
An Empirical Analysis of Compute-Optimal Inference for Problem-Solving with Language Models

Anonymous Author(s)

Affiliation

Address

email

Abstract

1 The optimal training configurations of large language models (LLMs) with respect
2 to model sizes and compute budgets have been extensively studied. But how to
3 optimally configure LLMs during inference has not been explored in sufficient
4 depth. We study *compute-optimal inference*: designing models and inference
5 strategies that optimally trade off additional inference-time compute for improved
6 performance. As a first step towards understanding and designing compute-optimal
7 inference methods, we assessed the effectiveness and computational efficiency
8 of multiple inference strategies such as Greedy Search, Majority Voting, Best-of-
9 N, Weighted Voting, and their variants on two different Tree Search algorithms,
10 involving different model sizes (e.g., 7B and 34B) and computational budgets. We
11 found that a smaller language model with a novel tree search algorithm typically
12 achieves a Pareto-optimal trade-off. These results highlight the potential benefits of
13 deploying smaller models equipped with more sophisticated decoding algorithms
14 in end-devices to enhance problem-solving accuracy. For instance, we show that
15 the Llemma-7B model can achieve competitive accuracy to a Llemma-34B model
16 on MATH500 while using $2\times$ less FLOPs. Our findings could potentially apply to
17 any generation task with a well-defined measure of success.

18 1 Introduction

19 Scaling laws of neural networks [Hestness et al., 2017, Rosenfeld et al., 2019] have been established
20 across a range of domains, including language modeling [Kaplan et al., 2020, Hoffmann et al., 2022,
21 OpenAI, 2023], image modeling [Henighan et al., 2020, Yu et al., 2022, Peebles and Xie, 2023],
22 video modeling [Brooks et al., 2024], reward modeling [Gao et al., 2023], and board games [Jones,
23 2021]. These studies have demonstrated how model performance is influenced by both the size of the
24 model and the amount of training computation. However, there is limited knowledge on how varying
25 the compute during inference affects model performance after the model has been trained.

26 To improve the task performance of large language models (LLMs), inference techniques typically
27 involve additional computation in a *performance maximization* step at inference time [Nye et al.,
28 2021, Wei et al., 2022, Wang et al., 2022b, Yao et al., 2023, Chen et al., 2024b]. This cost must be
29 taken into account for *compute-optimal inference*. For example, a Monte Carlo Tree Search (MCTS)
30 method [Jones, 2021] may improve task performance, but potentially cost much more than simply
31 sampling solutions multiple times. Generally speaking, we need a comprehensive understanding of
32 how various inference-time methods (e.g., Best-of-N, majority voting) trade off between performance
33 and cost. To improve our understanding, this paper presents a thorough empirical evaluation with
34 careful analysis over various configurations of representative LLMs and inference algorithms.

35 Specifically, we explore how to select an optimal model size (e.g., 7B or 34B) for the policy model
36 and an effective inference strategy (e.g., Greedy Search, Majority Voting, Best-of-N, Weighted Voting,

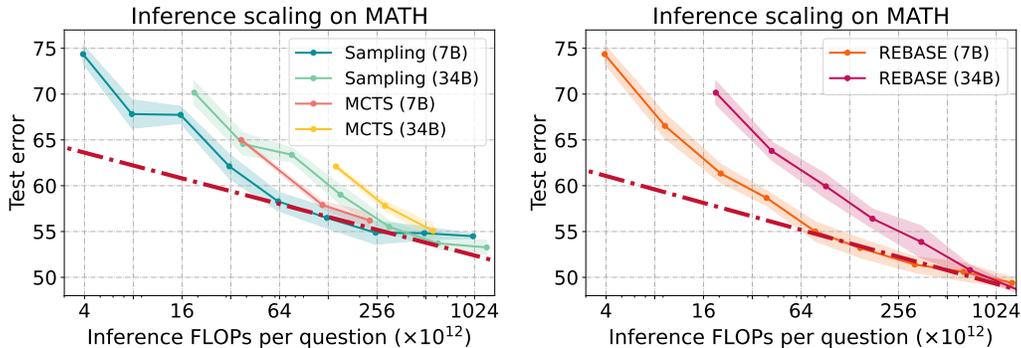


Figure 1: **The inference computation scaling laws** exhibited in error rate on the **MATH500** test set based on weighted majority voting, where the left figure shows sampling vs. MCTS, and the right figure shows our proposed REBASE. Clearly, the error rate decreases steadily when the computation increases, and REBASE exhibits a Pareto-optimal tradeoff during inference.

37 and their Tree Search variants) to maximize performance (i.e., accuracy) within a given compute
 38 budget. We manipulate the inference computation (FLOPs) of a fixed model by generating additional
 39 tokens through the policy model, sampling further candidate solutions, and ranking them with a
 40 reward model. We analyze the performance of a family of math-specialized LLMs (i.e., Llemma-7B
 41 and Llemma-34B [Azerbayev et al., 2023]) fine-tuned on the MetaMath dataset [Yu et al., 2023] and
 42 measure the error rate on the GSM8K test set [Cobbe et al., 2021a] and MATH500 test set [Hendrycks
 43 et al., 2021b, Lightman et al., 2023b].

44 Our analysis shows that voting-based methods generally outperform the strategy which selects the
 45 best solution (i.e., Best-of-N), and weighted voting has the most favorable results (Section 4.3,
 46 Figure 5 & 6). However, neither method shows a desirable behavior at high levels of compute. For
 47 instance, weighted voting saturates when sampling more than 128 solutions (Figure 1). We have also
 48 found that the commonly used MCTS method does not perform well with weighted voting, as it often
 49 yields many unfinished solutions, hence having less votes. To address this issue, we propose a novel
 50 tree search algorithm, *REward BALanced SEarch* (REBASE), which pairs well with weighted voting
 51 and improves the Pareto-optimal trade-off between accuracy and inference compute. The key idea of
 52 REBASE is to use a node-quality based reward to control the exploitation and pruning properties of
 53 tree search, while ensuring enough candidate solutions for voting or selection.

54 In our experiments, REBASE consistently outperforms sampling and MCTS methods across all
 55 settings, models, and tasks. Importantly, we find that REBASE with a *smaller* language model
 56 typically achieves a Pareto-optimal trade-off. For instance, we show that the Llemma-7B model can
 57 achieve competitive accuracy to a Llemma-34B model while using $2\times$ less FLOPs when evaluating
 58 on MATH500 (Figure 1) or GSM8K (Figure 4). These findings underscore the advantages of using
 59 smaller models with advanced inference-time algorithms on end-devices to improve problem-solving.

60 2 Related Works

61 **Mathematical Reasoning with LLMs.** Large language models have made significant progress
 62 in recent years, and have exhibited strong reasoning abilities [Brown et al., 2020, Hoffmann et al.,
 63 2022, Chowdhery et al., 2022, Lewkowycz et al., 2022]. Mathematical problem solving is a key task
 64 for measuring LLM reasoning abilities [Cobbe et al., 2021a, Hendrycks et al., 2021b]. [Ling et al.,
 65 2017] first developed the method of producing step by step solutions that lead to the final answer.
 66 Later, [Cobbe et al., 2021b] extended the work by finetuning the pre-trained language model on a
 67 large dataset to solve math word problems, a verifier is trained for evaluating solutions and ranking
 68 solutions. Nye et al. [2021] train models to use a scratchpad and improve their performance on
 69 algorithmic tasks. Wei et al. [2022] demonstrate that the reasoning ability of a language model can
 70 be elicited through the prompting. Subsequent research [Kojima et al., 2022, Lewkowycz et al., 2022,
 71 Zhou et al., 2022] in reasoning tasks has also highlighted the efficacy of rationale augmentation. We
 72 choose problem solving in mathematics as the task to study the compute-optimal strategy since it
 73 allows us to accurately evaluate the problem solving ability of LLMs.

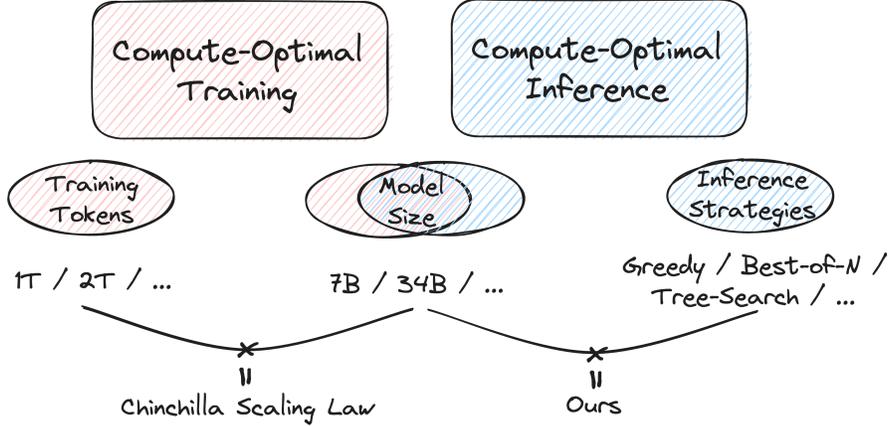


Figure 2: **Illustration of compute-optimal scaling laws in training and inference.** The Chinchilla scaling law shows how to choose a model size and number of training tokens under a training-compute budget, while ours shows how to choose a model size and an inference strategy under a inference-compute budget.

