

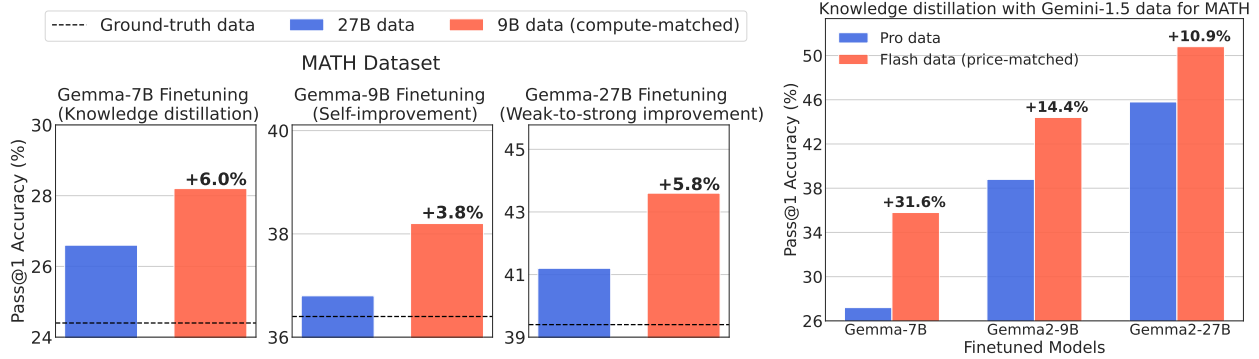
# Smaller, Weaker, Yet Better: Training LLM Reasoners via Compute-Optimal Sampling

Hritik Bansal<sup>1,2</sup>, Arian Hosseini<sup>1,3</sup>, Rishabh Agarwal<sup>1,3</sup>, Vinh Q. Tran<sup>1</sup> and Mehran Kazemi<sup>1</sup>

<sup>1</sup>Google DeepMind, <sup>2</sup>UCLA, <sup>3</sup>Mila

Training on high-quality synthetic data from strong language models (LMs) is a common strategy to improve the reasoning performance of LMs. In this work, we revisit whether this strategy is compute-optimal under a fixed inference budget (e.g., FLOPs). To do so, we investigate the trade-offs between generating synthetic data using a stronger but more expensive (SE) model versus a weaker but cheaper (WC) model. We evaluate the generated data across three key metrics: coverage, diversity, and false positive rate, and show that the data from WC models may have higher coverage and diversity, but also exhibit higher false positive rates. We then finetune LMs on data from SE and WC models in different settings: knowledge distillation, self-improvement, and a novel weak-to-strong improvement setup where a weaker LM teaches reasoning to a stronger LM. Our findings reveal that models finetuned on WC-generated data consistently outperform those trained on SE-generated data across multiple benchmarks and multiple choices of WC and SE models. These results challenge the prevailing practice of relying on SE models for synthetic data generation, suggesting that WC may be the compute-optimal approach for training advanced LM reasoners.

arXiv:2408.16737v1 [cs.CL] 29 Aug 2024



(a) Finetuning LMs with Gemma2 data.

(b) Finetuning LMs with Gemini 1.5 data.

**Figure 1 | Summary of the results.** (a) We finetune Gemma-7B, Gemma2-9B, and Gemma2-27B on the synthetic data collected from a stronger but more expensive LM (Gemma2-27B) and a weaker but cheaper LM (Gemma2-9B) in a compute-matched setup for the MATH dataset. We find that training with Gemma2-9B data is a more compute-optimal setting across diverse finetuning paradigms – knowledge distillation, self-improvement, and weak-to-strong improvement (i.e. using a weaker model to improve a stronger model). (b) We finetune Gemma models (7B/9B/27B) on synthetic data generated by the state-of-the-art LMs Gemini-1.5-Pro and Gemini-1.5-Flash in a price-matched setup. We find that finetuning with Flash-generated data consistently outperforms Pro-generated data.

