

# Revisiting the Test-Time Scaling of o1-like Models: Do they Truly Possess Test-Time Scaling Capabilities?

Anonymous ACL submission

## Abstract

The advent of test-time scaling in large language models (LLMs), exemplified by OpenAI’s o1 series, has advanced reasoning capabilities by scaling computational resource allocation during inference. While successors like QwQ, Deepseek-R1 (R1) and LIMO replicate these advancements, whether these models truly possess test-time scaling capabilities remains underexplored. This study found that longer CoTs of these o1-like models do not consistently enhance accuracy; in fact, correct solutions are often shorter than incorrect ones for the same questions. Further investigation shows this phenomenon is closely related to models’ self-revision capabilities - longer CoTs contain more self-revisions, which often lead to performance degradation. We then compare sequential and parallel scaling strategies on QwQ, R1 and LIMO, finding that parallel scaling achieves better coverage and scalability. Based on these insights, we propose Shortest Majority Vote, a method that combines parallel scaling strategies with CoT length characteristics, significantly improving models’ test-time scalability compared to conventional majority voting approaches.

## 1 Introduction

The release of the OpenAI o1 series models (OpenAI, 2024a,b) marked a pivotal advancement in the reasoning capabilities of Large Language Models (LLMs), introducing a novel scaling paradigm, test-time scaling, which allocates more compute resources during test time. The test-time scaling have two dimensions, sequential and parallel (Zeng et al., 2024). Sequential scaling increase test-time compute by scaling the length of Chain-of-Thought (CoT) (Wei et al., 2022), while parallel scaling parallelly samples multiple solutions and pick the best one.

Following o1’s success, models such as QwQ (Team, 2024), Deepseek-R1 (R1) (DeepSeek-AI

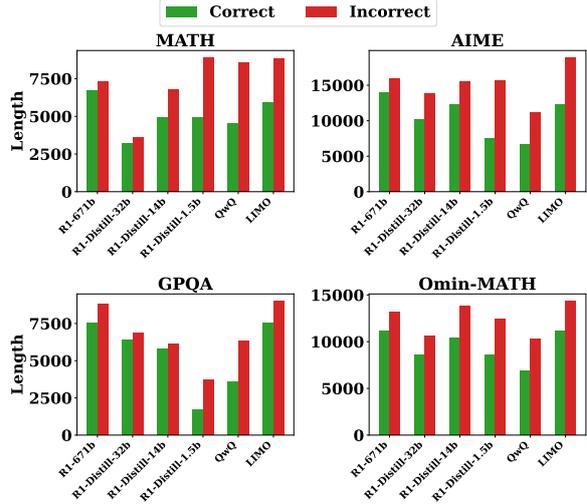


Figure 1: The average length of correct solutions versus incorrect solutions evaluated on the same questions. For each question, solution lengths were averaged separately for correct and incorrect responses, then averaged across all questions.

et al., 2025) and LIMO (Ye et al., 2025) have emerged as leading open-source successors, replicating o1’s achievements and demonstrating comparable reasoning abilities. Although both QwQ, R1 and LIMO demonstrate strong reasoning capabilities and the ability to generate lengthy CoT at test time, the existence of **true test-time scaling where performance consistently improves with longer CoTs** remains to be verified for these models.

To explore this question, we systematically investigate the relationship between CoT length and reasoning performance in QwQ, R1 and LIMO, challenging the conventional assumption that extended reasoning chains inherently lead to improved accuracy. Contrary to expectations, our analysis reveals that longer CoTs do not consistently improve accuracy of these o1-like models. Notably, we found that the average length of correct solutions is shorter than that of incorrect ones for the same ques-

tions, which is shown in Figure 1. This counterintuitive finding underscores the need for a deeper understanding of the test-time scaling of o1-like models.

To understand why the longer CoTs do not lead to the better performance, we compared the difference between long CoTs and short CoTs, finding that long CoTs contain more self-revisions (“Wait”, “Alternatively”) than the short CoTs, which is shown in Appendix E. Inspired by that, we iteratively prompted QwQ, R1 and LIMO for more self-revisions. Our observations revealed that QwQ and R1-Distill-1.5b exhibited performance degradation as the length of reflection increased. In contrast, R1-Distill-14b, R1-Distill-32b, and LIMO demonstrated initial performance improvements during early revisions, followed by oscillatory behavior in subsequent iterations. To further understand the limitations of sequential scaling, we evaluated the models’ capacity to revise incorrect answers. Our findings indicate that QwQ, R1 and LIMO all demonstrated limited ability to convert incorrect answers to correct ones during the revision process. Most revisions retained the original answers, and more concerning, both QwQ and R1-Distill-1.5b showed a higher propensity to change correct answers to incorrect ones rather than vice versa. These results reveal that **self-revision ability is a key factor in the effectiveness of sequential scaling for o1-like models.**

Given the limited effectiveness of sequential scaling, we explored an alternative test-time scaling strategie, parallel scaling. Our comparative analysis of sequential and parallel scaling revealed that parallel scaling not only achieves the better coverage (pass@k score) but also offers superior scalability compared to sequential scaling for QwQ and R1, which demonstrates that o1-like models have limited sequential-scaling capability, but strong parallel-scaling capability.

Building on these findings, we propose a novel test-time scaling method, **Shortest Majority Vote**, which incorporate parallel scaling approaches with our insight on sequential scaling. In particular, this method leverages the observation that shorter solutions tend to lead to better performance compared to longer ones. Shortest Majority Vote improves majority vote by prioritizing clusters that have both more solutions and shorter solution lengths. Experimental results demonstrate that Shortest Majority Vote substantially outperforms conventional Majority Vote, significantly improving the test-time

scalability of both QwQ and R1 models.

Our contributions are as follows:

- 1) We systematically investigate the test-time scaling capabilities of o1-like models QwQ, R1 and LIMO, and find that their performance can not be continuously improved through increasing CoT length.
- 2) We reveal that insufficient self-revision capability of o1-like models is the primary reason for their failure in sequential scaling.
- 3) We find that parallel scaling achieves better coverage and scalability than sequential revision for o1-like models.
- 4) Based on our insights into sequential and parallel scaling, we propose Shortest Majority Vote, a test-time scaling method that enhances majority voting by considering solution length, significantly outperforming traditional methods.

## 2 Related Work

The success of o1 has ushered in a new scaling paradigm, **test-time compute scaling**, which enables continuous improvements in model performance by increasing computational expenditure during inference (OpenAI, 2024a,b). Currently, scaling test-time compute can be approached in two dimensions: parallel scaling and sequential scaling (Snell et al., 2024; Zeng et al., 2024).

**Parallel Scaling** Parallel scaling typically samples multiple solutions in parallel and pick one according to some guidance signal like reward. Notable examples of parallel scaling include Best-of-N Search (Cobbe et al., 2021; Sun et al., 2024; Gui et al., 2024; Amini et al., 2024; Sessa et al., 2024), which is based on a reward model (Cobbe et al., 2021; Lightman et al., 2024), and Majority Vote (Wang et al., 2023), which exploits model uncertainty. The primary distinction between these approaches lies in the method used to select the final solution or answer after sampling multiple candidates. Both Best-of-N Search and Majority Vote are parallel scaling techniques at the solution level, while Tree-Search algorithms can be viewed as parallel scaling at the token or step level. Beam-Search (Qiu et al., 2024; Yu et al., 2024; Xie et al., 2023; Kool et al., 2019) and MCTS (Hao et al., 2023; Wan et al., 2024; Chen et al., 2024a; Zhang et al., 2023) are classic examples of Tree-Search algorithms. All parallel scaling methods rely on

163 guidance signals to select the optimal token, step,  
164 or solution from a set of candidates.