74 **Inference Strategies of LLM Problem Solving.** A variety of inference (also called decoding)
 75 strategies have been developed to generate sequences with a trained model. Deterministic methods
 76 such as greedy decoding and beam search [Teller, 2000, Graves, 2012] find highly probable sequences,
 77 often yielding high quality results but without diversity. Sampling algorithms (e.g., temperature
 78 sampling [Ackley et al., 1985]) can produce a diverse set of results which are then aggregated to
 79 achieve higher accuracy (e.g., via majority voting [Wang et al., 2022a]). Recent methods combine
 80 search algorithms with modern LLMs, including breadth-first or depth-first search [Yao et al., 2023],
 81 Monte-Carlo Tree Search (MCTS) [Zhang et al., 2023, Zhou et al., 2023, Liu et al., 2024, Choi et al.,
 82 2023], and Self-evaluation Guided Beam Search [Xie et al., 2023]. All of these methods show that
 83 using search at inference time can lead to performance gains in various tasks. However, the trade-off
 84 for the improved performance is the use of compute to perform the search. Analyzing the trade-off
 85 between compute budget and LLM inference performance remains understudied. In this paper, we
 86 systematically analyze the trade-off between compute budget and problem-solving performance, and
 87 propose a tree search method that is empirically Pareto-optimal.

88 **Process Reward Models.** Process reward models (PRMs) have emerged as a technique to improve
 89 the reasoning and problem-solving capabilities of LLMs. These models assign rewards to the
 90 intermediate steps of the LLM generated sequences. PRMs have been shown effective in selecting
 91 reasoning traces with a low error rate, and for providing rewards in reinforcement learning-style
 92 algorithms [Uesato et al., 2022, Polu and Sutskever, 2020, Gudibande et al., 2023]. Ma et al. [2023]
 93 applies the PRM to give rewards on the intermediate steps and guide the multi-step reasoning process.
 94 The PRM can be either trained on human labeled data [Lightman et al., 2023a] or model-labeled
 95 synthetic data [Wang et al., 2023]. In our work, we use the PRM as the reward model for selecting
 96 generated solutions, and for selecting which partial solutions to explore in tree search.

97 3 An Empirical Analysis of Compute-Optimal Inference for Problem-Solving

98 We explore the following question: *Given a fixed FLOPs budget, how should one select an optimal*
 99 *model size for the policy model, and an effective inference strategy to maximize performance (i.e.,*
 100 *accuracy)?* To address this, we represent the problem-solving error rate $E(N, T)$ as a function of the
 101 number of model parameters N and the number of generated tokens T . The computational budget C
 102 is a deterministic function $\text{FLOPs}(N, T)$, based on N and T . Our goal is to minimize E under the
 103 test-time compute constraint $\text{FLOPs}(N, T) = C$:

$$N_{opt}(C), T_{opt}(C) = \arg \min_{N, T \text{ s.t. } \text{FLOPs}(N, T) = C} E(N, T) \quad (1)$$

104 where $N_{opt}(C)$ and $T_{opt}(C)$ denote the optimal allocation of a computational budget C .

105 Here, the inference computation (FLOPs) for a fixed model can be modulated by generating more
106 tokens with the policy model, e.g., by sampling additional candidate solutions and subsequently
107 ranking them using a reward model. We primarily consider sampling and tree-search approaches
108 with reranking or majority voting as the means to consume more tokens, including Greedy Search,
109 Majority Voting, Best-of-N, Weighted Voting, and their variants on tree search methods.

110 3.1 Inference Strategies

111 3.1.1 Sampling

112 **Greedy Search.** This strategy generates tokens one at a time by selecting the highest probability token
113 at each step, without considering future steps. It is computationally efficient but often suboptimal in
114 terms of diversity.

115 **Best-of-n.** This strategy, also known as rejection sampling, samples multiple solutions and chooses
116 the one with the highest score given by the reward model.

117 **Majority Voting.** In this strategy, multiple model outputs are generated, and the final answer to the
118 problem is determined by the most frequently occurring answer in all the outputs.

119 **Weighted Majority Voting.** This strategy is a variant of majority voting in which the votes are
120 weighted based on the score given by the reward model.

121 3.1.2 Monte Carlo Tree Search (MCTS)

122 Monte Carlo Tree Search (MCTS) has proven effective in domains such as board games where
123 strategic decision-making is required [Silver et al., 2016, 2017, Jones, 2021]. Recent work has shown
124 that adapting MCTS to the context of LLMs can enhance the text generation process [Zhang et al.,
125 2023, Zhou et al., 2023, Liu et al., 2024, Choi et al., 2023, Chen et al., 2024a, Tian et al., 2024,
126 Chen et al., 2024a]. In this context, MCTS is often paired with a value model to score and guide the
127 exploration steps. For additional background, we provide a review of MCTS in Appendix B.

128 Recent work in MCTS or its variants (e.g., Tree of Thoughts [Yao et al., 2023]) mainly focus on
129 improving the performance (e.g., accuracy) on the studied tasks. However, generic comparisons of
130 MCTS with conventional methods like Best-of-n and Majority Voting in terms of computational
131 budget, measured in generated tokens or processing time, are either scarce or indicating inference-
132 time issues. For example, MCTS consumes substantially more resources, often requiring dozens of
133 times more generated tokens than simpler methods. Specifically, a significant portion of the paths
134 in the search tree are used to estimate and select nodes, and these paths do not necessarily become
135 a part of the final candidate solution, although MCTS ensures that the sampled solutions comprise
136 high-quality intermediate steps. In contrast, sampling methods generate multiple solutions in parallel
137 and independently, and all the generated sequences are included in the candidate solutions. However,
138 the intermediate steps in these sequences are not guaranteed to be of high quality, as there is no
139 mechanism for pruning poor steps or exploiting promising ones.

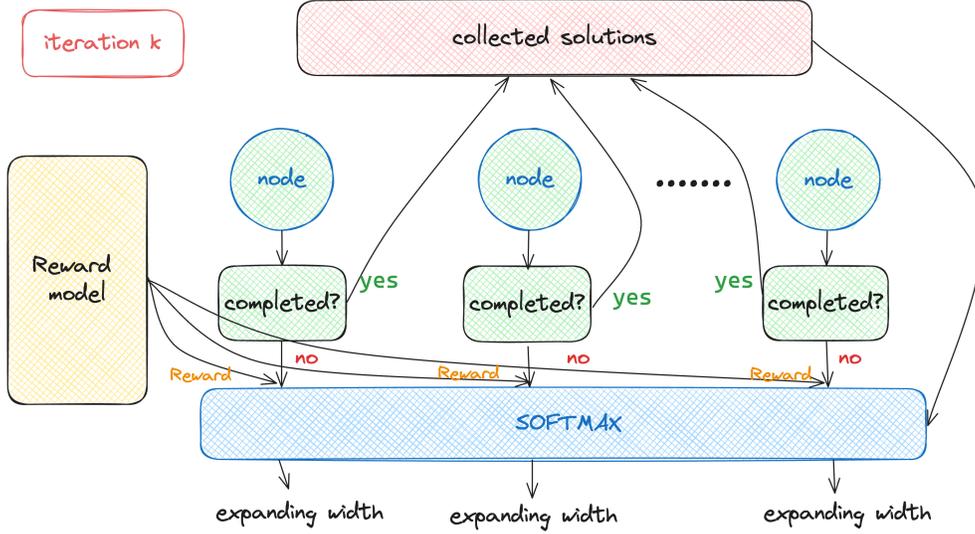
140 This highlights the need for developing a new tree search method that can achieve a comparable (or
141 better) performance as MCTS, and that is computationally less costly, just like weighted majority
142 voting and best-of-n. This need leads to the development of our new method named Reward Balanced
143 Search (REBASE), as introduced next.

144 3.1.3 Reward Balanced Search (REBASE)

145 The REBASE tree search method inherits the exploitation and pruning properties of tree search,
146 while using the reward model alone to estimate the nodes’ qualities without additional computation
147 for estimating values by sampling children. The efficiency is achieved by constraining the total
148 expansion width of the tree at a certain depth. REBASE balances the expansion width among the
149 nodes at the same depth based on the rewards given by the Process Reward Model (PRM). The details
150 are provided below:

151 **Notations.** We consider the fine-tuned LLM as a policy π_θ . Given a question q and the first k steps
152 of a solution x_1, \dots, x_k , the $(k + 1)$ -th step is produced by $\pi_\theta(x_{k+1}|q, x_1 \dots x_k)$. When generating
153 solutions using tree search, the root of the tree corresponds to the question q . The node corresponding

N = Computation budget (number of complete solutions) in tree-search
 C = Total number of collected solutions in the first k iterations.



Sum of expanding widths = Num of nodes in next iteration = $N - C$

Figure 3: Illustration of one iteration of REward BALanced SEarch (REBASE).

154 to x_{k+1} is the child of the node corresponding to x_k if it is sampled from $\pi_\theta(\cdot|q, x_1 \cdots, x_k)$. The
 155 reward of a node $n(x_k)$ is determined by the PRM as $R(n(x_k)) = R(q, x_1, \cdots, x_k)$.

156 **Initialization.** Given the question q , balance temperature T_b , and sampling number of solutions N ,
 157 we sample N instances of the first step for the question, yielding all the nodes of depth 1 in the search
 158 tree. We set the sampling budget of depth 0 $B_0 = N$ as initialization.

159 **Reward modeling and update.** In the i -th iteration, the PRM assigns the rewards to all the nodes
 160 at depth i . After that, the algorithm examines whether the solutions up to depth i are complete.
 161 Supposing there are C_i completed solutions, we update the sampling budget using $B_i \leftarrow B_{i-1} - C_i$.
 162 If $B_i = 0$, the process ends, and we obtain N solutions.

163 **Exploration balancing and expansion.** For all of the nodes n_j with reward $R(n_j)$ in the depth i of
 164 the tree, we calculate the expansion width of the n_j as:

$$W_j = \text{Round} \left(B_i \frac{\exp(R(n_j)/T_b)}{\sum_k \exp(R(n_k)/T_b)} \right). \quad (2)$$

165 Then we sample W_j children for n_j for all the nodes in depth i , and start the next iteration.

166 3.1.4 Theoretical Analysis

167 Before empirically studying the scaling effects of increasing the inference-time compute budget,
 168 we present two theorems which will help us understand the experimental results later. These two
 169 theorems give an upper bound on the performance of sampling when fixing the LLM generator.