## 1. Introduction

Language models (LMs) have demonstrated impressive capabilities in reasoning tasks, but their success heavily relies on being trained on vast amounts of (problem, solution) pairs. Collecting this data from humans is a costly and time-consuming process. Recent studies have demonstrated the

feasibility of synthetically generating this data using LMs themselves, offering a potentially more scalable and efficient approach to training data acquisition. One such widely-adopted approach is to sample multiple candidate solutions for a problem from an LM, filters them for final answer correctness, and finetune models on the correct solutions (Singh et al., 2023; Zelikman et al., 2022). Several works show that LMs trained with such synthetic solutions outperform those trained with human-written solutions (Pang et al., 2024; Singh et al., 2023; Yu et al., 2023; Yuan et al., 2023; Yue et al., 2023). Practitioners often sample solutions from strong LMs to ensure high quality (Mukherjee et al., 2023; Roziere et al., 2023; Teknium, 2023; Xu et al., 2023). However, sampling from strong LMs is expensive and resource-intensive, and limits the number of solutions that can be generated for practical sampling budgets.

In this paper, we explore an alternative sampling approach. Given a fixed compute budget, we investigate sampling from a **weaker but cheaper (WC)** model as opposed to the commonly-used approach of sampling from a **stronger but more expensive (SE)** model. We start by comparing data from WC vs SE across three axes that play crucial roles in the utility of such synthetic data: 1- *coverage*, the number of unique problems that are solved, 2- *diversity*, the average number of unique solutions we obtain per problem, and 3- *false positive rate (FPR)*, the percentage of problems that arrive at the correct final answer but with a wrong reasoning. We find that since we can generate more samples from the WC model compared to the SE model under a fixed budget, the data from WC may exhibit higher coverage and diversity. However, due to the lower quality of the WC model, it may also have a higher FPR. As a particular example for the Gemma2 family (Team et al., 2024a,b) on the MATH dataset (Hendrycks et al., 2021), Gemma2-9B achieves 11% higher coverage and 86% higher diversity, but also with 7% higher FPR compared to Gemma2-27B.

We then fine-tune models on data from SE and WC (see Figure 2) across diverse setups corresponding to three paradigms: 1) *knowledge distillation*, where a student LM learns from a teacher LM (Hinton et al., 2015); 2) *self-improvement*, where an LM learns from self-generated data (Huang et al., 2022); and 3) a new paradigm we introduce called *Weak-to-Strong Improvement*, where a strong student LM improves using synthetic data from a weaker teacher LM. Using two (WC, SE) model pairs, one from the Gemma2 family and another from the Gemini 1.5 family (Reid et al., 2024), we show on multiple benchmarks that training on WC-generated data consistently outperforms training on SE-generated data under the three setups, with relative gains of up to 31.6% percent (see Figure 1 for a summary of the results). Our results indicate that it is more compute-optimal to sample from a WC model as opposed to the common-practice of sampling from a SE model. With the performance gap between small and large LMs getting narrower over time (especially at larger scales), our results establish a solid foundation for training the next generation of LM reasoners.

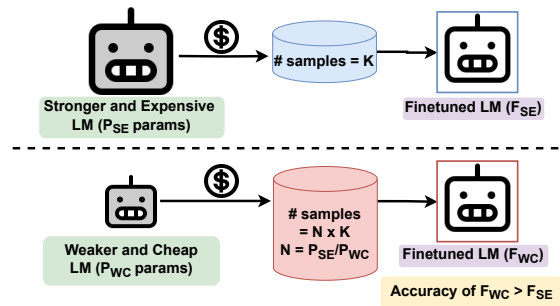


Figure 2 | **Illustration of the approach.** Given a fixed sampling budget, one can either generate fewer samples from a stronger but more expensive (SE) model or more samples from a weaker but cheaper (WC) model. The latter may lead to solving a wider range of problems and also more correct solutions per question. We compare the utility of these two synthetically generated datasets for training LM reasoners in various supervised finetuning setups and show that training with the data from WC consistently outperforms training on data from SE.

