{12}\right)^{1/4} \alpha \right\rfloor$. For a fixed $k \in \mathbb{Z}_{\geq 1}$ and real parameters τ_1, \dots, τ_k and $\alpha_1, \dots, \alpha_k$, define $i_\ell = i(\tau_\ell)$ and $u_\ell = u(\alpha_\ell)$. What is the limit as $n \rightarrow \infty$ of</p> $\det_{1 \leq \ell, \ell' \leq k} \left[\left(\frac{n}{12}\right)^{1/4} \left(\delta_{i_\ell, i_{\ell'}} \delta_{u_\ell, u_{\ell'}} - K^{\text{SSP}}(i_\ell, u_\ell; i_{\ell'}, u_{\ell'}) \right) \right]?$

Table 1: Examples of theorem-to-question conversion with arXiv IDs, paper titles, and domain classification tags with explanations.

101 3 Methodology

102 3.1 Synthesis Pipeline

103 **Retrieval and Theorem Identification** This synthesis pipeline is adapted based on the RealMath
104 Framework Zhang et al. [2025]. We retrieve arXiv papers directly from the arXiv API to obtain
105 comprehensive metadata including LaTeX source code, DOI, publication date, and arXiv domain
106 tags. Our synthesis pipeline utilizes OpenAI-GPT-4.1 and OpenAI-o4-mini OpenAI [2025] to
107 systematically parse LaTeX source and identify constructive theorems with fixed answers, excluding
108 theorems with multiple solutions or ambiguity to ensure automated verification while maintaining
109 mathematical rigor.

110 **QA Generation** We employ sophisticated multi-stage generation focusing on compositional multi-
111 step reasoning where solvers must synthesize multiple mathematical concepts, apply sequential
112 logical transformations, and maintain coherent reasoning across extended derivation sequences. Each
113 question undergoes automated analysis to verify solution paths requiring multiple interdependent
114 reasoning steps including concept application, algebraic manipulation, logical inference, and result
115 synthesis. We systematically filter problems requiring fewer than six distinct reasoning steps, ensuring
116 substantial cognitive effort that cannot be solved through pattern recognition.

117 **Filtering and Validation** Our filtering process includes comprehensive post-processing review
 118 of each generated question-answer pair, systematically removing low-quality samples with easily
 119 guessable answers or obvious solutions from context. We extend methodology beyond mathematics
 120 to Physics, creating problems requiring interdisciplinary reasoning combining mathematical formal-
 121 ism with physical intuition, dimensional analysis, and connections between theoretical predictions
 122 and measurable quantities. This multi-step reasoning validation ensures our benchmark captures
 123 sophisticated abilities distinguishing advanced models from pattern-matching approaches.

124 3.2 Dataset Construction

125 **Temporal Stratification** We synthesized questions with monthly resolution from May 2023 to
 126 June 2025, ensuring coverage of at least 6 months before and after models’ knowledge cutoffs.
 127 This 26-month window enables precise contamination detection across temporal boundaries while
 128 leveraging the natural chronological ordering of academic publications that prevents models from
 129 encountering evaluation content during training.

130 **Cross-Domain Coverage** Our dataset balances Mathematics and Physics domains, generating
 131 20 questions per month per domain. We focus on combinatorics and number theory domains
 132 providing suitable sources for constructive theorems with fixed answers. This cross-domain approach
 133 evaluates both domain-specific knowledge and general reasoning transfer capabilities across scientific
 134 disciplines.

135 **Quality Assurance** We performed manual expert inspection for technical accuracy and LaTeX
 136 formatting as quality filters, obtaining 1,643 questions equally spanning math and physics domains.
 137 This curation process ensures evaluation integrity while maintaining the scalability advantages of
 138 automated synthesis over manually curated benchmarks.

139 4 Experimental

140 We evaluated models from four distinct families (OpenAI OpenAI [2025], Gemini Deepmind [2025],
 141 Llama Grattafiori et al. [2024], and DeepSeek Guo et al. [2025]) where each family is represented
 142 by 2 models with different knowledge cutoff dates to assess the consistency of our findings. We
 143 summarize the models and their respective cutoff dates in Table 2. We evaluated all models by calling
 144 OpenRouter API with default settings and enabled thinking modes whenever possible.

Table 2: All evaluated models and their knowledge cutoff dates Hao00Wang [2025]

Model	Knowledge Cutoff
DeepSeek-R1-0528	2024.07
DeepSeek-R1	2023.10
OpenAI-o4-mini	2024.06
OpenAI-o3-mini	2023.10
Gemini-2.5-Flash	2025.01
Gemini-2.0-Flash	2024.08
Llama-4-Scout	2024.08
Llama-3.3-70B	2023.12

145 Figure 2 illustrates temporal performance across 26 months in both Mathematics and Physics
 146 domains. The monthly trends show no significant performance degradation around knowledge cutoff
 147 dates marked by dashed vertical lines. Rather than expected decay due to potential training data
 148 contamination, models demonstrate stable or improved performance on post-cutoff samples across
 149 both domains. This pattern suggests our multi-step reasoning synthesis pipeline creates contamination
 150 resistance by requiring genuine problem-solving rather than memorization-based pattern matching.

151 Figure 3 presents mean accuracy comparison across six months before and after knowledge cutoff
 152 dates for both mathematics and physics domains. Results demonstrate consistent performance
 153 patterns without significant degradation following cutoff boundaries, confirming temporal stability in
 154 synthesized evaluation content.