165 **Sequential Scaling** Sequential scaling enhances  
166 test-time computation by generating progressively  
167 longer solutions along the sequence dimension.  
168 The most prevalent method of sequential scaling  
169 is Self-Revision, where Madaan et al. (2023) first  
170 generate an initial response and then iteratively  
171 evaluate and refine it based on self-assessment. In  
172 contrast, Chen et al. (2024b); Gou et al. (2024)  
173 leverage external feedback—such as signals from  
174 a code execution environment—rather than self-  
175 evaluation to enhance solutions.

176 The effectiveness of sequential scaling with self-  
177 revision remains a contentious issue. Huang et al.  
178 (2024a); Kamoi et al. (2024) argue that models  
179 cannot achieve effective self-refinement without  
180 external feedback. Conversely, some researchers  
181 posit that evaluating a solution’s correctness is in-  
182 herently easier than generating a correct solution  
183 (Leike, 2022), suggesting that LLMs have the ca-  
184 pacity for self-evaluation. Kumar et al. (2024);  
185 Zhang et al. (2024) show that it is possible to teach  
186 LLM to self-refine through reinforcement learn-  
187 ing or supervised fine-tuning. Chen et al. (2024c)  
188 compared various test-time scaling algorithms and  
189 found that when feedback accuracy exceeds 90%,  
190 Self-Revision outperforms Best-of-N Search.

191 **o1-like Models** The release of o1 (OpenAI,  
192 2024a,b) has further underscored the significance  
193 of sequential scaling, as o1’s CoT length is sub-  
194 stantially greater than that of conventional models.  
195 The research community has made significant ef-  
196 forts to reproduce the capabilities of o1 (Qin et al.,  
197 2024; Huang et al., 2024b; Jiang et al., 2024; Min  
198 et al., 2024; Muennighoff et al., 2025), with QwQ  
199 (Team, 2024) and R1 (DeepSeek-AI et al., 2025)  
200 and LIMO (Ye et al., 2025) emerging as the most  
201 successful attempts. However, Our findings reveal  
202 that for R1 and QwQ, extending solution length  
203 does not necessarily yield better performance due  
204 to the models’ limited self-revision capabilities.  
205 Parallel findings by Wang et al. (2025) attribute this  
206 phenomenon to model underthinking, where mod-  
207 els initially reach correct intermediate solutions but  
208 subsequently deviate toward incorrect conclusions  
209 during extended reasoning.

### 3 Experiment Setting 210

211 **Models** Our experiments involved models from  
212 the QwQ (Team, 2024), LIMO (Ye et al., 2025)  
213 and Deepseek-R1 series (DeepSeek-AI et al.,  
214 2025), including Deepseek-R1, Deepseek-R1-  
215 Distill-Qwen-32b, Deepseek-R1-Distill-Qwen-14b,  
216 and Deepseek-R1-Distill-Qwen-1.5b. For simpli-  
217 city, we call these R1 models as R1-671b, R1-  
218 Distill-32b, R1-Distill-14b and R1-Distill-1.5b re-  
219 spectively. The models were run using SGLang  
220 framework (Zheng et al., 2024), with the sampling  
221 temperature set to 0.7 and the maximum generation  
222 length set to 32k. We show the system prompt and  
223 instructions used for evaluation in Appendix D.

224 **Benchmark** We conducted comprehensive evalu-  
225 ations across four benchmarks: MATH-500 (Light-  
226 man et al., 2024), AIME (AIMO, 2018), Omini-  
227 MATH (Gao et al., 2024), and GPQA (Rein et al.,  
228 2023). While MATH-500, AIME, and Omini-  
229 MATH focus on mathematical reasoning, GPQA  
230 encompasses broader scientific domains. For  
231 AIME evaluation, we utilized the AIMO validation  
232 set, comprising 90 questions from AIME 22, 23,  
233 and 24 (AIMO, 2018). Given the computational de-  
234 mands of evaluating the full Omini-MATH dataset  
235 (4.4K questions), we randomly sampled 500 ques-  
236 tions to maintain efficiency. For GPQA, we focused  
237 on the diamond subset containing 198 questions.  
238 To ensure robust evaluation of answer correctness,  
239 we employed both the OpenCompass (Contributors,  
240 2023) and Qwen Math (Yang et al., 2024) evalua-  
241 tors, considering an answer correct if validated by  
242 either evaluator.

### 4 The Failure of Sequential Scaling 243

#### 4.1 Invalid Scaling of CoT Length: Longer CoTs Do not Improve Performance 244

245 To investigate whether the accuracy of QwQ, R1  
246 and LIMO genuinely improves with increasing  
247 CoT length, we sampled each model five times  
248 on the same question and sorted the five solutions  
249 by length in ascending order. We grouped the so-  
250 lutions based on their rank in this sorted list, with  
251 the  $i$ -th ranked solutions forming a distinct group.  
252 For instance, all the longest solutions (rank 5) from  
253 different questions formed one group, while all the  
254 shortest solutions (rank 1) formed another, result-  
255 ing in 5 comprehensive solution groups for analy-  
256 sis.

257 We present the average lengths of the five groups  
258

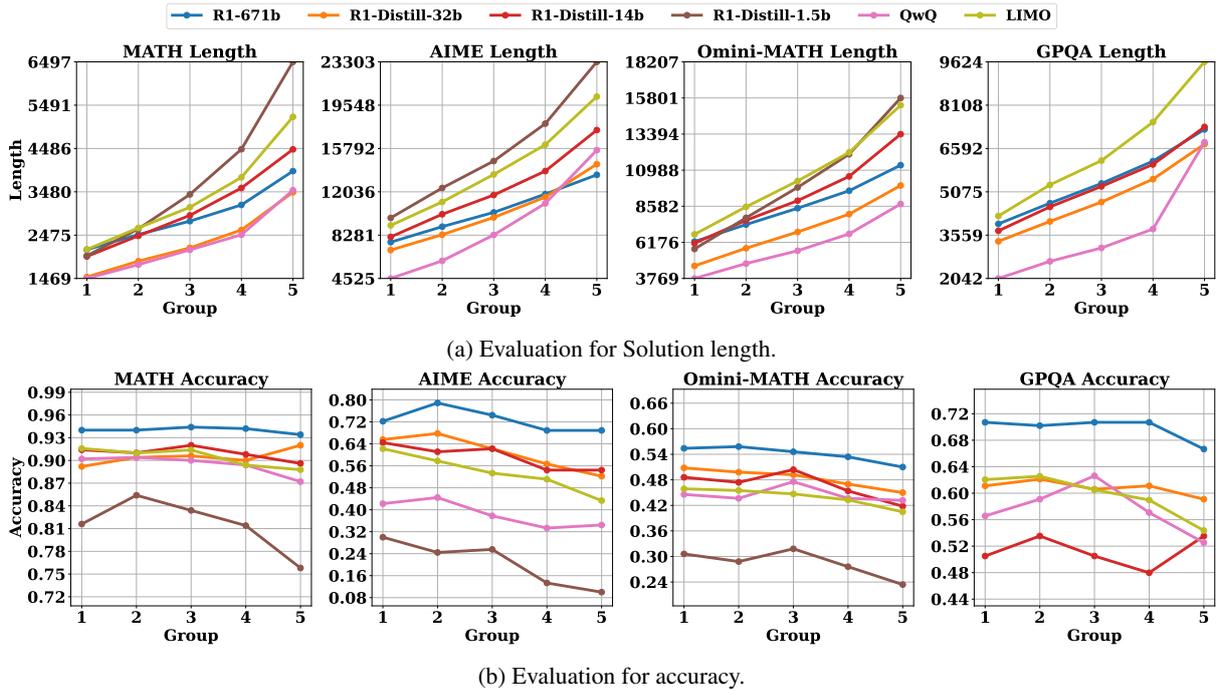


Figure 2: Solutions of QwQ and R1 were categorized into different groups according to their length and evaluated in terms of solution length (a) and accuracy (b). The categorization of solutions is progressed for each question independently, i.e., all groups of solutions are corresponding to the same questions.

of solutions in Figure 2a. Since the grouping of solutions is based on their lengths, the differences in length between the groups are pronounced. The average length of the longest solutions is approximately twice that of the shortest solutions. This indicates that long-chain-of-thought (CoT) models like QwQ, R1 and LIMO exhibit a high diversity in the lengths of the solutions they sample.

