
Mind the Gap: Examining the Self-Improvement Capabilities of Large Language Models

Anonymous Author(s)

Affiliation

Address

email

Abstract

1 Self-improvement is a mechanism in Large Language Model (LLM) pre-training,
2 post-training and test-time inference. We explore a framework where the model ver-
3 ifies its own outputs, filters or reweights data based on this verification, and distills
4 the filtered data. Despite several empirical successes, a fundamental understanding
5 is still lacking. In this work, we initiate a comprehensive, modular and controlled
6 study on LLM self-improvement. We provide a mathematical formulation for
7 self-improvement, which is largely governed by a quantity which we formalize as
8 the *generation-verification gap*. Through experiments with various model families
9 and tasks, we discover a scaling phenomenon of self-improvement – a variant of the
10 generation-verification gap scales monotonically with the model pre-training flops.
11 We also examine when self-improvement is possible, an iterative self-improvement
12 procedure, and ways to improve its performance. We believe our results have
13 several empirical implications, and our study leaves many exciting future directions
14 for understanding the potential and limits of LLM self-improvement.

15 1 Introduction

16 Recent work increasingly explores the use of synthetic data in training large language models (LLMs),
17 with applications in both pre-training and post-training (Bai et al., 2022; Meng et al., 2022; Li et al.,
18 2023b; Adler et al., 2024; Dubey et al., 2024; Yang et al., 2024; Hui et al., 2024; Li et al., 2024). While
19 synthetic data, often generated by LLMs, offers a valuable complement to human-generated data,
20 its misuse can harm performance. Bertrand et al. (2023) and Gerstgrasser et al. (2024) showed self-
21 training on model-generated data leads to degradation. To mitigate this, incorporating a “reliable” ver-
22 ifier to label data has shown promise in preventing such performance collapse (Gillman et al., 2024).

23 A straightforward verification mechanism is to train a reward model on human-annotated data to
24 assess the quality of synthetic data (Lightman et al., 2023; Wang et al., 2024a). However, this
25 approach can be prohibitively expensive and may offer no signal in domains where models exhibit
26 super-human performance. An alternative is to use a stronger model (Chang et al., 2023; Havrilla
27 et al., 2024) for annotation, but this becomes infeasible when the model is at the frontier of current
28 capabilities. A promising solution is to use the model to label its own generations. Motivated by the
29 intuition that “verification is easier than generation”, one can hypothesize that the model may act as a
30 better-than-random verifier of its own outputs, enabling *self-improvement* (Zelikman et al., 2022).

31 Most previous self-improvement algorithms can be summarized as follows: 1) make multiple
32 generations from the model, 2) use the same model to verify the generations, and 3) distill from the
33 reranked/filtered generation (Zelikman et al., 2022; Huang et al., 2022; Wang et al., 2022b; Yehudai
34 et al., 2024; Madaan et al., 2024; Yuan et al., 2024; Xu et al., 2024; Liang et al., 2024). With this
35 framework, self-improvement is also related to improving inference quality (Wang et al., 2022a;

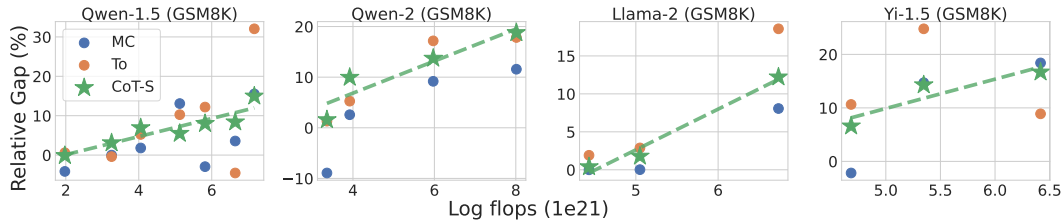


Figure 1: With proper verification method (e.g., CoT-S), the relative generation-verification gap (Definition 2.2) scales monotonically with respect to the pre-train flops. We conjecture that in this case, the relative gap is linear with respect to the log of the pre-train flops. MC denotes Multiple Choice verification, CoT-S denotes CoT-Score verification, and To denotes Tournament verification. The description of each verification can be found in Section 3.

36 Welleck et al., 2024) – if the model can verify its own generation, self-improvement can enhance
 37 test-time performance with additional computation towards more generations and updates.

38 Despite significant empirical progress and some impressive results, a fundamental understanding of
 39 LLM self-improvement remains limited. It is uncertain whether these results can be interpreted solely
 40 as an indication of the self-improvement capability of LLMs, given the potential for confounders at
 41 various stages of the process. Moreover, much of the existing research has concentrated on just one
 42 model family or a single verification mechanism, limiting the broader applicability of the findings.
 43 In this work we conduct a comprehensive study of the self-improvement capability of LLMs; our
 44 contribution is as follows:

45 • **Self-Improvement Framework:** Section 2 details the mathematical formulation of the self-
 46 improvement process, highlighting three critical desiderata. We propose the *generation-verification*
 47 *gap (GV-Gap)* (Definition 2.1) as the central metric for evaluation. GV-Gap captures the additional
 48 “precision” gained by using model verification over generations alone. It is defined as the perfor-
 49 mance improvement obtained by re-weighting generations by the model’s self-verification score
 50 (e.g., 0 or 1). Our empirical findings indicate that GV-Gap is a more accurate metric for measuring
 51 self-improvement versus the previous metric of the performance difference after a model update.

52 • **Scaling Properties:** In Section 4, we measure the generation-verification gap across multiple
 53 model families, verification mechanisms, and tasks. Certain verification methods induce a scaling
 54 phenomenon for self-improvement – relative GV-Gap (Definition 2.2) increases monotonically
 55 with pre-train flops – shown in Figure 1. We find that in cross-verification, the GV-gap increases
 56 with verifier capability and decreases with generator capability. Finally, we observe that most
 57 models do not achieve self-improvement in information retrieval and reasoning tasks that exceed
 58 their inherent capabilities. By studying multiple types of verification, our results indicate general
 59 patterns beyond just prompt engineering.

60 • **Iterative Self-Improvement:** In Section 5, we identify that **i)** GV-Gap saturates to 0 in handful
 61 rounds of iterative self-improvement; **ii)** the saturation rate is independent from the model capacity,
 62 **iii)** the effective diversity degrades during the iterative self-improvement.

63 • **Verification Mechanisms:** In Appendix A, we consider methods to enhance self-improvement
 64 through a fine-grained study on the verification methods. Key observations include: **i)** there is signif-
 65 icant non-overlap between different verification mechanisms, **ii)** GV-Gap is not positively correlated
 66 with generation accuracy, **iii)** using an ensemble of verification can improve self-improvement.

67 We believe our results provide an initial step towards a systematic understanding of the intriguing
 68 self-improvement framework in LLMs. While our observations provide several practical implications
 69 in pre-training, post-training, and test-time inferencing, we also leave several interesting future
 70 directions toward a more profound understanding of the mechanism of self-improvement.

71 Due to space constrain, we defer the related work section to Appendix B.

72 2 A Dissection of the Self-improvement Framework

73 In this section, we introduce the setup of the self-improvement framework considered in the paper,
74 which consists of the following three main components:

75 **Response Generation.** Let \mathcal{X} be the prompt/context space, and \mathcal{Y} be the response space. Let
76 $\mathcal{F} \subseteq \{f : \mathcal{X} \rightarrow \Delta(\mathcal{Y})\}$ be a class of generative models that maps a prompt to a distribution
77 over responses. A task, $\text{task} \in \mathcal{T}$ (e.g., math, code, trivia, puzzles, etc.), defines a distribution
78 $\mu_{\text{task}} \in \Delta(\mathcal{X})$ over the prompt space \mathcal{X} , and utility function $u_{\text{task}} : \mathcal{X} \times \mathcal{Y} \rightarrow [U_{\min}, U_{\max}]$, where
79 U_{\min}, U_{\max} denote bounds of the utility function.

80 The goal is to find a generative model that maximizes the expected utility: $f_{\text{task}}^* :=$
81 $\arg \max_{f \in \mathcal{F}} J_{\text{task}}(f)$, where $J_{\text{task}}(f) := \mathbb{E}_{x \sim \mu_{\text{task}}} [\mathbb{E}_{y \sim f(\cdot|x)} [u_{\text{task}}(x, y)]]$. We will often
82 drop the subscript task when it is clear from the context.

83 **Verify with Proxy Utility.** Without the access to the ground truth utility function u , we rely on a
84 proxy utility function constructed by some model f , denoted as $\hat{u}_f : \mathcal{X} \times \mathcal{Y} \rightarrow [U_{\min}, U_{\max}]$. For
85 example, let $Z = [10] := \{1, 2, \dots, 10\}$ be scores, the proxy utility $\hat{u}_f(x, y) = \mathbb{E}_{z \sim f(\cdot|x, y, \text{prom})} [z]$
86 rates the quality of the response with a score from 1 to 10 with an instruction prompt prom . In
87 [Section 3](#), we provide more proxy utility functions used in our experiments.

88 **Update via Reweighting.** Let $t \in [T]$ be the iteration index, f_t be the model at iteration t , and $\hat{u}_{f'}$
89 be the proxy utility defined by model f' . The reweighted distribution $f_t[w(\hat{u}_{f'})]$ is defined such that

$$f_t[w(\hat{u}_{f'})](y | x) \propto f_t(y | x) \cdot w(\hat{u}_{f'}(x, y)), \forall x, y \in \mathcal{X} \times \mathcal{Y}.$$

90 Here $w : \mathbb{R} \rightarrow \mathbb{R}_{\geq 0}$ is a weight function that maps a utility from the verification procedure to a
91 weight. The specific form of w is determined by the algorithm used for the model update step (we
92 provide two examples below). The objective of this update is to find a model $f_{t+1} \in \mathcal{F}$ such that
93 some distance $\ell(f_{t+1}, f_t[w(\hat{u}_{f'})])$ is small (i.e., ℓ -projecting $f_t[w(\hat{u}_{f'})]$ onto \mathcal{F}). Note that when
94 $f' = f_t$, we have self-improvement, though the framework can also be used for improving using
95 utilities produced by a different model as is studied in [Section 4.2](#).

96 **Example 2.1** (KL-regularized RL Update). *One can treat the proxy utility \hat{u} as reward, and perform*
97 *RLHF style RL update with a reverse KL constraint ([Christiano et al., 2017](#); [Ouyang et al., 2022](#)):*

$$f_{t+1} = \arg \max_{f \in \mathcal{F}} \mathbb{E}_{x, y \sim \mu \circ f} \left[\hat{u}_{f_t}(x, y) - \beta \log \left(\frac{f(y | x)}{f_t(y | x)} \right) \right].$$

98 *In this case, we have $w(s) = \exp(s/\beta)$, and ℓ can be the KL divergence between f_{t+1} and $f_t[w(\hat{u}_{f_t})]$*
99 *([Hazan, 2016](#)).* ◀

100 **Example 2.2** (Rejection Sampling). *In rejection sampling, we first filter the generation by a threshold*
101 *τ , and then fine-tune the model on the filtered data:*

$$f_{t+1} = \arg \max_{f \in \mathcal{F}} \mathbb{E}_{x, y \sim \mu \circ f_t} [\log(f(y | x)) \cdot \mathbb{1}[\hat{u}_{f_t}(x, y) \geq \tau]].$$

102 *In this case, we have $w(s) = \mathbb{1}[s \geq \tau]$, and ℓ can be the total variation distance between f_{t+1} and*
103 *$f_t[w(\hat{u}_{f_t})]$ ([Zhang, 2006](#)).* ◀

104 Finally, it is convenient to abuse the notation and allow w and \hat{u} to take batch input. For example, we
105 can allow w to take a list of score and then set the filtering threshold τ to the n quantile ($n \in [0, 1]$)
106 of the score. We denote this as top- n or quantile- n filtering.

107 2.1 Three Key Factors of Self-improvement

108 For any meaningful self-improvement, at iteration t , we would like to find f_{t+1} such that $J(f_{t+1}) >$
109 $J(f_t)$, where recall $J(f)$ is the expected utility under the model f . We identify the three key
110 conditions that may bottleneck improvement on model f :

111 **1. Improvable Generation.** Our framework involves reshaping the generation distribution towards
112 increased utility. In order for this to be useful, the utilities of generations must have variability. For
113 example, if generation were done with greedy decoding, no improvement in this process would be
114 possible. Fortunately, the improvable generation phenomenon has been well-observed in LLMs ([Li](#)
115 [et al., 2022](#); [Brown et al., 2024](#)) (see also [Figure 3](#)).

116 **2. Informative verification.** Recall that weight function $w(\hat{u}_g)$ is defined by the proxy utility function
 117 \hat{u}_g , which is constructed by a verifier model g . If the verification capability is limited, the weighting
 118 may not provide a useful signal for improvement. The following definition quantifies this intuition:

119 **Definition 2.1** (Generation-Verification Gap). *For a generator f and verifier g , we define the*
 120 *generation-verification gap (GV-Gap) between f and g as*

$$\text{gap}(f, g) := J(f[w(\hat{u}_g)]) - J(f).$$

121 This is the core metric for our analysis throughout the paper. When the generator f and verifier
 122 g are the same, we denote the shorthand $\text{gap}(f) := \text{gap}(f, f)$, which is the self-improvement
 123 generation-verification gap. However, the GV-Gap is an absolute quantity, which does not fully
 124 capture the various qualities of the generation. Imagine a generator that already achieves 99%
 125 accuracy on a task: first, the upper bound for $\text{gap}(f)$ is only 0.01; second, incorrect responses are
 126 likely to be very subtle, and thus any improvement in the reweighted distribution might require a
 127 very strong verification model. This motivates the relative GV-Gap:

128 **Definition 2.2** (Relative Generation Verification Gap). *For a generator f and verifier g , we define*
 129 *the relative generation-verification gap between f and g as*

$$\text{gap}_{\text{rel}}(f, g) := \mathbb{E}_{x \sim \mu} \left[\frac{\mathbb{E}_{y \sim f[w(\hat{u}_g)](\cdot|x)}[u(x, y)] - \mathbb{E}_{y \sim f(\cdot|x)}[u(x, y)]}{U_{\text{max}} - \mathbb{E}_{y \sim f(\cdot|x)}[u(x, y)]} \right],$$

130 That is, we weigh the gap of each prompt by its deficiency to the best possible utility. For simplicity,
 131 we will denote the self-GV-Gap as “gap” or gap and relative self-GV-Gap as “relative gap” or gap_{rel}
 132 when the context is clear. In domains where verification is easier than generation, $\text{gap} > 0$ likely
 133 holds, and indicates that there is additional signal that can be exploited. One can also check that,
 134 for all prompts $x \in \mathcal{X}$, if the weight function $w(\hat{u}_g)(x, \cdot)$ and $u(x, \cdot)$ is positively correlated¹ under
 135 the distribution of $y \sim f$, then we can always guarantee $\text{gap}(f, g) > 0$.

136 **3. High-fidelity model update.** The final condition is that the model f_{t+1} mimics/distills the
 137 performance of the reweighted distribution $f_t[w(\hat{u}_{f_t})]$, i.e., $|J(f_{t+1}) - J(f_t[w(\hat{u}_{f_t})])| \leq \varepsilon_{\text{update}}$,
 138 with some small ε . For example, if through MLE we bound the TV-distance between the two by ε' ,
 139 then by Holder’s inequality we have $\varepsilon_{\text{update}} \leq \varepsilon' U_{\text{max}}$. Combining it with the gap guarantee, we have:

$$J(f_{t+1}) - J(f_t) \geq \text{gap}(f_t, g) - \varepsilon_{\text{update}}.$$

140 With a sufficiently expressive LLM model class, it is often observed that the distillation error ε is small.

141 Note that sometimes we might observe $J(f_{t+1}) - J(f_t[w(\hat{u}_{f_t})]) > 0$. For instance, a benchmark
 142 may require outputs in a specific format; in such cases, the finetuned model f_{t+1} might outperform
 143 the reweighted distribution $f_t[w(\hat{u}_{f_t})]$ simply by aligning outputs with the required format, even
 144 if the underlying answers remain unchanged. Recent work (Dubois et al., 2024; Zhang et al.,
 145 2024c) highlight additional confounders, such as modifications to output length during finetuning,
 146 which may inflate perceived improvements without reflecting true model capabilities. Conversely,
 147 it is conceivable that intrinsic enhancements in the model’s reasoning might occur; for example,
 148 mastering simpler tasks could indirectly boost performance on more complex problems requiring
 149 similar skills. However, in our experiment we only observe the former scenario. That said, both cases
 150 emphasize the need for caution when interpreting such improvements, and this further emphasizes
 151 the importance of our modular approach in dissecting the components of self-improvement.