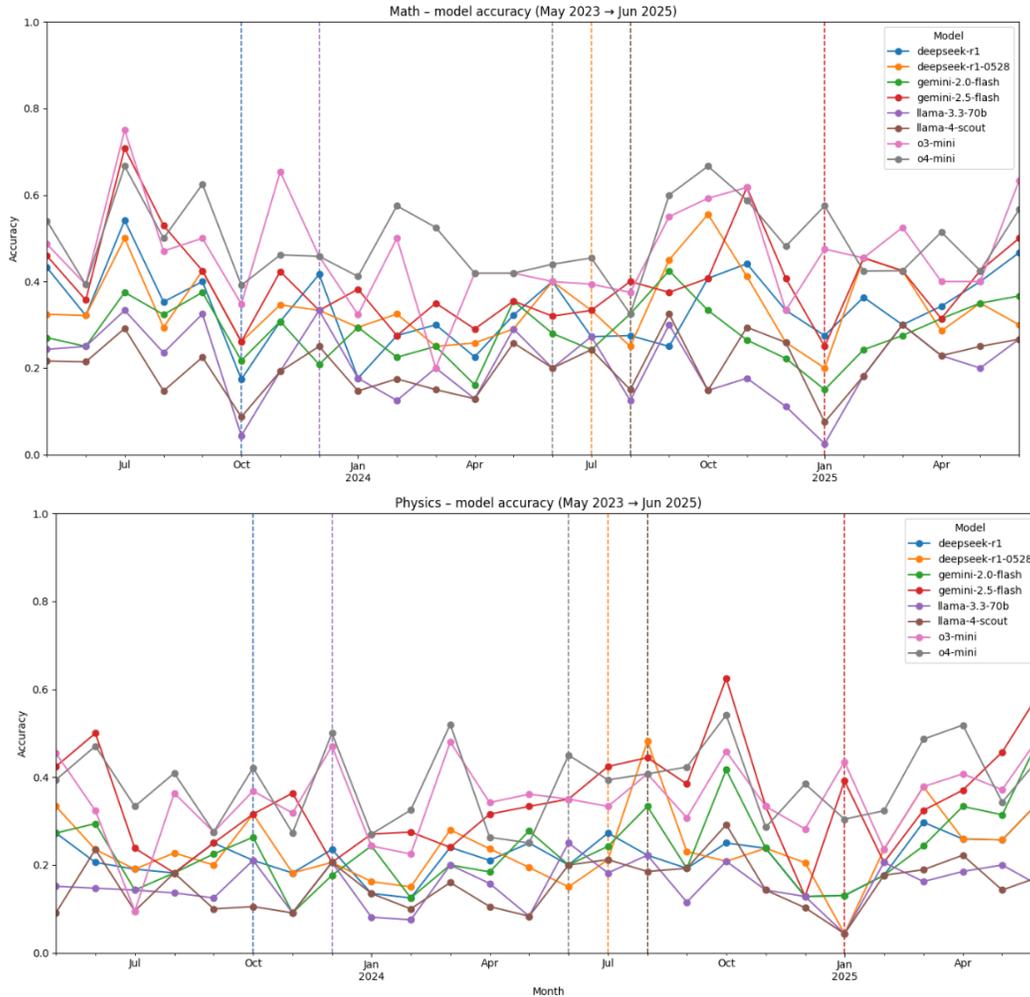


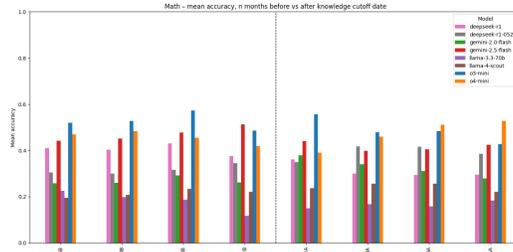
Figure 2: Model performance scores across 26 months in Math and Physics domains, with knowledge cutoff dates marked by dashed lines matching model legends. The figure demonstrates temporal performance stability without significant decay patterns around cutoff boundaries, suggesting contamination resistance through multi-step reasoning synthesis rather than performance degradation expected from contaminated benchmarks.

155 Figure 4 shows aggregated performance scores across temporal windows, further confirming stability
 156 of model capabilities on synthesized questions regardless of publication timing relative to training
 157 cutoffs. This consistency across all evaluated model families provides systematic evidence that
 158 reasoning-driven synthesis creates contamination resistance through cognitive complexity barriers.

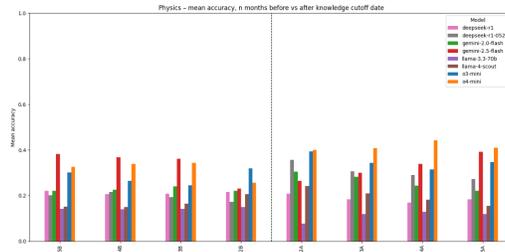
159 5 Discussion

160 We present a comparative analysis between retrieval-based and synthesis-based benchmarking ap-
 161 proaches through their temporal contamination patterns, providing empirical evidence that reason-
 162 ing-driven synthesis creates cognitive barriers against shallow memorization.

163 **Contamination Vulnerability in Retrieval-Based Approaches** Retrieval-based benchmarks con-
 164 sistently exhibit post-cutoff performance decay that reveals widespread contamination effects across
 165 multiple domains. Mathematical Olympiad evaluations demonstrate severe performance deterioration
 166 on problems published after training cutoffs, where models achieved artificially inflated scores on
 167 pre-cutoff content through memorized solution patterns from competition archives Roberts et al.



(a) Mean accuracy for mathematics-related QA pairs.



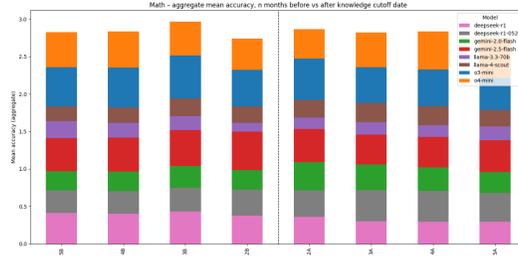
(b) Mean accuracy for physics-related QA pairs.