## 2. Preliminaries

Let  $\mathcal{D} = \{q_i, a_i\}_{i=1}^{i=n}$  be a training dataset of size  $n$  with reasoning questions  $q_i$  and final answers (aka labels)  $a_i$ . A successful approach to leverage such data to improve models for reasoning is as follows. We sample multiple solutions for each  $q_i$  at a non-zero temperature and create the synthetic data  $\mathcal{D}_G = \{q_i, \{(\hat{r}_{ij}, \hat{a}_{ij})_{j=1}^{j=k}\}\}$ , where  $k$  is the number of samples,  $\hat{r}_{ij}$  is the  $j$ -th reasoning chain (i.e. solution) generated by the model for  $q_i$ , and  $\hat{a}_{ij}$  is the model’s final answer for  $q_i$  in the  $j$ -th sample. Then, we filter the incorrect solutions by comparing  $\hat{a}_{ij}$  to  $a_i$  and removing the solutions whose final answer do not match that of the gold answer<sup>1</sup>. Finally, we supervise finetune a model on the remaining data  $\tilde{\mathcal{D}}_G$  to maximize  $J(\theta) = \mathbb{E}_{(q,r,a) \sim \tilde{\mathcal{D}}_G} [\log(p_\theta(r, a|q))]$ , i.e. the probability of generating the reasoning  $r$  and final answer  $a$  given the question  $q$ . This approach was first proposed in (Zelikman et al., 2022) and was then extended in multiple works including (Singh et al., 2023; Zelikman et al., 2024).

For a dataset  $\mathcal{D}_G$ , we compute *coverage@k* (aka *pass@k*) (Chen et al., 2021) as  $\mathbb{E}_{\mathcal{D}_G} [1 - \binom{M-c}{k} / \binom{M}{k}]$  where  $c$  is the number of solutions, out of  $M$ , with correct answers and  $\mathbb{E}_{\mathcal{D}_G} [\cdot]$  denotes the expectation over the problems and solutions in the generated dataset. Conceptually, *coverage@k* measures the fraction of *unique* questions that have at least one correct solution, assuming that we sample  $k$  solutions per question from the model. We also define *diversity@k* as the average number of unique correct solutions we obtain per question when we sample  $k$  solutions per question. Finally, we define *false positive rate (FPR)* as the percentage of solutions in  $\tilde{\mathcal{D}}_G$  where the reasoning is incorrect, despite the final answer being correct.

Different choices of the LM to sample solutions from and the LM to finetune lead to different setups. *Knowledge Distillation* (Hinton et al., 2015) corresponds to training a student LM on the synthetic data sampled from a stronger and larger LM. *Self-Improvement* (Huang et al., 2022) corresponds to training an LM on samples generated from itself.

## 3. Compute-Matched Sampling and Training

To generate a dataset  $\mathcal{D}_G$  with synthetic solutions from  $\mathcal{D}$ , one can leverage different models for generating solutions. Specifically, at a fixed sampling budget (FLOPs), one can generate more samples from a weaker but cheaper (WC) model or fewer samples from a stronger but more expensive (SE) model. Given a WC model with  $P_{WC}$  parameters and SE with  $P_{SE}$  parameters, we compute the sampling ratio at a fix budget for the two models, focusing on decoder-only transformer models (Vaswani, 2017). Following (Kaplan et al., 2020), we note that the FLOPs per inference token is  $2P$ , for a model with  $P$  parameters. As a result, the FLOPs for  $T$  inference tokens is  $2PT$ . Further, we assume that generating each solution requires an average of  $W$  inference tokens for both models<sup>2</sup>. Let  $S_{WC}$  and  $S_{SE}$  represent the number of samples we generate per question for the two models. The total cost of generating samples for the dataset  $\mathcal{D}$  will then be  $Cost_{WC} = n \times S_{WC} \times W \times (2P_{WC})$  and  $Cost_{SE} = n \times S_{SE} \times W \times (2P_{SE})$  for the cheap and expensive models, respectively. At a fixed sampling budget, we have:

$$n \times S_{WC} \times W \times (2P_{WC}) = n \times S_{SE} \times W \times (2P_{SE}) \quad \Rightarrow \quad S_{WC} = \frac{P_{SE}}{P_{WC}} S_{SE} \quad (1)$$