152 3 Experiment Setup

153 Our experiment is based on lm-evaluation-harness (Gao et al., 2024). For all tasks we use the
 154 following setup: for generations and verification, we use sampling parameters $p = 0.9, t = 0.7$, max
 155 length of 512 and 4-shot in-context samples. For each prompt, we sample 128 responses and sample
 156 1 verification for each response². In this work, we consider the following verification mechanisms (a
 157 formal description of each is provided in Appendix C along with example prompts in Appendix F):

¹Even under the case that $w(\hat{u}_g)(x, \cdot)$ and $u(x, \cdot)$ are negatively correlated, if we have a small set of holdout
 dataset with ground truth labels $u(x, y)$, we can always define $w(\alpha) := \exp(\alpha \cdot w(\hat{u}_g))$, use the holdout set
 to tune α , and use $w(\alpha)$ to reweight the distribution.

²Note that an ideal verification should be sampling multiple verifications per generation. We only sample
 one due to computational constraints and we leave multiple verifications along with understanding verification
 compute scaling for future work.

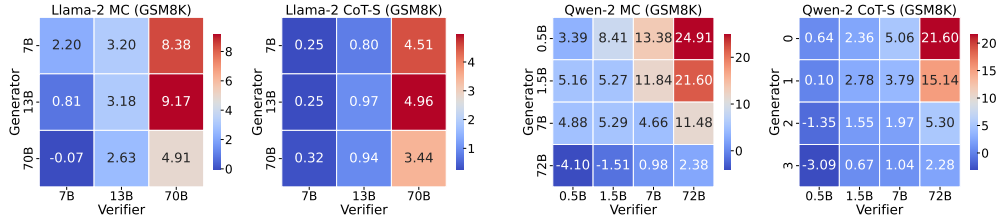


Figure 2: Gaps (%) in cross-improvement. For each row (a fixed generator), gaps increases as verifier capacity goes up. For each column (a fixed verifier), gap decreases as generator capacity goes up.

- 158 1. **Multiple Choice (MC)** (Li et al., 2023a; Gao et al., 2024) asks the LM to label responses as
159 “Correct” and “Incorrect” and uses the probability of “Correct” as a continuous score.
- 160 2. **Chain of Thought (CoT)** (Wei et al., 2022) asks the LM to score responses and to provide the
161 justification (i.e. CoT) and the score is parsed from the answer. Scores can be on a scale from
162 1 to 10 (**CoT-Score**) (Yuan et al., 2024; Liang et al., 2024) or binary (**CoT-Binary**).
- 163 3. **Tournament (To)** involves sampling a batch of generations and having a verifier compare
164 generation pairs in a single elimination tournament to produce a new generation distribution. We
165 can repeat the process until there is only one response left from the batch.

166 We consider the following models families: Qwen-1.5 (Bai et al., 2023), Qwen-2 (Yang et al., 2024),
167 Llama-2 (Touvron et al., 2023), Llama-3, Llama-3.1 (Dubey et al., 2024) and Yi-1.5 (Young et al.,
168 2024). To avoid the confounding effect of the post-training, all experiments in this paper are performed
169 on *base* models. Finally, all inference in this paper is performed with vLLM (Kwon et al., 2023).

170 4 Scaling Properties of Generation-verification Gap

171 In this section, we conduct a comprehensive study on measuring the scaling property of the generation-
172 verification gap due to its valuable practical guidance in both pre-training and downstream tasks
173 (Kaplan et al., 2020; Hernandez et al., 2021; Isik et al., 2024; Ruan et al., 2024).

174 4.1 Scaling Results

175 We start with the GSM8K benchmark (Cobbe et al., 2021), with 1320 questions on the test data split.
176 The ground truth utility $u(x, y) = 1$ if the end of response y is the correct answer to the question x ,
177 or $u(x, y) = 0$ otherwise. We compute $\text{gap}(f)$ for each model f and verification method, and we
178 record the full results in Tables 4 and 5. In particular, we observe the following phenomenon:

179 **Small Models can not Self-improve.** For small (in terms of pre-train flops) models such as Qwen-1.5
180 0.5B, Qwen-2 0.5B and Llama-2 7B, $\text{gap}(f)$ is non-positive for nearly all verification methods, even
181 though the models have non-trivial generation accuracy. We also observe this phenomenon in Pythia
182 (Biderman et al., 2023) and OPT (Zhang et al., 2022) model families. We believe this result indicates
183 that self-improvement requires a minimal level of instruction following and reasoning capabilities,
184 which is not present in these small models. We will further illustrate this point in Section 4.4.

185 **CoT Verification is More Stable than MC.** Some MC verification incurs non-positive gap even for
186 medium-sized models such as Qwen-1.5 14/32B and Llama-3/3.1 8B models, while CoT verification
187 always has a positive gap for medium/large-sized models. Our results align with recent studies
188 showing that MC evaluation might be unreliable, especially for small models (Dominguez-Olmedo
189 et al., 2024). We perform a more in-depth analysis on this point in Appendix A.

190 **$\text{gap}_{\text{rel}}(f)$ Scales with Pre-training Flops.** We observe that with certain verification methods (such as
191 CoT-Score), the relative gap grows monotonically with the pre-train flops, demonstrating a scaling
192 property. We visualize the scaling results in Figure 1, where we plot $\text{gap}_{\text{rel}}(f)$ with respect to the
193 log of pre-train flops. Specifically, we hypothesize that in the case where the verification elicits the
194 scaling property, $\text{gap}_{\text{rel}}(f)$ scales linearly with respect to the log of the pre-train flops. However, note
195 that we should not expect the slope for each model family to be the same. In Figure 7, we repeat the
196 same plot for $\text{gap}(f)$, but we do not observe a similar trend with the absolute gap.

Table 1: Gap (%) on Natural Question for Qwen-2 models. While all models have a non-trivial generation accuracy, all gaps are near 0, indicating that the task is unimprovable.

	0.5B	1.5B	7B	72B
Generation Accuracy	6.51	13.87	29.09	41.45
MC (top 0.8)	-0.06	0.04	0.79	0.28
MC ($\tau = 0.8$)	-0.05	0.02	-0.05	-0.05

Table 2: Generation accuracy, gap and relative gap (%) on Sudoku for Qwen-2 models. Only the 72B models can self-improve. For the 72B models, the improvement is around 200%.

	0.5B	1.5B	7B	72B
Generation Accuracy	0.66	0.62	2.09	8.82
gap	-0.09	0.04	-0.07	16.99
gap _{rel}	-0.15	-0.61	-0.01	20.81

197 4.2 Cross Verification

198 In self-improvement, both generator and verifier change when transitioning between different models.
 199 To better understand the relationship between generation/verification ability and model capacity, we
 200 perform a cross-verification study, where we only alter either the generator or the verifier at a time.
 201 We consider the Llama-2 and Qwen-2 model families, and the two most representative verification
 202 methods: MC with quantile threshold and CoT-Score. We present the results in Figure 2. We observe
 203 that the results are consistent with our intuition on the difficulty of verification: fix a generator model
 204 f , $\text{gap}(f, g)$ increases as the model capacity (defined by pre-train flops) of the verifier model g
 205 increases. On the other hand, fix a verifier model g , $\text{gap}(f, g)$ decreases as the model capacity of the
 206 generator model f increases, as the error of the generator model becomes more difficult to detect.

207 At first glance, the results seem to imply that selecting the largest model as the verifier, akin to a
 208 teacher-student setup, is advantageous. However, considering the computational costs associated
 209 with larger verifier models, this approach might be suboptimal. Alternatively, a weak-to-strong
 210 setup, where a smaller model verifies a larger one, might be more cost-effective, but our findings
 211 indicate that a positive gap cannot always be assured. We believe an interesting future direction is
 212 to explore the compute-optimal configuration for cross-verification. This, however, might require a
 213 combinatorially large number of experiments to pinpoint the optimal verifier for each generator.

Takeaway on scaling of self-improvement

LLMs demonstrate clear scaling trends in self- and cross-improvement:

- **Self-Improvement:** With stable verification, the relative gap grows monotonically with pre-training flops.
- **Cross-Improvement:** The gap scales **directly** with the verifier’s flops and **inversely** with the generator’s flops.

If the relative gap scales linearly with the logarithm of pre-training flops, this relationship could guide decisions on synthetic data generation strategies in self-improvement. Additionally, results from cross-verification suggest that a compute-optimal combination may exist to maximize efficiency in cross-improvement contexts.

214

215 4.3 Unimprovable Tasks

216 The primary objective of self-improvement is predicated on the assumption that “verification is easier
 217 than generation”. As such, it is also worthwhile to consider tasks where such intuition would not
 218 hold. One such scenario involves factual tasks that require generating a factually correct answer to a
 219 trivia question. We hypothesize that the capability to generate a correct answer is contingent solely
 220 on whether the model has been trained with the relevant factual knowledge, and verification would
 221 provide little additional signal. To test this, we measure $\text{gap}(f)$ on the Natural Question dataset
 222 (Kwiatkowski et al., 2019), where $u(x, y) = 1$ if y is one of the candidate answers to the question
 223 x , and $u(x, y) = 0$ otherwise. Our analysis on a test subset of 3610 questions, presented in Table 1,
 224 reveals that despite all models achieving non-trivial generation accuracy, the gap remains smaller
 225 than 1%, or is even negative, across all models. This suggests that certain tasks may not benefit from
 226 the current self-improvement framework. We include the full results in Table 6.

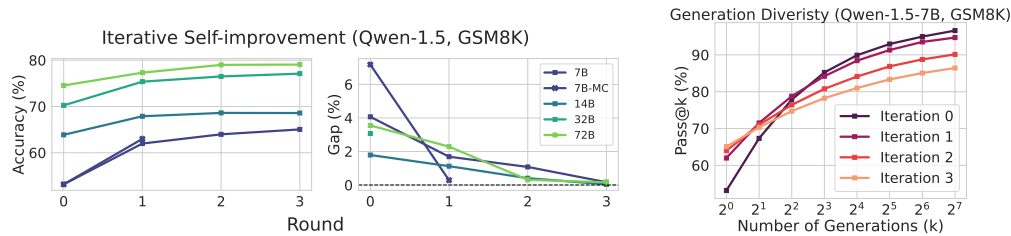


Figure 3: **Left:** The generation accuracy and gap along the iterative self-improvement process for Qwen-1.5 model family with CoT-Binary and MC verification. The horizontal line denotes 0.005. **Right:** The change of effective generation diversity along the iterative self-improvement process for Qwen-1.5 7B model, measured by pass@ k for different k .

227 4.4 Sudoku

228 Generalized sudoku is a canonical example where the generation (NP-hard) is harder than the
 229 verification (P) (Haythorpe, 2016). We consider 4 by 4 sudoku puzzles, each with a unique solution,
 230 with 288 puzzles in total. We task the models to use CoT reasoning for both generation and verification.
 231 The results, presented in Table 2 and detailed further in Appendix D.3, reveal a surprising pattern:
 232 only the largest models, such as Qwen-1.5/2 72B and Llama 3.1 70B, exhibit non-trivial gaps. For
 233 these models, the improvement is indeed more significant (50% – 300% improvement in accuracy)
 234 than the improvement in the math task.

235 While the second observation aligns with common intuition, the first may be unexpected, as most
 236 models demonstrate the ability to self-improve on tasks where the gap between generation and verifi-
 237 cation appears even narrower, such as in GSM8K. We hypothesize that, despite sudoku verification
 238 being simpler than generation, it still necessitates a certain level of reasoning and planning, even
 239 with explicit verification guidelines. This requirement is similar in mathematical tasks; however, it is
 240 likely that most models have been exposed to math verification during pre-training, unlike sudoku
 241 verification. Consequently, smaller models may lack the requisite reasoning capabilities to improve
 242 on sudoku tasks. Although our analysis is primarily post-hoc, an interesting avenue for future research
 243 would be to develop a metric to predict a model’s “self-improvability” on specific tasks.

Takeaway on improvable tasks

LLMs do not universally self-improve across all tasks:

- **Trivia Tasks:** There is no significant generation-verification gap, given the similarity in complexity between generation and verification.
- **Sudoku:** Despite the exponential complexity separation between generation and verification in generalized sudoku, most models fail to self-improve. When improvement occurs, it is notably significant.

These findings suggest that the crucial factor for general self-improvement is the model’s inherent reasoning and planning capabilities developed during pre-training.

245 5 Iterative Self-improvement

246 Building on our understanding of single-round self-improvement, a natural extension is to study
 247 iterative self-improvement. As there is no additional information introduced in the process, it is
 248 unrealistic to expect indefinite improvement. Thus in this section, we study the dynamics of the
 249 iterative self-improvement, and its relationship with model scales.

250 In our experiment, we perform iterative self-improvement on the Qwen-1.5 model family with
 251 CoT-Binary verification on GSM8K. We defer the finetuning hyperparameters to Appendix E. We
 252 present the results in Figure 3. We observe that 1) the gap diminishes nearly to zero within two or
 253 three rounds of self-improvement; this is consistent with the observation with previous works (Yuan
 254 et al., 2024; Liang et al., 2024). 2) The rate of saturation is similar across models with different
 255 capacities. 3) Notably, for the 7B and 14B models, the model accuracy at iteration 1 exceeds the sum

256 of the generation accuracy and the gap at iteration 0, i.e., $J(f_1) > J(f_0[w(\hat{u}_{f_0})])$. This increase is
257 attributed to improved adherence to the required answer format post-finetuning – the discrepancy
258 between “flexible match” and “exact match” (extract the answer from the required answer format)
259 disappears after the first round. As we argued in the previous section, this additional accuracy gain
260 is not due to the self-improvement capability of the model, and thus our modular study reduces the
261 confounding factors in understanding the self-improvement capability of the model.

262 To compare the dynamics between verification methods, in Figure 3 we also plot MC (top 0.7)
263 verification for the 7B model. We observe that the gap immediately drops to near 0 after the first
264 round of self-improvement, and thus multi-round self-improvement with MC verification is unlikely.
265 This rapid saturation is consistent across other thresholds for MC verification. We provide a more
266 detailed study on the cause of this phenomenon in Appendix A.1.

267 We also examine the “effective diversity” of generations throughout the iterative self-improvement
268 process using the metric $\text{pass}@k^3$. We present the results in Figure 3. We observe when k is small,
269 $\text{pass}@k$ increases with the number of rounds of self-improvement, validating the success of the
270 self-improvement process. However, when k is large, $\text{pass}@k$ decreases with the number of iterations,
271 indicating that the diversity of the generations is reduced through the self-improvement process. This
272 trend may result from the model’s inability to verify rare, yet correct, answers, potentially leading to
273 convergence on incorrect solutions during the self-improvement process.

Takeaway on iterative self-improvement

LLMs can perform iterative self-improvement with an effective verification method:

- **Saturation Limit:** Without new information, iterative self-improvement typically saturates after two or three rounds, regardless of the model’s capacity (measured in pre-training flops).
- **Cause of Saturation:** A potential cause for this saturation is a decrease in effective diversity, leading to convergence on incorrect answers for certain questions.

Addressing the reduction in diversity could potentially extend the duration and effectiveness of the self-improvement process.

274

275 6 Conclusion and Discussion

276 In this paper, we conduct a comprehensive and modular study on the LLM self-improvement frame-
277 work through multiple model families, tasks and verification mechanisms. We structure the mathemat-
278 ical framework of the self-improvement process and pinpoint the generation-verification gap as a criti-
279 cal metric. Our results reveal several intriguing properties such as the scaling properties of the relative
280 gap, saturation of iterative self-improvement and enhancement of verification via ensemble methods.
281 These insights are likely to have practical implications for improving pre-training, post-training, and
282 test-time inference. Additionally, our research opens several promising avenues for future exploration:

- 283 • While our scaling analysis is primarily observational (Ruan et al., 2024), pursuing a more extensive
284 scaling law study (Kaplan et al., 2020) based on our preliminary findings could provide robust
285 empirical guidelines.
- 286 • Our results hint at an inference-time scaling law (Wu et al., 2024) is possible for self-improvement
287 (or with cross-improvement (c.r. Section 4.2)). Identifying compute-optimal methods for
288 self-improvement across different tasks remains a critical challenge.
- 289 • The decline in the effective diversity of generations during iterative self-improvement presents a sig-
290 nificant obstacle. Developing strategies to mitigate this issue offers considerable empirical benefits.
- 291 • The distinct non-overlap property of verification mechanisms, despite their functional similarities,
292 suggests that combining compositional verification could significantly enhance self-improvement.
293 Exploring this potential further could yield fruitful results.