Figure 3: Mean accuracy for mathematics and physics QA pairs across six months before and after knowledge cutoff dates. Both domains show consistent performance without significant post-cutoff degradation, supporting our hypothesis that synthesis-based benchmarks resist contamination effects through cognitive complexity requirements that demand genuine reasoning rather than memorized pattern matching.

168 [2024]. Programming benchmarks reveal similarly concerning patterns, with coding problems show-
 169 ing sharp performance drops after knowledge boundaries while exhibiting suspicious correlation
 170 between problem popularity and model accuracy Huang et al. [2024]. LiveCodeBench documented
 171 systematic contamination through clear temporal performance boundaries in code generation tasks,
 172 demonstrating that models relied heavily on memorized solutions from competitive programming
 173 platforms and GitHub repositories rather than genuine algorithmic reasoning Jain et al. [2024]. These
 174 findings establish a consistent pattern where publicly sourced evaluation content becomes compro-
 175 mised through training data inclusion, creating evaluation frameworks that reward memorization over
 176 authentic problem-solving capabilities.

177 **Performance Stability in Synthesis-Based Approaches** In contrast, synthesis-based benchmarks
 178 performs transformation based on retrieved contents to vary the contexts and complexity of evaluation
 179 questions in question Zhu et al. [2024]. We hypothesize that the lack of post-cutoff performance
 180 decay may be attributed to the added transformation serving as an extra layer of cognitive distance
 181 that cognitively goes deeper than shallow memorization-based shortcuts. AntiLeakBench Wu et al.
 182 [2025] also highlights that data contamination before LLM cutoff dates can be effectively overcome
 183 by constructing samples with explicitly new knowledge after training cutoff instead of simply collect
 184 new data periodically. Our experimental results align well with Zhang et al. [2025] across diverse
 185 model families and extended temporal windows of 26 months, showcasing how our synthesis pipeline
 186 established contamination resistance through reasoning demands that cannot be satisfied through
 187 memorized associations alone.

188 **Reasoning-Driven Synthesis as Mitigation Strategy** We hypothesize the reason why there is a
 189 lack of significant post-cutoff performance decay in our results is because the transformation made
 190 in our synthesis pipeline goes deeper than shallow pattern-matching that enables memorization,
 191 creating effective contamination resistance through cognitive complexity requirements. Composi-
 192 tional reasoning tasks engage working memory integration, sequential logical processing, and
 193 abstract concept manipulation that operate through fundamentally different neural pathways than
 194 simple associative recall mechanisms established in cognitive science research Qiu et al. [2025].
 195 The Extended Edit Distance metric provides quantitative evidence that complex reasoning problems