There is no clear correlation between the length of solutions and the model’s size. For example, R1-Distill-1.5b produces the longest solutions while QwQ (32b) generates the shortest. A comparison of solution lengths across different datasets shows that solutions for simpler datasets, such as Math, are significantly shorter than those for more difficult datasets, like AIME. This suggests that the model adjusts the solution length based on the difficulty of the problem.

The accuracy of the five groups of solutions is presented in Figure 2b. Although there is a significant disparity in solution lengths across the groups, the differences in accuracy are much less pronounced. Notably, we do not observe a consistent improvement in accuracy for either QwQ or R1 as solution length increases. This trend holds true across all model variants as well as across all evaluated datasets. In some cases, we even

observe an inverse scaling phenomenon, where accuracy decreases with increasing CoT length, especially on more difficult datasets like AIME and Omini-MATH. These findings cast doubt on the presumed test-time scaling capabilities of o1-like models, challenging the assumption that extended reasoning chains inherently yield superior problem-solving performance.

To make the relationship between CoT length and accuracy more clear, we compared the lengths of correct and incorrect solutions for the same question. First, we identified questions that had both correct and incorrect answers. For each of these questions, we calculated the average length of correct and incorrect solutions. We then averaged these values across all questions to determine the overall average length for correct and incorrect solutions. The results are shown in Figure 1. We found that, for QwQ, R1 and LIMO, across all model sizes and datasets, the length of correct solutions is consistently shorter than that of incorrect solutions. This observation suggests that longer CoTs do not necessarily lead to better performance and may even be associated with lower accuracy. Moreover, we observed that for weaker models, such as QwQ and R1-Distill-1.5B, the gap in solution length between correct and incorrect solutions is

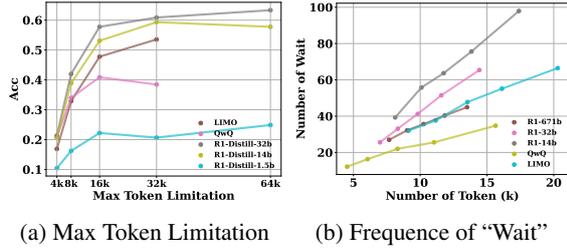


Figure 3: (a): The relationship between model accuracy and the generation parameter Max Token Limitation. (b): The relationship between solution length and the average number of “wait” occur in a solution.

significantly larger than for stronger models, such as R1-671b. This suggests that the invalid scaling phenomenon is more pronounced in the weaker models.

## 4.2 Explaining Invalid Scaling: The Key Factor is the Failure of Self-Revision

In Section 4.1, we observed the phenomenon that long solutions exhibit lower accuracy compared to short solutions. In this section, we investigate the underlying reasons for this phenomenon. We first analyzed how the maximum token limitation affects generation performance and confirmed that the observed invalid scaling phenomenon was not caused by constraints in the maximum token length. Next, we examined the differences between long and short solutions, finding that long solutions exhibit a higher frequency of self-revision. Moreover, our analysis suggests a strong correlation between self-revision, solution length, and accuracy.

**Max Token Limitation** The max token limitation parameter controls the maximum number of tokens a model can generate for a question, which plays a critical role in influencing model accuracy, especially when generating long solutions. To explore its impact, we tested several max token limitation values and compared the performance of QwQ, R1 and LIMO on the AIME benchmark. The results are shown in Figure 3a, which revealed that 16k is a key threshold: when the max token limitation is below this value, it significantly affects the model performance. However, increasing the max token limitation beyond 16k leads to diminishing returns, particularly for QwQ. In our other experiments, we set the max token limitation to 32k, suggesting that this parameter is not the main cause of invalid scaling.

**Difference between Short and Long CoT** To understand why long solutions of QwQ, R1 and LIMO is not better than short solutions, we analyzed their differences. We observed that QwQ, R1 and LIMO all primarily extend solution length through self-revision, characterized by markers such as “Wait” and “Alternatively”. We show some examples of that in Appendix E. To quantify this phenomenon, we counted the occurrences of “wait” in solutions of QwQ, R1 and LIMO in Figure 3b. The results demonstrate a strong linear correlation between solution length and the frequency of self-correction markers for all models. This suggests that the mechanisms of self-revision may play a significant role in generating longer solutions.

**Scaling Solution Length with Self-Revision** We have tried to investigate the revision behaviors inside the sampled solutions, however, it is difficult to extract the initial solution and the following revision exactly from QwQ, R1 and LIMO’s solutions. Alternatively to that, we prompted the models to continue thinking based on their sampled solutions.

QwQ, R1 and LIMO often conclude their solutions with phrases like “final answer: ...”, and R1 additionally outputs a ‘</think>’ tag followed by a final response. To facilitate smoother continuation of the reasoning process, we removed the “final answer” portion from the solutions. We then used the keyword “Wait” or “Alternatively” as the prompt to encourage self-revision. We calculated the probabilities of the model predicting the next token as “Wait” or “Alternatively” and selected the one with the higher probability as the prompt.

We prompted QwQ, R1 and LIMO to continue reasoning for 40 additional steps on the AIME benchmark. We show the results in Figure 4c, from which we observe that the solution length increase almost linearly with additional steps. After 40 steps, the solution length of QwQ and R1 is almost third as their original length.

We show the accuracy after sequential revision in Figure 4a and 4b. Our results reveal that the accuracy of QwQ and R1-Distill-1.5b decreases constantly as the number of reasoning steps increases, while the accuracy of R1-Distill-32b, R1-Distill-14b and LIMO initially improves and then oscillates with further reasoning steps. Further analysis in Appendix B reveal that the improvement on R1-Distill-32b, R1-Distill-14b and LIMO during revisions mainly comes from the revision on short solutions. These results corroborate our previous

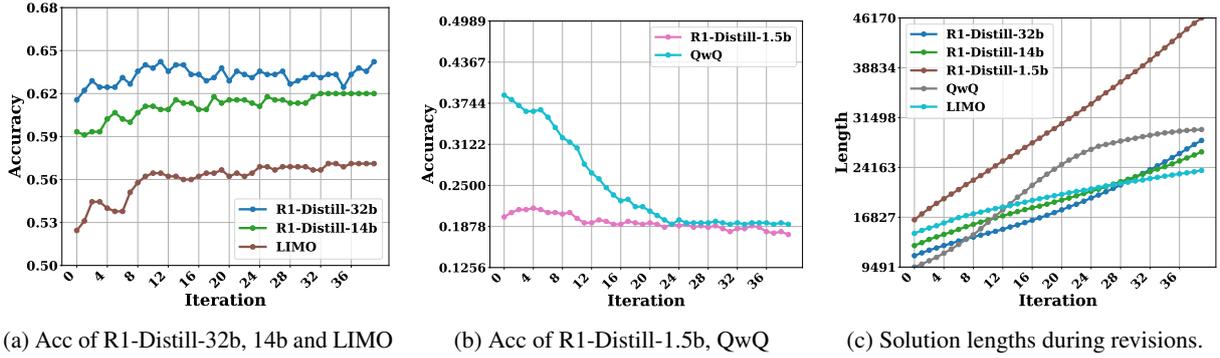


Figure 4: (a): Accuracy of R1-Distill-32b, R1-Distill-14b and LIMO during sequential revisions. (b): Accuracy of R1-Distill-1.5b and QwQ during sequential revisions. (c) Solution length increased with the more revision steps.