³Given a question, $\text{pass}@k$ is 1 if at least one of the k generations of the model is correct, or 0 otherwise.

294 References

- 295 Josh Achiam, Steven Adler, Sandhini Agarwal, Lama Ahmad, Ilge Akkaya, Florencia Leoni Aleman,
296 Diogo Almeida, Janko Altenschmidt, Sam Altman, Shyamal Anadkat, et al. Gpt-4 technical report.
297 *arXiv preprint arXiv:2303.08774*, 2023.
- 298 Bo Adler, Niket Agarwal, Ashwath Aithal, Dong H Anh, Pallab Bhattacharya, Annika Brundyn,
299 Jared Casper, Bryan Catanzaro, Sharon Clay, Jonathan Cohen, et al. Nemotron-4 340b technical
300 report. *arXiv preprint arXiv:2406.11704*, 2024.
- 301 Sina Alemohammad, Josue Casco-Rodriguez, Lorenzo Luzi, Ahmed Imtiaz Humayun, Hossein
302 Babaei, Daniel LeJeune, Ali Siahkoochi, and Richard G Baraniuk. Self-consuming generative
303 models go mad. *arXiv preprint arXiv:2307.01850*, 2023.
- 304 Zachary Ankner, Mansheej Paul, Brandon Cui, Jonathan D Chang, and Prithviraj Ammanabrolu.
305 Critique-out-loud reward models. *arXiv preprint arXiv:2408.11791*, 2024.
- 306 Jinze Bai, Shuai Bai, Yunfei Chu, Zeyu Cui, Kai Dang, Xiaodong Deng, Yang Fan, Wenbin Ge,
307 Yu Han, Fei Huang, et al. Qwen technical report. *arXiv preprint arXiv:2309.16609*, 2023.
- 308 Yuntao Bai, Saurav Kadavath, Sandipan Kundu, Amanda Askell, Jackson Kernion, Andy Jones, Anna
309 Chen, Anna Goldie, Azalia Mirhoseini, Cameron McKinnon, et al. Constitutional ai: Harmlessness
310 from ai feedback. *arXiv preprint arXiv:2212.08073*, 2022.
- 311 Hritik Bansal, Arian Hosseini, Rishabh Agarwal, Vinh Q Tran, and Mehran Kazemi. Smaller, weaker,
312 yet better: Training llm reasoners via compute-optimal sampling. *arXiv preprint arXiv:2408.16737*,
313 2024.
- 314 Quentin Bertrand, Avishek Joey Bose, Alexandre Duplessis, Marco Jiralerspong, and Gauthier Gidel.
315 On the stability of iterative retraining of generative models on their own data. *arXiv preprint*
316 *arXiv:2310.00429*, 2023.
- 317 Stella Biderman, Hailey Schoelkopf, Quentin Gregory Anthony, Herbie Bradley, Kyle O’Brien, Eric
318 Hallahan, Mohammad Aflah Khan, Shivanshu Purohit, USVSN Sai Prashanth, Edward Raff, et al.
319 Pythia: A suite for analyzing large language models across training and scaling. In *International*
320 *Conference on Machine Learning*, pp. 2397–2430. PMLR, 2023.
- 321 Martin Briesch, Dominik Sobania, and Franz Rothlauf. Large language models suffer from their own
322 output: An analysis of the self-consuming training loop. *arXiv preprint arXiv:2311.16822*, 2023.
- 323 Bradley Brown, Jordan Juravsky, Ryan Ehrlich, Ronald Clark, Quoc V Le, Christopher Ré, and
324 Azalia Mirhoseini. Large language monkeys: Scaling inference compute with repeated sampling.
325 *arXiv preprint arXiv:2407.21787*, 2024.
- 326 Sébastien Bubeck, Varun Chandrasekaran, Ronen Eldan, Johannes Gehrke, Eric Horvitz, Ece Kamar,
327 Peter Lee, Yin Tat Lee, Yuanzhi Li, Scott Lundberg, et al. Sparks of artificial general intelligence:
328 Early experiments with gpt-4. *arXiv preprint arXiv:2303.12712*, 2023.
- 329 Jonathan D Chang, Kianté Brantley, Rajkumar Ramamurthy, Dipendra Misra, and Wen Sun. Learning
330 to generate better than your llm. *arXiv preprint arXiv:2306.11816*, 2023.
- 331 Xinyun Chen, Maxwell Lin, Nathanael Schärli, and Denny Zhou. Teaching large language models to
332 self-debug. *arXiv preprint arXiv:2304.05128*, 2023.
- 333 Cheng-Han Chiang and Hung-yi Lee. Can large language models be an alternative to human
334 evaluations? *arXiv preprint arXiv:2305.01937*, 2023.
- 335 Wei-Lin Chiang, Zhuohan Li, Zi Lin, Ying Sheng, Zhanghao Wu, Hao Zhang, Lianmin Zheng,
336 Siyuan Zhuang, Yonghao Zhuang, Joseph E Gonzalez, et al. Vicuna: An open-source chatbot
337 impressing gpt-4 with 90%* chatgpt quality. See <https://vicuna.lmsys.org> (accessed 14 April
338 2023), 2(3):6, 2023.
- 339 Paul F Christiano, Jan Leike, Tom Brown, Miljan Martic, Shane Legg, and Dario Amodei. Deep
340 reinforcement learning from human preferences. *Advances in neural information processing*
341 *systems*, 30, 2017.

- 342 Karl Cobbe, Vineet Kosaraju, Mohammad Bavarian, Mark Chen, Heewoo Jun, Lukasz Kaiser,
343 Matthias Plappert, Jerry Tworek, Jacob Hilton, Reiichiro Nakano, et al. Training verifiers to solve
344 math word problems. *arXiv preprint arXiv:2110.14168*, 2021.
- 345 Michael Collins and Terry Koo. Discriminative reranking for natural language parsing. *Computational*
346 *Linguistics*, 31(1):25–70, 2005.
- 347 Elvis Dohmatob, Yunzhen Feng, and Julia Kempe. Model collapse demystified: The case of regression.
348 *arXiv preprint arXiv:2402.07712*, 2024.
- 349 Ricardo Dominguez-Olmedo, Florian E Dorner, and Moritz Hardt. Training on the test task confounds
350 evaluation and emergence. *arXiv preprint arXiv:2407.07890*, 2024.
- 351 Abhimanyu Dubey, Abhinav Jauhri, Abhinav Pandey, Abhishek Kadian, Ahmad Al-Dahle, Aiesha
352 Letman, Akhil Mathur, Alan Schelten, Amy Yang, Angela Fan, et al. The llama 3 herd of models.
353 *arXiv preprint arXiv:2407.21783*, 2024.
- 354 Yann Dubois, Balázs Galambosi, Percy Liang, and Tatsunori B Hashimoto. Length-controlled
355 alpacaeval: A simple way to debias automatic evaluators. *arXiv preprint arXiv:2404.04475*, 2024.
- 356 Leo Gao, Jonathan Tow, Baber Abbasi, Stella Biderman, Sid Black, Anthony DiPofi, Charles Foster,
357 Laurence Golding, Jeffrey Hsu, Alain Le Noac’h, Haonan Li, Kyle McDonell, Niklas Muennighoff,
358 Chris Ociepa, Jason Phang, Laria Reynolds, Hailey Schoelkopf, Aviya Skowron, Lintang Sutawika,
359 Eric Tang, Anish Thite, Ben Wang, Kevin Wang, and Andy Zou. A framework for few-shot
360 language model evaluation, 07 2024. URL <https://zenodo.org/records/12608602>.
- 361 Matthias Gerstgrasser, Rylan Schaeffer, Apratim Dey, Rafael Rafailov, Henry Sleight, John Hughes,
362 Tomasz Korbak, Rajashree Agrawal, Dhruv Pai, Andrey Gromov, et al. Is model collapse in-
363 evitable? breaking the curse of recursion by accumulating real and synthetic data. *arXiv preprint*
364 *arXiv:2404.01413*, 2024.
- 365 Nate Gillman, Michael Freeman, Daksh Aggarwal, Chia-Hong Hsu, Calvin Luo, Yonglong Tian, and
366 Chen Sun. Self-correcting self-consuming loops for generative model training. *arXiv preprint*
367 *arXiv:2402.07087*, 2024.
- 368 Suriya Gunasekar, Yi Zhang, Jyoti Aneja, Caio César Teodoro Mendes, Allie Del Giorno, Sivakanth
369 Gopi, Mojan Javaheripi, Piero Kauffmann, Gustavo de Rosa, Olli Saarikivi, et al. Textbooks are all
370 you need. *arXiv preprint arXiv:2306.11644*, 2023.
- 371 Ryuichiro Hataya, Han Bao, and Hiromi Arai. Will large-scale generative models corrupt future
372 datasets? In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pp.
373 20555–20565, 2023.
- 374 Alex Havrilla, Yuqing Du, Sharath Chandra Raparthy, Christoforos Nalmpantis, Jane Dwivedi-Yu,
375 Maksym Zhuravinskiy, Eric Hambro, Sainbayar Sukhbaatar, and Roberta Raileanu. Teaching large
376 language models to reason with reinforcement learning. *arXiv preprint arXiv:2403.04642*, 2024.
- 377 Michael Haythorpe. Reducing the generalised sudoku problem to the hamiltonian cycle problem.
378 *AKCE International Journal of Graphs and Combinatorics*, 13(3):272–282, 2016.
- 379 Elad Hazan. Introduction to online convex optimization. *Foundations and Trends® in Optimization*,
380 2(3-4):157–325, 2016.
- 381 Danny Hernandez, Jared Kaplan, Tom Henighan, and Sam McCandlish. Scaling laws for transfer.
382 *arXiv preprint arXiv:2102.01293*, 2021.
- 383 Jiaxin Huang, Shixiang Shane Gu, Le Hou, Yuexin Wu, Xuezhi Wang, Hongkun Yu, and Jiawei Han.
384 Large language models can self-improve. *arXiv preprint arXiv:2210.11610*, 2022.
- 385 Liang Huang and David Chiang. Forest rescoring: Faster decoding with integrated language models.
386 In *Proceedings of the 45th annual meeting of the association of computational linguistics*, pp.
387 144–151, 2007.

- 388 Binyuan Hui, Jian Yang, Zeyu Cui, Jiayi Yang, Dayiheng Liu, Lei Zhang, Tianyu Liu, Jiajun Zhang,
389 Bowen Yu, Kai Dang, et al. Qwen2. 5-coder technical report. *arXiv preprint arXiv:2409.12186*,
390 2024.
- 391 Berivan Isik, Natalia Ponomareva, Hussein Hazimeh, Dimitris Paparas, Sergei Vassilvitskii, and
392 Sanmi Koyejo. Scaling laws for downstream task performance of large language models. *arXiv*
393 *preprint arXiv:2402.04177*, 2024.
- 394 Jared Kaplan, Sam McCandlish, Tom Henighan, Tom B Brown, Benjamin Chess, Rewon Child, Scott
395 Gray, Alec Radford, Jeffrey Wu, and Dario Amodei. Scaling laws for neural language models.
396 *arXiv preprint arXiv:2001.08361*, 2020.
- 397 Kalpesh Krishna, Yapei Chang, John Wieting, and Mohit Iyyer. Rankgen: Improving text generation
398 with large ranking models. *arXiv preprint arXiv:2205.09726*, 2022.
- 399 Tom Kwiatkowski, Jennimaria Palomaki, Olivia Redfield, Michael Collins, Ankur Parikh, Chris
400 Alberti, Danielle Epstein, Illia Polosukhin, Jacob Devlin, Kenton Lee, et al. Natural questions: a
401 benchmark for question answering research. *Transactions of the Association for Computational*
402 *Linguistics*, 7:453–466, 2019.
- 403 Woosuk Kwon, Zhuohan Li, Siyuan Zhuang, Ying Sheng, Lianmin Zheng, Cody Hao Yu, Joseph
404 Gonzalez, Hao Zhang, and Ion Stoica. Efficient memory management for large language model
405 serving with pagedattention. In *Proceedings of the 29th Symposium on Operating Systems*
406 *Principles*, pp. 611–626, 2023.
- 407 Haoran Li, Qingxiu Dong, Zhengyang Tang, Chaojun Wang, Xingxing Zhang, Haoyang Huang, Shao-
408 han Huang, Xiaolong Huang, Zeqiang Huang, Dongdong Zhang, et al. Synthetic data (almost) from
409 scratch: Generalized instruction tuning for language models. *arXiv preprint arXiv:2402.13064*,
410 2024.
- 411 Xuechen Li, Tianyi Zhang, Yann Dubois, Rohan Taori, Ishaan Gulrajani, Carlos Guestrin, Percy
412 Liang, and Tatsunori B. Hashimoto. AlpacaEval: An automatic evaluator of instruction-following
413 models. https://github.com/tatsu-lab/alpaca_eval, 5 2023a.
- 414 Yuanzhi Li, Sébastien Bubeck, Ronen Eldan, Allie Del Giorno, Suriya Gunasekar, and Yin Tat Lee.
415 Textbooks are all you need ii: phi-1.5 technical report. *arXiv preprint arXiv:2309.05463*, 2023b.
- 416 Yujia Li, David Choi, Junyoung Chung, Nate Kushman, Julian Schrittwieser, Rémi Leblond, Tom
417 Eccles, James Keeling, Felix Gimeno, Agustin Dal Lago, et al. Competition-level code generation
418 with alphacode. *Science*, 378(6624):1092–1097, 2022.
- 419 Yiming Liang, Ge Zhang, Xingwei Qu, Tianyu Zheng, Jiawei Guo, Xinrun Du, Zhenzhu Yang,
420 Jiaheng Liu, Chenghua Lin, Lei Ma, et al. I-sheep: Self-alignment of llm from scratch through an
421 iterative self-enhancement paradigm. *arXiv preprint arXiv:2408.08072*, 2024.
- 422 Hunter Lightman, Vineet Kosaraju, Yura Burda, Harri Edwards, Bowen Baker, Teddy Lee, Jan
423 Leike, John Schulman, Ilya Sutskever, and Karl Cobbe. Let’s verify step by step. *arXiv preprint*
424 *arXiv:2305.20050*, 2023.
- 425 Liangchen Luo, Yinxiao Liu, Rosanne Liu, Samrat Phatale, Harsh Lara, Yunxuan Li, Lei Shu, Yun
426 Zhu, Lei Meng, Jiao Sun, et al. Improve mathematical reasoning in language models by automated
427 process supervision. *arXiv preprint arXiv:2406.06592*, 2024.
- 428 Aman Madaan, Niket Tandon, Prakhar Gupta, Skyler Hallinan, Luyu Gao, Sarah Wiegrefe, Uri
429 Alon, Nouha Dziri, Shrimai Prabhumoye, Yiming Yang, et al. Self-refine: Iterative refinement
430 with self-feedback. *Advances in Neural Information Processing Systems*, 36, 2024.
- 431 Gonzalo Martínez, Lauren Watson, Pedro Reviriego, José Alberto Hernández, Marc Juárez, and
432 Rik Sarkar. Combining generative artificial intelligence (ai) and the internet: Heading towards
433 evolution or degradation? *arXiv preprint arXiv:2303.01255*, 2023.
- 434 Yu Meng, Jiaxin Huang, Yu Zhang, and Jiawei Han. Generating training data with language models:
435 Towards zero-shot language understanding. *Advances in Neural Information Processing Systems*,
436 35:462–477, 2022.