221 References

- 222 Akari Ajith, Nathan Peng, Ronald Weston, Neel Elhage, Deep Ganguli, Danny Hernandez, Tom
223 Jones, Amanda Lovitt, Nova DasSarma, Peter Hase, et al. Openscholar: Synthesizing scientific
224 literature with retrieval-augmented lms. *arXiv preprint arXiv:2411.14199*, 2024. URL <https://arxiv.org/abs/2411.14199>.
225
- 226 Anonymous. Spiqqa: A dataset for multimodal question answering on scientific papers, 2024. URL
227 <https://neurips.cc/virtual/2024/events/datasets-benchmarks>. 270K multimodal
228 questions from scientific papers.
- 229 Anonymous. Curie: Evaluating llms on multitask scientific long context understanding and reasoning,
230 2025. URL <https://arxiv.org/abs/2503.13517>. To appear at ICLR 2025.
- 231 Simin Chen, Pranav Pusarla, and Baishakhi Ray. Dynamic benchmarking of reasoning capabilities in
232 code large language models under data contamination, 2025. URL [https://arxiv.org/abs/](https://arxiv.org/abs/2503.04149)
233 [2503.04149](https://arxiv.org/abs/2503.04149).
- 234 Deepmind. Gemini flash system card, 2025. URL [https://ai.google.dev/gemini-api/docs/](https://ai.google.dev/gemini-api/docs/models#gemini-2.5-flash)
235 [models#gemini-2.5-flash](https://ai.google.dev/gemini-api/docs/models#gemini-2.5-flash).
- 236 Cheng Ding, Jingfeng Zhang, and Shervin Malmasi. Benchmark data contamination of large language
237 models: A survey. *arXiv preprint arXiv:2406.04244*, 2024. URL [https://arxiv.org/abs/](https://arxiv.org/abs/2406.04244)
238 [2406.04244](https://arxiv.org/abs/2406.04244).
- 239 Yihong Dong, Xue Jiang, Huanyu Liu, Zhi Jin, and Ge Li. Generalization or memorization: Data
240 contamination and trustworthy evaluation for large language models. *ArXiv*, abs/2402.15938, 2024.
241 URL <https://arxiv.org/pdf/2402.15938.pdf>.
- 242 Yihang Fan, Yujun Hou, Jie Zhang, Qin Liu, Xuming Hu, and Xiaolin Wang. Generalization or
243 memorization: Data contamination and trustworthy evaluation for large language models. *arXiv*
244 *preprint arXiv:2402.15938*, 2024. URL <https://arxiv.org/abs/2402.15938>.
- 245 Elliot Glazer, Ege Erdil, Tamay Besiroglu, Diego Chicharro, Evan Chen, Alex Gunning, Caroline Falk-
246 man Olsson, Jean-Stanislas Denain, Anson Ho, Emily de Oliveira Santos, et al. Frontiermath: A
247 benchmark for evaluating advanced mathematical reasoning in ai. *arXiv preprint arXiv:2411.04872*,
248 2024. URL <https://arxiv.org/abs/2411.04872>.
- 249 Aaron Grattafiori, Abhimanyu Dubey, Abhinav Jauhri, Abhinav Pandey, Abhishek Kadian, Ahmad
250 Al-Dahle, Aiesha Letman, Akhil Mathur, Alan Schelten, Alex Vaughan, Amy Yang, Angela Fan,
251 Anirudh Goyal, Anthony Hartshorn, Aobo Yang, Archi Mitra, Archie Sravankumar, Artem Korenev,
252 Arthur Hinsvark, Arun Rao, Aston Zhang, Aurelien Rodriguez, Austen Gregerson, Ava Spataru,
253 Baptiste Roziere, Bethany Biron, Binh Tang, Bobbie Chern, Charlotte Caucheteux, Chaya Nayak,
254 Chloe Bi, Chris Marra, Chris McConnell, Christian Keller, Christophe Touret, Chunyang Wu,
255 Corinne Wong, Cristian Canton Ferrer, Cyrus Nikolaidis, Damien Allonsius, Daniel Song, Danielle
256 Pintz, Danny Livshits, Danny Wyatt, David Esiobu, Dhruv Choudhary, Dhruv Mahajan, Diego
257 Garcia-Olano, Diego Perino, Dieuwke Hupkes, Egor Lakomkin, Ehab AlBadawy, Elina Lobanova,
258 Emily Dinan, Eric Michael Smith, Filip Radenovic, Francisco Guzmán, Frank Zhang, Gabriel
259 Synnaeve, Gabrielle Lee, Georgia Lewis Anderson, Govind Thattai, Graeme Nail, Gregoire Mialon,
260 Guan Pang, Guillem Cucurell, Hailey Nguyen, Hannah Korevaar, Hu Xu, Hugo Touvron, Iliyan
261 Zarov, Imanol Arrieta Ibarra, Isabel Kloumann, Ishan Misra, Ivan Evtimov, Jack Zhang, Jade Copet,
262 Jaewon Lee, Jan Geffert, Jana Vranes, Jason Park, Jay Mahadeokar, Jeet Shah, Jelmer van der Linde,
263 Jennifer Billock, Jenny Hong, Jenya Lee, Jeremy Fu, Jianfeng Chi, Jianyu Huang, Jiawen Liu, Jie
264 Wang, Jiecao Yu, Joanna Bitton, Joe Spisak, Jongsoo Park, Joseph Rocca, Joshua Johnstun, Joshua
265 Saxe, Junteng Jia, Kalyan Vasuden Alwala, Karthik Prasad, Kartikeya Upasani, Kate Plawiak,
266 Ke Li, Kenneth Heafield, Kevin Stone, Khalid El-Arini, Krithika Iyer, Kshitiz Malik, Kuenley
267 Chiu, Kunal Bhalla, Kushal Lakhotia, Lauren Rantala-Yearly, Laurens van der Maaten, Lawrence
268 Chen, Liang Tan, Liz Jenkins, Louis Martin, Lovish Madaan, Lubo Malo, Lukas Blecher, Lukas
269 Landzaat, Luke de Oliveira, Madeline Muzzi, Mahesh Pasupuleti, Mannat Singh, Manohar Paluri,
270 Marcin Kardas, Maria Tsimpoukelli, Mathew Oldham, Mathieu Rita, Maya Pavlova, Melanie
271 Kambadur, Mike Lewis, Min Si, Mitesh Kumar Singh, Mona Hassan, Naman Goyal, Narjes
272 Torabi, Nikolay Bashlykov, Nikolay Bogoychev, Niladri Chatterji, Ning Zhang, Olivier Duchenne,