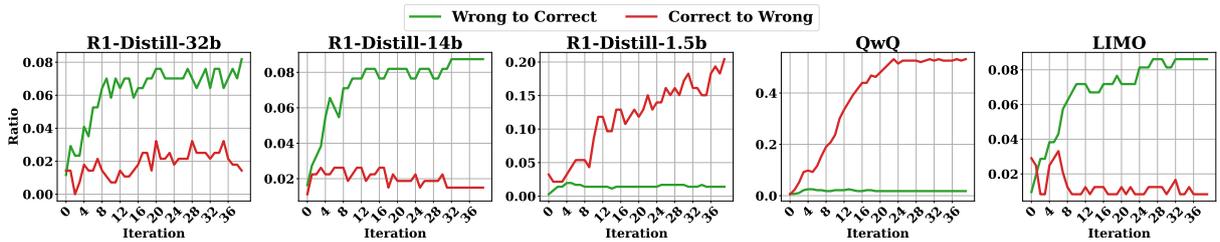


Figure 5: The ratio of turning an initial correct answer to incorrect one (correct to wrong) and an initial incorrect answer to a correct one (wrong to correct) during sequential scaling.

experimental findings, suggesting that longer solutions do not improve performance, especially for weaker models such as QwQ and R1-Distill-1.5b. These findings suggest that the reason why longer solutions do not consistently lead to better performance in QwQ, R1 and LIMO may lie in the failure of self-revision.

**Investigating Self-Revision Behavior** To further investigate the effectiveness of self-revision, we analyzed the proportion of cases where the model corrected an initial incorrect answer to a correct one versus changing an initial correct answer to an incorrect one during scaling solution length. We found that, the proportions of changing a incorrect answer to an correct one is extremely low, always below 10%. Notably, for QwQ and R1-Distill-1.5b, the proportion of changing a correct answer to an incorrect one was even higher than that of correcting an incorrect answer to a correct one. This observation helps explain why prompting QwQ and R1-Distill-1.5b to continue reasoning led to a decrease in accuracy. For simplicity, we call the proportions of changing a incorrect answer to an correct one as the successful-revision rate, while the reverse as the failed-revision rate.

Although R1-Distill-32b, R1-Distill-14b and LIMO exhibit a higher successful-revision rate

R1-32b	R1-14b	R1-1.5b	QwQ	LIMO
72%	70%	58%	32%	54%

Table 1: The proportion of the revisions that models stick to the original wrong answers.

than failed-revision rate, the increase of successful-revision rate plateaus after approximately 10 steps, with further revisions providing no additional benefits. This observation explains why their accuracy during sequential scaling initially increases with multiple rounds of revision but later stabilizes with fluctuations.

The successful-revision rate of QwQ, R1 and LIMO are all below 10%, what is the outcome of the model’s self-revision in unsuccessful cases? We hypothesize that, in most instances, the model simply keeps its original answer unchanged. To validate that, we computed the proportion of instances where the model persists with its original answer, even when it is incorrect, and the results were as expected. As shown in Figure 5, when the original answer is wrong, both R1-Distill-32b and R1-Distill-14b maintain the original answer in over 70% of cases. Although retaining the original answer does not reduce accuracy, it also makes the scaling solution length ineffective. This phenomenon suggests

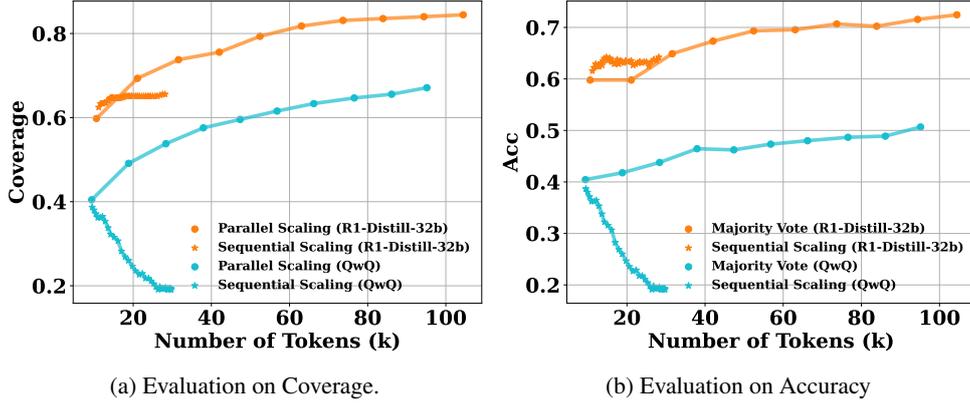


Figure 6: (a): the coverage of sequential scaling and parallel scaling on AIME. (b): the accuracy of sequential revision and majority vote on AIME.

that the model’s ability to early stop may also be a critical factor influencing whether its performance improves with an increasing solution length.

The above analysis indicates that the key factor determining whether o1-like models’ performance improve with an increase in solution length is their ability to self-revise. The model’s accuracy increases with the more incorrect answers revised to correct and vice versa.

## 5 Sequential Scaling vs. Parallel Scaling

Based on our experimental findings presented in Section 4.2, sequential scaling demonstrates limited effectiveness for QwQ, R1 and LIMO. An alternative approach to scaling test-time compute is parallel scaling, which generates multiple solutions in parallel and selects the best one as the final answer.

We compared the performance of sequential scaling and parallel scaling in terms of the coverage (pass@k score) and accuracy of QwQ and R1, which are shown in Figure 6a and 6b respectively. For sequential scaling, we iteratively prompt models to self-revise for 40 steps. While for parallel scaling, we parallelly sample 10 solutions. The coverage is evaluated by counting the proportion of whether multiple candidate answers contain a correct one. In parallel scaling, coverage increases by one if at least one sampled solution is correct. Similarly, in sequential scaling, coverage increases by one if at least one revision iteration succeeds.

Our findings show that, for the same number of generated tokens, parallel scaling provides a significantly larger improvement in coverage compared to sequential scaling, for both R1-Distill-32b and QwQ. However, a practical parallel scaling method

must select a final answer from a set of candidate answers. We implement parallel scaling using majority vote (Wang et al., 2023) and sequential scaling by taking the answer from the last revision as the final answer. Since majority voting requires at least three solutions to be effective, it does not provide any benefit when scaling the number of solutions from 1 to 2. In contrast, sequential revision is effective for R1-Distill-32b when scaling the number of tokens to 10k, but further scaling does not yield additional benefits. Additionally, because sequential scaling involves attention over a longer context, its computational cost is much higher than that of parallel scaling when generating the same number of tokens.

## 6 Application of Our Findings: Shortest Majority Vote

Given the limitation of sequential scaling of the current o1-like models, we turn to parallel scaling techniques and incorporate it with our insight on sequential scaling. Specifically, we propose a new Parallel Scaling algorithm: Shortest Majority Vote. Shortest Majority Vote is an extension of Majority Vote, but it accounts for the length of the solutions generated by the model. In the original Majority Vote, solutions with the same answer are grouped into a single category, and the number of solutions in each category is counted, with the answer corresponding to the category with the most solutions selected as the final answer. In contrast, Shortest Majority Vote not only counts the number of solutions in each category, but also computes the average length of the solutions in each category. Let the number of solutions in the  $i$ -th category be  $c_i$  and the average solution length in that category

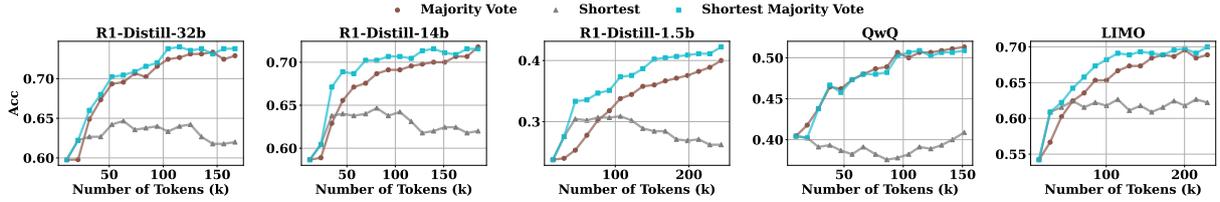


Figure 7: Parallel-scaling performance of Majority Vote, Shortest and Shortest Majority Vote on AIME.