- 437 Long Ouyang, Jeffrey Wu, Xu Jiang, Diogo Almeida, Carroll Wainwright, Pamela Mishkin, Chong
438 Zhang, Sandhini Agarwal, Katarina Slama, Alex Ray, et al. Training language models to follow
439 instructions with human feedback. *Advances in neural information processing systems*, 35:27730–
440 27744, 2022.
- 441 Yangjun Ruan, Chris J Maddison, and Tatsunori Hashimoto. Observational scaling laws and the
442 predictability of language model performance. *arXiv preprint arXiv:2405.10938*, 2024.
- 443 Noah Shinn, Federico Cassano, Ashwin Gopinath, Karthik Narasimhan, and Shunyu Yao. Reflexion:
444 Language agents with verbal reinforcement learning. *Advances in Neural Information Processing
445 Systems*, 36, 2024.
- 446 Avi Singh, John D Co-Reyes, Rishabh Agarwal, Ankesh Anand, Piyush Patil, Peter J Liu, James
447 Harrison, Jaehoon Lee, Kelvin Xu, Aaron Parisi, et al. Beyond human data: Scaling self-training
448 for problem-solving with language models. *arXiv preprint arXiv:2312.06585*, 2023.
- 449 Charlie Snell, Jaehoon Lee, Kelvin Xu, and Aviral Kumar. Scaling llm test-time compute optimally
450 can be more effective than scaling model parameters. *arXiv preprint arXiv:2408.03314*, 2024.
- 451 Nisan Stiennon, Long Ouyang, Jeffrey Wu, Daniel Ziegler, Ryan Lowe, Chelsea Voss, Alec Radford,
452 Dario Amodei, and Paul F Christiano. Learning to summarize with human feedback. *Advances in
453 Neural Information Processing Systems*, 33:3008–3021, 2020.
- 454 Tianxiang Sun, Xiaotian Zhang, Zhengfu He, Peng Li, Qinyuan Cheng, Hang Yan, Xiangyang Liu,
455 Yunfan Shao, Qiong Tang, Xingjian Zhao, et al. Moss: Training conversational language models
456 from synthetic data. *arXiv preprint arXiv:2307.15020*, 7:3, 2023.
- 457 Rohan Taori and Tatsunori Hashimoto. Data feedback loops: Model-driven amplification of dataset
458 biases. In *International Conference on Machine Learning*, pp. 33883–33920. PMLR, 2023.
- 459 Rohan Taori, Ishaan Gulrajani, Tianyi Zhang, Yann Dubois, Xuechen Li, Carlos Guestrin, Percy
460 Liang, and Tatsunori B Hashimoto. Stanford alpaca: An instruction-following llama model, 2023.
- 461 Gemini Team, Rohan Anil, Sebastian Borgeaud, Yonghui Wu, Jean-Baptiste Alayrac, Jiahui Yu, Radu
462 Soricut, Johan Schalkwyk, Andrew M Dai, Anja Hauth, et al. Gemini: a family of highly capable
463 multimodal models. *arXiv preprint arXiv:2312.11805*, 2023.
- 464 Hugo Touvron, Louis Martin, Kevin Stone, Peter Albert, Amjad Almahairi, Yasmine Babaei, Nikolay
465 Bashlykov, Soumya Batra, Prajjwal Bhargava, Shruti Bhosale, et al. Llama 2: Open foundation
466 and fine-tuned chat models. *arXiv preprint arXiv:2307.09288*, 2023.
- 467 Peiyi Wang, Lei Li, Zhihong Shao, Runxin Xu, Damai Dai, Yifei Li, Deli Chen, Yu Wu, and Zhifang
468 Sui. Math-shepherd: Verify and reinforce llms step-by-step without human annotations. In
469 *Proceedings of the 62nd Annual Meeting of the Association for Computational Linguistics (Volume
470 1: Long Papers)*, pp. 9426–9439, 2024a.
- 471 Tianlu Wang, Ilia Kulikov, Olga Golovneva, Ping Yu, Weizhe Yuan, Jane Dwivedi-Yu,
472 Richard Yuanzhe Pang, Maryam Fazel-Zarandi, Jason Weston, and Xian Li. Self-taught evaluators.
473 *arXiv preprint arXiv:2408.02666*, 2024b.
- 474 Xuezhi Wang, Jason Wei, Dale Schuurmans, Quoc Le, Ed Chi, Sharan Narang, Aakanksha Chowdh-
475 ery, and Denny Zhou. Self-consistency improves chain of thought reasoning in language models.
476 *arXiv preprint arXiv:2203.11171*, 2022a.
- 477 Yizhong Wang, Yeganeh Kordi, Swaroop Mishra, Alisa Liu, Noah A Smith, Daniel Khashabi, and
478 Hannaneh Hajishirzi. Self-instruct: Aligning language models with self-generated instructions.
479 *arXiv preprint arXiv:2212.10560*, 2022b.
- 480 Jason Wei, Xuezhi Wang, Dale Schuurmans, Maarten Bosma, Fei Xia, Ed Chi, Quoc V Le, Denny
481 Zhou, et al. Chain-of-thought prompting elicits reasoning in large language models. *Advances in
482 neural information processing systems*, 35:24824–24837, 2022.
- 483 Yuxiang Wei, Zhe Wang, Jiawei Liu, Yifeng Ding, and Lingming Zhang. Magicoder: Empowering
484 code generation with oss-instruct. In *Forty-first International Conference on Machine Learning*,
485 2024.

- 486 Sean Welleck, Amanda Bertsch, Matthew Finlayson, Hailey Schoelkopf, Alex Xie, Graham Neubig,
487 Ilya Kulikov, and Zaid Harchaoui. From decoding to meta-generation: Inference-time algorithms
488 for large language models. *arXiv preprint arXiv:2406.16838*, 2024.
- 489 Yangzhen Wu, Zhiqing Sun, Shanda Li, Sean Welleck, and Yiming Yang. An empirical analy-
490 sis of compute-optimal inference for problem-solving with language models. *arXiv preprint*
491 *arXiv:2408.00724*, 2024.
- 492 Wenda Xu, Guanglei Zhu, Xuandong Zhao, Liangming Pan, Lei Li, and William Wang. Pride and
493 prejudice: Llm amplifies self-bias in self-refinement. In *Proceedings of the 62nd Annual Meeting*
494 *of the Association for Computational Linguistics (Volume 1: Long Papers)*, pp. 15474–15492,
495 2024.
- 496 An Yang, Baosong Yang, Binyuan Hui, Bo Zheng, Bowen Yu, Chang Zhou, Chengpeng Li,
497 Chengyuan Li, Dayiheng Liu, Fei Huang, et al. Qwen2 technical report. *arXiv preprint*
498 *arXiv:2407.10671*, 2024.
- 499 Asaf Yehudai, Boaz Carmeli, Yosi Mass, Ofir Arviv, Nathaniel Mills, Assaf Toledo, Eyal Shnarch,
500 and Leshem Choshen. Genie: Achieving human parity in content-grounded datasets generation.
501 *arXiv preprint arXiv:2401.14367*, 2024.
- 502 Alex Young, Bei Chen, Chao Li, Chengen Huang, Ge Zhang, Guanwei Zhang, Heng Li, Jiangcheng
503 Zhu, Jianqun Chen, Jing Chang, et al. Yi: Open foundation models by 01. ai. *arXiv preprint*
504 *arXiv:2403.04652*, 2024.
- 505 Weizhe Yuan, Richard Yuanzhe Pang, Kyunghyun Cho, Sainbayar Sukhbaatar, Jing Xu, and Jason
506 Weston. Self-rewarding language models. *arXiv preprint arXiv:2401.10020*, 2024.
- 507 Eric Zelikman, Yuhuai Wu, Jesse Mu, and Noah Goodman. Star: Bootstrapping reasoning with
508 reasoning. *Advances in Neural Information Processing Systems*, 35:15476–15488, 2022.
- 509 Eric Zelikman, Georges Harik, Yijia Shao, Varuna Jayasiri, Nick Haber, and Noah D Goodman.
510 Quiet-star: Language models can teach themselves to think before speaking. *arXiv preprint*
511 *arXiv:2403.09629*, 2024.
- 512 Di Zhang, Jiatong Li, Xiaoshui Huang, Dongzhan Zhou, Yuqiang Li, and Wanli Ouyang. Accessing
513 gpt-4 level mathematical olympiad solutions via monte carlo tree self-refine with llama-3 8b. *arXiv*
514 *preprint arXiv:2406.07394*, 2024a.
- 515 Lunjun Zhang, Arian Hosseini, Hritik Bansal, Mehran Kazemi, Aviral Kumar, and Rishabh Agarwal.
516 Generative verifiers: Reward modeling as next-token prediction, 2024b. URL [https://arxiv.
517 org/abs/2408.15240](https://arxiv.org/abs/2408.15240).
- 518 Susan Zhang, Stephen Roller, Naman Goyal, Mikel Artetxe, Moya Chen, Shuohui Chen, Christopher
519 Dewan, Mona Diab, Xian Li, Xi Victoria Lin, et al. Opt: Open pre-trained transformer language
520 models. *arXiv preprint arXiv:2205.01068*, 2022.
- 521 Tong Zhang. From ϵ -entropy to kl-entropy: Analysis of minimum information complexity density
522 estimation. *The Annals of Statistics*, pp. 2180–2210, 2006.
- 523 Xuanchang Zhang, Wei Xiong, Lichang Chen, Tianyi Zhou, Heng Huang, and Tong Zhang. From
524 lists to emojis: How format bias affects model alignment. *arXiv preprint arXiv:2409.11704*, 2024c.
- 525 Lianmin Zheng, Wei-Lin Chiang, Ying Sheng, Siyuan Zhuang, Zhanghao Wu, Yonghao Zhuang,
526 Zi Lin, Zhuohan Li, Dacheng Li, Eric Xing, et al. Judging llm-as-a-judge with mt-bench and
527 chatbot arena. *Advances in Neural Information Processing Systems*, 36:46595–46623, 2023.
- 528 Chunting Zhou, Pengfei Liu, Puxin Xu, Srinivasan Iyer, Jiao Sun, Yuning Mao, Xuezhe Ma, Avia
529 Efrat, Ping Yu, Lili Yu, et al. Lima: Less is more for alignment. *Advances in Neural Information*
530 *Processing Systems*, 36, 2024.
- 531 Banghua Zhu, Evan Frick, Tianhao Wu, Hanlin Zhu, and Jiantao Jiao. Starling-7b: Improving llm
532 helpfulness & harmlessness with rlaif, 2023.

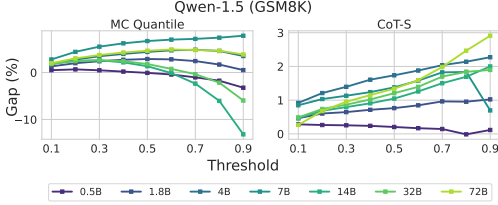


Figure 4: Change in the gap as we vary the threshold for each verification method. We only present the results for MC with quantile threshold and CoT-Score, because CoT-Binary’s gap does not change as we change the threshold.

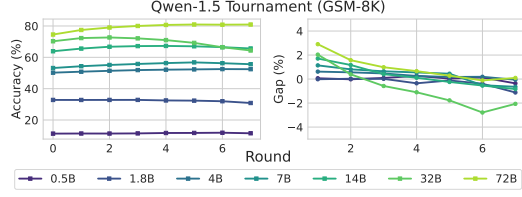


Figure 5: Change in the generation accuracy of the filtered dataset (left) and gap (right) with respect to the round of tournament. The right figure plots the gap with respect to the accuracy from the previous round instead of the base accuracy.

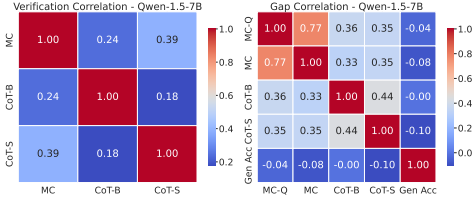
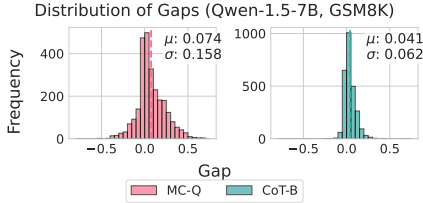


Figure 6: **Left:** The empirical distribution of gaps of MC Quantile and CoT-Binary of Qwen-1.5-7B on GSM8K. We cluster gaps in bins of intervals with width 0.005. We label the mean (μ) and standard deviation (σ) of each distribution. **Right:** the correlation plot of the output of each verification \hat{u} and the correlation plot of the gap from each verification and generation accuracy.

533 A A Fine-grained Study on Verification

534 Among the three components of self-improvement, the verification step offers the most flexibility,
 535 whereas generation and update follow more fixed procedures. Therefore, this section presents a
 536 detailed examination of the verification mechanisms. Through this focused study, we aim to uncover
 537 practical ways to enhance the overall self-improvement process.

538 A.1 Generalization of Verification Thresholds

539 In the rejection sampling framework, selecting an appropriate threshold for filtering generations based
 540 on verification is a crucial practical concern. We explore verification methods adaptable to various
 541 thresholds, including MC, CoT-Score, and Tournament. Our analysis focuses on how the performance
 542 gap changes with different thresholds for these methods. We present results for the Qwen-1.5
 543 model family using MC, MC-quantile and CoT-Score in Figure 4, and results using Tournament in Figure 5.
 544 For tournament verification, the threshold is defined as the number of rounds of tournament. We defer
 545 the results of other models to Appendix D.4. In Tournament, we note that the gap with respect to the
 546 accuracy in the previous iteration generally decreases monotonically; this trend occurs as the verifica-
 547 tion error is more likely to be exploited by the remaining generations at later stages of the tournament.

548 We observe that the relationship between the gap and the threshold is consistent across most models
 549 when using any fixed verification method. For MC and Tournament, the gap follows a concave curve
 550 relative to the threshold, while for CoT-Score, it increases monotonically. In addition, the majority of
 551 models agrees on the optimal thresholds for each verification across model families, and we use the
 552 optimal thresholds to report our results for this paper. However, in general, one should not expect the
 553 optimal threshold transfers between different tasks. That said, the consistency in threshold effects
 554 suggests a practical approach: if determining the optimal threshold for a large model is costly, one
 555 might first establish it for a smaller model and then apply it to the larger model.

556 A.2 Correlation Studies between Verification Methods

557 As the verification methods are functionally similar, one might question the necessity to study multiple
 558 verifications. To address this, We start by comparing the distribution of gaps induced by different

Table 3: Gaps of combining verification with AND operation on Qwen-2 model family. For MC, we use a quantile threshold of 0.8, for CoT-S, we use the global threshold of 8, as they are the best-performing thresholds from the analysis in the previous sections.

	MC	CoT-B	CoT-S	MC+CoT-B	MC+CoT-S	CoT-B+CoT-S	All
0.5B	3.39	0.21	0.64	3.12	3.01	0.72	3.11
1.5B	5.27	1.26	2.78	5.08	7.68	3.15	7.61
7B	4.66	2.36	1.97	5.24	4.96	3.35	5.46
72B	2.38	2.51	2.08	3.21	3.13	3.35	3.65

559 verification methods. We use Qwen-1.5 7B model as an example and we present the bar plot of the
 560 gap distribution in Figure 6. Notably, there are significant discrepancies, especially between MC
 561 and CoT methods – the variance in the gap is considerably larger with MC. This aligns with our
 562 previous findings in Section 5 where iterative self-improvement with MC verification saturates more
 563 quickly than with CoT. While CoT slightly improves the accuracy on most questions, MC will drive
 564 the accuracy to the extreme in one round of self-improvement. We observe that this pattern holds
 565 across all models, as detailed in Figure 10.

566 To further compare the verification methods, we calculate the Pearson correlation coefficient between
 567 the outputs of the proxy utility \hat{u} the gaps, shown in Figure 6. We use Qwen-1.5 7B as an example
 568 and defer full results to Figs. 11 and 12. We observe that the correlations between \hat{u} are generally low,
 569 suggesting potential benefits in combining different verification methods. Notably, the correlation
 570 between MC verifications and the correlation between CoT verifications are generally the highest,
 571 and larger models tend to have a higher correlation between the gaps. Surprisingly, the gaps of any
 572 verification method do not positively correlate with generation accuracy, reinforcing the idea that the
 573 relative gap may be a more appropriate metric for measuring self-improvement capability.

574 A.3 Improvement via Ensemble

575 The non-overlap property of different verification methods suggests the potential for enhanced
 576 verification performance through their combination. We again focus on the rejection sampling setup.
 577 In the rejection sampling framework, we employ a logical AND operation, keeping samples only if
 578 they pass all verification filters. We provide an example result from Qwen-2 model family in Table 8,
 579 and we defer the full result to Appendix D.6. We observe that combining any verifications with
 580 non-trivial gaps improves the verification performance (with the exception of CoT for 0.5B model
 581 with near 0 gap). This promising outcome indicates that despite functional similarities, different
 582 verification mechanisms can still be combined to improve self-improvement efficacy. The consistent
 583 improvements across different model sizes also suggest that strategies developed using smaller
 584 models can be effectively applied to larger ones, if all verifications are valid.