170 We assume the vocabulary is limited and the sequence length is constrained, thus the number of
 171 possible solutions and answers are finite. The proofs are provided in the Appendix A.

172 **Theorem 1.** Given a test dataset \mathcal{D} and a LLM π . $|\mathcal{A}|$ is the finite set of all possible answers given
 173 by LLM, the ground truth function g maps test data d to the true answer. Denote the accuracy of the
 174 LLM on this dataset with majority over N samples as $ACC_{MV}(\pi, \mathcal{D}, N)$. The accuracy of majority
 175 voting on the LLM will eventually saturate, i.e.

$$\lim_{N \rightarrow \infty} ACC_{MV}(\pi, \mathcal{D}, N) = \frac{\sum_{d \in \mathcal{D}} \mathbb{I}((g(d) = \arg \max_{a \in \mathcal{A}} \pi(a|d))}{|\mathcal{D}|}. \quad (3)$$

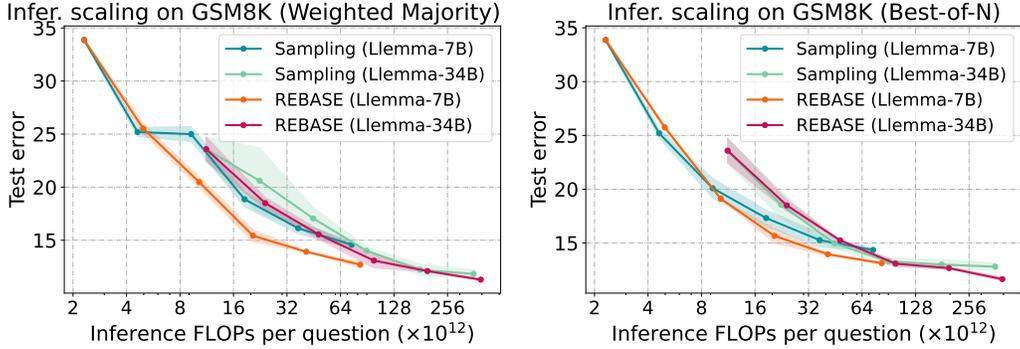


Figure 4: **The inference computation scaling comparisons across different model sizes.** The left/right panel shows the GSM8K problem-solving error rate on GSM8K based on Weighted Majority/Best-of-N.

176 where $\pi(x|d)$ denotes the probability that the LLM answers x given input d and \mathbb{I} is the indicator
 177 function.

178 **Theorem 2.** Assume the reward model assigns an expected reward of $R(a)$ to $a \in \mathcal{A}$ among the
 179 different solutions generated by LLM that yields a . Given a test dataset \mathcal{D} and a LLM π , $|\mathcal{A}|$ is the
 180 finite set of all possible answers given by LLM, the ground truth function g maps test data d to the
 181 true answer. Denote the accuracy of the LLM on this dataset with weighted majority over N samples
 182 as $ACC_{WV}(\pi, \mathcal{D}, N, R)$. The accuracy of weighted majority voting on the LLM will eventually
 183 saturate, i.e.

$$\lim_{N \rightarrow \infty} ACC_{WV}(\pi, \mathcal{D}, N, R) = \frac{\sum_{d \in \mathcal{D}} \mathbb{I}((g(d) = \arg \max_{a \in \mathcal{A}} R(a)\pi(a|d)))}{|\mathcal{D}|}. \quad (4)$$

184 where $\pi(x|d)$ denotes the probability that the LLM answers x given input d and \mathbb{I} denotes the
 185 indicator function.

186 Theorem 2 shows that as long as the reward model assigns higher rewards than the policy for correct an-
 187 swers versus other answers in expectation, the upper bound of Weighted Majority Voting will be higher
 188 than Majority Voting since $\mathbb{I}((g(d) = \arg \max_{a \in \mathcal{A}} R(a)\pi(a|d))) > \mathbb{I}((g(d) = \arg \max_{a \in \mathcal{A}} \pi(a|d)))$.
 189 We put the figures comparing BoN and Weighted Majority Voting in the main paper and leave the
 190 Majority Voting figures in Appendix D since Majority Voting is dominated by Weighted Majority
 191 Voting.

192 4 Experiments

193 4.1 Setup

194 **Datasets.** We conduct experiments on two mathematical problem-solving datasets to investigate
 195 the scaling effects of computation and our REBASE method for both challenging and simpler
 196 problems. Specifically, MATH [Hendrycks et al., 2021a] and GSM8K[Cobbe et al., 2021b] are
 197 datasets containing high school mathematics competition-level problems and grade-school level
 198 mathematical reasoning problems, respectively. Following [Lightman et al., 2023b, Wang et al., 2024,
 199 Sun et al., 2024], we use the MATH500 subset as our test set.

200 **Generators.** We use Llemma-7B and Llemma-34B [Azerbaiyev et al., 2024] as our base models and
 201 finetune them on the MetaMath dataset [Yu et al., 2024] using full parameter supervised fine-tuning
 202 (Full-SFT). The detailed finetuning configuration is given in the Appendix. Additionally, we test the
 203 Mistral-7B model to expand our findings across different models.

204 **Reward Model.** All of the experiments use the same Llemma-34B reward model, which we
 205 finetuned on the synthetic process reward modeling dataset, Math-Shepherd [Wang et al., 2024]. We
 206 added a reward head to make the model, enabling it to output a scalar reward at the end of each step.

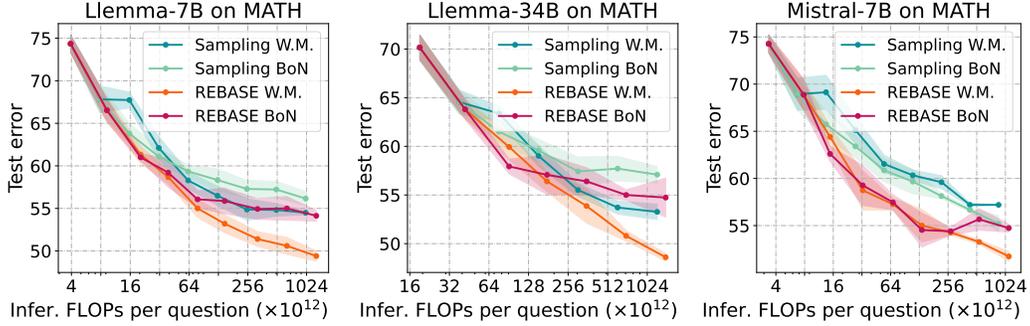


Figure 5: **The inference computation scaling laws** of different models for the problem-solving error rate on **MATH500** test set. The tested models are Llemma-7B (left), Llemma-34B (middle), & Mistral-7B (right). In the legend, W.M. and BoN refer to Weighted Majority and Best-of-N, respectively.

207 **Inference Configuration.** For the MATH dataset, we sample 1, 2, 4, 8, 16, 32, 64, 128, and 256
 208 solutions for the 7B models, and 1 to 64 solutions for the 34B Llemma model. For the GSM8K dataset,
 209 we sample 1 to 32 solutions, as it is relatively easier. We use sampling and REBASE to generate
 210 these samples and select the answer through Best-of-N, Majority Voting, and Weighted Voting.
 211 Each configuration is run multiple times to calculate the mean and variance, thereby mitigating the
 212 randomness and ensuring the reliability of our conclusions.

213 4.2 Main Results of Compute-Optimal Inference

214 In order to compare the compute budgets of 7B and 34B models, we plot the figures with the number
 215 of FLOPs used per question during inference. We compute the inference FLOPs based on the standard
 216 formula from [Kaplan et al., 2020].

217 **Llemma-7B model achieves competitive accuracy to Llemma-34B model with lower compute**
 218 **budget.** Figures 1 and 4 show the curves of error rates versus total number of inference FLOPs per
 219 question. Inference methods with different model sizes are plotted in the same diagram. We found
 220 that Llemma-7B costs approximately 2x less total FLOPs than Llemma-34B under the same method
 221 (Sampling, MCTS, REBASE) and task (MATH, GSM8K) while achieving competitive accuracy.
 222 This result suggests that, with the same training dataset and model family, training and inference with
 223 a smaller model could be more favorable in terms of compute budget if multiple sampling or search
 224 methods are employed.

225 **All inference configurations will saturate eventually.** This is expected as Theorem 1 and Theorem
 226 2 show. Also illustrated in Figures 5 and 6, the slope of the error rate curves start large, then decreases
 227 and the curves finally become nearly flat as the number of samples scales, showing the effect of
 228 saturation.

229 **Scaling law of compute-optimal inference.** The findings in our experiments are consistent with
 230 the Theorem 1 and 2, After a threshold the accuracy of sampling more solutions saturate, we should
 231 scale the model size. We interpolate the smoothed test error rate curve in Figure 1 and Figure 4,
 232 and fit power laws to estimate the optimal model size N and number of generated tokens T for any
 233 given amount of compute. We obtained a relationship $N_{opt} \propto C^a$ and $T_{opt} \propto C^b$, where $a = 1.0$
 234 and $b = 0.0$ for both sampling-based weighted voting and our tree-search method REBASE. Our
 235 fitted curves indicate that the optimal inference strategy is invariant to the amount of compute (e.g.,
 236 re-ranking with 32 sampled solutions or REBASE tree search with a compute budget of 64 for
 237 MATH), and the optimal model size grows linearly with the increased compute budget.