273 Onur Çelebi, Patrick Alrassy, Pengchuan Zhang, Pengwei Li, Petar Vasic, Peter Weng, Prajjwal
274 Bhargava, Pratik Dubal, Praveen Krishnan, Punit Singh Koura, Puxin Xu, Qing He, Qingxiao Dong,
275 Ragavan Srinivasan, Raj Ganapathy, Ramon Calderer, Ricardo Silveira Cabral, Robert Stojnic,
276 Roberta Raileanu, Rohan Maheswari, Rohit Girdhar, Rohit Patel, Romain Sauvestre, Ronnie
277 Polidoro, Roshan Sumbaly, Ross Taylor, Ruan Silva, Rui Hou, Rui Wang, Saghar Hosseini, Sahana
278 Chennabasappa, Sanjay Singh, Sean Bell, Seohyun Sonia Kim, Sergey Edunov, Shaoliang Nie,
279 Sharan Narang, Sharath Rapparthi, Sheng Shen, Shengye Wan, Shruti Bhosale, Shun Zhang, Simon
280 Vandenhende, Soumya Batra, Spencer Whitman, Sten Sootla, Stephane Collot, Suchin Gururangan,
281 Sydney Borodinsky, Tamar Herman, Tara Fowler, Tarek Sheasha, Thomas Georgiou, Thomas
282 Scialom, Tobias Speckbacher, Todor Mihaylov, Tong Xiao, Ujjwal Karn, Vedanuj Goswami,
283 Vibhor Gupta, Vignesh Ramanathan, Viktor Kerkez, Vincent Gonguet, Virginie Do, Vish Vogeti,
284 Vítor Albiero, Vladan Petrovic, Weiwei Chu, Wenhan Xiong, Wenyin Fu, Whitney Meers, Xavier
285 Martinet, Xiaodong Wang, Xiaofang Wang, Xiaoqing Ellen Tan, Xide Xia, Xinfeng Xie, Xuchao
286 Jia, Xuwei Wang, Yaelle Goldschlag, Yashesh Gaur, Yasmine Babaei, Yi Wen, Yiwen Song,
287 Yuchen Zhang, Yue Li, Yuning Mao, Zacharie Delpierre Coudert, Zheng Yan, Zhengxing Chen, Zoe
288 Papakipos, Aaditya Singh, Aayushi Srivastava, Abha Jain, Adam Kelsey, Adam Shajnfeld, Adithya
289 Gangidi, Adolfo Victoria, Ahuva Goldstand, Ajay Menon, Ajay Sharma, Alex Boesenberg, Alexei
290 Baevski, Allie Feinstein, Amanda Kallet, Amit Sangani, Amos Teo, Anam Yunus, Andrei Lupu,
291 Ram Alvarado, Andrew Caples, Andrew Gu, Andrew Ho, Andrew Poulton, Andrew Ryan, Ankit
292 Ramchandani, Annie Dong, Annie Franco, Anuj Goyal, Aparajita Saraf, Arkabandhu Chowdhury,
293 Ashley Gabriel, Ashwin Bharambe, Assaf Eisenman, Azadeh Yazdan, Beau James, Ben Maurer,
294 Benjamin Leonhardi, Bernie Huang, Beth Loyd, Beto De Paola, Bhargavi Paranjape, Bing Liu,
295 Bo Wu, Boyu Ni, Braden Hancock, Bram Wasti, Brandon Spence, Brani Stojkovic, Brian Gamido,
296 Britt Montalvo, Carl Parker, Carly Burton, Catalina Mejia, Ce Liu, Changhan Wang, Changkyu
297 Kim, Chao Zhou, Chester Hu, Ching-Hsiang Chu, Chris Cai, Chris Tindal, Christoph Feichtenhofer,
298 Cynthia Gao, Damon Civin, Dana Beaty, Daniel Kreymer, Daniel Li, David Adkins, David Xu,
299 Davide Testuggine, Delia David, Devi Parikh, Diana Liskovich, Didem Foss, Dingkang Wang, Duc
300 Le, Dustin Holland, Edward Dowling, Eissa Jamil, Elaine Montgomery, Eleonora Presani, Emily
301 Hahn, Emily Wood, Eric-Tuan Le, Erik Brinkman, Esteban Arcaute, Evan Dunbar, Evan Smothers,
302 Fei Sun, Felix Kreuk, Feng Tian, Filippos Kokkinos, Firat Ozgenel, Francesco Caggioni, Frank
303 Kanayet, Frank Seide, Gabriela Medina Florez, Gabriella Schwarz, Gada Badeer, Georgia Swee,
304 Gil Halpern, Grant Herman, Grigory Sizov, Guangyi, Zhang, Guna Lakshminarayanan, Hakan Inan,
305 Hamid Shojanazeri, Han Zou, Hannah Wang, Hanwen Zha, Haroun Habeeb, Harrison Rudolph,
306 Helen Suk, Henry Aspegren, Hunter Goldman, Hongyuan Zhan, Ibrahim Damlaj, Igor Molybog,
307 Igor Tufanov, Ilias Leontiadis, Irina-Elena Veliche, Itai Gat, Jake Weissman, James Geboski, James
308 Kohli, Janice Lam, Japhet Asher, Jean-Baptiste Gaya, Jeff Marcus, Jeff Tang, Jennifer Chan, Jenny
309 Zhen, Jeremy Reizenstein, Jeremy Teboul, Jessica Zhong, Jian Jin, Jingyi Yang, Joe Cummings,
310 Jon Carvill, Jon Shepard, Jonathan McPhie, Jonathan Torres, Josh Ginsburg, Junjie Wang, Kai
311 Wu, Kam Hou U, Karan Saxena, Kartikay Khandelwal, Katayoun Zand, Kathy Matosich, Kaushik
312 Veeraraghavan, Kelly Michelena, Keqian Li, Kiran Jagadeesh, Kun Huang, Kunal Chawla, Kyle
313 Huang, Lailin Chen, Lakshya Garg, Lavender A, Leandro Silva, Lee Bell, Lei Zhang, Liangpeng
314 Guo, Licheng Yu, Liron Moshkovich, Luca Wehrstedt, Madian Khabsa, Manav Avalani, Manish
315 Bhatt, Martynas Mankus, Matan Hasson, Matthew Lennie, Matthias Reso, Maxim Groshev, Maxim
316 Naumov, Maya Lathi, Meghan Keneally, Miao Liu, Michael L. Seltzer, Michal Valko, Michelle
317 Restrepo, Mihir Patel, Mik Vyatskov, Mikayel Samvelyan, Mike Clark, Mike Macey, Mike Wang,
318 Miquel Jubert Hermoso, Mo Metanat, Mohammad Rastegari, Munish Bansal, Nandhini Santhanam,
319 Natascha Parks, Natasha White, Navyata Bawa, Nayan Singhal, Nick Egebo, Nicolas Usunier,
320 Nikhil Mehta, Nikolay Pavlovich Laptev, Ning Dong, Norman Cheng, Oleg Chernoguz, Olivia
321 Hart, Omkar Salpekar, Ozlem Kalinli, Parkin Kent, Parth Parekh, Paul Saab, Pavan Balaji, Pedro
322 Rittner, Philip Bontrager, Pierre Roux, Piotr Dollar, Polina Zvyagina, Prashant Ratanchandani,
323 Pritish Yuvraj, Qian Liang, Rachad Alao, Rachel Rodriguez, Rafi Ayub, Raghotham Murthy,
324 Raghu Nayani, Rahul Mitra, Rangaprabhu Parthasarathy, Raymond Li, Rebekkah Hogan, Robin
325 Battey, Rocky Wang, Russ Howes, Rutu Rinott, Sachin Mehta, Sachin Siby, Sai Jayesh Bondu,
326 Samyak Datta, Sara Chugh, Sara Hunt, Sargun Dhillon, Sasha Sidorov, Satadru Pan, Saurabh
327 Mahajan, Saurabh Verma, Seiji Yamamoto, Sharadh Ramaswamy, Shaun Lindsay, Shaun Lindsay,
328 Sheng Feng, Shenghao Lin, Shengxin Cindy Zha, Shishir Patil, Shiva Shankar, Shuqiang Zhang,
329 Shuqiang Zhang, Sinong Wang, Sneha Agarwal, Soji Sajuyigbe, Soumith Chintala, Stephanie
330 Max, Stephen Chen, Steve Kehoe, Steve Satterfield, Sudarshan Govindaprasad, Sumit Gupta,
331 Summer Deng, Sungmin Cho, Sunny Virk, Suraj Subramanian, Sy Choudhury, Sydney Goldman,