Equation 1 indicates that at a fixed sampling budget, for each question we can generate  $P_{SE}/P_{WC}$

<sup>1</sup>While it is possible to use more sophisticated approaches for filtering (e.g., process-based or outcome-based reward model (Uesato et al., 2022)), in this work we focus on final answer correctness for filtering as it has shown to be strong.

<sup>2</sup>This is a reasonable assumption given that the solution to a question is expected to be model-agnostic. We note, however, that it is possible for some questions that one model solves a question using a more optimal way compared to the other model thus producing a smaller solution.

Data (↓) / Finetuning setup (→)	Student-LM	WC-LM	SE-LM
<b>WC (Compute-matched)</b>	Knowledge distillation	Self-improvement	Weak-to-strong improvement
<b>SE</b>	Knowledge distillation	Knowledge distillation	Self-improvement

Table 1 | **Summary of the supervised finetuning setups.** We finetuned the language models under three setups: (a) Student LM, (b) Weak-Cheap (WC) LM, and (c) Strong-Expensive (SE) LM. For each setup, we employed different finetuning paradigms based on the source of the synthetic data. For example, training a separate student LM with data from both WC and SE models falls under the knowledge distillation paradigm. In contrast, training a WC model with its own samples is self-improvement. Finally, we also introduce a new paradigm, weak-to-strong improvement, where the samples from the WC model is used to improve the reasoning capabilities of the SE model at the fixed compute budget.

more samples from WC; the ratio scales linearly with the model parameters ratio<sup>3</sup>. Sampling more solutions from WC may increase the likelihood of correctly solving a larger subset of the problems (high coverage) and obtaining more correct solutions per question (high diversity).

Given a fixed budget, we can either generate fewer samples from a SE model or more samples from a WC model, and then finetune models for a fixed number of steps on the data from each of these models to measure and compare the utility of the data from each model. Specifically, we generate  $P_{SE}/P_{WC}$  more samples from the WC model compared to the SE model. We consider three finetuning setups that consists of diverse finetuning paradigms. The paradigms include the widely used knowledge distillation, the emerging framework of self-improvement, and a novel weak-to-strong improvement paradigm we introduce in this work. We define weak-to-strong improvement (W2S-I) as *enhancing* the reasoning capabilities of a strong model using samples generated from a weaker model. The three setups are as follows:

- **Student-LM finetuning:** Conventionally, the supervised finetuning data for training student LM is acquired from SE models to ensure high-quality (Teknium, 2023). However, we aim to understand whether WC models can replace SE models for distillation at the fixed sampling budget. To do so, we finetune a student LM separate from the WC and SE models on the WC and SE data, which corresponds to distillation in both the cases.
- **WC-LM finetuning:** Prior work (Singh et al., 2023) has shown that finetuning a WC model through self-generated data lags behind distillation from SE data. However, their setup spends a higher sampling budget (FLOPs) on collecting data from the SE model than collecting SI data from the WC model. In this work, we revisit this finetuning setup under the fixed sampling budget and finetune the WC model on the WC and SE data at a fixed budget for both. Note that training the WC model on its own data corresponds to self-improvement whereas training WC on the data from SE corresponds to distillation. Hence, this setup compares self-improvement on WC data with distillation from SE data.
- **SE-LM finetuning:** It is commonly believed that to improve a SE model, we either need synthetic data from the SE model itself or from an even stronger (and perhaps more expensive) model than SE. Here, we test an alternative approach to understand whether the synthetic data from the WC model can improve the SE model. To this end, we finetune the SE model on the WC and SE data. Training SE on data from WC corresponds to W2S-I and training SE on data from SE corresponds to self-improvement. Overall, this setup compares W2S-I by WC data with self-improvement by SE data.