Model	Solutions	AIME			GPQA		
		MV	Shortest	Shortest MV	MV	Shortest	Shortest MV
R1-Distill-32b		59.77	<b>62.22</b>	<b>62.22</b>	61.41	<b>62.52</b>	<b>62.52</b>
R1-Distill-14b		58.88	<b>60.44</b>	<b>60.44</b>	51.21	<b>52.32</b>	<b>52.32</b>
R1-Distill-1.5b	2	24	<b>27.55</b>	<b>27.55</b>	15.25	<b>15.35</b>	<b>15.35</b>
QwQ		<b>41.77</b>	40.22	40.22	<b>58.05</b>	57.02	57.02
LIMO		56.66	<b>60.88</b>	<b>60.88</b>	50.46	<b>54.56</b>	<b>54.56</b>
R1-Distill-32b		72.88	61.99	<b>73.77</b>	63.33	61.21	<b>63.53</b>
R1-Distill-14b		<b>71.77</b>	62.00	71.55	56.16	<b>56.66</b>	56.46
R1-Distill-1.5b	16	40.00	26.22	<b>42.22</b>	29.59	27.77	<b>30.20</b>
QwQ		<b>51.33</b>	40.88	50.88	<b>62.25</b>	56.82	<b>62.25</b>
LIMO		68.88	62.22	<b>70.00</b>	55.58	50.15	<b>55.89</b>

Table 2: Performance comparison between Majority Vote (MV), Shortest and Shortest Majority Vote (Shortest MV) on AIME and GPQA, when there are 2 and 16 solutions sampled.

be  $l_i$ . The score for category  $i$  in Shortest Majority Vote is computed as:

$$s_i = \frac{c_i}{\log l_i} \quad (1)$$

and the final answer is chosen from the category with the highest score. The score  $s_i$  is designed with the assumption that the correct answer is more likely to appear in categories with a larger number of solutions and shorter solution lengths. Shortest Majority Vote offers two key advantages: first, it is particularly effective for some o1-like models, where performance deteriorates with increasing solution length; second, it enables the use of solution length as a guidance signal for identifying superior solutions when candidate solutions are limited, especially in cases where conventional Majority Vote becomes ineffective due to having only two candidate solutions.

We evaluated the performance of Shortest Majority Vote and Majority Vote through experiments on the AIME and GPQA benchmarks, sampling 16 solutions from QwQ, R1 and LIMO models. We implemented a simple baseline approach, denoted as "Shortest," which selects the answer from the solution with the minimal length. The experimental results are presented in Table 2 and Figure 7.

Table 2 demonstrates that Shortest Majority Vote significantly outperforms both Majority Vote and Shortest methods, particularly on the AIME benchmark. Figure 7 illustrates the parallel-scaling performance of these three methods, showing that as the number of generated tokens increases, Shortest Majority Vote maintains superior performance over both alternatives on AIME. The corresponding parallel-scaling results for GPQA are provided in Appendix C. Notably, while Shortest performs better than Majority Vote when only two solutions are sampled, it exhibits inferior performance in all other scenarios. These empirical findings strongly support the effectiveness of the Shortest Majority Vote approach.

## 7 Conclusion

In this study, we challenged the assumption that o1-like models like QwQ and R1 have test-time scaling capability. We found that the longer solutions not necessarily yield better performance, and that sequential scaling through self-revision has limited effectiveness. Based on these insights, we developed Shortest Majority Vote, a parallel scaling method that considers solution length, which significantly outperformed traditional majority vote.

## 568 Limitations

- 569 1. Given the considerable cost of R1-671b, eval-  
570 uation on it was limited to the experiments  
571 in Figures 1 and 2, whereas distilled R1 was  
572 utilized for all subsequent.
- 573 2. Our experimental framework was limited to  
574 static model checkpoints. Future research  
575 should investigate test-time scaling behavior  
576 using dynamic checkpoints in reinforcement  
577 learning settings.
- 578 3. While the proposed shortest majority method  
579 may have limited applicability for models  
580 with strong sequential-scaling capabilities, so-  
581 lution length remains a valuable guidance sig-  
582 nal for candidate selection in parallel scaling  
583 scenarios. The method can be adapted to a  
584 Longest Majority Vote variant for such cases.

## 585 Ethics Statement

586 This paper honors the ACL Code of Ethics. The  
587 dataset used in the paper does not contain any pri-  
588 vate information. All data and tools used in this  
589 study comply with their respective licenses and  
590 terms of use.

## 591 References

- 592 AIMO. 2018. Dataset card for aimo valida-  
593 tion aime. [https://huggingface.co/datasets/  
594 AI-M0/aimo-validation-aime](https://huggingface.co/datasets/AI-M0/aimo-validation-aime).
- 595 Afra Amini, Tim Vieira, and Ryan Cotterell. 2024. [Vari-  
596 ational best-of-n alignment](#). *CoRR*, abs/2407.06057.
- 597 Guoxin Chen, Minpeng Liao, Chengxi Li, and Kai Fan.  
598 2024a. [Alphamath almost zero: process supervision  
599 without process](#). *CoRR*, abs/2405.03553.
- 600 Xinyun Chen, Maxwell Lin, Nathanael Schärli, and  
601 Denny Zhou. 2024b. [Teaching large language mod-  
602 els to self-debug](#). In *The Twelfth International Con-  
603 ference on Learning Representations, ICLR 2024,  
604 Vienna, Austria, May 7-11, 2024*. OpenReview.net.
- 605 Ziru Chen, Michael White, Raymond J. Mooney, Ali  
606 Payani, Yu Su, and Huan Sun. 2024c. [When is  
607 tree search useful for LLM planning? it depends  
608 on the discriminator](#). In *Proceedings of the 62nd  
609 Annual Meeting of the Association for Computa-  
610 tional Linguistics (Volume 1: Long Papers), ACL  
611 2024, Bangkok, Thailand, August 11-16, 2024*, pages  
612 13659–13678. Association for Computational Lin-  
613 guistics.