332 Tal Remez, Tamar Glaser, Tamara Best, Thilo Koehler, Thomas Robinson, Tianhe Li, Tianjun
333 Zhang, Tim Matthews, Timothy Chou, Tzook Shaked, Varun Vontimitta, Victoria Ajayi, Victoria
334 Montanez, Vijai Mohan, Vinay Satish Kumar, Vishal Mangla, Vlad Ionescu, Vlad Poenaru,
335 Vlad Tiberiu Mihailescu, Vladimir Ivanov, Wei Li, Wenchen Wang, Wenwen Jiang, Wes Bouaziz,
336 Will Constable, Xiaocheng Tang, Xiaojian Wu, Xiaolan Wang, Xilun Wu, Xinbo Gao, Yaniv
337 Kleinman, Yanjun Chen, Ye Hu, Ye Jia, Ye Qi, Yenda Li, Yilin Zhang, Ying Zhang, Yossi Adi,
338 Youngjin Nam, Yu, Wang, Yu Zhao, Yuchen Hao, Yundi Qian, Yunlu Li, Yuzi He, Zach Rait,
339 Zachary DeVito, Zef Rosnbrick, Zhaoduo Wen, Zhenyu Yang, Zhiwei Zhao, and Zhiyu Ma. The
340 llama 3 herd of models, 2024. URL <https://arxiv.org/abs/2407.21783>.

341 Daya Guo, Dejian Yang, Haowei Zhang, Junxiao Song, Ruoyu Zhang, Runxin Xu, Qihao Zhu,
342 Shirong Ma, Peiyi Wang, Xiao Bi, et al. Deepseek-r1: Incentivizing reasoning capability in llms
343 via reinforcement learning. *arXiv preprint arXiv:2501.12948*, 2025.

344 HaooWang. llm-knowledge-cutoff-dates: This repository contains a summary of knowledge cut-
345 off dates for various large language models (llms), such as gpt, claude, gemini, llama, and
346 more. <https://github.com/HaooWang/llm-knowledge-cutoff-dates>, 2025. GitHub
347 repository.

348 Dan Hendrycks, Collin Burns, Saurav Kadavath, Akul Arora, Steven Basart, Eric Tang, Dawn Song,
349 and Jacob Steinhardt. Measuring mathematical problem solving with the MATH dataset. In Joaquin
350 Vanschoren and Sai-Kit Yeung, editors, *Proceedings of the Neural Information Processing Systems
351 Track on Datasets and Benchmarks 1, NeurIPS Datasets and Benchmarks 2021, December 2021,
352 virtual*, 2021. URL [https://datasets-benchmarks-proceedings.neurips.cc/paper/
353 2021/hash/be83ab3ecd0db773eb2dc1b0a17836a1-Abstract-round2.html](https://datasets-benchmarks-proceedings.neurips.cc/paper/2021/hash/be83ab3ecd0db773eb2dc1b0a17836a1-Abstract-round2.html).

354 Yiming Huang, Zhenghao Lin, Xiao Liu, Yeyun Gong, Shuai Lu, Fangyu Lei, Yaobo Liang, Yelong
355 Shen, Chen Lin, Nan Duan, and Weizhu Chen. Competition-level problems are effective LLM
356 evaluators. In Lun-Wei Ku, Andre Martins, and Vivek Srikumar, editors, *Findings of the Association
357 for Computational Linguistics, ACL 2024, Bangkok, Thailand and virtual meeting, August 11-
358 16, 2024*, pages 13526–13544. Association for Computational Linguistics, 2024. doi: 10.18653/
359 V1/2024.FINDINGS-ACL.803. URL [https://doi.org/10.18653/v1/2024.findings-acl.
360 803](https://doi.org/10.18653/v1/2024.findings-acl.803).

361 Naman Jain, King Han, Alex Gu, Wen-Ding Li, Fanjia Yan, Tianjun Zhang, Sida Wang, Armando
362 Solar-Lezama, Koushik Sen, and Ion Stoica. Livecodebench: Holistic and contamination free
363 evaluation of large language models for code. *ArXiv*, abs/2403.07974, 2024. URL [https:
364 //api.semanticscholar.org/CorpusId:268379413](https://api.semanticscholar.org/CorpusId:268379413).

365 Sachin Jain, Rusheb Batra, Subbarao Kambhampati, Ellie Pavlick, William Merrill, Danqi Chen,
366 Graham Neubig, et al. Humanity’s last exam. *arXiv preprint arXiv:2501.14249*, 2025. URL
367 <https://arxiv.org/abs/2501.14249>.

368 Tanishq Lala, Odhran Pradhan, Adrian Akiki, Carolyn Wang, Noel Iskander, Amalie Trewartha,
369 Alexander Dunn, et al. Paperqa: Retrieval-augmented generative agent for scientific research.
370 *arXiv preprint arXiv:2312.07559*, 2023. URL <https://arxiv.org/abs/2312.07559>.

371 Yucheng Li, Tianshi Wang, Zhaowei Xu, Chenghua Zhang, Bo Wang, et al. A survey on data
372 contamination for large language models. *arXiv preprint arXiv:2502.14425*, 2025. URL [https:
373 //arxiv.org/abs/2502.14425](https://arxiv.org/abs/2502.14425).