Takeaway on verification mechanisms

A fine-grained study on verification reveals several implications for practice:

- **Verification Consistency:** The distribution of the gaps and optimal verification threshold typically generalize across models.
- **Verification Distinction:** Despite functional similarities, the outputs and gaps of verification methods show non-trivial differences among each other.
- **Ensemble Heuristic:** Simple verification ensemble heuristics can improve the performance.

The consistency result suggests that configurations from smaller models can be applied to larger ones to avoid the costs associated with tuning big models. The discovery that simple ensemble techniques can enhance performance highlights the potential for more sophisticated algorithms to further improve self-improvement strategies.

585

586 B Related Works

587 **Synthetic Data and Self-Training.** Training LLMs with a mixture of “real” data (generated by
 588 human) and synthetic data has been the standard protocol nowadays given the limited number of
 589 human data and extensive amount of data required as we scale up the LLM training. Initial studies
 590 generated synthetic data from more powerful models (Gunasekar et al., 2023; Li et al., 2023b; Team

591 et al., 2023; Sun et al., 2023; Taori et al., 2023; Zhu et al., 2023; Wei et al., 2024), while recent
592 approaches involve models training on their own outputs (Achiam et al., 2023; Adler et al., 2024;
593 Dubey et al., 2024; Yang et al., 2024; Hui et al., 2024).

594 On the theoretical front, extensive research has explored the phenomenon of model collapse during
595 self-training and strategies to counter this degenerate behavior (Hataya et al., 2023; Martínez et al.,
596 2023; Bertrand et al., 2023; Briesch et al., 2023; Taori & Hashimoto, 2023; Alemohammad et al.,
597 2023; Dohmatob et al., 2024; Gillman et al., 2024).

598 **LLM Self-improvement.** One of the most effective strategies to prevent model collapse during
599 self-training is the use of a reliable verifier (Gillman et al., 2024). In the absence of additional
600 resources like labeled data or an external oracle, models can utilize their own verification capabilities.
601 This is particularly effective if the model is more proficient at verification than generation. Numerous
602 studies have proposed variations of self-improvement algorithms based on this principle, resulting
603 in significant practical achievements (Zelikman et al., 2022; Wang et al., 2022b; Huang et al., 2022;
604 Singh et al., 2023; Chen et al., 2023; Madaan et al., 2024; Xu et al., 2024; Yuan et al., 2024; Liang et al.,
605 2024; Wang et al., 2024b; Shinn et al., 2024; Zelikman et al., 2024). Previous research, however, often
606 relied on additional data to enhance verification, used surrogate metrics for improvement, or limited
607 their focus to a small number of models. In this work, instead of proposing any new algorithm, we
608 aim to rigorously analyze the self-improvement phenomenon in a controlled, comprehensive manner.

609 **Improving Test-time Inference with Additional Computation.** Recent research has demonstrated
610 that the performance of models can be enhanced by allocating more computational resources to
611 inference (Welleck et al., 2024). This typically leverages the observation that LLMs can make diverse
612 generations, and with a small probability it can generate high-quality responses (Li et al., 2022;
613 Brown et al., 2024; Bansal et al., 2024). Thus with oracle verifier, or with training a high-quality
614 reward model, model performance can be improved by simply making multiple generations and
615 selecting the best ones according to the oracle or the reward model (Cobbe et al., 2021). There are
616 also works on training process-based reward models (Lightman et al., 2023) to improve the model’s
617 reasoning results (Luo et al., 2024; Wang et al., 2024a; Zhang et al., 2024a).

618 Concurrently there are also works on the test-time scaling law which investigates the computational
619 trade-off between the model size (which determines the number of generations given a computation
620 budget) and final accuracy combined with reward model or oracle (Wu et al., 2024; Snell et al., 2024).
621 The results provide the compute-optimal solution for test-time inferencing with a fixed compute
622 budget and a fixed verifier. We believe a better understanding of self-improvement can also lead
623 to a test-time scaling law without an external verifier.

624 **LLM-as-a-Judge.** LLM-as-a-judge refers to using an LLM to verify the generation of some other (or
625 the same) LLM (Chiang et al., 2023; Zheng et al., 2023; Bubeck et al., 2023; Chiang & Lee, 2023;
626 Zhou et al., 2024). Recently the same idea has also been applied to train a generative reward model
627 (Ankner et al., 2024; Zhang et al., 2024b). Having a model that can verify its own generation is one
628 of the key components of self-improvement, and in this work, we perform a fine-grained study on
629 various types of LLM verification mechanisms.

630 **Reranking Algorithms.** The self-improvement framework we study in this paper relies on reweight-
631 ing the generation distribution. Prior to self-improvement, the reranking algorithm has already been
632 widely applied in various NLP applications (Collins & Koo, 2005; Huang & Chiang, 2007; Stiennon
633 et al., 2020; Cobbe et al., 2021; Krishna et al., 2022; Lightman et al., 2023).

634 C Verification Mechanisms

635 In this section, we provide a more complete description of the verification mechanism we use
636 throughout the paper.

- 637 • **Multiple Choice (MC):** Multiple choice verification asks the LM to label responses as “Correct”
638 and “Incorrect”. Let prom_{mc} be a verification prompt and denote $\hat{u}_f^{\text{mc}}(x, y)$ a utility derived from
639 the verifier generating a single token t^+, t^- , representing the word “Correct” and “Incorrect”
640 respectively. The score uses the logits from these tokens to find the probability of “Correct”

641 conditioned on the next token being “Correct” or “Incorrect”:

$$\hat{u}_f^{\text{mc}}(x, y) := \frac{f(t^+ | x, y, \text{prom}_{\text{mc}})}{f(t^+ | x, y, \text{prom}_{\text{mc}}) + f(t^- | x, y, \text{prom}_{\text{mc}})}.$$

642 • **Chain of Thought (CoT):** CoT verification asks the LM to score responses and to provide the CoT.
 643 Denote by $\mathcal{S} \subset \mathbb{R}$ the set of verification scores and by $\text{prom}_{\mathcal{S}}$ a verification prompt. We can define
 644 a utility

$$\hat{u}_f^{\mathcal{S}}(x, y) := \mathbb{E}_{s, z \sim f(\cdot | x, y, \text{prom}_{\mathcal{S}})}[s(z)],$$

645 where z is the verification CoT, and $s(z) \in \mathcal{S}$ is the score extracted from the CoT. In our experiments
 646 we consider two versions, CoT-Score with $\mathcal{S} = [10]$ and CoT-Binary with $\mathcal{S} = \{0, 1\}$.

647 • **Tournament (To):** The tournament verification does not directly fit the utility framework described
 648 in Section 2. Rather, this verification procedure involves comparisons of a batch of generations to
 649 provide a modified distribution⁴. Given a comparison prompt prom_{com} , we perform a tournament-
 650 style elimination over a batch of 2^r generations by comparing disjoint pairs in each round until a
 651 single generation remains. Let $\mathcal{Y}^{(0)} = y_1, y_2, \dots, y_{2^r}$ be the initial set of generations. At round k ,
 652 the set $\mathcal{Y}^{(k)}$ contains 2^{r-k} remaining generations. These are split into disjoint pairs $(y_i, y_j) \in \mathcal{Y}^{(k)}$.
 653 Each pair is compared using the prompt prom_{com} , and the verifier’s output $s \in A, B$ indicates the
 654 preferred generation:

$$y_{\text{win}} = \begin{cases} y_i & \text{if } f(\cdot | y_i, y_j, \text{prom}_{\text{com}}) = A, \\ y_j & \text{if } f(\cdot | y_i, y_j, \text{prom}_{\text{com}}) = B. \end{cases}$$

655 where y_{win} is the winner of the pairwise comparison and advances to the next round. After each
 656 round k , the set of winners $\mathcal{Y}^{(k+1)}$ contains half the number of generations. This process is repeated
 657 until $k = r$, leaving only one generation, the lone element in \mathcal{Y}^r .

658 D Additional Results

659 D.1 Additional Results for Section 4.1

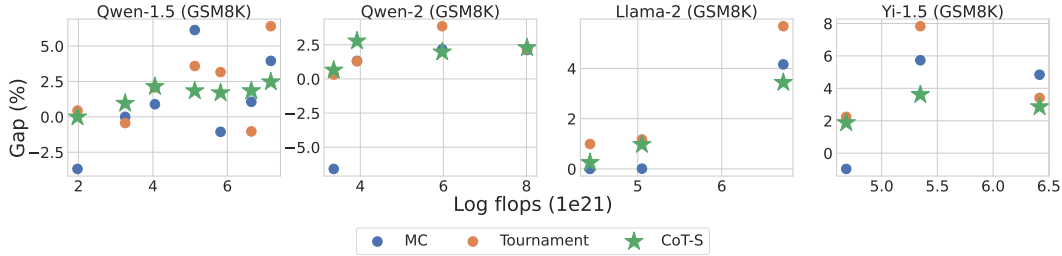


Figure 7: Scaling phenomenon of the gap with respect to the pretrain flops. The x-axis is the log of pretrain flops, and the y-axis is the relative gap. MC denotes Multiple Choice verification with quantile threshold 0.8 (top 0.8), CoT-S denotes CoT-Score verification with global threshold 8, and To denotes Tournament verification with 5 rounds.

⁴This batch-style distribution weighting also applies to strategies like top-k wherein we take the highest k utility generations for a particular question.

Table 4: Gap on GSM-8K for all models. For each verification, “top n ” denotes taking the threshold as the n quantile of the proximity utility for each prompt, and $\tau = n$ denotes taking the threshold as n for all prompts. All numbers denote the percentage.

Name	Size	Accuracy	MC				CoT			To round 5
			top 0.7	top 0.8	$\tau = 0.7$	$\tau = 0.8$	Bin	S ($\tau = 8$)	S ($\tau = 9$)	
Qwen-1.5	0.5B	11.31	-0.62	-1.42	-2.27	-3.68	-0.02	-0.01	0.12	0.44
	1.8B	32.81	3.06	2.70	-0.01	0.00	1.58	0.95	1.02	-0.44
	4B	50.21	4.85	4.63	0.48	0.89	2.36	2.14	2.27	2.08
	7B	53.17	7.18	7.42	3.68	6.13	4.07	1.84	0.70	3.59
	14B	63.87	-2.04	-5.61	2.36	-1.06	1.79	1.69	2.00	3.16
	32B	70.25	-0.22	-1.98	1.72	1.06	3.07	1.84	1.90	-1.03
	72B	74.55	4.84	4.75	2.00	3.95	3.55	2.47	2.91	6.40
Qwen-2	0.5B	26.19	4.59	3.39	3.81	-6.59	0.21	0.64	0.72	0.32
	1.5B	48.82	6.09	5.27	5.02	1.32	1.26	2.78	2.80	1.29
	7B	76.42	4.69	4.66	3.90	2.17	2.36	1.97	2.08	3.86
	72B	81.69	2.39	2.38	0.89	2.12	2.51	2.08	2.45	2.20
Llama-2	7B	11.64	2.33	2.20	0.10	0.00	0.16	0.25	0.23	0.99
	13B	21.57	3.45	3.18	-0.13	0.01	1.13	0.97	1.01	1.17
	70B	48.42	5.14	4.91	4.77	4.16	3.98	3.44	3.45	5.68
Llama-3	8B	45.66	5.34	5.33	-0.50	0.44	2.67	2.10	2.13	4.08
	70B	74.19	5.06	4.52	4.68	1.35	3.89	2.59	2.72	2.00
Llama-3.1	8B	49.31	4.68	4.57	-2.25	-0.24	3.37	2.09	2.05	4.78
	70B	71.71	6.88	6.77	6.00	0.93	3.29	2.71	2.88	-0.40
Yi-1.5	6B	55.53	4.40	4.27	2.98	-0.97	2.01	1.88	1.95	2.24
	9B	61.04	7.74	7.50	5.61	5.73	2.32	3.61	3.72	7.83
	34B	73.71	6.29	6.23	3.34	4.84	2.49	2.86	2.95	3.41

660 **D.2 Additional Results for Section 4.3**

Model	Qwen-1.5						
	0.5B	1.8B	4B	7B	14B	32B	72B
661 Generation Accuracy	6.20	11.40	17.16	21.20	26.83	35.79	39.97
MC (top 0.2)	-0.62	-0.02	-1.20	-1.15	-1.99	-0.43	-0.75
MC ($\tau = 0.8$)	-0.51	-0.32	-0.72	-1.07	-1.21	-0.10	0.32

Model	Qwen-2			
	0.5B	1.5B	7B	72B
662 Generation Accuracy	6.51	13.87	29.09	41.45
MC (top 0.2)	-0.06	0.04	0.79	0.28
MC ($\tau = 0.8$)	-0.05	0.02	-0.05	-0.05

Model	Llama-2		
	7B	13B	70B
663 Generation Accuracy	25.52	41.00	29.09
MC (top 0.2)	-0.96	0.76	0.30
MC ($\tau = 0.8$)	-0.81	-2.31	-0.44

Model	Llama-3	
	8B	70B
664 Generation Accuracy	30.40	45.59
MC (top 0.2)	0.27	0.32
MC ($\tau = 0.8$)	-0.23	-0.41

Table 5: Relative gaps on GSM-8K for all models. For each verification, “top n ” denotes taking the threshold as the n quantile of the proximity utility for each prompt, and $\tau = n$ denotes taking the threshold as n for all prompts. All numbers denote the percentage

Name	Size	Accuracy	MC				CoT			To round 5
			top 0.7	top 0.8	$\tau = 0.7$	$\tau = 0.8$	Bin	S ($\tau = 8$)	S ($\tau = 9$)	
Qwen-1.5	0.5B	11.31	-0.70	-1.60	-2.56	-4.15	-0.03	-0.13	0.04	0.59
	1.8B	32.81	4.55	4.01	-0.01	-0.01	4.45	3.10	3.15	-0.39
	4B	50.21	9.75	9.31	0.96	1.80	8.22	6.92	8.12	5.24
	7B	53.17	15.34	15.84	7.87	13.08	14.06	5.50	-0.35	10.26
	14B	63.87	-5.64	-15.54	6.52	-2.94	7.77	8.00	9.02	12.19
	32B	70.25	-0.75	-6.67	5.78	3.57	15.31	8.40	8.49	-4.55
	72B	74.55	19.01	18.65	7.87	15.52	21.82	14.96	16.97	32.06
Qwen-2	0.5B	26.19	6.22	4.59	5.16	-8.93	0.44	1.58	1.79	1.31
	1.5B	48.82	11.90	10.30	9.82	2.58	4.45	9.96	10.11	5.26
	7B	76.42	19.89	19.74	16.53	9.18	16.46	13.73	14.34	17.17
	72B	81.69	13.07	13.00	4.87	11.57	20.16	14.34	19.92	17.84
Llama-2	7B	11.64	2.64	2.49	0.11	0.00	0.25	0.39	0.37	1.91
	13B	21.57	4.40	4.05	-0.16	0.02	2.45	1.79	1.82	2.88
	70B	48.42	9.97	9.52	9.25	8.07	13.77	12.22	12.07	18.57
Llama-3	8B	45.66	9.62	9.60	-0.90	0.80	8.51	6.16	6.70	12.55
	70B	74.19	18.08	16.13	16.71	4.84	19.09	13.24	13.74	11.28
Llama-3.1	8B	49.31	8.97	8.77	-4.31	-0.47	11.74	7.02	6.89	14.97
	70B	71.71	22.83	22.47	19.91	3.09	16.68	12.65	13.59	-1.03
Yi-1.5	6B	55.53	9.89	9.60	6.70	-2.17	9.10	6.60	6.69	10.64
	9B	61.04	19.86	19.25	14.40	14.70	10.09	14.29	14.35	24.78
	34B	73.71	23.94	23.68	12.70	18.40	15.02	16.69	16.96	8.89

665

Model	Llama-3.1	
	8B	70B
Generation Accuracy	27.75	45.13
MC (top 0.2)	0.42	0.44
MC ($\tau = 0.8$)	-0.59	-0.37

666

Model	Yi-1.5		
	6B	9B	34B
Generation Accuracy	22.82	25.94	35.31
MC (top 0.2)	0.09	0.21	0.24
MC ($\tau = 0.8$)	-0.07	0.61	0.30

Table 6: Gap on Natural Question for all models. With non-trivial generation accuracy, all gaps are near 0, indicating that the task is non-improvable.

