
Disentangling Test-Time and Parameter Scaling for Cost-Efficient Accuracy Improvements in Agentic Evaluation

Anonymous Author(s)

Affiliation

Address

email

Abstract

1 Large language models (LLMs) offer two primary levers for improving accuracy
2 in agentic systems: *test-time scaling* (e.g., Chain-of-Thought reasoning) and *pa-*
3 *rameter scaling* (upgrading to larger models). Despite widespread adoption, the
4 field lacks principled evaluation of the accuracy-cost-latency trade-offs under controlled
5 conditions. We present a comprehensive evaluation framework and conduct
6 experiments on GSM8K (1,319 items) and PopQA (2,000-item subset) to establish
7 these trade-offs. Our key findings reveal that: (i) on mathematical reasoning tasks,
8 Chain-of-Thought is highly effective for smaller models but becomes redundant
9 when internal reasoning capabilities are available; (ii) on knowledge-intensive QA,
10 performance is primarily capacity-bound, with Chain-of-Thought often increasing
11 costs without improving accuracy; (iii) for models with advanced reasoning capa-
12 bilities, external Chain-of-Thought becomes largely redundant and can even harm
13 performance while increasing costs. We formalize Pareto frontiers and cost-per-
14 point metrics that translate into actionable deployment policies for more efficient
15 agentic systems.

16 1 Introduction

17 Modern agentic LLM systems increasingly handle complex scientific workflows, from hypothesis
18 generation to experimental design and manuscript drafting. As these systems evolve from assistants
19 to autonomous actors, practitioners face a fundamental optimization question: when performance
20 is insufficient, should they apply Chain-of-Thought (CoT) [13] prompting to enhance reasoning, or
21 upgrade to a more capable model?

22 The former approach increases inference-time computation through explicit step-by-step reasoning
23 but inflates token usage and latency. The latter improves underlying capacity but typically increases
24 per-token costs. While teams routinely make these trade-offs, systematic evidence comparing their
25 cost-effectiveness across different task domains remains limited.

26 This paper advocates for a domain-aware, budget-conscious approach to model selection. We
27 hypothesize that different task types fail for fundamentally different reasons. Multi-step mathematical
28 problems are often *reasoning-limited*, where errors stem from computational mistakes or inadequate
29 problem decomposition—issues that explicit reasoning chains can address. Conversely, open-domain
30 factual QA tasks are typically *knowledge-limited*, where failures arise from missing or imprecise
31 information that longer reasoning chains cannot remedy [8].

32 We develop a controlled evaluation framework that systematically separates external Chain-of-
33 Thought prompting from internal reasoning capabilities (via model-native controls) while simultane-
34 ously measuring accuracy, monetary cost, and latency. Our experiments span two complementary

35 domains: GSM8K for mathematical reasoning [3] and PopQA for the retrieval of factual knowl-
36 edge [9].

37 **Research Hypotheses.** **H1:** On mathematical tasks, Chain-of-Thought provides substantial benefits
38 for smaller models but offers diminishing returns when internal reasoning is available. **H2:** On
39 knowledge-intensive tasks, parameter scaling dominates Chain-of-Thought in both accuracy gains
40 and cost efficiency. **H3:** For models with internal reasoning capabilities, external Chain-of-Thought
41 becomes largely redundant and can even harm performance while increasing costs.

42 **Contributions.**

- 43 • **Methodological framework:** A controlled comparison isolating external Chain-of-Thought
44 from internal reasoning to identify when they are substitutable versus complementary.
- 45 • **Cost-latency analysis:** Comprehensive logging of per-sample costs and latencies alongside
46 accuracy, with metrics for cost-per-percentage-point improvements and Pareto frontier
47 analysis.
- 48 • **Deployment guidance:** Evidence-based recommendations for practitioners: prioritize
49 parameter scaling initially, then apply reasoning techniques selectively based on task charac-
50 teristics.

51 **2 Related Work**

52 **2.1 Test-Time Scaling Techniques**

53 Chain-of-Thought prompting [13] elicits intermediate reasoning steps through few-shot examples
54 or explicit instructions, demonstrating particular effectiveness on arithmetic and logical reason-
55 ing tasks [7]. Extensions include Program-of-Thoughts [2], which leverages code execution for
56 mathematical reasoning, and Tree-of-Thoughts [14], which explores multiple reasoning branches.
57 Self-consistency approaches [12] sample multiple reasoning paths and select the most frequent
58 answer.

59 These techniques consistently improve performance on algorithmic tasks but incur computational
60 overhead through increased token generation, longer inference times, and potential verbosity that can
61 harm exact-match extraction in retrieval-oriented tasks.

62 **2.2 Parameter Scaling and Model Capacity**

63 The scaling hypothesis suggests that larger models with more parameters generally achieve better
64 performance across diverse tasks [6, 5]. Recent work has explored the trade-offs between model
65 size and inference efficiency [11], particularly in deployment scenarios with strict latency or cost
66 constraints.

67 However, systematic comparisons of parameter scaling versus test-time techniques under matched
68 experimental conditions remain limited, particularly when controlling for domain-specific task
69 characteristics.

70 **2.3 Internal Reasoning and Thinking Tokens**

71 Recent model architectures incorporate internal reasoning capabilities that allocate additional compu-
72 tation without generating observable tokens [15]. These systems can perform latent reasoning similar
73 to explicit Chain-of-Thought but with different cost and latency profiles [4]. Our work leverages
74 these capabilities where available to isolate the effects of reasoning from token generation overhead.