238 4.3 Comparing REBASE to Other Baselines

239 **REBASE is Pareto-optimal.** While MCTS underperforms Sampling (Fig. 1), from Fig. 1, 4, 5,
 240 and 6, we found that REBASE consistently outperforms the Sampling method in all settings, when

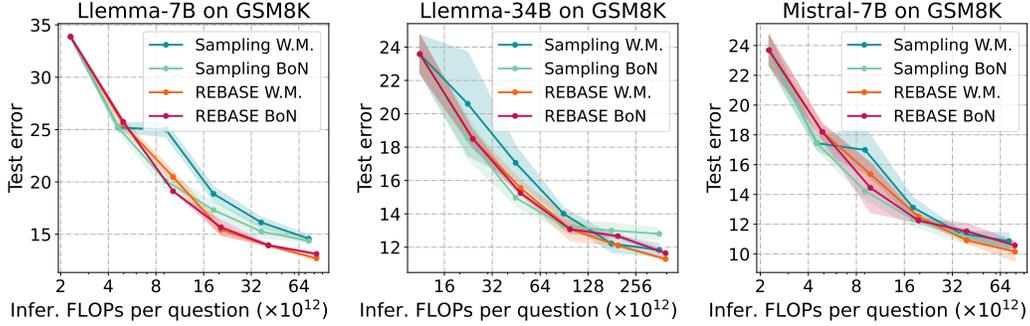


Figure 6: **The inference computation scaling laws** of different models for the problem-solving error rate on **GSM8K** test set. The tested models are Llemma-7B (left), Llemma-34B (middle), & Mistral-7B (right). In the legend, W.M. and BoN refer to Weighted Majority and Best-of-N, respectively.

Table 1: REBASE with lower compute budget has competitive accuracy against Sampling with higher compute budget. We use weighted voting to aggregate the candidate solutions in both Sampling and REBASE.

	# SAMPLES	FLOPS	MATH500
MISTRAL-7B			
SAMPLING	256	8.7×10^{14}	42.8
REBASE	32	1.36×10^{14}	45.0
LLEMMA-7B			
SAMPLING	256	10.0×10^{14}	45.5
REBASE	32	1.48×10^{14}	46.8
LLEMMA-34B			
SAMPLING	64	12.1×10^{14}	46.7
REBASE	32	7.08×10^{14}	49.2

241 fixing the model and the evaluation task. Table 1 shows that REBASE can achieve competitive
 242 accuracy with even a lower compute budget than the sampling method. This finding is novel, and
 243 differs from previous tree search works which typically improve the performance at the cost of higher
 244 computational expense compared to sampling [Chen et al., 2024a, Xie et al., 2023]. Table 2 shows
 245 that given the same compute budget (sampling 32 solutions for the 7B model and 8 solutions for 34B
 246 model), using REBASE yields higher accuracy than sampling.

247 **Weaker models gain more from Tree Search.** From Fig. 2, we saw that compared with sampling,
 248 Mistral-7B, Llemma-7B, Llemma-34B increase 5.3%, 3.3%, 2.6% in MATH and 0.7%, 1.9%, 0.9%
 249 in GSM8K. The order of accuracy increase is inversely related to the model’s corresponding greedy
 250 search on those datasets. This suggests that weaker models, as indicated by their lower greedy search
 251 accuracy, benefit more from tree search methods like REBASE.

252 **REBASE saturates later than sampling with higher accuracy.** From Figure 5 and Figure 6, we
 253 observe that both sampling and REBASE saturate early in GSM8K and relatively late in MATH,
 254 which we attribute to the difference of the difficulty level. This can be explained through the LLM
 255 may assign high probability to the true answer in easy problems than those of harder problem, as
 256 suggested by Theorem 1 and 2 with their proofs A. On MATH (Figure 5), we see that REBASE
 257 finally saturates with a higher accuracy than sampling. We hypothesize the reason is that REBASE
 258 samples the truth answer with higher probability than sampling. And as demonstrated by Theorem 1
 259 and 2, the upper bound becomes higher.

Table 2: Accuracy of different inference configurations under a specific compute budget. MV, BoN and WV denote Majority Voting, Best-of-N and Weighted Voting, respectively.

	# SAMPLES	MATH FLOPS	GSM8K FLOPS	MATH500	GSM8K
MISTRAL-7B					
GREEDY	1	3.4×10^{12}	2.3×10^{12}	28.6	77.9
SAMPLING + MV	32	109.2×10^{12}	72.6×10^{12}	36.1	85.7
SAMPLING + BoN	32	109.2×10^{12}	72.6×10^{12}	40.3	89.4
SAMPLING + WV	32	109.2×10^{12}	72.6×10^{12}	39.7	89.1
REBASE + MV	32	136.2×10^{12}	78.9×10^{12}	44.1	88.8
REBASE + BoN	32	136.2×10^{12}	78.9×10^{12}	45.4	89.4
REBASE + WV	32	136.2×10^{12}	78.9×10^{12}	45.0	89.8
LLEMMA-7B					
GREEDY	1	3.92×10^{12}	2.3×10^{12}	30.0	68.5
SAMPLING + MV	32	125.4×10^{12}	73.9×10^{12}	41.0	80.0
SAMPLING + BoN	32	125.4×10^{12}	73.9×10^{12}	41.7	85.6
SAMPLING + WV	32	125.4×10^{12}	73.9×10^{12}	43.5	85.4
REBASE + MV	32	148.0×10^{12}	82.6×10^{12}	46.1	86.1
REBASE + BoN	32	148.0×10^{12}	82.6×10^{12}	44.1	86.9
REBASE + WV	32	148.0×10^{12}	82.6×10^{12}	46.8	87.3
LLEMMA-34B					
GREEDY	1	19.0×10^{12}	11.2×10^{12}	33.0	78.4
SAMPLING + MV	8	152.3×10^{12}	89.7×10^{12}	39.9	84.3
SAMPLING + BoN	8	152.3×10^{12}	89.7×10^{12}	40.4	86.7
SAMPLING + WV	8	152.3×10^{12}	89.7×10^{12}	41.0	86.0
REBASE + MV	8	176.8×10^{12}	98.7×10^{12}	43.9	86.1
REBASE + BoN	8	176.8×10^{12}	98.7×10^{12}	43.6	86.9
REBASE + WV	8	176.8×10^{12}	98.7×10^{12}	42.9	86.9

260 5 Conclusion & Limitations

261 In this work, we have conducted a comprehensive empirical analysis of compute-optimal inference
 262 for problem-solving with language models. Our study has revealed several key findings. First, with
 263 an optimal inference configuration, a small language model can achieve competitive accuracy to a
 264 $4\times$ larger model while using approximately $2\times$ less total FLOPs under the same inference method
 265 (Sampling, MCTS, REBASE) and task (MATH, GSM8K), suggesting that training and inference
 266 with smaller models could be more favorable in terms of compute budget when combined with
 267 multiple sampling or search strategies. Second, our new REBASE tree-search method consistently
 268 outperforms sampling (and MCTS) across all settings, models, and tasks, achieving competitive
 269 accuracy with even lower compute budget compared to sampling. Our findings highlight the potential
 270 of deploying smaller models equipped with more sophisticated inference strategies like REBASE to
 271 enhance problem-solving accuracy while maintaining computational efficiency.

272 **Limitations** First, our experiments focused specifically on mathematical problem-solving tasks,
 273 and the generalizability of our findings to other domains remains to be explored. Second, we only
 274 investigated a limited range of model scales, primarily focusing on 7B and 34B models. Future
 275 research could extend our analysis to a wider range of model sizes to gain a more comprehensive
 276 understanding of the scaling laws for compute-optimal inference. Finally, our experiments mainly
 277 utilized the MetaMath dataset for training the models. It would be valuable to explore the impact of
 278 different training datasets on the performance and efficiency of compute-optimal inference strategies
 279 for mathematical problem-solving.

280 References

- 281 David H Ackley, Geoffrey E Hinton, and Terrence J Sejnowski. A learning algorithm for boltzmann
282 machines. *Cognitive science*, 9(1):147–169, 1985.
- 283 Zhangir Azerbayev, Hailey Schoelkopf, Keiran Paster, Marco Dos Santos, Stephen McAleer, Albert Q
284 Jiang, Jia Deng, Stella Biderman, and Sean Welleck. Llemma: An open language model for
285 mathematics. *arXiv preprint arXiv:2310.10631*, 2023.
- 286 Zhangir Azerbayev, Hailey Schoelkopf, Keiran Paster, Marco Dos Santos, Stephen McAleer, Albert Q.
287 Jiang, Jia Deng, Stella Biderman, and Sean Welleck. Llemma: An open language model for
288 mathematics, 2024.
- 289 Tim Brooks, Bill Peebles, Connor Holmes, Will DePue, Yufei Guo, Li Jing, David Schnurr,
290 Joe Taylor, Troy Luhman, Eric Luhman, Clarence Ng, Ricky Wang, and Aditya Ramesh.
291 Video generation models as world simulators. 2024. URL [https://openai.com/research/
292 video-generation-models-as-world-simulators](https://openai.com/research/video-generation-models-as-world-simulators).
- 293 Tom Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared D. Kaplan, Prafulla Dhariwal,
294 Arvind Neelakantan, Pranav Shyam, Girish Sastry, Amanda Askell, et al. Language models are
295 few-shot learners. *Advances in Neural Information Processing Systems*, 33:1877–1901, 2020.
- 296 Guoxin Chen, Minpeng Liao, Chengxi Li, and Kai Fan. Alphamath almost zero: process supervision
297 without process, 2024a.
- 298 Ziru Chen, Michael White, Raymond Mooney, Ali Payani, Yu Su, and Huan Sun. When is tree search
299 useful for llm planning? it depends on the discriminator. *arXiv preprint arXiv:2402.10890*, 2024b.
- 300 Sehyun Choi, Tianqing Fang, Zhaowei Wang, and Yangqiu Song. Kcts: Knowledge-constrained tree
301 search decoding with token-level hallucination detection, 2023.
- 302 Aakanksha Chowdhery, Sharan Narang, Jacob Devlin, Maarten Bosma, Gaurav Mishra, Adam
303 Roberts, Paul Barham, Hyung Won Chung, Charles Sutton, Sebastian Gehrmann, et al. PaLM:
304 Scaling language modeling with pathways. *arXiv preprint arXiv:2204.02311*, 2022.
- 305 Karl Cobbe, Vineet Kosaraju, Mohammad Bavarian, Mark Chen, Heewoo Jun, Lukasz Kaiser,
306 Matthias Plappert, Jerry Tworek, Jacob Hilton, Reiichiro Nakano, Christopher Hesse, and John
307 Schulman. Training verifiers to solve math word problems. *arXiv preprint arXiv:2110.14168*,
308 2021a.
- 309 Karl Cobbe, Vineet Kosaraju, Mohammad Bavarian, Mark Chen, Heewoo Jun, Lukasz Kaiser,
310 Matthias Plappert, Jerry Tworek, Jacob Hilton, Reiichiro Nakano, et al. Training verifiers to solve
311 math word problems. *arXiv preprint arXiv:2110.14168*, 2021b.
- 312 Leo Gao, John Schulman, and Jacob Hilton. Scaling laws for reward model overoptimization. In
313 *International Conference on Machine Learning*, pages 10835–10866. PMLR, 2023.
- 314 Alex Graves. Sequence transduction with recurrent neural networks, 2012.
- 315 Arnav Gudibande, Eric Wallace, Charlie Snell, Xinyang Geng, Hao Liu, Pieter Abbeel, Sergey Levine,
316 and Dawn Song. The false promise of imitating proprietary llms. *arXiv preprint arXiv:2305.15717*,
317 2023.
- 318 Dan Hendrycks, Steven Basart, Saurav Kadavath, Mantas Mazeika, Akul Arora, Ethan Guo, Collin
319 Burns, Samir Puranik, Horace He, Dawn Song, et al. Measuring coding challenge competence
320 with apps. *arXiv preprint arXiv:2105.09938*, 2021a.
- 321 Dan Hendrycks, Collin Burns, Saurav Kadavath, Akul Arora, Steven Basart, Eric Tang, Dawn
322 Song, and Jacob Steinhardt. Measuring mathematical problem solving with the math dataset. In
323 *Thirty-fifth Conference on Neural Information Processing Systems Datasets and Benchmarks Track*
324 *(Round 2)*, 2021b.
- 325 Tom Henighan, Jared Kaplan, Mor Katz, Mark Chen, Christopher Hesse, Jacob Jackson, Heewoo Jun,
326 Tom B. Brown, Prafulla Dhariwal, Scott Gray, et al. Scaling laws for autoregressive generative
327 modeling. *arXiv preprint arXiv:2010.14701*, 2020.