374 OpenAI. Openai o3 and o4-mini system card, 2025. URL [https://openai.com/index/
375 introducing-o3-and-o4-mini/](https://openai.com/index/introducing-o3-and-o4-mini/).

376 Shi Qiu, Shaoyang Guo, Zhuo-Yang Song, Yunbo Sun, Zeyu Cai, Jiashen Wei, Tianyu Luo, Yixuan
377 Yin, Haoxu Zhang, Yi Hu, Chenyang Wang, Chencheng Tang, Haoling Chang, Qi Liu, Ziheng
378 Zhou, Tianyu Zhang, Jingtian Zhang, Zhangyi Liu, Minghao Li, Yuku Zhang, Boxuan Jing,
379 Xianqi Yin, Yutong Ren, Zizhuo Fu, Jiaming Ji, Weike Wang, Xudong Tian, Anqi Lv, Laifu
380 Man, Jianxiang Li, Feiyu Tao, Qihua Sun, Zhou Liang, Yushu Mu, Zhongxuan Li, Jing-Jun
381 Zhang, Shutao Zhang, Xiaotian Li, Xingqi Xia, Jiawei Lin, Zheyu Shen, Jiahang Chen, Qiuha
382 Xiong, Binran Wang, Fengyuan Wang, Ziyang Ni, Bohan Zhang, Fan Cui, Changkun Shao,
383 Qing-Hong Cao, Ming xing Luo, Yaodong Yang, Muhan Zhang, and Hua Xing Zhu. Phybench:

- 384 Holistic evaluation of physical perception and reasoning in large language models, 2025. URL
385 <https://arxiv.org/abs/2504.16074>.
- 386 David Rein, Betty Li Hou, Asa Cooper Stickland, Jackson Petty, Richard Yuanzhe Pang, Julien Dirani,
387 Julian Michael, and Samuel R. Bowman. Gpqa: A graduate-level google-proof q&a benchmark,
388 2023. URL <https://arxiv.org/abs/2311.12022>.
- 389 Manley Roberts, Himanshu Thakur, Christine Herlihy, Colin White, and Samuel Dooley. To the
390 cutoff... and beyond? a longitudinal perspective on LLM data contamination. In *The Twelfth
391 International Conference on Learning Representations*, 2024. URL [https://openreview.net/
392 forum?id=m2NVG4Htxs](https://openreview.net/forum?id=m2NVG4Htxs).
- 393 Nimrod Shabtay, Felipe Maia Polo, Sivan Doveh, Wei Lin, Muhammad Jehanzeb Mirza, Leshem
394 Choshen, Mikhail Yurochkin, Yuekai Sun, Assaf Arbel, Leonid Karlinsky, and Raja Giryes.
395 Livexiv - A multi-modal live benchmark based on arxiv papers content. In *The Thirteenth
396 International Conference on Learning Representations, ICLR 2025, Singapore, April 24-28, 2025*.
397 OpenReview.net, 2025. URL <https://openreview.net/forum?id=Su1RfnEVK4>.
- 398 Michael Skarlinski, Sam Chiappa, Shayne Longpre, Nelson Elhage, and Jacob Steinhardt. Di-
399 gestables: Synthesizing scientific literature into tables using language models. *arXiv preprint
400 arXiv:2410.22360*, 2024. URL <https://arxiv.org/abs/2410.22360>.
- 401 William Thomson. *Popular Lectures and Addresses, Vol. 1: Electrical Units of Mea-
402 surement*. Macmillan and Co., London, 1889. URL [https://archive.org/details/
403 popularlecturesa01kelvuoft/page/73/mode/1up](https://archive.org/details/popularlecturesa01kelvuoft/page/73/mode/1up). Lecture delivered 3 May 1883.
- 404 Yuanhan Wang, Jianxi Li, Yifan Wang, Junyan Li, Dahua Lin, and Yu Qiao. Sciqag: A framework for
405 auto-generated scientific question answering dataset with fine-grained evaluation. *arXiv preprint
406 arXiv:2405.09939*, 2024. URL <https://arxiv.org/abs/2405.09939>.
- 407 Colin White, Samuel Dooley, Manley Roberts, Arka Pal, Benjamin Feuer, Siddhartha Jain, Ravid
408 Shwartz-Ziv, Neel Jain, Khalid Saifullah, Sreemanti Dey, Shubh-Agrawal, Sandeep Singh Sandha,
409 Siddhartha V. Naidu, Chinmay Hegde, Yann LeCun, Tom Goldstein, Willie Neiswanger, and Micah
410 Goldblum. Livebench: A challenging, contamination-limited LLM benchmark. In *The Thirteenth
411 International Conference on Learning Representations, ICLR 2025, Singapore, April 24-28, 2025*.
412 OpenReview.net, 2025. URL <https://openreview.net/forum?id=sKYHBTaxVa>.
- 413 Xiaobao Wu, Liangming Pan, Yuxi Xie, Ruiwen Zhou, Shuai Zhao, Yubo Ma, Mingzhe Du, Rui
414 Mao, Anh Tuan Luu, and William Yang Wang. Antileakbench: Preventing data contamination
415 by automatically constructing benchmarks with updated real-world knowledge. In Wanxiang
416 Che, Joyce Nabende, Ekaterina Shutova, and Mohammad Taher Pilehvar, editors, *Proceedings
417 of the 63rd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long
418 Papers)*, ACL 2025, Vienna, Austria, July 27 - August 1, 2025, pages 18403–18419. Association for
419 Computational Linguistics, 2025. URL <https://aclanthology.org/2025.acl-long.901/>.
- 420 Zhaowei Xu, Yapei Wang, Mengzhou Li, Jiaheng Chen, Qingxiu Dong, Hongyi Wang, Tao Gui,
421 Qi Zhang, and Xuanjing Huang. Unveiling the spectrum of data contamination in language
422 models: A survey from detection to remediation. *arXiv preprint arXiv:2406.14644*, 2024. URL
423 <https://arxiv.org/abs/2406.14644>.
- 424 Shuo Yang, Wei-Lin Chiang, Lianmin Zheng, Joseph Gonzalez, and Ion Stoica. Rethinking bench-
425 mark and contamination for language models with rephrased samples. *ArXiv*, abs/2311.04850,
426 2023. URL <https://arxiv.org/pdf/2311.04850.pdf>.
- 427 Jie Zhang, Cezara Petru, Kristina Nikolić, and Florian Tramèr. Realmath: A continuous benchmark
428 for evaluating language models on research-level mathematics, 2025. URL [https://arxiv.org/
429 abs/2505.12575](https://arxiv.org/abs/2505.12575).
- 430 Kaijie Zhu, Jiaao Chen, Jindong Wang, Neil Zhenqiang Gong, Diyi Yang, and Xing Xie. Dyval:
431 Dynamic evaluation of large language models for reasoning tasks. In *The Twelfth International
432 Conference on Learning Representations, ICLR 2024, Vienna, Austria, May 7-11, 2024*. OpenRe-
433 view.net, 2024. URL <https://openreview.net/forum?id=gjf0L9z5Xr>.