75 **2.4 Cost-Aware LLM Evaluation**

76 While most benchmark evaluations focus primarily on accuracy metrics, recent work has begun
77 incorporating cost and latency considerations for practical deployment scenarios [1, 10]. Our approach
78 extends this line of work by providing systematic cost-per-improvement metrics and Pareto frontier
79 analysis for strategic decision-making.

80 **3 Methodology**

81 **3.1 Experimental Framework**

82 For a dataset $D = \{(x_i, a_i^*)\}_{i=1}^N$ with inputs x_i and ground truth answers a_i^* , we define strategy
83 $s \in \{\text{No-CoT}, \text{CoT}\}$, model family f , and capacity level m . Our evaluation metrics are:

$$A(f, m, s) = \frac{1}{N} \sum_{i=1}^N \mathbf{1}[g(y_i(f, m, s)) = a_i^*] \quad (1)$$

$$C(f, m, s) = \frac{1}{N} \sum_{i=1}^N \text{cost}(y_i) \quad (2)$$

$$L(f, m, s) = \frac{1}{N} \sum_{i=1}^N \text{latency}(y_i) \quad (3)$$

84 where $g(\cdot)$ is an extraction function, y_i is the model output, and we measure accuracy A , cost C ,
85 and latency L . We report improvements relative to baseline configurations and calculate cost-per-
86 percentage-point (CPP) as $\Delta C/\Delta A$.

87 **3.2 Task-Specific Prompting**

88 **GSM8K (Mathematical Reasoning).** We employ an 8-shot Chain-of-Thought template that en-
89 forces deterministic answer extraction:

```
90 Question: Janet's ducks lay 16 eggs per day...  
91 Reason step by step, then give the final numeric answer.  
92 Final line must be: #### <number>  
93 Answer: Step 1: Calculate... Step 2: Therefore...  
94 #### 18
```

95 The No-CoT variant requests only the final numeric answer in the specified format.

96 **PopQA (Knowledge Retrieval).** The Chain-of-Thought template elicits brief reasoning while
97 ensuring clean answer extraction:

```
98 Question: Who wrote "Pride and Prejudice"?  
99 Answer step by step in 1-3 sentences, then provide:  
100 Final Answer: <entity>
```

101 The No-CoT variant requests only the direct answer. We evaluate using exact match (EM) and F1
102 scores with standard normalization.

103 **3.3 Internal Reasoning Controls**

104 For models supporting internal reasoning controls (specifically Gemini Flash), we manipulate the
105 thinking budget parameter: setting it to zero disables internal reasoning, while omitting the parameter
106 enables it. This allows us to isolate capacity effects from reasoning effects while maintaining identical
107 prompts and decoding parameters.

108 **3.4 Experimental Controls**

109 All experiments use consistent random seeds (42), temperature 0.7, top-p 1.0, and appropriate
110 timeouts. We log input/output tokens, end-to-end latency, and per-sample costs. Self-consistency
111 sampling is disabled due to API reliability constraints, which we discuss in our limitations.

Table 1: Complete GSM8K results across all model configurations.

Model	Reasoning	CoT	Accuracy	Cost/Sample (\$)	Total Cost (\$)	Latency/Sample (s)
GPT-4.1-mini	N/A	No	45.19%	0.000042	0.055	0.65
GPT-4.1-mini	N/A	Yes	95.15%	0.000763	1.006	2.29
GPT-4.1	N/A	No	57.01%	0.000208	0.275	0.75
GPT-4.1	N/A	Yes	94.69%	0.003889	5.130	2.55
Gemini 2.5 Flash-Lite	N/A	No	36.67%	0.000025	0.033	0.80
Gemini 2.5 Flash-Lite	N/A	Yes	93.85%	0.000250	0.330	1.45
Gemini 2.5 Flash	Disabled	No	55.19%	0.000038	0.050	0.58
Gemini 2.5 Flash	Disabled	Yes	95.60%	0.000944	1.246	1.45
Gemini 2.5 Flash	Enabled	No	95.36%	0.000038	0.049	2.07
Gemini 2.5 Flash	Enabled	Yes	95.27%	0.000889	1.172	2.99
Gemini 2.5 Pro	Enabled	No	96.18%	0.000154	0.203	7.57
Gemini 2.5 Pro	Enabled	Yes	96.41%	0.003867	5.100	10.88

Table 2: Complete PopQA results across all model configurations.

Model	Reasoning	CoT	Accuracy	Cost/Sample (\$)	Total Cost (\$)	Latency/Sample (s)
GPT-4.1-mini	N/A	No	36.05%	0.000022	0.045	0.64
GPT-4.1-mini	N/A	Yes	40.85%	0.000201	0.402	1.99
GPT-4.1	N/A	No	49.55%	0.000110	0.220	0.74
GPT-4.1	N/A	Yes	44.55%	0.000911	1.822	2.40
Gemini 2.5 Flash-Lite	N/A	No	29.50%	0.000005	0.009	0.70
Gemini 2.5 Flash-Lite	N/A	Yes	29.59%	0.000046	0.093	1.13
Gemini 2.5 Flash	Disabled	No	33.50%	0.000018	0.035	0.56
Gemini 2.5 Flash	Disabled	Yes	38.95%	0.000208	0.416	0.94
Gemini 2.5 Flash	Enabled	No	40.13%	0.000018	0.037	2.48
Gemini 2.5 Flash	Enabled	Yes	38.35%	0.000192	0.384	2.65
Gemini 2.5 Pro	Enabled	No	45.60%	0.000080	0.160	9.34
Gemini 2.5 Pro	Enabled	Yes	40.33%	0.001172	2.345	12.84

112 4 Experiments and Results

113 4.1 Experimental Setup

114 We evaluate on the complete GSM8K test set (1,319 items) and a uniformly sampled subset of
 115 PopQA (2,000 items, seed=42). Our baselines are the smallest models in each family without Chain-
 116 of-Thought: GPT-4.1-mini and Gemini 2.5 Flash-Lite. We compare two improvement strategies:
 117 (i) adding Chain-of-Thought to the baseline model, and (ii) upgrading to a larger model while
 118 maintaining reasoning settings.