- 328 Joel Hestness, Sharan Narang, Newsha Ardalani, Gregory Diamos, Heewoo Jun, Hassan Kianinejad,
329 Md Patwary, Mostofa Ali, Yang Yang, and Yanqi Zhou. Deep learning scaling is predictable,
330 empirically. *arXiv preprint arXiv:1712.00409*, 2017.
- 331 Jordan Hoffmann, Sebastian Borgeaud, Arthur Mensch, Elena Buchatskaya, Trevor Cai, Eliza
332 Rutherford, Diego de Las Casas, Lisa Anne Hendricks, Johannes Welbl, Aidan Clark, et al.
333 Training compute-optimal large language models. *arXiv preprint arXiv:2203.15556*, 2022.
- 334 Andy L Jones. Scaling scaling laws with board games. *arXiv preprint arXiv:2104.03113*, 2021.
- 335 Jared Kaplan, Sam McCandlish, Tom Henighan, Tom B. Brown, Benjamin Chess, Rewon Child,
336 Scott Gray, Alec Radford, Jeffrey Wu, and Dario Amodei. Scaling laws for neural language models.
337 *arXiv preprint arXiv:2001.08361*, 2020.
- 338 Takeshi Kojima, Shixiang Shane Gu, Machel Reid, Yutaka Matsuo, and Yusuke Iwasawa. Large
339 language models are zero-shot reasoners. *arXiv preprint arXiv:2205.11916*, 2022.
- 340 Aitor Lewkowycz, Anders Andreassen, David Dohan, Ethan Dyer, Henryk Michalewski, Vinay Ra-
341 masesh, Ambrose Slone, Cem Anil, Imanol Schlag, Theo Gutman-Solo, et al. Solving quantitative
342 reasoning problems with language models. *arXiv preprint arXiv:2206.14858*, 2022.
- 343 Hunter Lightman, Vineet Kosaraju, Yura Burda, Harri Edwards, Bowen Baker, Teddy Lee, Jan
344 Leike, John Schulman, Ilya Sutskever, and Karl Cobbe. Let’s verify step by step. *arXiv preprint*
345 *arXiv:2305.20050*, 2023a.
- 346 Hunter Lightman, Vineet Kosaraju, Yura Burda, Harri Edwards, Bowen Baker, Teddy Lee, Jan Leike,
347 John Schulman, Ilya Sutskever, and Karl Cobbe. Let’s verify step by step, 2023b.
- 348 Wang Ling, Dani Yogatama, Chris Dyer, and Phil Blunsom. Program induction by rationale gen-
349 eration: Learning to solve and explain algebraic word problems. In *Proceedings of the 55th*
350 *Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages
351 158–167, 2017.
- 352 Jiacheng Liu, Andrew Cohen, Ramakanth Pasunuru, Yejin Choi, Hannaneh Hajishirzi, and Asli
353 Celikyilmaz. Don’t throw away your value model! generating more preferable text with value-
354 guided monte-carlo tree search decoding, 2024.
- 355 Qianli Ma, Haotian Zhou, Tingkai Liu, Jianbo Yuan, Pengfei Liu, Yang You, and Hongxia Yang.
356 Let’s reward step by step: Step-level reward model as the navigators for reasoning, 2023.
- 357 Maxwell Nye, Anders Johan Andreassen, Guy Gur-Ari, Henryk Michalewski, Jacob Austin, David
358 Bieber, David Dohan, Aitor Lewkowycz, Maarten Bosma, David Luan, et al. Show your work:
359 Scratchpads for intermediate computation with language models. *arXiv preprint arXiv:2112.00114*,
360 2021.
- 361 OpenAI. Gpt-4 technical report, 2023.
- 362 William Peebles and Saining Xie. Scalable diffusion models with transformers. In *Proceedings of*
363 *the IEEE/CVF International Conference on Computer Vision*, pages 4195–4205, 2023.
- 364 Stanislas Polu and Ilya Sutskever. Generative language modeling for automated theorem proving.
365 *arXiv preprint arXiv:2009.03393*, 2020.
- 366 Jonathan S Rosenfeld, Amir Rosenfeld, Yonatan Belinkov, and Nir Shavit. A constructive prediction
367 of the generalization error across scales. *arXiv preprint arXiv:1909.12673*, 2019.
- 368 David Silver, Aja Huang, Chris J Maddison, Arthur Guez, Laurent Sifre, George Van Den Driessche,
369 Julian Schrittwieser, Ioannis Antonoglou, Veda Panneershelvam, Marc Lanctot, et al. Mastering
370 the game of Go with deep neural networks and tree search. *Nature*, 529(7587):484–489, 2016.
- 371 David Silver, Julian Schrittwieser, Karen Simonyan, Ioannis Antonoglou, Aja Huang, Arthur Guez,
372 Thomas Hubert, Lucas Baker, Matthew Lai, Adrian Bolton, et al. Mastering the game of go without
373 human knowledge. *nature*, 550(7676):354–359, 2017.

- 374 Zhiqing Sun, Longhui Yu, Yikang Shen, Weiyang Liu, Yiming Yang, Sean Welleck, and Chuang
375 Gan. Easy-to-hard generalization: Scalable alignment beyond human supervision. *arXiv preprint*
376 *arXiv:2403.09472*, 2024.
- 377 Virginia Teller. Speech and language processing: an introduction to natural language processing,
378 computational linguistics, and speech recognition, 2000.
- 379 Ye Tian, Baolin Peng, Linfeng Song, Lifeng Jin, Dian Yu, Haitao Mi, and Dong Yu. Toward self-
380 improvement of llms via imagination, searching, and criticizing. *arXiv preprint arXiv:2404.12253*,
381 2024.
- 382 Jonathan Uesato, Nate Kushman, Ramana Kumar, Francis Song, Noah Siegel, Lisa Wang, Antonia
383 Creswell, Geoffrey Irving, and Irina Higgins. Solving math word problems with process- and
384 outcome-based feedback. *arXiv preprint arXiv:2211.14275*, 2022.
- 385 Peiyi Wang, Lei Li, Zhihong Shao, RX Xu, Damai Dai, Yifei Li, Deli Chen, Y Wu, and Zhifang
386 Sui. Math-shepherd: Verify and reinforce llms step-by-step without human annotations. *CoRR*,
387 *abs/2312.08935*, 2023.
- 388 Peiyi Wang, Lei Li, Zhihong Shao, R. X. Xu, Damai Dai, Yifei Li, Deli Chen, Y. Wu, and Zhifang
389 Sui. Math-shepherd: Verify and reinforce llms step-by-step without human annotations, 2024.
- 390 Xuezhi Wang, Jason Wei, Dale Schuurmans, Quoc Le, Ed Chi, Sharan Narang, Aakanksha Chowdh-
391 ery, and Denny Zhou. Self-consistency improves chain of thought reasoning in language models.
392 *International Conference on Learning Representations (ICLR 2023)*, 2022a.
- 393 Yizhong Wang, Yeganeh Kordi, Swaroop Mishra, Alisa Liu, Noah A Smith, Daniel Khashabi, and
394 Hannaneh Hajishirzi. Self-instruct: Aligning language model with self generated instructions.
395 *arXiv preprint arXiv:2212.10560*, 2022b.
- 396 Jason Wei, Xuezhi Wang, Dale Schuurmans, Maarten Bosma, Ed Chi, Quoc Le, and Denny Zhou.
397 Chain-of-thought prompting elicits reasoning in large language models. *NeurIPS*, 2022.
- 398 Yuxi Xie, Kenji Kawaguchi, Yiran Zhao, Xu Zhao, Min-Yen Kan, Junxian He, and Qizhe Xie.
399 Self-evaluation guided beam search for reasoning, 2023.
- 400 Shunyu Yao, Dian Yu, Jeffrey Zhao, Izhak Shafran, Thomas L Griffiths, Yuan Cao, and Karthik
401 Narasimhan. Tree of thoughts: Deliberate problem solving with large language models. *arXiv*
402 *preprint arXiv:2305.10601*, 2023.
- 403 Jiahui Yu, Yuanzhong Xu, Jing Yu Koh, Thang Luong, Gunjan Baid, Zirui Wang, Vijay Vasudevan,
404 Alexander Ku, Yinfei Yang, Burcu Karagol Ayan, et al. Scaling autoregressive models for content-
405 rich text-to-image generation. *arXiv preprint arXiv:2206.10789*, 2(3):5, 2022.
- 406 Longhui Yu, Weisen Jiang, Han Shi, Jincheng Yu, Zhengying Liu, Yu Zhang, James T Kwok, Zhenguo
407 Li, Adrian Weller, and Weiyang Liu. Metamath: Bootstrap your own mathematical questions for
408 large language models. *arXiv preprint arXiv:2309.12284*, 2023.
- 409 Longhui Yu, Weisen Jiang, Han Shi, Jincheng Yu, Zhengying Liu, Yu Zhang, James T. Kwok,
410 Zhenguo Li, Adrian Weller, and Weiyang Liu. Metamath: Bootstrap your own mathematical
411 questions for large language models, 2024.
- 412 Shun Zhang, Zhenfang Chen, Yikang Shen, Mingyu Ding, Joshua B. Tenenbaum, and Chuang Gan.
413 Planning with large language models for code generation, 2023.
- 414 Andy Zhou, Kai Yan, Michal Shlapentokh-Rothman, Haohan Wang, and Yu-Xiong Wang. Language
415 agent tree search unifies reasoning acting and planning in language models, 2023.
- 416 Yongchao Zhou, Andrei Ioan Muresanu, Ziwen Han, Keiran Paster, Silviu Pitit, Harris Chan,
417 and Jimmy Ba. Large language models are human-level prompt engineers. *arXiv preprint*
418 *arXiv:2211.01910*, 2022.