- Karl Cobbe, Vineet Kosaraju, Mohammad Bavarian, 614  
Mark Chen, Heewoo Jun, Lukasz Kaiser, Matthias 615  
Plappert, Jerry Tworek, Jacob Hilton, Reiichiro 616  
Nakano, Christopher Hesse, and John Schulman. 617  
2021. [Training verifiers to solve math word prob-  
618 lems](#). *CoRR*, abs/2110.14168. 619
- OpenCompass Contributors. 2023. Opencompass: 620  
A universal evaluation platform for foundation 621  
models. [https://github.com/open-compass/  
622 opencompass](https://github.com/open-compass/opencompass). 623
- DeepSeek-AI, Daya Guo, Dejian Yang, Haowei Zhang, 624  
Junxiao Song, Ruoyu Zhang, Runxin Xu, Qihao Zhu, 625  
Shirong Ma, Peiyi Wang, Xiao Bi, Xiaokang Zhang, 626  
Xingkai Yu, Yu Wu, Z. F. Wu, Zhibin Gou, Zhi- 627  
hong Shao, Zhuoshu Li, Ziyi Gao, and 181 others. 628  
2025. [Deepseek-r1: Incentivizing reasoning capa-  
629 bility in llms via reinforcement learning](#). *Preprint*,  
630 arXiv:2501.12948. 631
- Bofei Gao, Feifan Song, Zhe Yang, Zefan Cai, Yibo 632  
Miao, Qingxiu Dong, Lei Li, Chenghao Ma, Liang 633  
Chen, Runxin Xu, Zhengyang Tang, Benyou Wang, 634  
Daoguang Zan, Shanghaoran Quan, Ge Zhang, Lei 635  
Sha, Yichang Zhang, Xuancheng Ren, Tianyu Liu, 636  
and Baobao Chang. 2024. [Omni-math: A univer-  
637 sal olympiad level mathematic benchmark for large  
638 language models](#). *CoRR*, abs/2410.07985. 639
- Zhibin Gou, Zhihong Shao, Yeyun Gong, Yelong Shen, 640  
Yujia Yang, Nan Duan, and Weizhu Chen. 2024. [CRITIC: large language models can self-correct with  
641 tool-interactive critiquing](#). In *The Twelfth Inter-  
642 national Conference on Learning Representations,  
643 ICLR 2024, Vienna, Austria, May 7-11, 2024*. Open-  
644 Review.net. 645
- Lin Gui, Cristina Gârbacea, and Victor Veitch. 2024. 647  
[Bonbon alignment for large language models and  
648 the sweetness of best-of-n sampling](#). *CoRR*,  
649 abs/2406.00832. 650
- Shibo Hao, Yi Gu, Haodi Ma, Joshua Jiahua Hong, Zhen 651  
Wang, Daisy Zhe Wang, and Zhiting Hu. 2023. [Reason-  
652 ing with language model is planning with world  
653 model](#). In *Proceedings of the 2023 Conference on  
654 Empirical Methods in Natural Language Process-  
655 ing, EMNLP 2023, Singapore, December 6-10, 2023*,  
656 pages 8154–8173. Association for Computational  
657 Linguistics. 658
- Jie Huang, Xinyun Chen, Swaroop Mishra, 659  
Huaixiu Steven Zheng, Adams Wei Yu, Xiny- 660  
ing Song, and Denny Zhou. 2024a. [Large language  
661 models cannot self-correct reasoning yet](#). In *The  
662 Twelfth International Conference on Learning  
663 Representations, ICLR 2024, Vienna, Austria, May  
664 7-11, 2024*. OpenReview.net. 665
- Zhen Huang, Haoyang Zou, Xuefeng Li, Yixiu Liu, 666  
Yuxiang Zheng, Ethan Chern, Shijie Xia, Yiwei Qin, 667  
Weizhe Yuan, and Pengfei Liu. 2024b. [O1 replica-  
668 tion journey – part 2: Surpassing o1-preview through  
669 simple distillation, big progress or bitter lesson?](#)  
670 *Preprint*, arXiv:2411.16489. 671

672	Jinhao Jiang, Zhipeng Chen, Yingqian Min, Jie Chen,	OpenAI. 2024a. <a href="#">Learning to reason with llms.</a>	729
673	Xiaoxue Cheng, Jiapeng Wang, Yiru Tang, Hao-	OpenAI. 2024b. <a href="#">Openai o1 system card.</a>	730
674	xiang Sun, Jia Deng, Wayne Xin Zhao, and 1 oth-	Yiwei Qin, Xuefeng Li, Haoyang Zou, Yixiu Liu, Shijie	731
675	ers. 2024. Technical report: Enhancing llm reason-	Xia, Zhen Huang, Yixin Ye, Weizhe Yuan, Hector	732
676	ing with reward-guided tree search. <i>arXiv preprint</i>	Liu, Yuanzhi Li, and Pengfei Liu. 2024. <a href="#">O1 repli-</a>	733
677	<i>arXiv:2411.11694.</i>	<a href="#">cation journey: A strategic progress report - part 1.</a>	734
678	Ryo Kamoi, Yusen Zhang, Nan Zhang, Jiawei Han,	<i>CoRR</i> , abs/2410.18982.	735
679	and Rui Zhang. 2024. <a href="#">When can llms actually cor-</a>	Jiahao Qiu, Yifu Lu, Yifan Zeng, Jiacheng Guo, Ji-	736
680	<a href="#">rect their own mistakes? A critical survey of self-</a>	ayi Geng, Huazheng Wang, Kaixuan Huang, Yue	737
681	<a href="#">correction of llms.</a> <i>CoRR</i> , abs/2406.01297.	Wu, and Mengdi Wang. 2024. <a href="#">Treebon: Enhancing</a>	738
682	Wouter Kool, Herke van Hoof, and Max Welling. 2019.	<a href="#">inference-time alignment with speculative tree-search</a>	739
683	<a href="#">Stochastic beams and where to find them: The</a>	<a href="#">and best-of-n sampling.</a> <i>CoRR</i> , abs/2410.16033.	740
684	<a href="#">gumbel-top-k trick for sampling sequences without</a>	David Rein, Betty Li Hou, Asa Cooper Stickland,	741
685	<a href="#">replacement.</a> In <i>Proceedings of the 36th Interna-</i>	Jackson Petty, Richard Yuanzhe Pang, Julien Di-	742
686	<i>tional Conference on Machine Learning, ICML 2019,</i>	rani, Julian Michael, and Samuel R. Bowman. 2023.	743
687	<i>9-15 June 2019, Long Beach, California, USA,</i> vol-	<a href="#">GPQA: A graduate-level google-proof q&amp;a bench-</a>	744
688	ume 97 of <i>Proceedings of Machine Learning Re-</i>	<a href="#">mark.</a> <i>CoRR</i> , abs/2311.12022.	745
689	<i>search</i> , pages 3499–3508. PMLR.	Pier Giuseppe Sessa, Robert Dadashi, Léonard	746
690	Aviral Kumar, Vincent Zhuang, Rishabh Agarwal, Yi Su,	Husseinot, Johan Ferret, Nino Vieillard, Alexandre	747
691	John D. Co-Reyes, Avi Singh, Kate Baumli, Shariq	Ramé, Bobak Shahriari, Sarah Perrin, Abe Friesen,	748
692	Iqbal, Colton Bishop, Rebecca Roelofs, Lei M.	Geoffrey Cideron, Sertan Girgin, Piotr Stanczyk,	749
693	Zhang, Kay McKinney, Disha Shrivastava, Cosmin	Andrea Michi, Danila Sinopalnikov, Sabela Ramos,	750
694	Paduraru, George Tucker, Doina Precup, Feryal M. P.	Amélie Héliou, Aliaksei Severyn, Matt Hoffman,	751
695	Behbahani, and Aleksandra Faust. 2024. <a href="#">Training</a>	Nikola Momchev, and Olivier Bachem. 2024.	752
696	<a href="#">language models to self-correct via reinforcement</a>	<a href="#">BOND: aligning llms with best-of-n distillation.</a>	753
697	<a href="#">learning.</a> <i>CoRR</i> , abs/2409.12917.	<i>CoRR</i> , abs/2407.14622.	754
698	Jan Leike. 2022. <a href="#">Why i’m excited about ai-assisted</a>	Charlie Snell, Jaehoon Lee, Kelvin Xu, and Aviral Ku-	755
699	<a href="#">human feedback.</a>	mar. 2024. <a href="#">Scaling LLM test-time compute optimally</a>	756
700	Hunter Lightman, Vineet Kosaraju, Yuri Burda, Harri-	<a href="#">can be more effective than scaling model parameters.</a>	757
701	son Edwards, Bowen Baker, Teddy Lee, Jan Leike,	<i>CoRR</i> , abs/2408.03314.	758
702	John Schulman, Ilya Sutskever, and Karl Cobbe.	Hanshi Sun, Momin Haider, Ruiqi Zhang, Huitao	759
703	2024. <a href="#">Let’s verify step by step.</a> In <i>The Twelfth In-</i>	Yang, Jiahao Qiu, Ming Yin, Mengdi Wang, Peter	760
704	<i>ternational Conference on Learning Representations,</i>	L. Bartlett, and Andrea Zanette. 2024. <a href="#">Fast</a>	761
705	<i>ICLR 2024, Vienna, Austria, May 7-11, 2024.</i> Open-	<a href="#">best-of-n decoding via speculative rejection.</a> <i>CoRR</i> ,	762
706	Review.net.	abs/2410.20290.	763
707	Aman Madaan, Niket Tandon, Prakhar Gupta, Skyler	Qwen Team. 2024. <a href="#">Qwq: Reflect deeply on the bound-</a>	764
708	Hallinan, Luyu Gao, Sarah Wiegrefe, Uri Alon,	<a href="#">aries of the unknown.</a>	765
709	Nouha Dziri, Shrimai Prabhunoye, Yiming Yang,	Ziyu Wan, Xidong Feng, Muning Wen, Stephen Marcus	766
710	Shashank Gupta, Bodhisattwa Prasad Majumder,	McAleer, Ying Wen, Weinan Zhang, and Jun Wang.	767
711	Katherine Hermann, Sean Welleck, Amir Yazdan-	2024. <a href="#">Alphazero-like tree-search can guide large</a>	768
712	bakhsh, and Peter Clark. 2023. <a href="#">Self-refine: Itera-</a>	<a href="#">language model decoding and training.</a> In <i>Forty-</i>	769
713	<a href="#">tive refinement with self-feedback.</a> In <i>Advances in</i>	<i>first International Conference on Machine Learning,</i>	770
714	<i>Neural Information Processing Systems 36: Annual</i>	<i>ICML 2024, Vienna, Austria, July 21-27, 2024.</i> Open-	771
715	<i>Conference on Neural Information Processing Sys-</i>	Review.net.	772
716	<i>tems 2023, NeurIPS 2023, New Orleans, LA, USA,</i>	Xuezhi Wang, Jason Wei, Dale Schuurmans, Quoc V.	773
717	<i>December 10 - 16, 2023.</i>	Le, Ed H. Chi, Sharan Narang, Aakanksha Chowd-	774
718	Yingqian Min, Zhipeng Chen, Jinhao Jiang, Jie Chen,	hery, and Denny Zhou. 2023. <a href="#">Self-consistency</a>	775
719	Jia Deng, Yiwen Hu, Yiru Tang, Jiapeng Wang, Xi-	<a href="#">improves chain of thought reasoning in language</a>	776
720	xiaoxue Cheng, Huatong Song, and 1 others. 2024.	<a href="#">models.</a> In <i>The Eleventh International Conference</i>	777
721	<a href="#">Imitate, explore, and self-improve: A reproduction</a>	<i>on Learning Representations, ICLR 2023, Kigali,</i>	778
722	<a href="#">report on slow-thinking reasoning systems.</a> <i>arXiv</i>	<i>Rwanda, May 1-5, 2023.</i> OpenReview.net.	779
723	<i>preprint arXiv:2412.09413.</i>	Yue Wang, Qiuzhi Liu, Jiahao Xu, Tian Liang, Xingyu	780
724	Niklas Muennighoff, Zitong Yang, Weijia Shi, Xi-	Chen, Zhiwei He, Linfeng Song, Dian Yu, Juntao Li,	781
725	ang Lisa Li, Li Fei-Fei, Hannaneh Hajishirzi, Luke	Zhuosheng Zhang, Rui Wang, Zhaopeng Tu, Haitao	782
726	Zettlemoyer, Percy Liang, Emmanuel Candès, and	Mi, and Dong Yu. 2025. <a href="#">Thoughts are all over the</a>	783
727	Tatsunori Hashimoto. 2025. <a href="#">s1: Simple test-time</a>		
728	<a href="#">scaling.</a> <i>Preprint</i> , arXiv:2501.19393.		