119 4.2 Complete Performance Results

120 Tables 1 and 2 present the complete experimental results, showing absolute accuracy, cost, and latency
 121 metrics for all model configurations. These raw results reveal striking patterns that become clearer
 122 when examined alongside the relative improvements.

123 4.3 GSM8K Results

124 Table 3 summarizes performance improvements over baseline configurations. For OpenAI models,
 125 both Chain-of-Thought and parameter scaling provide substantial accuracy gains, with Chain-of-
 126 Thought achieving larger absolute improvements (+49.96 pp vs +11.82 pp) but at higher cost and
 127 latency.

128 For Gemini models, parameter scaling with reasoning disabled achieves substantial gains (+18.52
 129 pp) while actually reducing latency and maintaining very low costs. Most notably, when internal
 130 reasoning is enabled (Gemini 2.5 Flash with reasoning), the model achieves 95.36% accuracy without
 131 Chain-of-Thought, while adding Chain-of-Thought yields 95.27%—a negligible improvement that
 132 comes with substantial cost and latency penalties.

Table 3: GSM8K performance improvements over small No-CoT baseline by family.

Family	Strategy	Δ Acc (pp)	Δ Cost (\$)	Δ Latency (s)
OpenAI	Small + CoT	+49.96	+0.000721	+1.64
OpenAI	Upgrade (No-CoT)	+11.82	+0.000166	+0.10
Gemini	Small + CoT	+57.18	+0.000225	+0.65
Gemini	Upgrade (No reasoning)	+18.52	+0.000013	-0.22

Table 4: PopQA performance improvements over small No-CoT baseline by family.

Family	Strategy	Δ Acc (pp)	Δ Cost (\$)	Δ Latency (s)
OpenAI	Small + CoT	+4.80	+0.000179	+1.35
OpenAI	Upgrade (No-CoT)	+13.50	+0.000088	+0.10
Gemini	Small + CoT	+0.09	+0.000041	+0.43
Gemini	Upgrade (No reasoning)	+4.00	+0.000013	-0.14

133 4.4 PopQA Results

134 Table 4 shows markedly different patterns for knowledge retrieval tasks. Parameter scaling consistently outperforms Chain-of-Thought in both accuracy gains and cost efficiency. More strikingly, for 135 models with reasoning capabilities (GPT-4.1 and Gemini models with reasoning enabled), Chain- 136 of-Thought actually *decreases* accuracy: GPT-4.1 drops from 49.55% to 44.55%, Gemini 2.5 Flash 137 (enabled) drops from 40.13% to 38.35%, and Gemini 2.5 Pro drops from 45.60% to 40.33%. 138

139 4.5 Cost Efficiency Analysis

140 Table 5 presents cost-per-percentage-point metrics, revealing clear efficiency patterns. Parameter 141 scaling consistently provides better cost efficiency than Chain-of-Thought, particularly for Gemini 142 models where the difference is substantial (5.6 \times more efficient for GSM8K).

143 4.6 Pareto Frontier Analysis

144 Figures 1 and 2 visualize the cost-accuracy Pareto frontiers for both domains. The mathematical 145 reasoning domain shows multiple viable paths to high accuracy, while the knowledge domain 146 demonstrates clear dominance of parameter scaling approaches.

147 5 Discussion

148 5.1 Chain-of-Thought Redundancy in Advanced Models

149 Our results reveal a critical insight: **Chain-of-Thought becomes largely redundant or even** 150 **counterproductive when internal reasoning capabilities are available.** This pattern is evident 151 across both domains but manifests differently:

152 **Mathematical Reasoning with Internal Capabilities.** Gemini 2.5 Flash with reasoning en- 153 abled achieves 95.36% accuracy without Chain-of-Thought, while adding Chain-of-Thought yields 154 95.27%—a negligible difference that comes with substantial cost penalties (23 \times higher cost per 155 sample). Similarly, Gemini 2.5 Pro shows minimal improvement (96.18% to 96.41%) at dramatically 156 higher cost (25 \times increase).

157 **Knowledge Retrieval with Advanced Models.** The pattern is even more pronounced in knowledge 158 tasks, where Chain-of-Thought consistently degrades performance for capable models. This suggests 159 that external reasoning chains can interfere with internal knowledge retrieval processes and introduce 160 paraphrase drift that harms exact-match scoring.

Table 5: Cost efficiency comparison: cost per +1 percentage point (lower is better).