A summary of the three setups and the finetuning paradigms that each case corresponds to can be found in Table 1.

<sup>3</sup>Note that this may also depend on the available hardware, which we ignore in this work.

## 4. Experimental Setup

**Datasets:** We utilize MATH (Hendrycks et al., 2021) and GSM-8K (Cobbe et al., 2021) as the reasoning datasets due to their wide adoption for mathematical problem solving. Specifically, MATH consists of competition level problems with various levels of difficulty (Level 1-5), and GSM-8K comprises of grade school level math problems. Each dataset contains 7500 math problems in their training split. We evaluate the models on 500 problems from the MATH test split (Lightman et al., 2023) and 1319 problems from the GSM-8K test split. Further, we use 500 problems from the MATH test split and 500 problems from GSM-8K as the validation dataset. We also use the Functional MATH dataset (Srivastava et al., 2024) for a transfer study. Further, we present the results for a coding task in Appendix A.

**Data Generation:** We use Gemma2 models for synthetic data generation, with pretrained Gemma2-9B and Gemma2-27B acting as the WC and SE models respectively. We generate the solutions for the problems in the MATH using a 4-shot prompt and for GSM-8K using an 8-shot prompt. Since the 9B model is roughly 3 times smaller than the 27B model, at a fixed sampling compute budget we can sample  $3\times$  more sample solutions per problem for Gemma2-9B. For our experiments, we consider two sampling budgets: a *low budget*, where we generate 1 and 3 candidate solutions per problem from Gemma2-27B and Gemma2-9B, respectively, and a *high budget*, where we generate 10 and 30 candidate solutions per problem. Further, we study the transfer of the reasoning capabilities for the models trained on MATH at the high sampling budget on the Functional MATH dataset.

**Model Finetuning:** We summarize the details for our finetuning setups in the Table 1. In the Student-LM finetuning setup, we finetune the Gemma-7B model (Team et al., 2024a) on data from Gemma2-9B (WC) and Gemma2-27B (SE). In addition, we use Gemma2-9B and Gemma2-27B for the WC-LM and SE-LM finetuning setups, respectively. Further, we train the LMs across different setups with the human-written solutions as a ground-truth baseline. We provide the finetuning details in Appendix F.

**Synthetic Data Evaluation:** To assess the quality of the synthetic data from the SE and WC models, we measure the false positive rate, as well as *coverage* and *diversity* at a fixed cost. From Equation 1, we know that sampling one solution from SE takes the same FLOPs as sampling  $P_{SE}/P_{WC}$  solutions from WC. Therefore, we compare *coverage@k* for SE to *coverage@( $\frac{P_{SE}}{P_{WC}}k$ )* for WC to allow a similar budget to both models. Specifically, we compare *coverage@k* and *coverage@3k* for our SE and WC models. Similarly we compare *diversity@k* and *diversity@3k* for our SE and WC models. Since FPR cannot be computed automatically, we compute it using two proxies: 1- a human evaluation on a subset of the data, where 50 solutions from each model were selected randomly and rated for reasoning correctness by the authors, and 2- automatic evaluation where we sampled 500 solutions and prompted Gemini-Pro-1.5 (Reid et al., 2024) to rate the correctness of the reasoning paths. To sample solutions, for the MATH dataset we selected uniformly from each diversity level. In our experiments, we find that the FPR estimates are close to each other for the human and automatic evaluation. We provide a few qualitative examples for the false positive instances in Appendix B.