434 A Supplementary Materials and Instructions for Reproduction

435 The README file in the supplementary materials contains detailed instructions for how to reproduce
436 the results, where we also provide the full synthetic dataset for evaluation as well as the codebase
437 used for both QA synthesis and model evaluation, which ensures the reproducibility of this study.
438 After peer review, we plan to open-source both the code and data in their entirety.

439 The dataset includes 1,643 QA pairs, with 856 QA pairs from mathematics papers and 787 from
440 physics papers. These QA pairs were synthesized from 20,277 arXiv papers published between May
441 2023 and June 2025. Along with the datasets, we also provide the prompts that were used for question
442 generation and evaluation.

443 To reproduce our results, the papers should be retrieved from the ArXiv API, and the LaTeX code
444 should be extracted from these papers. Then, the provided prompts should be run using a large
445 language model (LLM) to standardize the LaTeX, extract theorems, check their quality, and generate
446 QA pairs. Another LLM is used to judge the quality of the generated QA pairs.

447 B Prompting Context

448 In this section, we detail all the prompts used for our QA generation and evaluation using LLM-as-a-
449 Judge paradigm as adapted from RealMath Zhang et al. [2025].

450 **System Prompt:** You are an expert research scientist designing clear question-answer pairs that requires at least 6 steps of scientific reasoning from research papers. The questions should have a unique numerical or analytical answer. Some common examples: - "If and only if condition A holds, then we can get X.", then we can ask "what condition must hold for X to be true?". This is also a unique answer. - Existence and Uniqueness Theorems: e.g., "There exists a unique X that satisfies A.", then we can ask "what is the unique solution that satisfies A?". This is also a unique answer. - Exact Formula Calculations: e.g., "The answer of formula (1) is 10", then we can ask "what is the value of formula (1)?". This is also a unique answer. - Unique Maximum/Minimum Points: e.g., "The maximum value of function f is 10 at point x=1", then we can ask "what is the maximum value of function f?". This is also a unique answer. - Exact Complexity Results in Computational Complexity: e.g., "The time complexity of algorithm A is exactly $\Theta(n^2)$ " (not $\Omega(n^2)$ or $O(n^2)$, because big-O and big-omega are not exact), then we can ask "what is the exact time complexity of algorithm A?". This is also a fixed answer.

If the theorem does not have a unique answer, you can skip this theorem and return empty result.

If the theorem is a good candidate, your questions should: - clearly state the context of this theorem, and clearly define all quantities to make the question statement clear and self-contained - requires at least 6 steps of scientific reasoning. - never reveal the answer in the question statement - never ask yes or no question, never ask questions that are easy to answer without any reasoning. - if the theorem says "There exists an X that satisfies A" but the numerical value of X is not unique, skip the theorem - if the conditions A under which we can get X are not unique (i.e. necessary and sufficient), skip the theorem - re-define in the question the quantities from the theorem statement (without revealing the answer) so that the question can be solved in a self-contained manner.

If the theorem is a good candidate, your answers should have: - a unique numerical or analytical answer, easy to verify without ambiguity; - if there's any approximation, the condition must be specified in the question body (e.g. to 2 decimal places)

Standardize LaTeX Prompt: You are an expert in LaTeX. Your task is to review contents from a scientific paper and ensure it can be directly rendered in standard LaTeX without requiring custom command definitions. We should only use usepackage: amsmath, amssymb, enumerate, amsfonts, mathrsfs, mathtools, logicproof. For any commonly used commands, you should not change them, e.g., mathbb, sum, prod, int, lim, frac, sin, cos, tan, ln, exp, log, etc. But if you find some words are similar to the custom command definitions but hard to parse, you can change them to the standard latex command, e.g., 'mathbb' should be changed to 'mathbb', because 'mathbb' is meaningless.

451 For any custom commands used in the content, please replace them with standard LaTeX notation. Make sure to check if for each begin command, there is a corresponding end command and viceversa. Moreover, make sure that \$ is not missing and insert it when needed.

I will compile the latex content into a pdf.

IMPORTANT: You must not change the mathematical meaning of the content. Focus only on syntax corrections.

Return the standardized content in this exact JSON format:

"theorem": "the well-formatted theorem in latex format without any custom commands", "changes": "explanation of what changes were made to the theorem, don't change the theorem content"

452