Domain	Strategy	Cost per +1pp (\$)
GSM8K	OpenAI: Small + CoT	1.44×10^{-5}
GSM8K	OpenAI: Upgrade	1.40×10^{-5}
GSM8K	Gemini: Small + CoT	3.94×10^{-6}
GSM8K	Gemini: Upgrade	7.02×10^{-7}
PopQA	OpenAI: Small + CoT	3.73×10^{-5}
PopQA	OpenAI: Upgrade	6.52×10^{-6}
PopQA	Gemini: Small + CoT	4.56×10^{-4}
PopQA	Gemini: Upgrade	3.25×10^{-6}

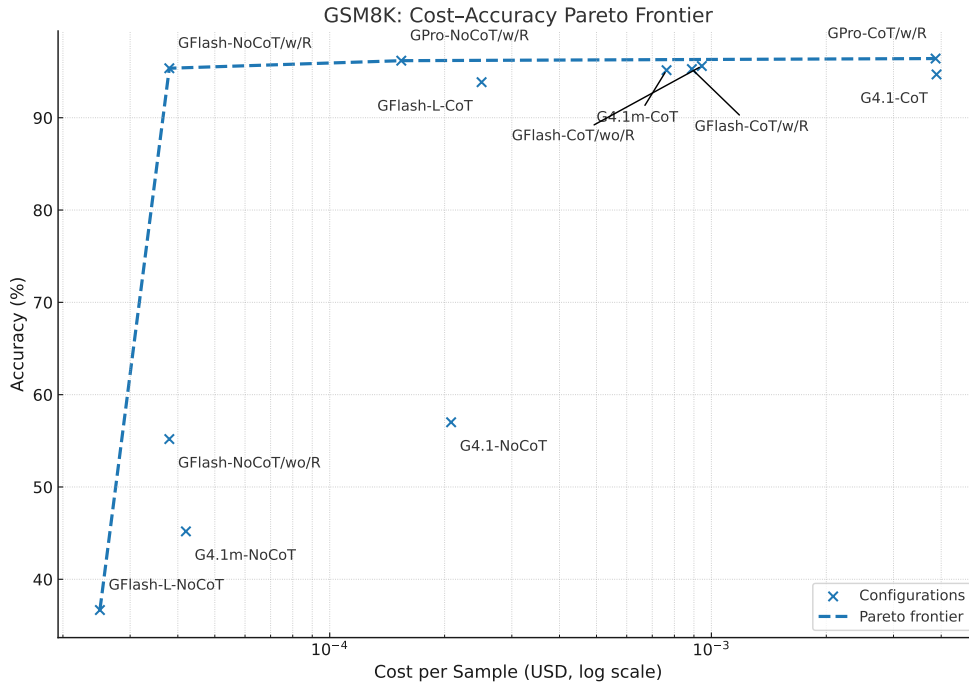


Figure 1: GSM8K cost-accuracy Pareto frontier. Each point represents a (model, CoT, reasoning) configuration plotted by average cost per sample (log scale) versus accuracy on the full test set. The dashed line connects non-dominated configurations. Abbreviations: G4.1m=GPT-4.1-mini, G4.1=GPT-4.1, GFlash-L=Gemini 2.5 Flash-Lite, GFlash=Gemini 2.5 Flash, GPro=Gemini 2.5 Pro, w/R=with internal reasoning, wo/R=without internal reasoning.

161 **5.2 Domain-Specific Optimization Strategies**

162 Beyond the redundancy of external reasoning for capable models, our results support distinct opti-
 163 mization strategies:

164 **Mathematical Reasoning.** Both Chain-of-Thought and internal reasoning address the core problem
 165 of computational errors through structured intermediate steps. Chain-of-Thought materializes these
 166 steps as tokens (external scratchpad), while internal reasoning processes them in hidden states (latent
 167 scratchpad). When tasks admit algorithmic decomposition, both approaches approximate similar
 168 computational patterns, explaining their effectiveness on GSM8K.

169 **Knowledge Retrieval.** Factual QA tasks fail primarily due to missing or imprecise information
 170 rather than reasoning errors. Extended reasoning chains can introduce paraphrase drift, where models

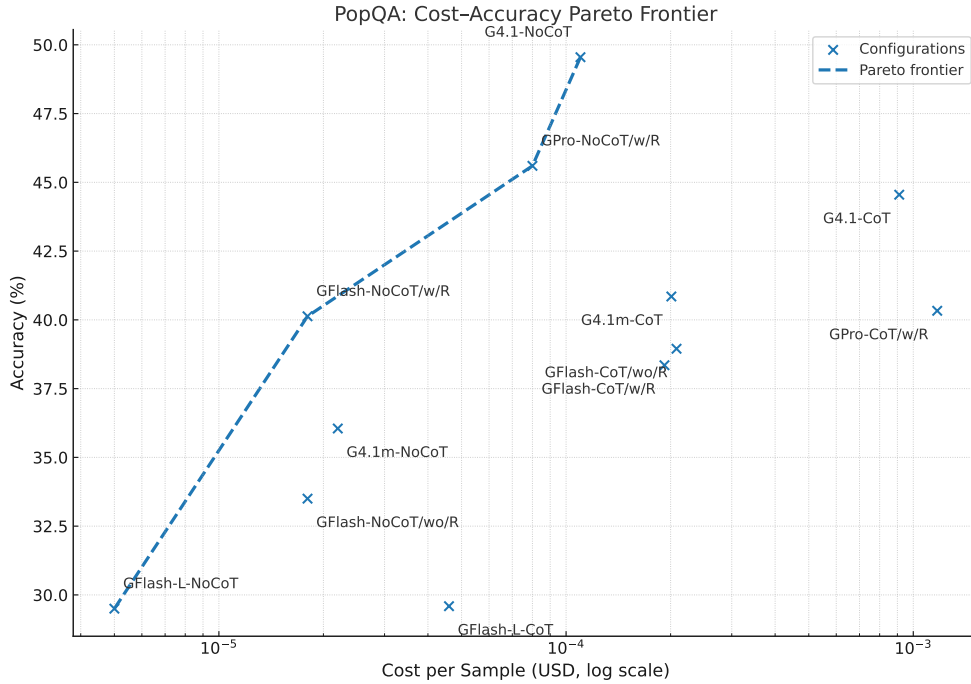


Figure 2: PopQA cost-accuracy Pareto frontier. Parameter scaling configurations dominate the efficient frontier for knowledge retrieval tasks, with Chain-of-Thought approaches generally falling below the optimal cost-accuracy trade-off. Same abbreviations as Figure 1.