**Evaluating Finetuned Models:** We use *pass@1* accuracy to evaluate the performance of the finetuned LMs. Specifically, we generate a single solution for the problem (zero-shot) from the test split, using a sampling temperature of 0.0 (greedy decoding) for the fine-tuned LM and measure the percentage of problems that where the final answer matches the golden final answer. We also report *maj@k* ( $k = 1, 4, 8, 16$ ) for part of our experiments, where we generate  $k$  solutions per problem at a sampling temperature of 0.7 and select the final answer that appears most among the  $k$  samples.

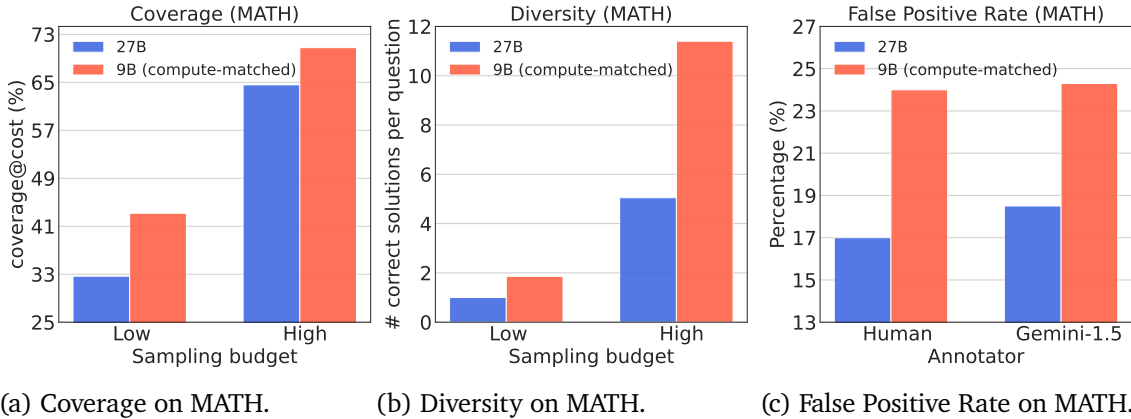


Figure 3 | **Synthetic data analysis for MATH dataset.** The (a) coverage, (b) diversity, and (c) false positive rates for Gemma2-27B and Gemma2-9B on the MATH dataset, at two sampling budgets.

## 5. Experiments and Results

We compare data from WC and SE models along several axes. First, we analyze the data along various quality metrics (§5.1). Subsequently, we present the supervised finetuning results for the different setups (§5.2). Finally, we perform ablation studies to study the impact of dataset size, sampling strategy, and the role of quality dimensions in the model performance (§5.3).

### 5.1. Synthetic Data Analysis

**Coverage:** Here, we aim to understand the pros and cons of generating solutions from the WC and SE models at a fixed sampling budget. We present the coverage, diversity, and FPR for the MATH at the low and high sampling budgets in Figure 3. The results for GSM-8K are presented in the Appendix – Figure 15. We find that in terms of coverage, the data from Gemma2-9B (WC) outperforms Gemma2-27B (SE) by 11% and 6% at the low and high sampling budgets, respectively, for the MATH dataset, and 8% and 1% for GSM-8K. This highlights that the higher number of samples for the WC model aids in solving more unique problems for both the reasoning datasets. We provide the coverage trends for diverse sampling budgets in Appendix C. In addition, we observe that the coverage of the WC model increases across various difficulty levels in the MATH dataset for the high sampling budget (see Appendix – Figure 16). This highlights that synthetic data from the WC model can solve more unique questions at various difficulty levels compare to the SE model, at a fixed sampling budget (Tong et al., 2024). Further, we provide a qualitative example that gets solved by repeated sampling from Gemma2-9B but remains unsolved by Gemma2-27B at the fixed high sampling budget (Table 5).

**Diversity:** The diversity for the data from Gemma2-9B is higher than Gemma2-27B by 86% and 125% at the low and high sampling budgets for the MATH dataset, and 134% and 158% at for the GSM-8K dataset. This implies that many unique reasoning chains in the synthetic data from the WC model lead to the correct solutions. We also observe that the absolute diversity scores are lower for MATH compared to GSM-8K at high sampling budget, indicating that models generate fewer correct solutions for the more challenging datasets when using repeated sampling.