667

D.3 Additional Results for Section 4.4

668

Model	Qwen-1.5						
	0.5B	1.8B	4B	7B	14B	32B	72B
Generation Accuracy	0.43	1.00	0.88	0.95	1.57	2.67	2.02
Gap	0.02	-0.03	-0.15	-0.64	0.22	0.07	1.23
Relative Gap	-0.10	-2.80	-1.39	-3.06	0.67	-1.25	1.14

Model	Qwen-2					
	0.5B	1.5B	7B	72B	7B-Instruct	72B-Instruct
669 Generation Accuracy	0.66	0.62	2.09	8.82	2.16	8.15
Gap	-0.09	0.04	-0.07	16.99	0.13	22.97
Relative Gap	-0.15	-0.61	-0.01	20.81	0.20	26.40

Model	Llama-2		
	7B	13B	70B
670 Generation Accuracy	0.82	0.89	0.86
Gap	-0.13	-0.63	-0.86
Relative Gap	0.45	-2.02	-3.57

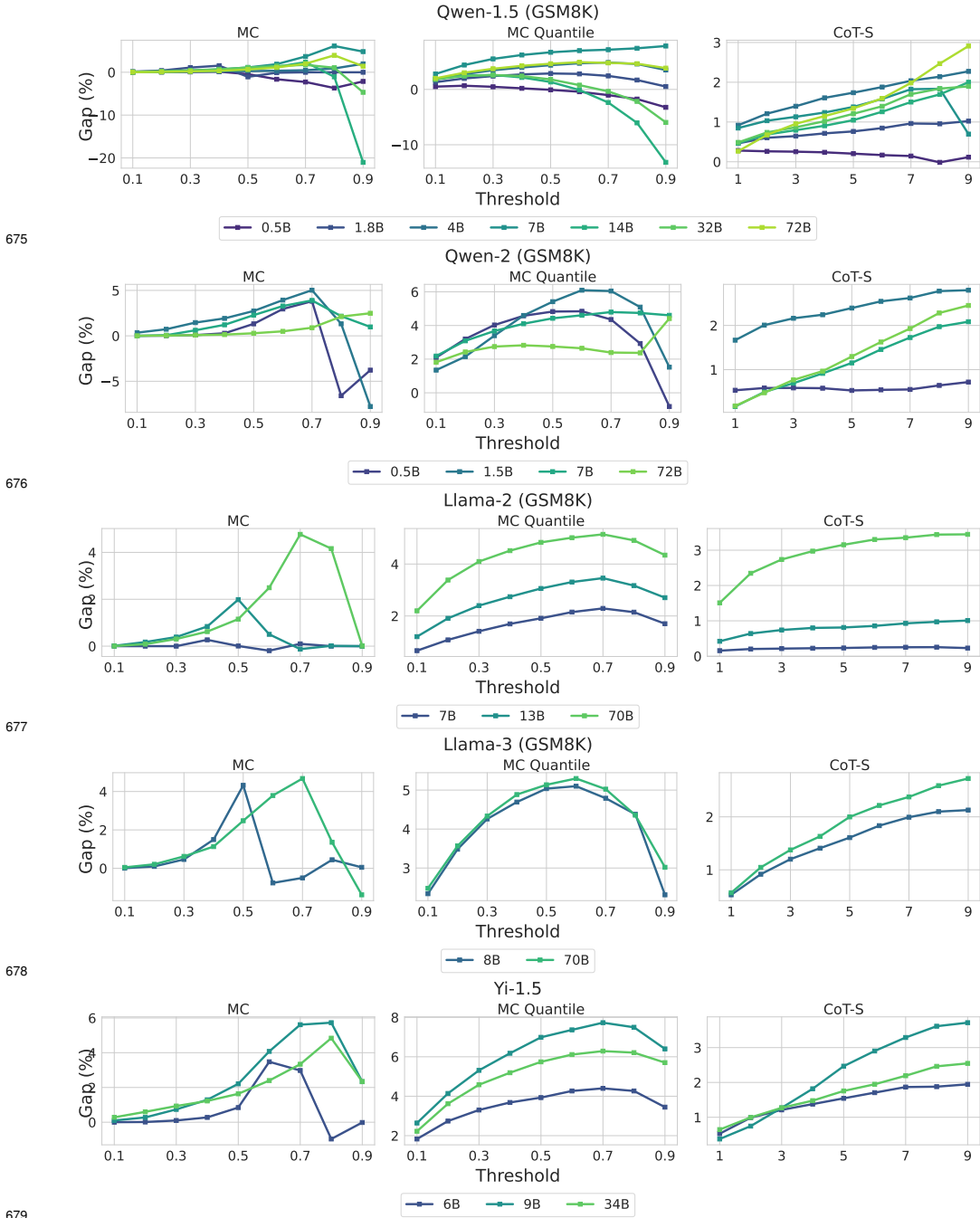
Model	Llama-3	
	8B	70B
671 Generation Accuracy	1.39	1.63
Gap	-1.10	-0.84
Relative Gap	-15.4	-36.12

Model	Llama-3.1	
	8B	70B
672 Generation Accuracy	1.11	1.68
Gap	-0.19	5.5
Relative Gap	-4.52	6.87

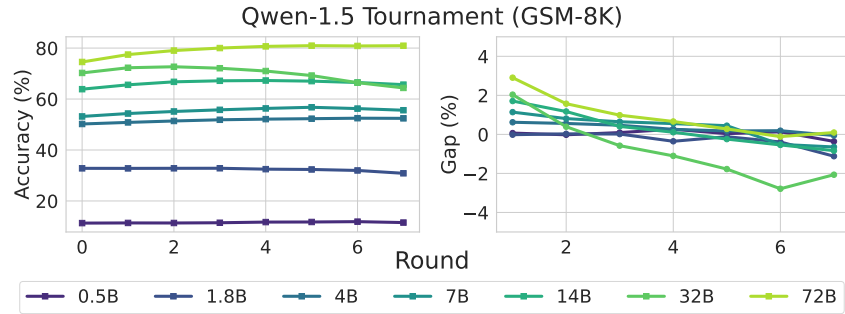
Model	Yi-1.5		
	6B	9B	34B
673 Generation Accuracy	0.59	1.29	4.48
Gap	-0.60	0.22	-1.75
Relative Gap	-0.94	0.43	-0.77

Table 7: Generation accuracy, gap and relative gap on Sudoku for all models.

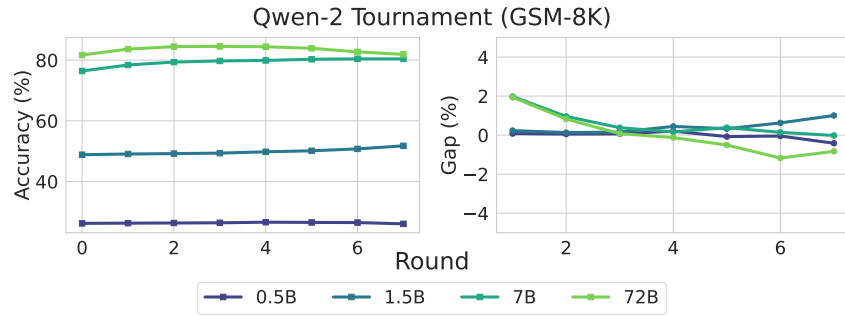
674 **D.4 Additional Results for Appendix A.1**



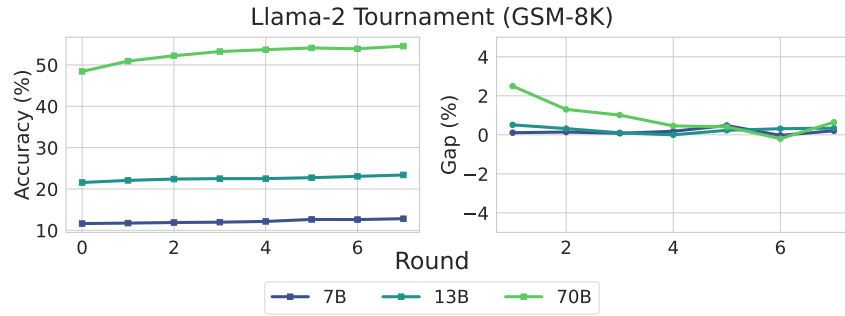
675
676
677
678
679
Figure 8: Change in the gap as we vary the threshold for each verification method. We only present the results for MC with global threshold, quantile threshold and CoT-Score, because CoT-Binary’s gap does not change as we change the threshold.



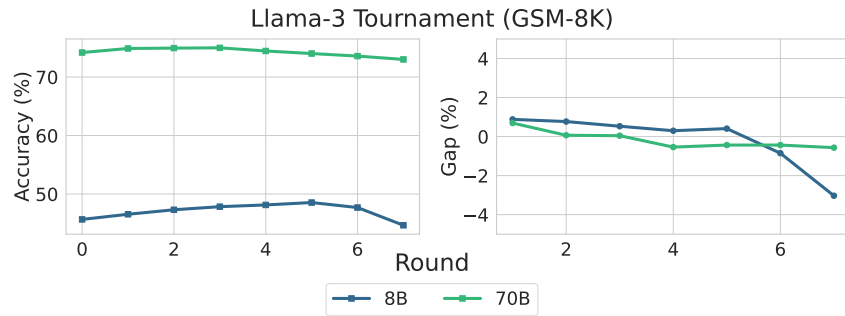
680



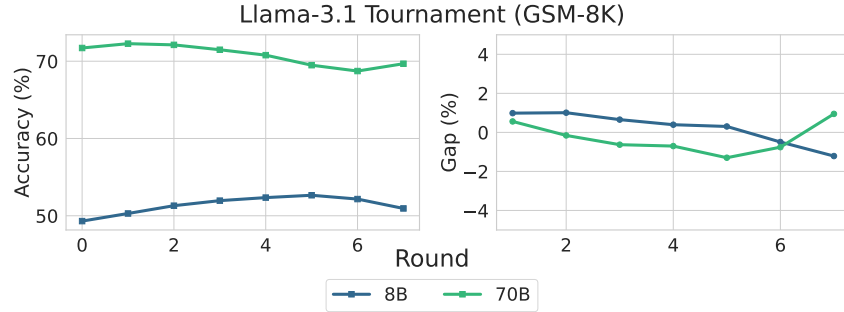
681



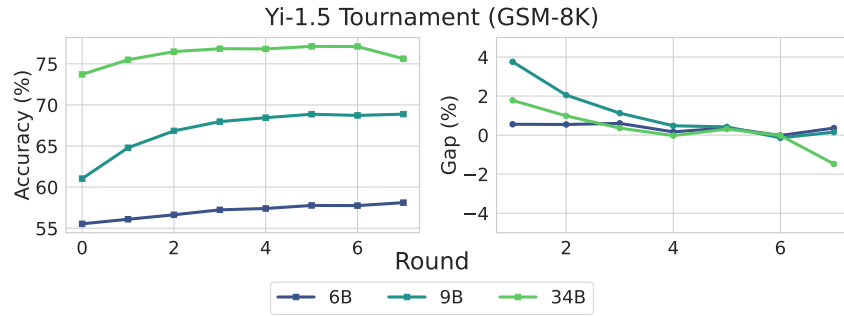
682



683



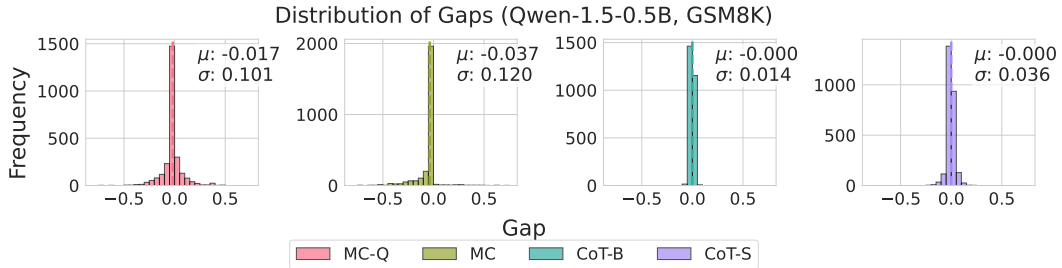
684



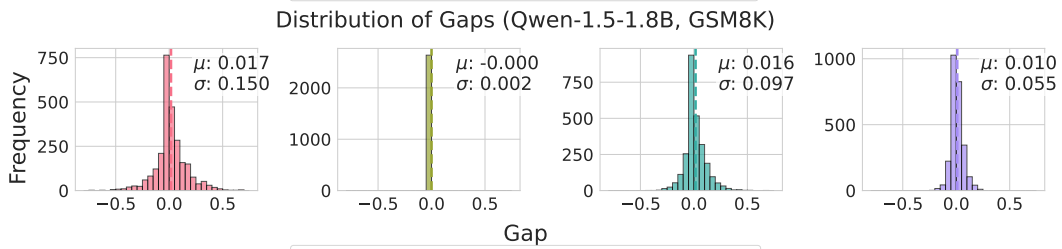
685

Figure 9: Tournament

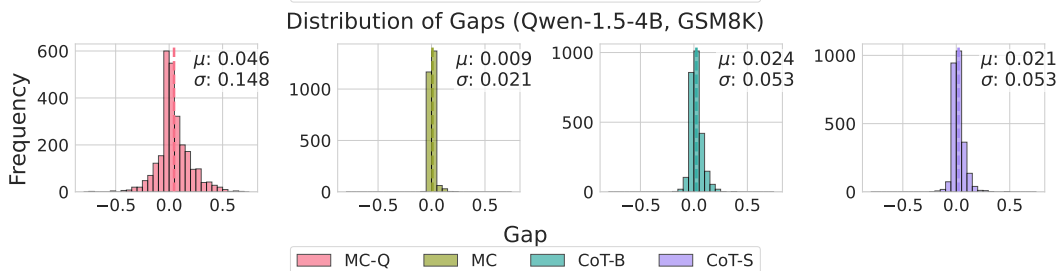
686 **D.5 Additional Results for Appendix A.2**