171 subtly alter answer phrasing in ways that harm exact-match scoring. Parameter scaling addresses
 172 the root cause by improving factual knowledge, while Chain-of-Thought often adds cost without
 173 corresponding benefits.

174 5.3 Practical Deployment Guidelines

175 Based on our findings, we recommend the following deployment strategy:

- 176 1. **Initialize with small models:** Begin with the most cost-effective baseline (small model, no
 177 Chain-of-Thought).
- 178 2. **Scale parameters first:** Upgrade model capacity before adding reasoning techniques, as
 179 this typically provides better cost efficiency and can reduce latency.
- 180 3. **Avoid Chain-of-Thought for advanced models:** For models with internal reasoning
 181 capabilities, external Chain-of-Thought is typically redundant and can harm performance
 182 while increasing costs.
- 183 4. **Apply reasoning selectively:** Use Chain-of-Thought primarily for mathematical/logical
 184 tasks with smaller models, or when pursuing maximum accuracy regardless of cost.
- 185 5. **Consider domain routing:** Implement task-aware routing that applies different strategies
 186 based on predicted task characteristics.

187 5.4 Limitations

188 Several factors limit the generalizability of our findings:

189 **Dataset scope:** Our PopQA evaluation uses a 2,000-item subset, which may not fully represent the
 190 diversity of the complete dataset, particularly for long-tail entities.

191 **Self-consistency exclusion:** API reliability issues prevented evaluation of self-consistency techniques,
 192 which could improve Chain-of-Thought performance, particularly on mathematical tasks.

193 **Provider-specific features:** Internal reasoning controls are available only for certain models (Gemini
194 in our study), limiting cross-provider comparisons.

195 **Prompt design:** Our deterministic answer extraction may favor certain approaches; alternative
196 prompt designs could shift absolute performance while preserving relative trends.

197 6 Conclusion

198 This work provides systematic evidence for optimizing the accuracy-cost-latency trade-offs in agentic
199 LLM systems. Our key insight is that optimal strategies depend critically on both task domain and
200 model capabilities: mathematical reasoning benefits from parameter scaling and reasoning techniques
201 for smaller models, while knowledge retrieval tasks favor parameter scaling alone. Most importantly,
202 for models with internal reasoning capabilities, external Chain-of-Thought becomes largely redundant
203 and can even harm performance.

204 The practical recommendation is straightforward: begin with cost-effective baselines, prioritize
205 parameter scaling for initial improvements, and avoid external reasoning techniques for advanced
206 models with internal reasoning capabilities. This approach can significantly reduce deployment costs
207 while maintaining or improving accuracy across diverse agentic workflows.

208 Future work should expand evaluation to additional domains, incorporate self-consistency techniques
209 under stable conditions, and develop automated routing systems that adapt strategies based on
210 real-time task classification and model capability assessment.

211 References

212 [1] Lingjiao Chen, Matei Zaharia, and James Zou. Frugalgpt: How to use large language models
213 while reducing cost and improving performance. *arXiv preprint arXiv:2305.05176*, 2023.

214 [2] Wenhu Chen, Xueguang Ma, Xinyi Wang, and William W Cohen. Program of thoughts prompt-
215 ing: Disentangling computation from reasoning for numerical reasoning tasks. *Transactions on*
216 *Machine Learning Research*, 2022.

217 [3] Karl Cobbe, Vineet Kosaraju, Mohammad Bavarian, Mark Chen, Heewoo Jun, Lukasz Kaiser,
218 Matthias Plappert, Jerry Tworek, Jacob Hilton, Reiichiro Nakano, et al. Training verifiers to
219 solve math word problems. *arXiv preprint arXiv:2110.14168*, 2021.

220 [4] Shibo Hao, Sainbayar Sukhbaatar, DiJia Su, Xian Li, Zhiting Hu, Jason Weston, and Yuandong
221 Tian. Training large language models to reason in a continuous latent space. *arXiv preprint*
222 *arXiv:2412.06769*, 2024.

223 [5] Jordan Hoffmann, Sebastian Borgeaud, Arthur Mensch, Elena Buchatskaya, Trevor Cai, Eliza
224 Rutherford, Diego de Las Casas, Lisa Anne Hendricks, Johannes Welbl, Aidan Clark, et al.
225 Training compute-optimal large language models. *arXiv preprint arXiv:2203.15556*, 2022.

226 [6] Jared Kaplan, Sam McCandlish, Tom Henighan, Tom B Brown, Benjamin Chess, Rewon Child,
227 Scott Gray, Alec Radford, Jeffrey Wu, and Dario Amodei. Scaling laws for neural language
228 models. *arXiv preprint arXiv:2001.08361*, 2020.

229 [7] Takeshi Kojima, Shixiang Shane Gu, Machel Reid, Yutaka Matsuo, and Yusuke Iwasawa. Large
230 language models are zero-shot reasoners. In *Advances in neural information processing systems*,
231 volume 35, pages 22199–22213, 2022.