**FPR:** Since we utilize the final answer correctness for filtering the synthetic data, it does not remove the solutions with incorrect intermediate reasoning steps. Our human evaluations suggest that the FPR for the WC-generated solutions is 7% and 2% higher than SE-generated solutions on the MATH and GSM-8K, respectively. The trends from the automatic evaluation are similar to that of human evaluation. Due to the differences in the difficulty of the problems, we note that the absolute

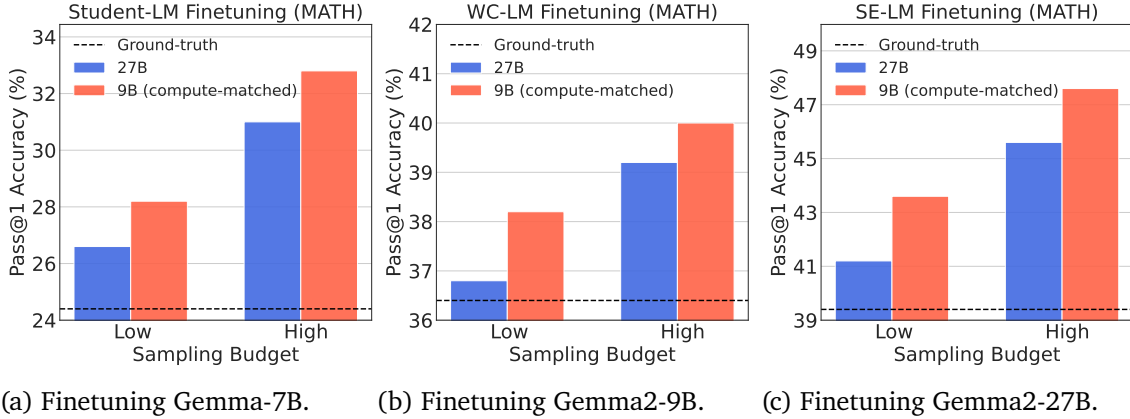


Figure 4 | **Supervised-finetuning results (MATH)**. The results for finetuning various LMs on the MATH synthetic data from the WC (Gemma2-9B) and SE (Gemma2-27B) models, at a fixed sampling budget. We observe that training with the samples from the WC model consistently outperforms training with SE data.

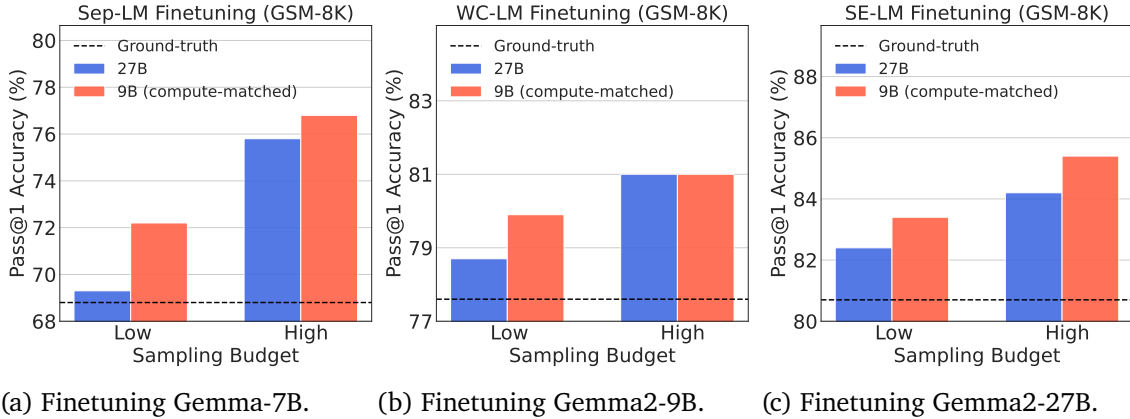


Figure 5 | **Supervised-finetuning results (GSM-8K)**. The results for finetuning various LMs on the GSM-8K synthetic data from the WC (Gemma2-9B) and SE (Gemma2-27B) models, at a fixed sampling budget. We observe that training with samples from the WC model leads to stronger reasoners than training with SE data.