784 [place: On the underthinking of o1-like llms.](#) *Preprint*,  
785 arXiv:2501.18585. 841

786 Jason Wei, Xuezhi Wang, Dale Schuurmans, Maarten  
787 Bosma, Brian Ichter, Fei Xia, Ed H. Chi, Quoc V. Le,  
788 and Denny Zhou. 2022. [Chain-of-thought prompting](#)  
789 [elicits reasoning in large language models.](#) In *Ad-*  
790 *vances in Neural Information Processing Systems 35:*  
791 *Annual Conference on Neural Information Process-*  
792 *ing Systems 2022, NeurIPS 2022, New Orleans, LA,*  
793 *USA, November 28 - December 9, 2022.* 842

794 Yuxi Xie, Kenji Kawaguchi, Yiran Zhao, James Xu  
795 Zhao, Min-Yen Kan, Junxian He, and Michael Qizhe  
796 Xie. 2023. [Self-evaluation guided beam search for](#)  
797 [reasoning.](#) In *Advances in Neural Information Pro-*  
798 *cessing Systems 36: Annual Conference on Neural*  
799 *Information Processing Systems 2023, NeurIPS 2023,*  
800 *New Orleans, LA, USA, December 10 - 16, 2023.* 843

801 An Yang, Beichen Zhang, Binyuan Hui, Bofei Gao,  
802 Bowen Yu, Chengpeng Li, Dayiheng Liu, Jian-  
803 hong Tu, Jingren Zhou, Junyang Lin, Keming Lu,  
804 Mingfeng Xue, Runji Lin, Tianyu Liu, Xingzhang  
805 Ren, and Zhenru Zhang. 2024. [Qwen2.5-math technical](#)  
806 [report: Toward mathematical expert model via](#)  
807 [self-improvement.](#) *CoRR*, abs/2409.12122. 844

808 Yixin Ye, Zhen Huang, Yang Xiao, Ethan Chern, Shijie  
809 Xia, and Pengfei Liu. 2025. [Limo: Less is more for](#)  
810 [reasoning.](#) *Preprint*, arXiv:2502.03387. 845

811 Fei Yu, Anningzhe Gao, and Benyou Wang. 2024. [Ovm,](#)  
812 [outcome-supervised value models for planning in](#)  
813 [mathematical reasoning.](#) In *Findings of the Associ-*  
814 *ation for Computational Linguistics: NAACL 2024,*  
815 *Mexico City, Mexico, June 16-21, 2024*, pages 858–  
816 875. Association for Computational Linguistics.

817 Zhiyuan Zeng, Qinyuan Cheng, Zhangyue Yin,  
818 Bo Wang, Shimin Li, Yunhua Zhou, Qipeng Guo,  
819 Xuanjing Huang, and Xipeng Qiu. 2024. [Scaling](#)  
820 [of search and learning: A roadmap to reproduce](#)  
821 [o1 from reinforcement learning perspective.](#) *CoRR*,  
822 abs/2412.14135.

823 Shun Zhang, Zhenfang Chen, Yikang Shen, Mingyu  
824 Ding, Joshua B. Tenenbaum, and Chuang Gan. 2023.  
825 [Planning with large language models for code](#)  
826 [generation.](#) In *The Eleventh International Conference*  
827 *on Learning Representations, ICLR 2023, Kigali,*  
828 *Rwanda, May 1-5, 2023.* OpenReview.net.

829 Yunxiang Zhang, Muhammad Khalifa, Lajanugen Lo-  
830 geswaran, Jaekyeom Kim, Moontae Lee, Honglak  
831 Lee, and Lu Wang. 2024. [Small language models](#)  
832 [need strong verifiers to self-correct reasoning.](#) In  
833 *Findings of the Association for Computational Lin-*  
834 *guistics, ACL 2024, Bangkok, Thailand and virtual*  
835 *meeting, August 11-16, 2024*, pages 15637–15653.  
836 Association for Computational Linguistics.