232 [8] Patrick Lewis, Ethan Perez, Aleksandra Piktus, Fabio Petroni, Vladimir Karpukhin, Naman
233 Goyal, Heinrich Küttler, Mike Lewis, Wen-tau Yih, Tim Rocktäschel, et al. Retrieval-augmented
234 generation for knowledge-intensive nlp tasks. *Advances in neural information processing*
235 *systems*, 33:9459–9474, 2020.

236 [9] Alex Mallen, Akari Asai, Victor Zhong, Rajarshi Das, Daniel Khashabi, and Hannaneh Ha-
237 jishirzi. When not to trust language models: Investigating effectiveness of parametric and
238 non-parametric memories. *arXiv preprint arXiv:2212.10511*, 2022.

- 239 [10] Yujia Qin, Shihao Liang, Yining Ye, Kunlun Zhu, Lan Yan, Yaxi Lu, Yankai Lin, Xin Cong,
240 Xiangru Tang, Bill Qian, et al. Toolllm: Facilitating large language models to master 16000+
241 real-world apis. *arXiv preprint arXiv:2307.16789*, 2023.
- 242 [11] Hugo Touvron, Louis Martin, Kevin Stone, Peter Albert, Amjad Almahairi, Yasmine Babaei,
243 Nikolay Bashlykov, Soumya Batra, Prajjwal Bhargava, Shruti Bhosale, et al. Llama 2: Open
244 foundation and fine-tuned chat models. *arXiv preprint arXiv:2307.09288*, 2023.
- 245 [12] Xuezhi Wang, Jason Wei, Dale Schuurmans, Quoc Le, Ed Chi, Sharan Narang, Aakanksha
246 Chowdhery, and Denny Zhou. Self-consistency improves chain of thought reasoning in language
247 models. *arXiv preprint arXiv:2203.11171*, 2022.
- 248 [13] Jason Wei, Xuezhi Wang, Dale Schuurmans, Maarten Bosma, Fei Xia, Ed Chi, Quoc V Le,
249 Denny Zhou, et al. Chain-of-thought prompting elicits reasoning in large language models. In
250 *Advances in neural information processing systems*, volume 35, pages 24824–24837, 2022.
- 251 [14] Shunyu Yao, Dian Yu, Jeffrey Zhao, Izhak Shafran, Tom Griffiths, Yuan Cao, and Karthik
252 Narasimhan. Tree of thoughts: Deliberate problem solving with large language models. *Ad-
253 vances in neural information processing systems*, 36:11809–11822, 2023.
- 254 [15] Eric Zelikman, Georges Harik, Yijia Shao, Varuna Jayasiri, Nick Haber, and Noah D Goodman.
255 Quiet-star: Language models can teach themselves to think before speaking. *arXiv preprint
256 arXiv:2403.09629*, 2024.

257 A Implementation Details

258 **Software Environment.** Experiments use Python 3.11 with official API clients for OpenAI and
259 Google. Core dependencies include `pandas`, `numpy`, `tqdm`, and `httpx` for HTTP requests. Complete
260 dependency specifications and environment setup instructions are provided with our code release.

261 **Reproducibility.** Each experimental run generates detailed JSONL logs containing per-sample
262 inputs, outputs, token counts, latencies, and costs. Aggregated results are exported to CSV format for
263 analysis. All random seeds, API parameters, and reasoning toggles are centrally configured to ensure
264 consistency across conditions.

265 **Data Availability.** GSM8K uses the standard test split available through HuggingFace datasets.
266 PopQA subset selection uses `numpy` random sampling with `seed=42` for reproducibility. Specific
267 item indices used in our evaluation will be released with our code.

268 **Code Availability.** Our implementation is available as open source at [https://anonymous.4open.
269 science/r/agent4science-2D92](https://anonymous.4open.science/r/agent4science-2D92).

270 Agents4Science AI Involvement Checklist

271 1. **Hypothesis development:** Hypothesis development includes the process by which you
272 came to explore this research topic and research question. This can involve the background
273 research performed by either researchers or by AI. This can also involve whether the idea
274 was proposed by researchers or by AI.

275 Answer: [D]

276 Explanation: I used Liner’s "Hypothesis Generator" agent to propose LLM-related hypothe-
277 ses that could be executed by AI across the full research pipeline. From several candidates, I
278 selected one and lightly refined it with my own perspective. I then evaluated it with Liner’s
279 "Hypothesis Evaluator" agent and incorporated its feedback. Through this iteration, with
280 minimal human steering but significant AI ideation and critique, the final hypothesis used in
281 the paper was produced.

282 2. **Experimental design and implementation:** This category includes design of experiments
283 that are used to test the hypotheses, coding and implementation of computational methods,
284 and the execution of these experiments.

285 Answer: [C]

286 Explanation: I used cursor with the claude sonnet-4 model to generate the experiment code.
287 The initial plan was to run open-source models on GPU instances, but the AI-generated
288 code yielded implausible results; my manual “vibe coding” attempts did not fix them. I
289 pivoted to LLM API calls to simplify execution. Even then, many errors remained, so I
290 read the code, diagnosed issues, and guided the coding agent on where and how to patch
291 them. AI produced most of the code, while I performed validation, debugging, and design
292 corrections—hence a rating of C.

293 3. **Analysis of data and interpretation of results:** This category encompasses any process to
294 organize and process data for the experiments in the paper. It also includes interpretations of
295 the results of the study.

296 Answer: [B]

297 Explanation: Data analysis and interpretation were primarily done by human, with cursor
298 and chat-gpt-5 assisting. Human verified whether results were valid, flagged anomalies, and
299 decided when to proceed or halt analyses. chat-gpt-5 helped summarize findings, suggest
300 checks, and organize tables/figures, but the key judgments—accepting results, revising
301 analyses, or updating conclusions—were mine. This is why I rate this category as B.