419 A Proofs of Theorem 1 and 2

420 A.1 Proof of Theorem 1

421 *Proof.* Suppose the possible answers of the LLM are $x_1, x_2, x_3, \dots, x_{|\mathcal{A}|}$, with $\pi(x_1|d) >$
 422 $\pi(x_2|d) > \dots > \pi(x_{|\mathcal{A}|}|d)$. After sampling N solutions from the LLM, we denote the occurrence of
 423 x_i as $f(x_i)$, the probability that x_1 is not the most frequent output is

$$P(f(x_1) \neq \arg \max_x f(x)) \quad (5)$$

424 With Union bound, we get

$$P(x_1 \neq \arg \max_x f(x)) \quad (6)$$

$$\leq \sum_{i=2}^{|\mathcal{A}|} P(f(x_1) \leq f(x_i)) \quad (7)$$

$$\leq |\mathcal{A}| P(f(x_1) \leq f(x_2)) \quad (8)$$

$$= |\mathcal{A}| (1 - P(f(x_1) \geq f(x_2))) \quad (9)$$

$$\leq |\mathcal{A}| \left(1 - P \left(f(x_1) \geq \frac{\pi(x_1|d) + \pi(x_2|d)}{2} N \right) P \left(f(x_2) \leq \frac{\pi(x_1|d) + \pi(x_2|d)}{2} N \right) \right) \quad (10)$$

$$\leq |\mathcal{A}| \left(1 - \left(1 - e^{-\frac{\delta_1^2}{2} \pi(x_1|d) N} \right) \left(1 - e^{-\frac{\delta_2^2}{2} \pi(x_2|d) N} \right) \right) \quad (11)$$

$$\leq |\mathcal{A}| C^N \quad \text{for some } C < 1. \quad (12)$$

425 Where (11) is by Chernoff Bound, $\delta_1 = \frac{\pi(x_1|d) - \pi(x_2|d)}{2\pi(x_1|d)}$ and $\delta_2 = \frac{\pi(x_1|d) - \pi(x_2|d)}{2\pi(x_2|d)}$. As $N \rightarrow \infty$, we
 426 have

$$f(x) = \begin{cases} M(x|N) = 1 & \text{if } x = \arg \max_{a \in \mathcal{A}} \pi(a|d) \\ M(x|N) = 0 & \text{otherwise.} \end{cases} \quad (13)$$

427 Where $M(x|N)$ denotes the probability that majority voting over N sampled solutions gives x . The
 428 proof of original theorem is automatically completed by (13). \square

429 A.2 Proof of Theorem 2

430 *Proof.* The proof is similar to the Theorem 1, We rank $x_1, x_2, \dots, x_{|\mathcal{A}|}$ with $R(x_1)f(x_1) > \dots >$
 431 $R(x_{|\mathcal{A}|})f(x_{|\mathcal{A}|})$. Denotes $w(x_i)$ as the the total weights of answer x_i after sampling N solutions. As
 432 $N \rightarrow \infty$, $w(x_i) \rightarrow R(x_i)f(x_i)$. Same as proof in theorem 1, we have

$$P(x_1 \neq \arg \max_x f(x)) \quad (14)$$

$$\leq |\mathcal{A}| P(w(x_1) \leq w(x_2)) \quad (15)$$

$$= |\mathcal{A}| (1 - P(w(x_1) \geq w(x_2))) \quad (16)$$

$$\leq |\mathcal{A}| \left(1 - P \left(w(x_1) \geq \frac{v(x_1) + v(x_2)}{2} N \right) P \left(w(x_2) \leq \frac{v(x_1) + v(x_2)}{2} N \right) \right). \quad (17)$$

433 Where $v(x) = R(x)\pi(x|d)$, the remaining proof completely follows Theorem 1. \square

434 B MCTS Details

435 The MCTS process can be represented as the following steps:

436 **Selection** The process begins at the root node. Here, the algorithm recursively selects the child
 437 node that offers the highest Upper Confidence Bound applied to Trees (UCT) value, continuing until
 438 a node is reached that has not been expanded. The UCT is calculated using the formula

$$UCT(s) = Q(s) + C \sqrt{\frac{\ln(N(\text{Parent}(s)))}{N(s)}}. \quad (18)$$

Table 3: Fine-tuning Hyper-parameters: LR refers to the learning rate, BS refers to the batch size. Llemma-7B and Llemma-34B are the generators we use in our experiments, RM is short for Reward Model.

Model	# Epoch	Dataset	BS	LR	Max Seq Length	Dtype
Llemma-7B	1	MetaMath	128	8E-6	1024	FP32
Llemma-34B	1	MetaMath	128	8E-6	768	FP32
Llemma-34B RM	2	Math-Shepherd	128	1E-5	768	BF16

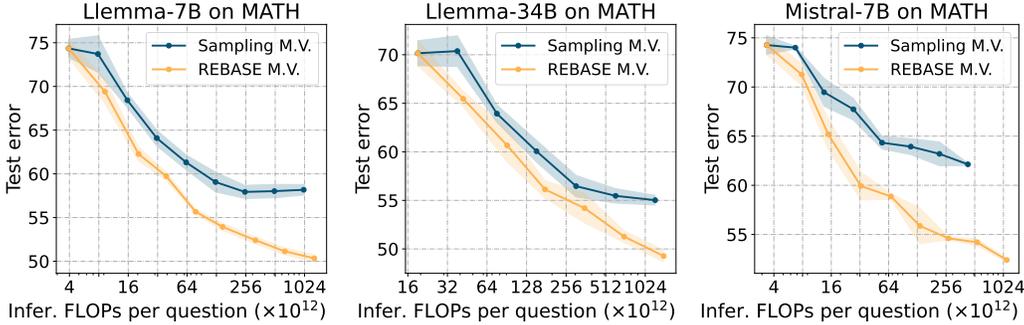


Figure 7: **The inference computation scaling laws** of different models for the problem-solving error rate on **MATH** test set. The tested models are Llemma-7B (left), Llemma-34B (middle), & Mistral-7B (right). In the legend, M.V. refer to Majority Voting.

439 Where $Q(s)$ represents the quality score of node s , $N(s)$ is the number of visits to node s , and C is a
 440 constant determining the level of exploration.

441 **Expansion and evaluation** Upon reaching a non-terminal node s , the node is expanded by gener-
 442 ating multiple child nodes. Each child node c is then evaluated using a value function $V(c)$, which
 443 predicts the potential quality of continuing the sequence from node c .

444 **Backpropagation** After evaluation, the algorithm updates the UCT values and the visit counts for
 445 all nodes along the path from the selected node back to the root. For any node n in this path, the
 446 updates are made as follows:

$$N(n) \leftarrow N(n) + 1,$$

$$Q(n) \leftarrow \frac{(N(n) - 1)Q(n) + V(s)}{N(n)}.$$

447 C Hyper-parameters

448 **Finetuning** We put all the hyperparameters of fine-tuned models in the table 3. We preprocess the
 449 MetaMath Dataset to make the solutions in a stepwise format.

450 **Inference** For all the inference strategies, the temperature of the LLM is set to 1.0. Max tokens
 451 for the output is 1024 and max tokens for one step is 256. For REBASE, we chose the balance
 452 temperature (the softmax temperature in the REBASE algorithm) as $T_b = 0.1$. For MCTS, we set
 453 C in the UCT value as 1 and we expand 4, 8, 16 children for the root, 2 children for other selected
 454 nodes with total 32, 64, 128 expansions respectively.