837 Lianmin Zheng, Liangsheng Yin, Zhiqiang Xie, Chuyue  
838 Sun, Jeff Huang, Cody Hao Yu, Shiyi Cao, Christos  
839 Kozyrakis, Ion Stoica, Joseph E. Gonzalez, Clark W.  
840 Barrett, and Ying Sheng. 2024. [Sglang: Efficient](#)

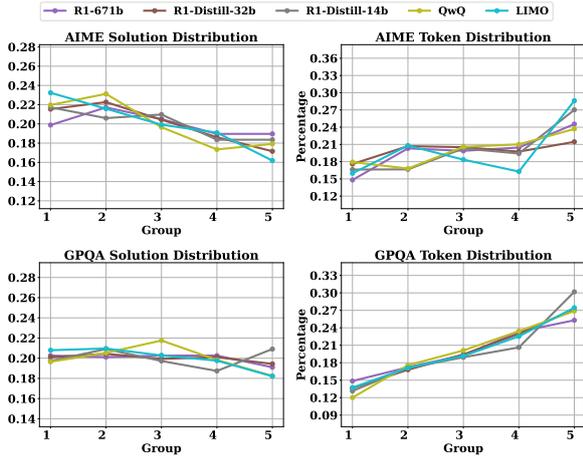


Figure 8: The number of correct solutions and tokens distributed across groups of different lengths.

## A Is Invalid Scaling Phenomenon Conflict to Findings of R1 technique Report?

The training objective of R1 aims to improve model accuracy, yet we observe that correct solutions tend to be shorter than incorrect ones. This raises an intriguing question: Why does R1’s reinforcement learning (RL) training consistently produce longer solutions?

To investigate this phenomenon, we analyzed five solutions per question, organizing them into groups by length in ascending order. Figure 8 illustrates the distribution of correct solutions across these groups.

Our analysis revealed that correct solutions predominantly appear in shorter-length groups, particularly in the AIME dataset. However, when examining the token distribution, we found that correct solution tokens are concentrated in longer-solution groups. This apparent contradiction arises because the total token count is determined by both the number of solutions and the average tokens per solution. As shown in Figure 2a, solutions in the longest group contain nearly twice as many tokens as those in the shortest group. This explains why, despite having fewer individual solutions, longer solutions account for a greater share of the total tokens.

We hypothesize that this discrepancy explains why RL training tends to produce longer solutions: the training process may favor generating longer solutions, even if they are less accurate, because they contribute more tokens to the gradient.

## B Further analysis on Sequential Scaling on R1-Distill-14b, R1-Distill-32b and LIMO

In Section 4.2, we observed that R1-Distill-14b, R1-Distill-32b and LIMO demonstrated some performance improvements after multiple rounds of self-revision, followed by stabilization. Furthermore, in Section 4.1, we found that the correct solutions generated by R1-Distill-14b, R1-Distill-32b and LIMO were generally shorter than incorrect solutions. To reconcile these seemingly contradictory findings and further analyze how R1-Distill-14b, R1-Distill-32b and LIMO benefit from self-revision, we conducted a detailed analysis of self-revision outcomes on both long and short solutions. Our methodology for collecting long and short solutions involved sampling five solutions for each question, ordering them by length, and then segregating the longest and shortest solutions into separate groups. The results of self-revision on both short and long solutions are presented in Figure 9. Our analysis reveals that short solutions exhibited significant performance improvements following self-revision, while this trend was less pronounced for long solutions. Therefore, the performance improvements we observed through self-revision in R1-Distill-14b, R1-Distill-32b and LIMO primarily stem from the self-revision on short solutions. This suggests that the relationship between accuracy and solution length for these models is complex, demonstrating neither a strictly positive nor negative correlation with length.

## C Parallel Scaling of Shortest Majority Vote on GPQA

In Section 6, we demonstrated that our proposed Shortest Majority Vote achieves superior test-time scaling performance compared to the other two methods on the AIME benchmark. In this section, we present the parallel-scaling results on GPQA in Figure 10. While Shortest Majority Vote consistently outperforms the Shortest method on GPQA, it does not exhibit significantly better parallel scaling performance compared to Majority Vote on this benchmark. This phenomenon might be attributed to the smaller performance gap between short and long solutions on GPQA compared to AIME, suggesting that solution length plays a less critical role in determining solution quality on the GPQA benchmark, which can be observed from Figure 2b

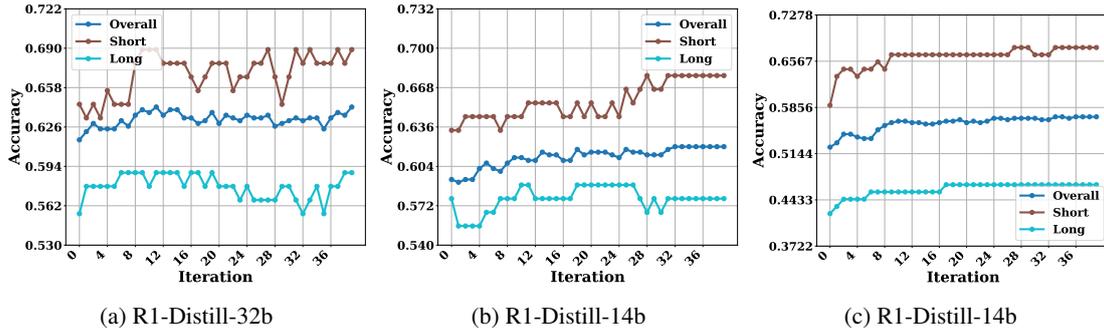


Figure 9: Accuracy of short solutions and long solutions of R1-Distill-14b (a) and R1-Distill-32b (b) during sequential revision.

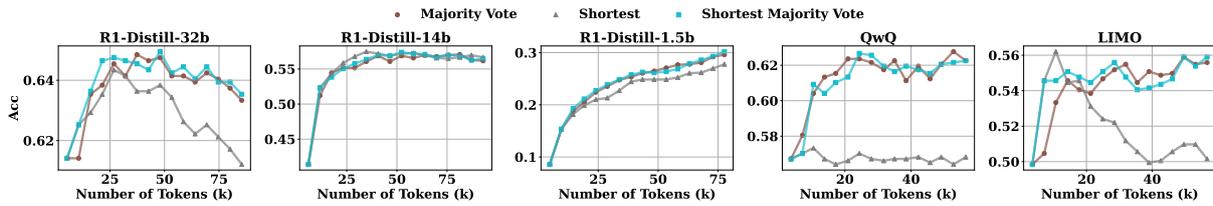


Figure 10: Performance Comparison between Majority Vote and Shortest Majority Vote on GPQA.

## D Prompt

System prompt:

**System prompt**

You are a helpful and harmless assistant.  
You should think step-by-step.

Instruction for MATH-500, AIME and Omini-MATH:

**Instruction**

Answer the question and enclose the final answer in boxed{ }

Instruction for GPQA:

**Instruction**

Select the best answer from the following options. Output only the letter corresponding to the correct answer, enclosed in boxed{ }.

## E Examples of self-revision

### Examples

Wait, let me verify that again ...

Wait, but that seems straightforward, but let me check if I got the constants right ...

Wait, but let me verify this to ensure I didn't make a mistake ...

Wait, so is the answer 756? But let me check if this is consistent ...

Wait, but in 3D space, the centers might not be coplanar? ...

Alternatively, try to find a general formula ...

Alternatively, consider that  $m$  is such that  $m$  divides  $k$  where  $k$  is from 1 to 999 ...

Alternatively, maybe we can use modulo 8 to get constraints ...

Alternatively, perhaps there's a smarter approach ...

Alternatively, another way to think about this problem is to recognize that  $w$  and  $z$  are roots of unity ...