302 4. **Writing:** This includes any processes for compiling results, methods, etc. into the final
303 paper form. This can involve not only writing of the main text but also figure-making,
304 improving layout of the manuscript, and formulation of narrative.

305 Answer: [D]

306 Explanation: For the initial draft, I used chat-gpt-5 and claude sonnet-4, supplying the
307 hypothesis, experimental plans, results, and cursor prompts. To refine the draft, I used
308 Liner’s Peer Review Agent, incorporated its feedback with claude sonnet-4, and repeated
309 this feedback loop about three times while checking for content drift. I stopped around
310 the fourth iteration when forced changes began to induce hallucinations. Figures were
311 generated by chat-gpt-5 from the results, and citations were suggested by Liner’s "Citation
312 Recommender" and integrated into \LaTeX . Overall, writing relied heavily on AI agents, with
313 human oversight for correctness and coherence.

314 5. **Observed AI Limitations:** What limitations have you found when using AI as a partner or
315 lead author?

316 Description: The hardest stage was coding and running experiments in cursor. The AI often
317 showed overconfidence—treating incomplete runs as “finished,” missing global context, or
318 producing plausible but incorrect outputs. When the code failed semantically (no crash,
319 wrong results), the agent struggled to localize faults. I had to perform root-cause analysis and
320 propose concrete fixes, then direct the agent to implement them. Using AI-generated (rather
321 than human-written) code increased verification overhead. In short: limited end-to-end
322 verification and insufficient epistemic humility were the main pain points.

323 Agents4Science Paper Checklist

324 1. Claims

325 Question: Do the main claims made in the abstract and introduction accurately reflect the
326 paper's contributions and scope?

327 Answer: [Yes]

328 Justification: The abstract and Introduction state the central claim—comparing test-time
329 scaling (CoT) vs. parameter scaling under cost/latency—formalize H1–H3, and these are
330 evaluated in Experiments and Results sections with supporting analyses in Cost Efficiency
331 Analysis and Discussion section. The scope and limitations are also explicitly acknowledged
332 in Limitations section.

333 Guidelines:

- 334 • The answer NA means that the abstract and introduction do not include the claims
335 made in the paper.
- 336 • The abstract and/or introduction should clearly state the claims made, including the
337 contributions made in the paper and important assumptions and limitations. A No or
338 NA answer to this question will not be perceived well by the reviewers.
- 339 • The claims made should match theoretical and experimental results, and reflect how
340 much the results can be expected to generalize to other settings.
- 341 • It is fine to include aspirational goals as motivation as long as it is clear that these goals
342 are not attained by the paper.

343 2. Limitations

344 Question: Does the paper discuss the limitations of the work performed by the authors?

345 Answer: [Yes]

346 Justification: Limitations section discusses dataset scope (PopQA subset), exclusion of self-
347 consistency due to API reliability, provider-specific reasoning controls, and prompt-design
348 sensitivities, clarifying how these factors affect generalization.

349 Guidelines:

- 350 • The answer NA means that the paper has no limitation while the answer No means that
351 the paper has limitations, but those are not discussed in the paper.
- 352 • The authors are encouraged to create a separate "Limitations" section in their paper.
- 353 • The paper should point out any strong assumptions and how robust the results are to
354 violations of these assumptions (e.g., independence assumptions, noiseless settings,
355 model well-specification, asymptotic approximations only holding locally). The authors
356 should reflect on how these assumptions might be violated in practice and what the
357 implications would be.
- 358 • The authors should reflect on the scope of the claims made, e.g., if the approach was
359 only tested on a few datasets or with a few runs. In general, empirical results often
360 depend on implicit assumptions, which should be articulated.
- 361 • The authors should reflect on the factors that influence the performance of the approach.
362 For example, a facial recognition algorithm may perform poorly when image resolution
363 is low or images are taken in low lighting.
- 364 • The authors should discuss the computational efficiency of the proposed algorithms
365 and how they scale with dataset size.
- 366 • If applicable, the authors should discuss possible limitations of their approach to
367 address problems of privacy and fairness.
- 368 • While the authors might fear that complete honesty about limitations might be used by
369 reviewers as grounds for rejection, a worse outcome might be that reviewers discover
370 limitations that aren't acknowledged in the paper. Reviewers will be specifically
371 instructed to not penalize honesty concerning limitations.

372 3. Theory assumptions and proofs

373 Question: For each theoretical result, does the paper provide the full set of assumptions and
374 a complete (and correct) proof?

375
376
377
378
379
380
381
382
383
384
385
386
387
388
389
390
391
392
393
394
395
396
397
398
399
400
401
402
403
404
405
406
407
408
409
410
411
412
413
414
415
416
417
418
419
420
421
422
423
424
425
426

Answer: [NA]

Justification: The paper is empirical and does not present new theorems or formal proofs; therefore this item is not applicable.

Guidelines:

- The answer NA means that the paper does not include theoretical results.
- All the theorems, formulas, and proofs in the paper should be numbered and cross-referenced.
- All assumptions should be clearly stated or referenced in the statement of any theorems.
- The proofs can either appear in the main paper or the supplemental material, but if they appear in the supplemental material, the authors are encouraged to provide a short proof sketch to provide intuition.

4. Experimental result reproducibility

Question: Does the paper fully disclose all the information needed to reproduce the main experimental results of the paper to the extent that it affects the main claims and/or conclusions of the paper (regardless of whether the code and data are provided or not)?