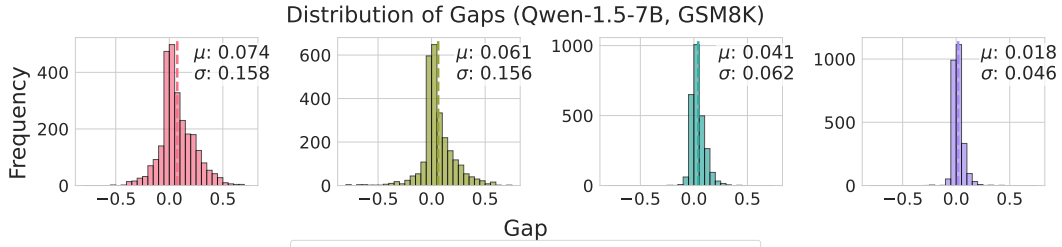
687



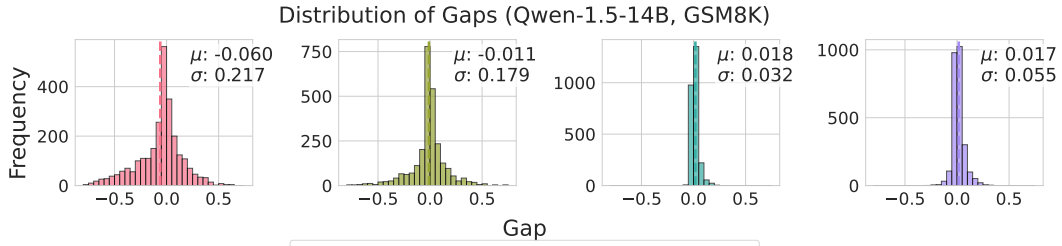
688



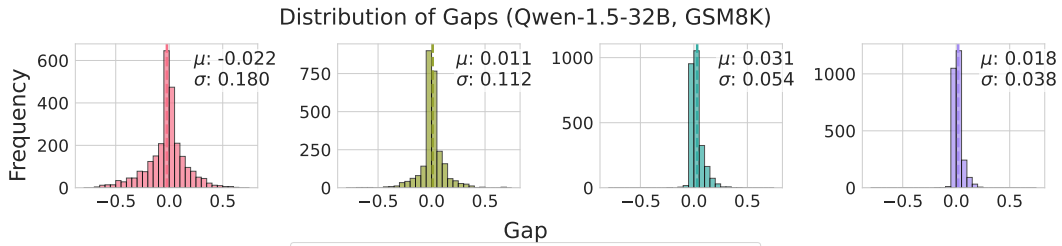
689



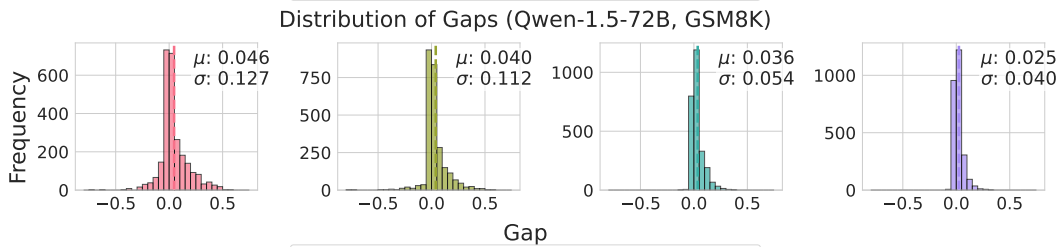
690



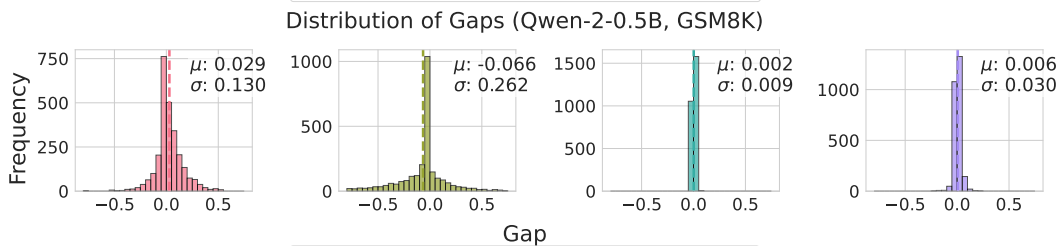
691



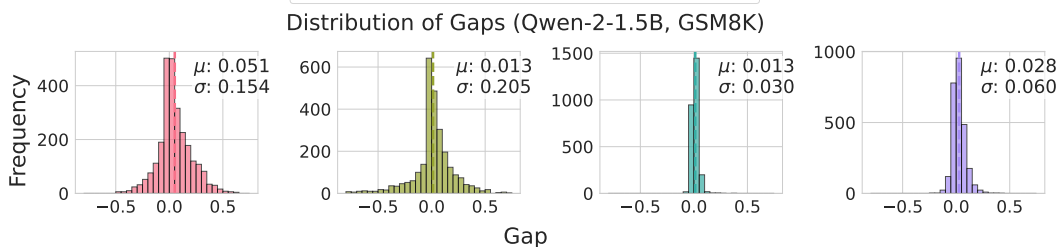
692



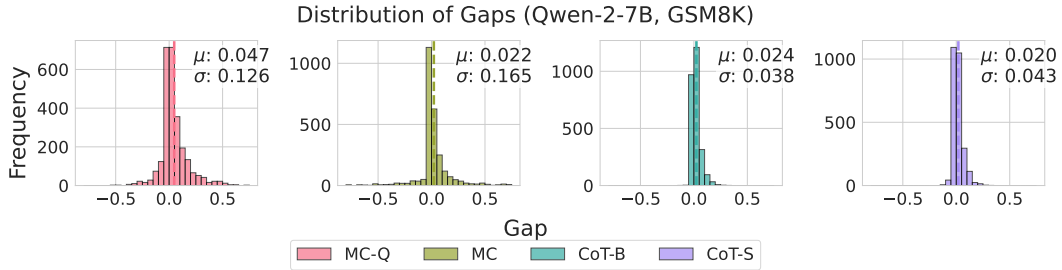
693



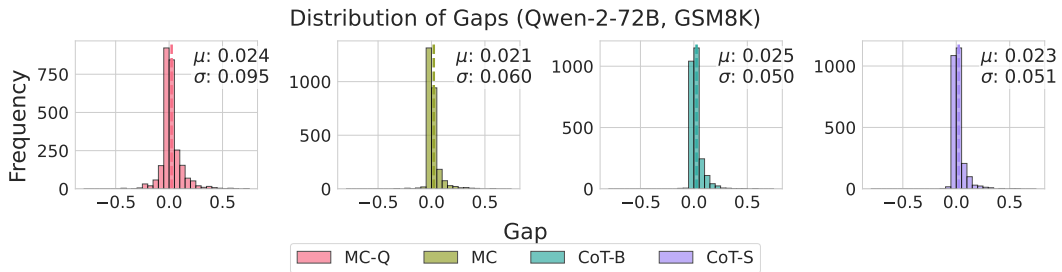
694



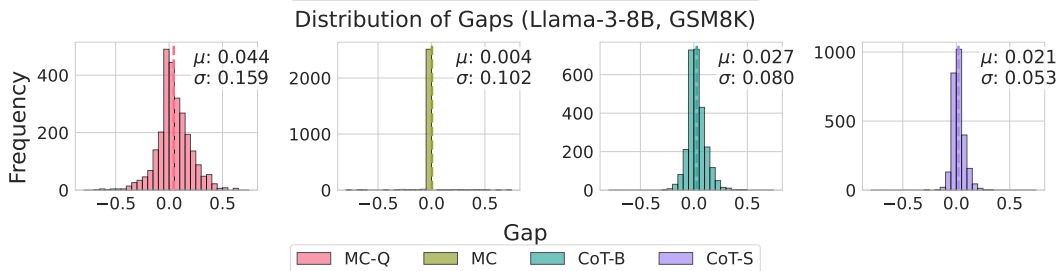
695



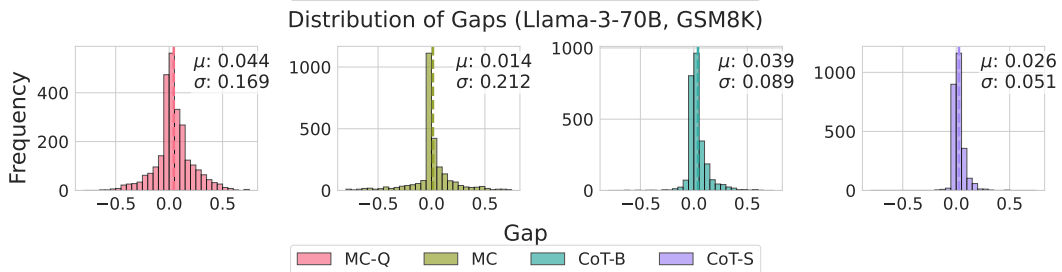
696



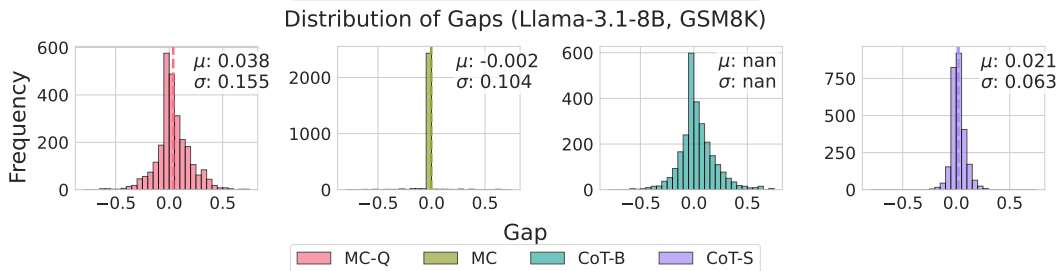
697



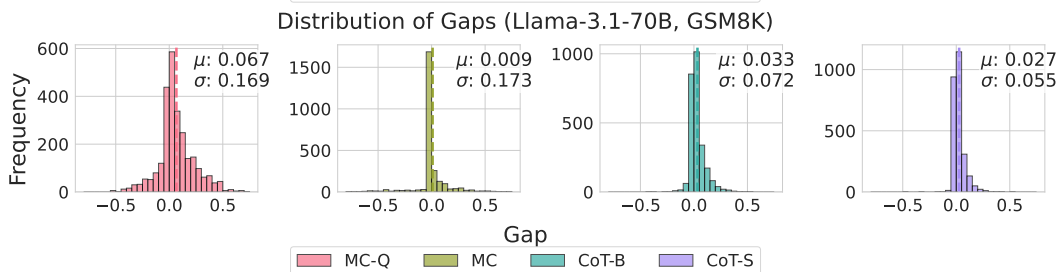
698



699



700



701

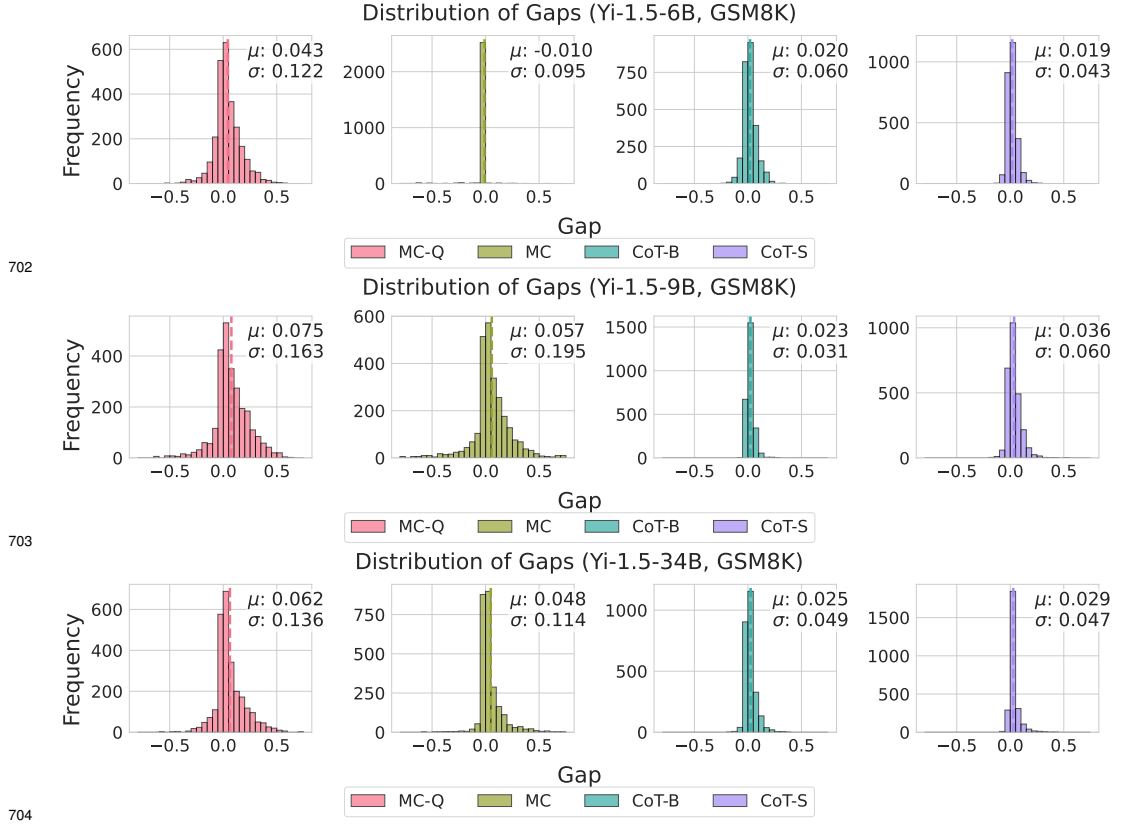
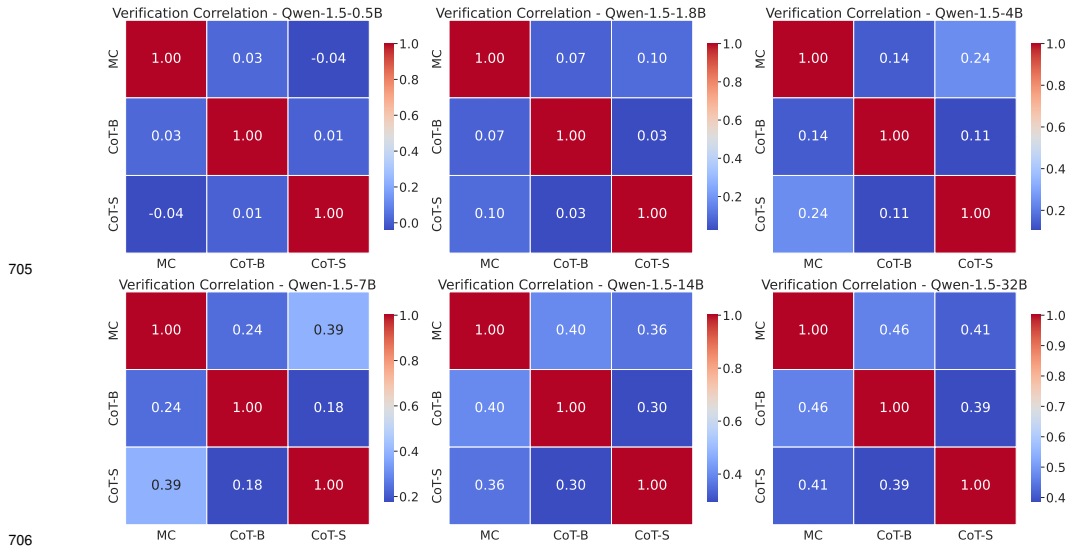


Figure 10: The empirical distribution of gaps of each verification method of each model on GSM8K. We cluster gaps in bins of intervals with width 0.005. We label the mean (μ) and standard deviation (σ) of each distribution.