455 D Supplementary Figures

456 In the main part of paper, there isn't enough space for showing the scaling effects of Majority Voting,
 457 we append the figures about Majority Voting and Majority Voting v.s. Weighted Majority Voting

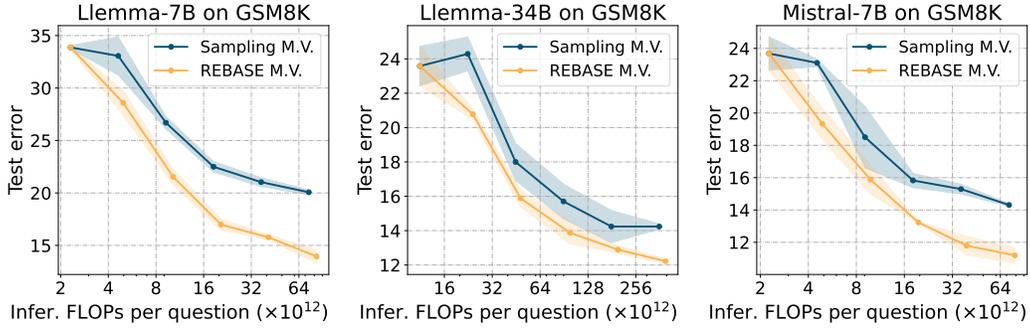


Figure 8: **The inference computation scaling laws** of different models for the problem-solving error rate on **GSM8K** test set. The tested models are Llemma-7B (left), Llemma-34B (middle), & Mistral-7B (right). In the legend, M.V. refer to Majority Voting.

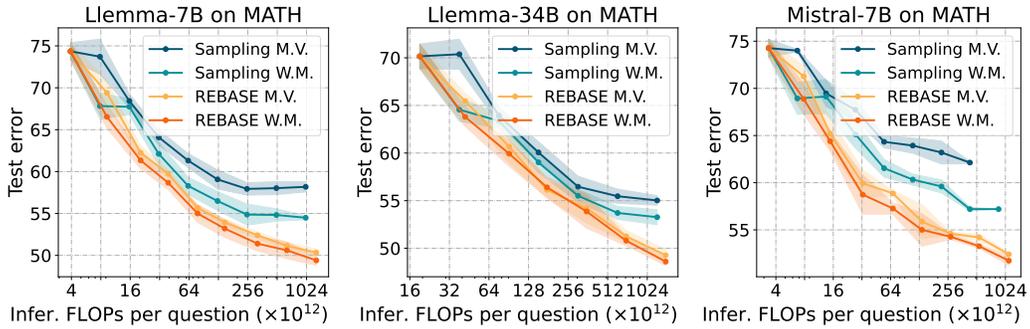


Figure 9: **The inference computation scaling laws** of different models for the problem-solving error rate on **MATH** test set. The tested models are Llemma-7B (left), Llemma-34B (middle), & Mistral-7B (right). In the legend, M.V. and W.M. refer to Majority Voting and Weighted Majority, respectively.

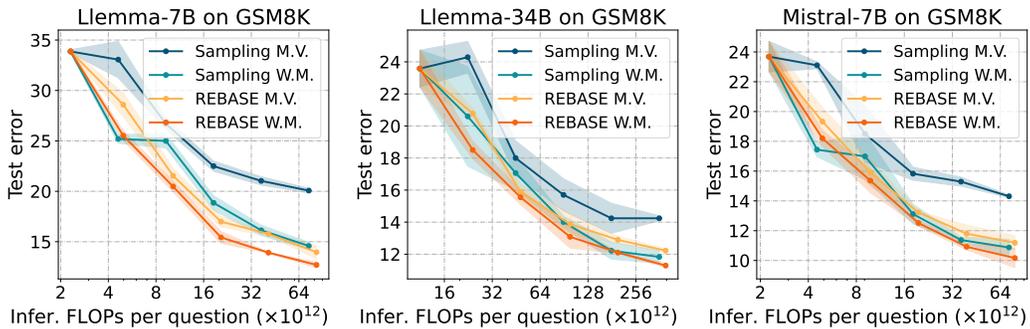


Figure 10: **The inference computation scaling laws** of different models for the problem-solving error rate on **GSM8K** test set. The tested models are Llemma-7B (left), Llemma-34B (middle), & Mistral-7B (right). In the legend, M.V. and W.M. refer to Majority Voting and Weighted Majority, respectively.

458 (Fig. 7, 8, 9, 10) in this appendix. The experiments show that although the gap between Majority
 459 Voting and Weighted Majority Voting on sampling is huge. This gap becomes much smaller if we
 460 apply REBASE. This phenomenon can be caused by the selection ability of tree search like REBASE.
 461 Once REBASE already samples solutions with high rewards, conducting weighted majority voting
 462 gains less since the sampled solutions may all have relatively high and stable rewards compared with
 463 those of sampling.

464 **NeurIPS Paper Checklist**

465 **1. Claims**

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

468 Answer: [\[Yes\]](#)

469 Justification: In Abstract and Introduction, we claim that we investigate the compute-optimal
470 inference: designing models and inference strategies that optimally trade off additional
471 inference-time compute for improved performance.

472 Guidelines:

- 473 • The answer NA means that the abstract and introduction do not include the claims
474 made in the paper.
- 475 • The abstract and/or introduction should clearly state the claims made, including the
476 contributions made in the paper and important assumptions and limitations. A No or
477 NA answer to this question will not be perceived well by the reviewers.
- 478 • The claims made should match theoretical and experimental results, and reflect how
479 much the results can be expected to generalize to other settings.
- 480 • It is fine to include aspirational goals as motivation as long as it is clear that these goals
481 are not attained by the paper.

482 **2. Limitations**

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

484 Answer: [\[Yes\]](#)

485 Justification: The discussion is in the last section of the main paper.

486 Guidelines:

- 487 • The answer NA means that the paper has no limitation while the answer No means that
488 the paper has limitations, but those are not discussed in the paper.
- 489 • The authors are encouraged to create a separate "Limitations" section in their paper.
- 490 • The paper should point out any strong assumptions and how robust the results are to
491 violations of these assumptions (e.g., independence assumptions, noiseless settings,
492 model well-specification, asymptotic approximations only holding locally). The authors
493 should reflect on how these assumptions might be violated in practice and what the
494 implications would be.
- 495 • The authors should reflect on the scope of the claims made, e.g., if the approach was
496 only tested on a few datasets or with a few runs. In general, empirical results often
497 depend on implicit assumptions, which should be articulated.
- 498 • The authors should reflect on the factors that influence the performance of the approach.
499 For example, a facial recognition algorithm may perform poorly when image resolution
500 is low or images are taken in low lighting. Or a speech-to-text system might not be
501 used reliably to provide closed captions for online lectures because it fails to handle
502 technical jargon.
- 503 • The authors should discuss the computational efficiency of the proposed algorithms
504 and how they scale with dataset size.
- 505 • If applicable, the authors should discuss possible limitations of their approach to
506 address problems of privacy and fairness.
- 507 • While the authors might fear that complete honesty about limitations might be used by
508 reviewers as grounds for rejection, a worse outcome might be that reviewers discover
509 limitations that aren't acknowledged in the paper. The authors should use their best
510 judgment and recognize that individual actions in favor of transparency play an impor-
511 tant role in developing norms that preserve the integrity of the community. Reviewers
512 will be specifically instructed to not penalize honesty concerning limitations.

513 **3. Theory Assumptions and Proofs**

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

516
517
518
519
520
521
522
523
524
525
526
527
528
529
530
531
532
533
534
535
536
537
538
539
540
541
542
543
544
545
546
547
548
549
550
551
552
553
554
555
556
557
558
559
560
561
562
563
564
565
566
567
568

Answer: [Yes]

Justification: It's in Appendix.

Guidelines:

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

4. Experimental Result Reproducibility

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

Answer: [Yes]

Justification: We introduce our method in Section 3 and the hyperparameters are introduced in the Appendix.

Guidelines:

- The answer NA means that the paper does not include experiments.
- If the paper includes experiments, a No answer to this question will not be perceived well by the reviewers: Making the paper reproducible is important, regardless of whether the code and data are provided or not.
- If the contribution is a dataset and/or model, the authors should describe the steps taken to make their results reproducible or verifiable.
- Depending on the contribution, reproducibility can be accomplished in various ways. For example, if the contribution is a novel architecture, describing the architecture fully might suffice, or if the contribution is a specific model and empirical evaluation, it may be necessary to either make it possible for others to replicate the model with the same dataset, or provide access to the model. In general, releasing code and data is often one good way to accomplish this, but reproducibility can also be provided via detailed instructions for how to replicate the results, access to a hosted model (e.g., in the case of a large language model), releasing of a model checkpoint, or other means that are appropriate to the research performed.
- While NeurIPS does not require releasing code, the conference does require all submissions to provide some reasonable avenue for reproducibility, which may depend on the nature of the contribution. For example
 - (a) If the contribution is primarily a new algorithm, the paper should make it clear how to reproduce that algorithm.
 - (b) If the contribution is primarily a new model architecture, the paper should describe the architecture clearly and fully.
 - (c) If the contribution is a new model (e.g., a large language model), then there should either be a way to access this model for reproducing the results or a way to reproduce the model (e.g., with an open-source dataset or instructions for how to construct the dataset).
 - (d) We recognize that reproducibility may be tricky in some cases, in which case authors are welcome to describe the particular way they provide for reproducibility. In the case of closed-source models, it may be that access to the model is limited in some way (e.g., to registered users), but it should be possible for other researchers to have some path to reproducing or verifying the results.

5. Open access to data and code

569 Question: Does the paper provide open access to the data and code, with sufficient instruc-
570 tions to faithfully reproduce the main experimental results, as described in supplemental
571 material?