Answer: [Yes]

Justification: Experimental Controls and Methodology sections specify seeds (42), decoding parameters, prompts, reasoning toggles, and evaluation metrics; Implementation Details section describes logging of tokens/latency/cost and release of PopQA indices, enabling reproduction with the same APIs.

Guidelines:

- The answer NA means that the paper does not include experiments.
- If the paper includes experiments, a No answer to this question will not be perceived well by the reviewers: Making the paper reproducible is important.
- If the contribution is a dataset and/or model, the authors should describe the steps taken to make their results reproducible or verifiable.
- We recognize that reproducibility may be tricky in some cases, in which case authors are welcome to describe the particular way they provide for reproducibility. In the case of closed-source models, it may be that access to the model is limited in some way (e.g., to registered users), but it should be possible for other researchers to have some path to reproducing or verifying the results.

5. Open access to data and code

Question: Does the paper provide open access to the data and code, with sufficient instructions to faithfully reproduce the main experimental results, as described in supplemental material?

Answer: [Yes]

Justification: An anonymized archive (supplement) includes scripts, configs, and run instructions; public repository and PopQA subset indices will be released upon acceptance. All datasets used (GSM8K test split, PopQA) are public.

Guidelines:

- The answer NA means that paper does not include experiments requiring code.
- Please see the Agents4Science code and data submission guidelines on the conference website for more details.
- While we encourage the release of code and data, we understand that this might not be possible, so “No” is an acceptable answer. Papers cannot be rejected simply for not including code, unless this is central to the contribution (e.g., for a new open-source benchmark).
- The instructions should contain the exact command and environment needed to run to reproduce the results.
- At submission time, to preserve anonymity, the authors should release anonymized versions (if applicable).

6. Experimental setting/details

427 Question: Does the paper specify all the training and test details (e.g., data splits, hyper-
428 parameters, how they were chosen, type of optimizer, etc.) necessary to understand the
429 results?

430 Answer: [Yes]

431 Justification: Methodology section details datasets/splits, prompt templates, decoding pa-
432 rameters, and reasoning controls; Experimental Setup and Implementation Details sections
433 provide additional configuration and logging specifics.

434 Guidelines:

- 435 • The answer NA means that the paper does not include experiments.
- 436 • The experimental setting should be presented in the core of the paper to a level of detail
437 that is necessary to appreciate the results and make sense of them.
- 438 • The full details can be provided either with the code, in appendix, or as supplemental
439 material.

440 7. Experiment statistical significance

441 Question: Does the paper report error bars suitably and correctly defined or other appropriate
442 information about the statistical significance of the experiments?

443 Answer: [No]

444 Justification: We report point estimates (e.g., Exact Match accuracy) without confidence
445 intervals or hypothesis tests; we plan to add bootstrap CIs in the appendix for key results in
446 a revised version.

447 Guidelines:

- 448 • The answer NA means that the paper does not include experiments.
- 449 • The authors should answer "Yes" if the results are accompanied by error bars, confi-
450 dence intervals, or statistical significance tests, at least for the experiments that support
451 the main claims of the paper.
- 452 • The factors of variability that the error bars are capturing should be clearly stated
453 (for example, train/test split, initialization, or overall run with given experimental
454 conditions).

455 8. Experiments compute resources

456 Question: For each experiment, does the paper provide sufficient information on the com-
457 puter resources (type of compute workers, memory, time of execution) needed to reproduce
458 the experiments?

459 Answer: [No]

460 Justification: Inference was executed via hosted APIs (OpenAI, Google), so provider-side
461 hardware is not user-controllable; we report per-sample latency and monetary cost, but do
462 not enumerate provider compute specs. We can add client machine specs and wall-clock
463 durations in the supplement.

464 Guidelines:

- 465 • The answer NA means that the paper does not include experiments.
- 466 • The paper should indicate the type of compute workers CPU or GPU, internal cluster,
467 or cloud provider, including relevant memory and storage.
- 468 • The paper should provide the amount of compute required for each of the individual
469 experimental runs as well as estimate the total compute.

470 9. Code of ethics

471 Question: Does the research conducted in the paper conform, in every respect, with the
472 Agents4Science Code of Ethics (see conference website)?

473 Answer: [Yes]

474 Justification: The work uses publicly available datasets and model APIs, contains no human-
475 subjects or personal data, and follows dataset licenses and provider terms; we see no
476 deviations from the Code of Ethics.

477 Guidelines:

- 478 • The answer NA means that the authors have not reviewed the Agents4Science Code of
479 Ethics.
480 • If the authors answer No, they should explain the special circumstances that require a
481 deviation from the Code of Ethics.

482 **10. Broader impacts**

483 Question: Does the paper discuss both potential positive societal impacts and negative
484 societal impacts of the work performed?

485 Answer: [No]

486 Justification: The current draft does not include a dedicated broader-impacts discussion; we
487 will add a brief section covering efficiency benefits, risks of over-automation/misuse, and
488 mitigation strategies in a revised version.

489 Guidelines:

- 490 • The answer NA means that there is no societal impact of the work performed.
491 • If the authors answer NA or No, they should explain why their work has no societal
492 impact or why the paper does not address societal impact.
493 • Examples of negative societal impacts include potential malicious or unintended uses
494 (e.g., disinformation, generating fake profiles, surveillance), fairness considerations,
495 privacy considerations, and security considerations.
496 • If there are negative societal impacts, the authors could also discuss possible mitigation
497 strategies.