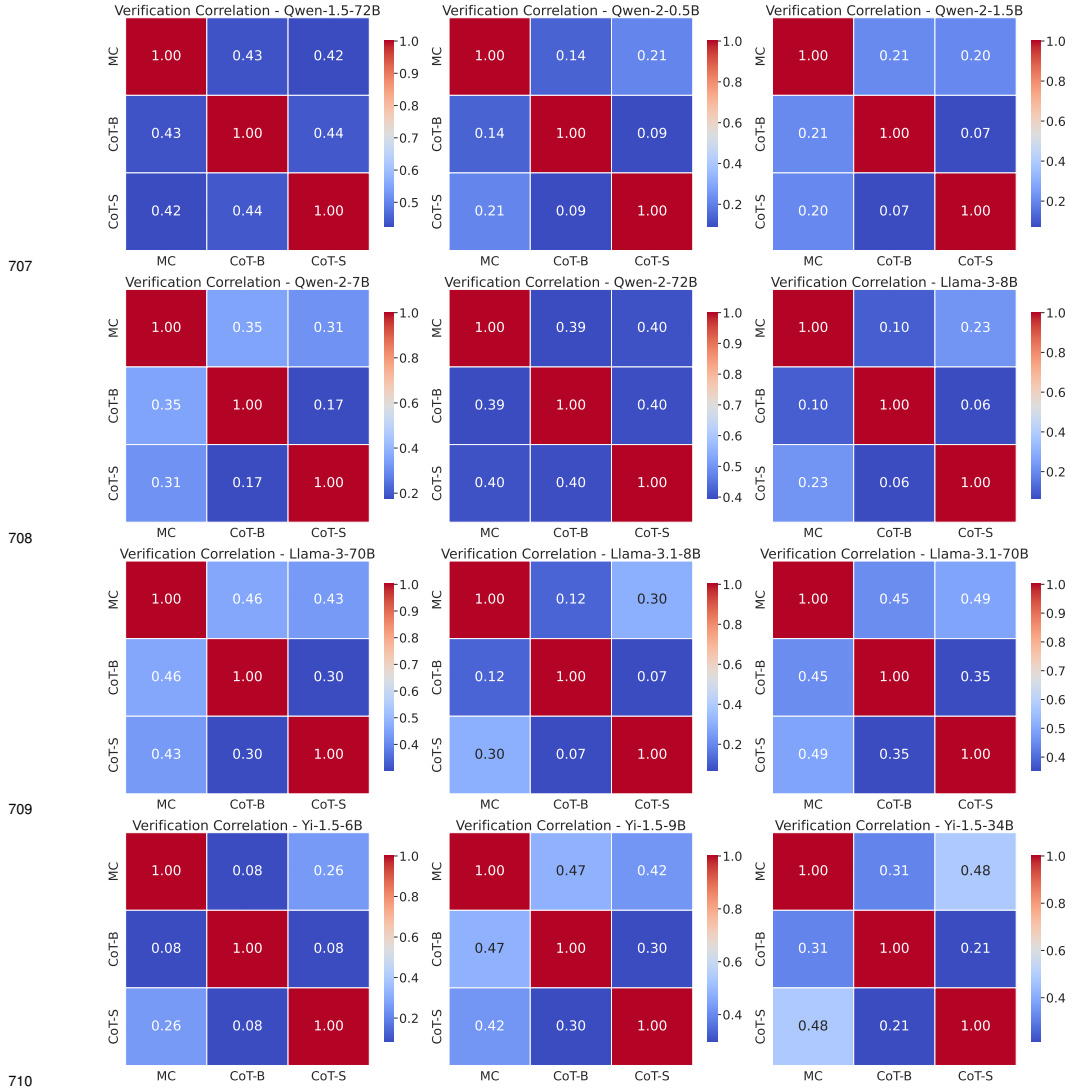
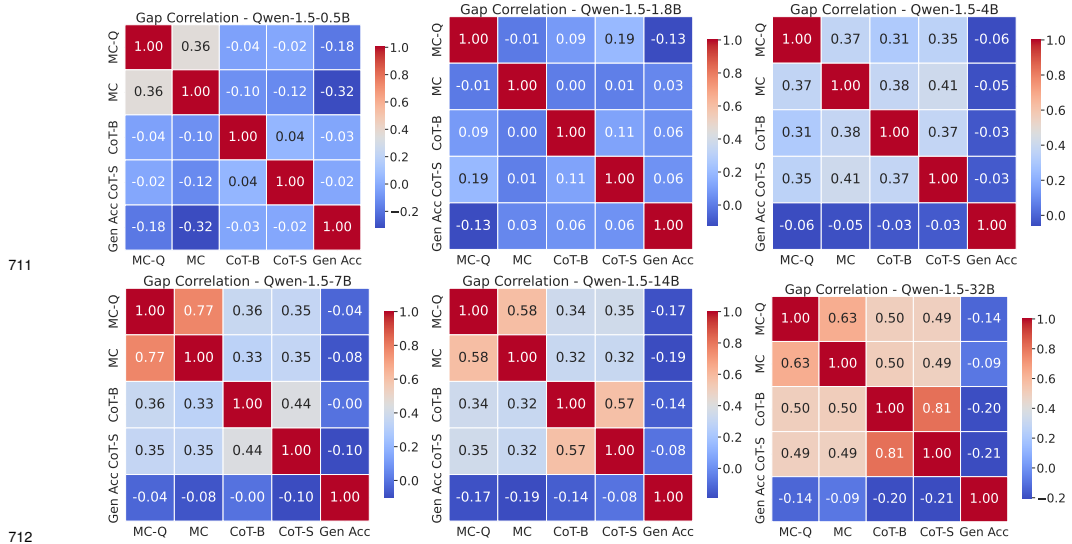


Figure 11: The correlation plot of the output of each verification \hat{u} .



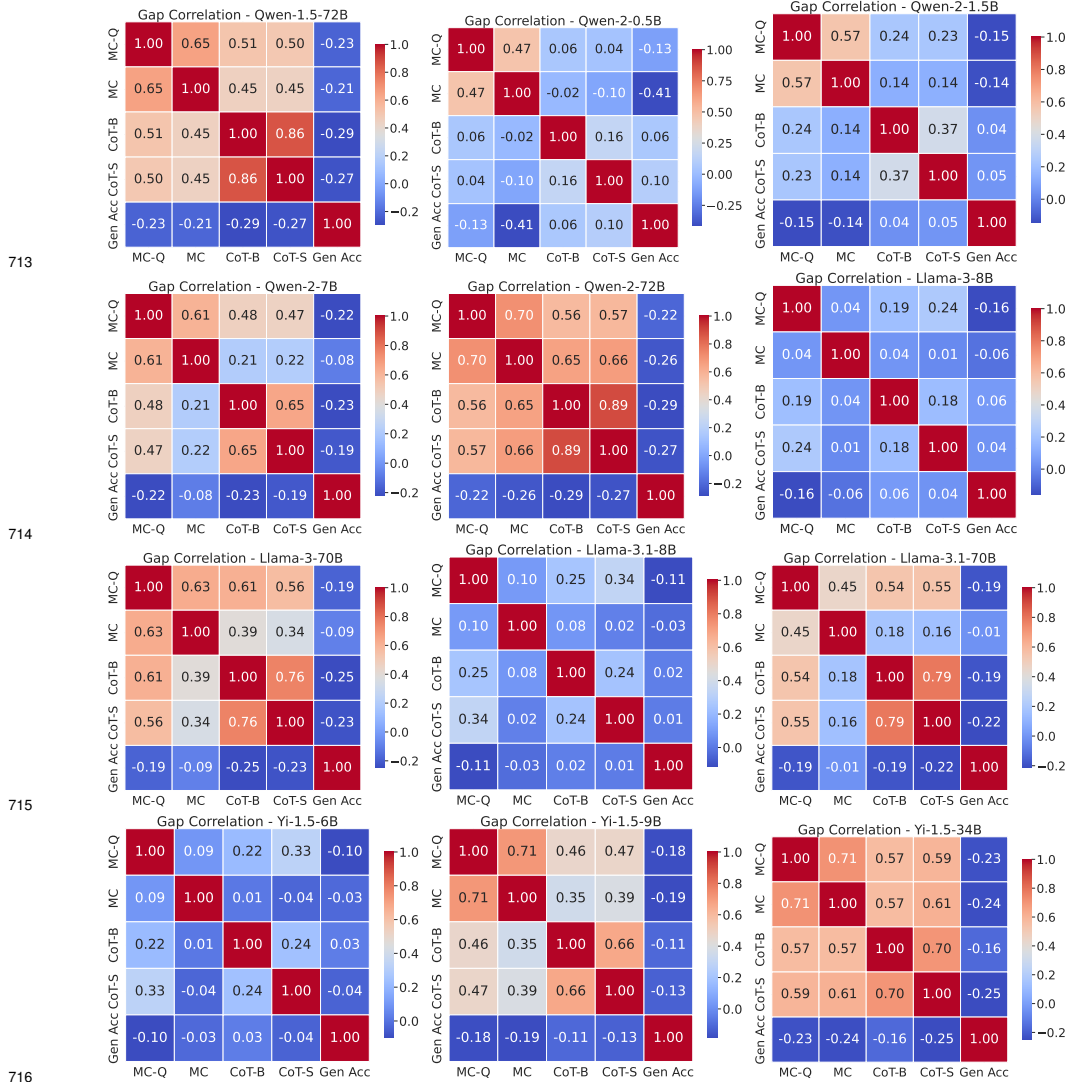


Figure 12: The correlation plot of the gap from each verification and generation accuracy.

Table 8: Relative gaps on GSM-8K for all models. For each verification, “top n ” denotes taking the threshold as the n quantile of the proxy utility for each prompt, and $\tau = n$ denotes taking the threshold as n for all prompts. All numbers denote the percentage

Name	Size	MC	CoT-B	CoT-S	MC+CoT-B	MC+CoT-S	CoT-B+CoT-S	All
Qwen-1.5	0.5B	-1.62	-0.07	-0.01	-1.65	-1.36	-0.08	-1.42
	1.8B	2.36	1.04	0.95	2.16	2.22	1.18	2.14
	4B	4.55	2.26	2.14	4.81	4.77	3.50	4.90
	7B	7.46	3.98	1.84	9.02	7.88	4.88	9.12
	14B	-6.09	1.79	1.69	-5.73	-5.85	2.45	-5.60
	32B	-2.30	3.07	1.84	-1.87	-2.21	3.62	-2.00
	72B	4.61	3.43	2.47	5.70	5.25	4.46	6.07
Qwen-2	0.5B	3.01	0.18	0.64	3.01	3.13	0.72	3.11
	1.5B	5.04	0.86	2.78	5.08	7.68	3.15	7.61
	7B	4.68	2.31	1.97	5.25	4.96	3.40	5.46
	72B	2.44	2.50	2.28	3.21	3.13	3.35	3.65
Llama-2	7B	2.17	0.13	0.25	2.21	2.42	0.37	2.44
	13B	3.19	0.91	0.97	3.47	3.38	1.76	3.21
	70B	4.78	3.73	3.44	5.52	5.61	5.79	4.90
Llama-3	8B	4.98	2.59	2.10	5.88	5.81	4.09	5.83
	70B	4.37	3.88	2.59	4.91	4.71	4.45	4.98
Llama-3.1	8B	4.34	3.50	2.09	4.47	4.57	3.27	2.96
	70B	6.72	3.29	2.71	7.08	7.03	3.94	7.21
Yi-1.5	6B	4.27	2.01	1.88	4.83	4.72	3.16	4.95
	9B	7.50	2.32	3.61	7.79	7.87	4.80	8.11
	34B	6.23	2.49	2.86	6.32	6.35	3.76	6.41

718 **E Hyperparameters**

Table 9: Hyperparameter for Iterative Self-improvement

Minibatch size	64
Learning rate	1e-6
Optimizer	AdamW
Gradient step	2000
Max Sequence Length	2048
Data Type	bf16

719 **F Generation and Verification Prompts**

Multiple Choice Verification Prompt (GSM8K / nq_open)

Judge the correctness of the following solution of the problem. Answer with either Correct or Incorrect. Problem: {problem}
Solution: {generation}
Judge:

720

Chain of Thoughts Binary Prompt (GSM8K)

Review the following math problem and the attempted solution and verify the correctness of the attempted solution, with a judgement of <correct> or <incorrect>. Your judgement should follow each criterion below: - The final ANSWER is after the phrase "The answer is ANSWER", and verify if the answer is correct with respect to the problem. If there is no such phrase, treat the answer as incorrect. - Each solution contains a derivation before the final answer, check the soundness of the derivation as well. - Your final judgement should reflect solely on the correctness of the final answer, but if there are issues in the derivation, please mention them in your justification.

Problem: {problem}

Attempted Solution: {generation}

After examining the problem and the attempted solution: - Briefly justify your judgement, up to 50 words. - Conclude with the judgement using the format: "Correctness: <correct> or <incorrect>".

Remember to assess from the math verifier perspective and be critical and verify carefully.

Judgement:

721

Chain of Thoughts Score Prompt (GSM8K)

Review the following math problem and the attempted solution and give a score from 1 to 10 to the attempted solution. The final ANSWER is after the phrase "Final Answer: The final answer is ANSWER". Give the answer a 1 if there is no such phrase or ANSWER is wrong, and give the answer a 10 if both the answer and the derivation are correct.

Problem: {problem}

Attempted Solution: {generation}

After examining the problem and the attempted solution: - Briefly justify your score, up to 50 words. - Conclude with the score using the format: "Score: <score>".

Remember to assess from the math verifier perspective and be critical and verify carefully.

Judgement:

722

Tournament Prompt (GSM8K)

Review the following math problem and two attempted solutions. Your task is to determine the better solution between the two. Your judgement should follow each criterion below: - The final ANSWER is after the phrase "The answer is ANSWER", and you should always prefer correct answers over incorrect answer. - Always prefer solutions with the phrase "The answer is ANSWER" over ones without it. - If both answers are correct or incorrect, you should prefer the one with better reasonings. Problem: {problem}

Solution A: {generation1}

Solution B: {generation2}

After examining the problem and the attempted solutions: - Briefly justify your judgement, up to 50 words. - Conclude with the judgement using the format: "Preferred solution: <A> or ".

Remember to assess from the math verifier perspective and be critical and verify carefully.

Judgement:

723

Sudoku Generation Prompt

You are a Sudoku solver specialized in 4x4 puzzles. You will be given a string of 16 digits representing an initial 4x4 Sudoku puzzle, where 0 represents an empty cell. Your task is to solve the puzzle and provide the complete solution.

Rules for solving a 4x4 Sudoku:

Each row must contain the numbers 1-4 without repetition.

Each column must contain the numbers 1-4 without repetition.

Each 2x2 quadrant must contain the numbers 1-4 without repetition.

The solution must maintain all the initial non-zero numbers from the puzzle in their original position.

Input format:

ab cd

ef gh

ij kl

mn op

Where each letter represents a single digit (0-4), and 0 indicates an empty cell. Your tasks:

Analyze the given puzzle.

Fill in all empty cells (0s) with appropriate numbers following Sudoku rules.

Provide your solution in the same format as the input, maintaining the original structure.

Puzzle: {problem}

Answer:

724

Sudoku Generation Example

Puzzle:

03 21

00 30

04 00

21 00

Answer: I filled A1 with 4 to complete Row 1. In Quadrant 1, the missing numbers 1 and 2 were assigned to B1 and B2 based on Column constraints. B4 was set to 4 to complete Row 2. In Quadrant 3, C1 was filled with 3. Row 3 missing numbers 1 and 2 were placed in C3 and C4. D4 was assigned 3 to complete Column 4, and D3 was set to 4 to finish Row 4.

Solution:

43 21

12 34

34 12

21 43

725

Sudoku Verification Prompt

You are a Sudoku validator for 4x4 puzzles. You will be given two strings of 16 numbers each: the first represents the initial Sudoku puzzle, and the second represents the proposed solution. Your task is to determine if the solution is correct.

Rules for a valid 4x4 Sudoku solution:

Each row must contain the numbers 1-4 without repetition.

Each column must contain the numbers 1-4 without repetition.

Each 2x2 quadrant must contain the numbers 1-4 without repetition.

The solution must maintain all the initial numbers from the puzzle.

Input format:

Question: Input format:

ab cd

ef gh

ij kl

mn op

Where each letter represents a single digit (0-4), and 0 indicates an empty cell. Answer: [same format but there should be no 0]

Your task:

Check if the solution maintains all initial numbers from the puzzle.

Verify that each row, column, and 2x2 quadrant contains 1-4 without repetition.

Respond with "Correct" if the solution is valid, or "Incorrect" if it's not.

Question: {problem}

Answer: {generation}

Judgment:

726

Sudoku Verification Example

problem:

03 00

21 00

00 01

10 43

generation:

34 12

21 31

34 21

12 44

correctness: Checking if the solution maintains all initial numbers from the puzzle:

Position (Row 1, Column 2): Puzzle has 3, Solution has 3 — Correct

Position (Row 2, Column 1): Puzzle has 2, Solution has 2 — Correct

Position (Row 2, Column 2): Puzzle has 1, Solution has 1 — Correct

Position (Row 3, Column 4): Puzzle has 1, Solution has 1 — Correct

Position (Row 4, Column 1): Puzzle has 1, Solution has 1 — Correct

Position (Row 4, Column 3): Puzzle has 4, Solution has 4 — Correct

Position (Row 4, Column 4): Puzzle has 3, Solution has 3 — Correct

All initial numbers are maintained.

Checking the rows:

Row 1: Numbers are 4, 3, 1, 2 (No repetition) - Correct

Row 2: Numbers are 2, 1, 3, 1 (Number 1 repeats) - Invalid

Row 3: Numbers are 3, 4, 2, 1 (No repetition) - Correct

Row 4: Numbers are 1, 2, 4, 3 (No repetition) - Correct

Checking the columns:

Column 1: Numbers are 4, 2, 3, 1 (No repetition) - Correct

Column 2: Numbers are 3, 1, 4, 2 (No repetition) - Correct

Column 3: Numbers are 1, 3, 2, 4 (No repetition) - Correct

Column 4: Numbers are 2, 1, 1, 3 (Number 1 repeats) - Invalid

727

Checking the 2x2 quadrants:

Quadrant 1 (Rows 1-2, Columns 1-2): Numbers are 4, 3, 2, 1 (No repetition) - Correct

Quadrant 2 (Rows 1-2, Columns 3-4): Numbers are 1, 2, 3, 1 (Number 1 repeats) - Invalid

Quadrant 3 (Rows 3-4, Columns 1-2): Numbers are 3, 4, 1, 2 (No repetition) - Correct

Quadrant 4 (Rows 3-4, Columns 3-4): Numbers are 2, 1, 4, 3 (No repetition) - Correct

Conclusion: There are repetitions in Row 2, Column 4, and Quadrant 2, the solution violates the Sudoku rules.

Therefore, the response is: Incorrect

728