572 Answer: [Yes]

573 Justification: We only used open-source models in this work. The code will be released.

574 Guidelines:

- 575 • The answer NA means that paper does not include experiments requiring code.
- 576 • Please see the NeurIPS code and data submission guidelines ([https://nips.cc/
577 public/guides/CodeSubmissionPolicy](https://nips.cc/public/guides/CodeSubmissionPolicy)) for more details.
- 578 • While we encourage the release of code and data, we understand that this might not be
579 possible, so “No” is an acceptable answer. Papers cannot be rejected simply for not
580 including code, unless this is central to the contribution (e.g., for a new open-source
581 benchmark).
- 582 • The instructions should contain the exact command and environment needed to run to
583 reproduce the results. See the NeurIPS code and data submission guidelines ([https:
584 //nips.cc/public/guides/CodeSubmissionPolicy](https://nips.cc/public/guides/CodeSubmissionPolicy)) for more details.
- 585 • The authors should provide instructions on data access and preparation, including how
586 to access the raw data, preprocessed data, intermediate data, and generated data, etc.
- 587 • The authors should provide scripts to reproduce all experimental results for the new
588 proposed method and baselines. If only a subset of experiments are reproducible, they
589 should state which ones are omitted from the script and why.
- 590 • At submission time, to preserve anonymity, the authors should release anonymized
591 versions (if applicable).
- 592 • Providing as much information as possible in supplemental material (appended to the
593 paper) is recommended, but including URLs to data and code is permitted.

594 6. Experimental Setting/Details

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

598 Answer: [Yes]

599 Justification: We used the standard training and test splits or MATH and GSM8K and report
600 the hyperparameters in the appendix.

601 Guidelines:

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

607 7. Experiment Statistical Significance

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

610 Answer: [Yes]

611 Justification: The error bars are included in our figures.

612 Guidelines:

- 613 • The answer NA means that the paper does not include experiments.
- 614 • The authors should answer "Yes" if the results are accompanied by error bars, confi-
615 dence intervals, or statistical significance tests, at least for the experiments that support
616 the main claims of the paper.
- 617 • The factors of variability that the error bars are capturing should be clearly stated (for
618 example, train/test split, initialization, random drawing of some parameter, or overall
619 run with given experimental conditions).

- 620 • The method for calculating the error bars should be explained (closed form formula,
621 call to a library function, bootstrap, etc.)
- 622 • The assumptions made should be given (e.g., Normally distributed errors).
- 623 • It should be clear whether the error bar is the standard deviation or the standard error
624 of the mean.
- 625 • It is OK to report 1-sigma error bars, but one should state it. The authors should
626 preferably report a 2-sigma error bar than state that they have a 96% CI, if the hypothesis
627 of Normality of errors is not verified.
- 628 • For asymmetric distributions, the authors should be careful not to show in tables or
629 figures symmetric error bars that would yield results that are out of range (e.g. negative
630 error rates).
- 631 • If error bars are reported in tables or plots, The authors should explain in the text how
632 they were calculated and reference the corresponding figures or tables in the text.

633 8. Experiments Compute Resources

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

637 Answer: [Yes]

638 Justification: All the experiments are conducted on $8 \times$ H100 GPUs.

639 Guidelines:

- 640 • The answer NA means that the paper does not include experiments.
- 641 • The paper should indicate the type of compute workers CPU or GPU, internal cluster,
642 or cloud provider, including relevant memory and storage.
- 643 • The paper should provide the amount of compute required for each of the individual
644 experimental runs as well as estimate the total compute.
- 645 • The paper should disclose whether the full research project required more compute
646 than the experiments reported in the paper (e.g., preliminary or failed experiments that
647 didn't make it into the paper).

648 9. Code Of Ethics

649 Question: Does the research conducted in the paper conform, in every respect, with the
650 NeurIPS Code of Ethics <https://neurips.cc/public/EthicsGuidelines?>

651 Answer: [Yes]

652 Justification: Yes, we conform with the NeurIPS Code of Ethics.

653 Guidelines:

- 654 • The answer NA means that the authors have not reviewed the NeurIPS Code of Ethics.
- 655 • If the authors answer No, they should explain the special circumstances that require a
656 deviation from the Code of Ethics.
- 657 • The authors should make sure to preserve anonymity (e.g., if there is a special consid-
658 eration due to laws or regulations in their jurisdiction).

659 10. Broader Impacts

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

662 Answer: [NA]

663 Justification: We do not find significant positive societal impacts and negative societal
664 impacts of our work.

665 Guidelines:

- 666 • The answer NA means that there is no societal impact of the work performed.
- 667 • If the authors answer NA or No, they should explain why their work has no societal
668 impact or why the paper does not address societal impact.

- 669
- 670
- 671
- 672
- 673
- 674
- 675
- 676
- 677
- 678
- 679
- 680
- 681
- 682
- 683
- 684
- 685
- 686
- 687
- Examples of negative societal impacts include potential malicious or unintended uses (e.g., disinformation, generating fake profiles, surveillance), fairness considerations (e.g., deployment of technologies that could make decisions that unfairly impact specific groups), privacy considerations, and security considerations.
 - The conference expects that many papers will be foundational research and not tied to particular applications, let alone deployments. However, if there is a direct path to any negative applications, the authors should point it out. For example, it is legitimate to point out that an improvement in the quality of generative models could be used to generate deepfakes for disinformation. On the other hand, it is not needed to point out that a generic algorithm for optimizing neural networks could enable people to train models that generate Deepfakes faster.
 - The authors should consider possible harms that could arise when the technology is being used as intended and functioning correctly, harms that could arise when the technology is being used as intended but gives incorrect results, and harms following from (intentional or unintentional) misuse of the technology.
 - If there are negative societal impacts, the authors could also discuss possible mitigation strategies (e.g., gated release of models, providing defenses in addition to attacks, mechanisms for monitoring misuse, mechanisms to monitor how a system learns from feedback over time, improving the efficiency and accessibility of ML).

688 11. Safeguards

689 Question: Does the paper describe safeguards that have been put in place for responsible
690 release of data or models that have a high risk for misuse (e.g., pretrained language models,
691 image generators, or scraped datasets)?

692 Answer: [NA]

693 Justification:

694 Guidelines:

- 695
- 696
- 697
- 698
- 699
- 700
- 701
- 702
- 703
- 704
- The answer NA means that the paper poses no such risks.
 - Released models that have a high risk for misuse or dual-use should be released with necessary safeguards to allow for controlled use of the model, for example by requiring that users adhere to usage guidelines or restrictions to access the model or implementing safety filters.
 - Datasets that have been scraped from the Internet could pose safety risks. The authors should describe how they avoided releasing unsafe images.
 - We recognize that providing effective safeguards is challenging, and many papers do not require this, but we encourage authors to take this into account and make a best faith effort.

705 12. Licenses for existing assets

706 Question: Are the creators or original owners of assets (e.g., code, data, models), used in
707 the paper, properly credited and are the license and terms of use explicitly mentioned and
708 properly respected?

709 Answer: [Yes]

710 Justification: We use the proper citations.

711 Guidelines:

- 712
- 713
- 714
- 715
- 716
- 717
- 718
- 719
- 720
- 721
- 722
- The answer NA means that the paper does not use existing assets.
 - The authors should cite the original paper that produced the code package or dataset.
 - The authors should state which version of the asset is used and, if possible, include a URL.
 - The name of the license (e.g., CC-BY 4.0) should be included for each asset.
 - For scraped data from a particular source (e.g., website), the copyright and terms of service of that source should be provided.
 - If assets are released, the license, copyright information, and terms of use in the package should be provided. For popular datasets, paperswithcode.com/datasets has curated licenses for some datasets. Their licensing guide can help determine the license of a dataset.

- 723
- 724
- 725
- 726
- For existing datasets that are re-packaged, both the original license and the license of the derived asset (if it has changed) should be provided.
 - If this information is not available online, the authors are encouraged to reach out to the asset’s creators.

727 **13. New Assets**

728 Question: Are new assets introduced in the paper well documented and is the documentation
729 provided alongside the assets?

730 Answer: [Yes]

731 Justification: We use the proper citations.

732 Guidelines:

- 733
- 734
- 735
- 736
- 737
- 738
- 739
- 740
- The answer NA means that the paper does not release new assets.
 - Researchers should communicate the details of the dataset/code/model as part of their submissions via structured templates. This includes details about training, license, limitations, etc.
 - The paper should discuss whether and how consent was obtained from people whose asset is used.
 - At submission time, remember to anonymize your assets (if applicable). You can either create an anonymized URL or include an anonymized zip file.

741 **14. Crowdsourcing and Research with Human Subjects**

742 Question: For crowdsourcing experiments and research with human subjects, does the paper
743 include the full text of instructions given to participants and screenshots, if applicable, as
744 well as details about compensation (if any)?

745 Answer: [NA]

746 Justification: No crowdsourcing experiments are used.

747 Guidelines:

- 748
- 749
- 750
- 751
- 752
- 753
- 754
- 755
- The answer NA means that the paper does not involve crowdsourcing nor research with human subjects.
 - Including this information in the supplemental material is fine, but if the main contribution of the paper involves human subjects, then as much detail as possible should be included in the main paper.
 - According to the NeurIPS Code of Ethics, workers involved in data collection, curation, or other labor should be paid at least the minimum wage in the country of the data collector.

756 **15. Institutional Review Board (IRB) Approvals or Equivalent for Research with Human
757 Subjects**

758 Question: Does the paper describe potential risks incurred by study participants, whether
759 such risks were disclosed to the subjects, and whether Institutional Review Board (IRB)
760 approvals (or an equivalent approval/review based on the requirements of your country or
761 institution) were obtained?

762 Answer: [NA]

763 Justification: No human subjects are involved.

764 Guidelines:

- 765
- 766
- 767
- 768
- 769
- 770
- 771
- 772
- 773
- 774
- The answer NA means that the paper does not involve crowdsourcing nor research with human subjects.
 - Depending on the country in which research is conducted, IRB approval (or equivalent) may be required for any human subjects research. If you obtained IRB approval, you should clearly state this in the paper.
 - We recognize that the procedures for this may vary significantly between institutions and locations, and we expect authors to adhere to the NeurIPS Code of Ethics and the guidelines for their institution.
 - For initial submissions, do not include any information that would break anonymity (if applicable), such as the institution conducting the review.