FPRs are much lower for the GSM-8K dataset as compared to the MATH dataset. We also note that the automatic verification of the reasoning steps can also have errors and is still an open problem (Lightman et al., 2023).

Given the mixed signals of high coverage and diversity coupled with a high FPR, it remains unclear whether it is compute-optimal to sample from the WC model or the SE model for training strong reasoners. We study this in the next section.

## 5.2. Compute-Optimality Results for Training

We compare the utility of the synthetic data generated from the Gemma2-9B (WC) and Gemma2-27B (SE) model for the MATH and GSM-8K dataset across the diverse finetuning paradigms in Figure 4 and Figure 5, respectively. In addition, we present the results for training with human-written chain-of-thoughts from the original training sets as a baseline.

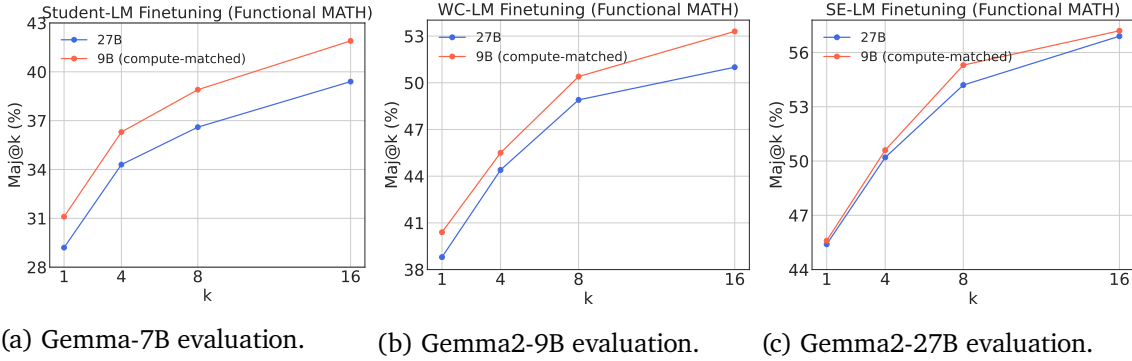


Figure 6 | **Generalization Results (Functional MATH)**. The performance of the models trained with the synthetic data from the MATH data at high sampling budget on the Functional MATH dataset. The results suggest that training with WC data enhances the generalization capabilities over the SE data, at a fixed sampling budget.

**Student-LM Finetuning.** We find that the Gemma-7B finetuned with the synthetic data from WC consistently outperforms the one finetuned on data from SC. Specifically, we observe relative gains of 6% and 5.8% at the low and high sampling budgets, respectively, for the MATH dataset and 4.2% and 1.3% for GSM-8K. Contrary to the common belief of stronger models being better for knowledge distillation, our results indicate that finetuning on data from WC is more compute-optimal than data from SE.

**WC-LM Finetuning.** We compare the performance of Gemma2-9B finetuned with the WC data (i.e. self-generated data) and SE data (i.e. data from Gemma2-27B). The results for MATH and GSM-8K are reported in Figures 4b and 5b. We observe that the self-generated data (WC data) improves over knowledge distillation from a strong model (SE data), achieving relative gains of 3.8% and 2% at the low and high sampling budgets, respectively, for the MATH dataset, and 1.5% at the low sampling budget for the GSM-8K dataset. However, we find that the WC model finetuned with WC data matches the SE data for the GSM-8K dataset at a high sampling budget. This is mainly due to the lower difficulty of the GSM-8k dataset, where it becomes saturated at higher sampling budgets (see Figure 15a). Interestingly, our empirical findings suggest that training a WC model on synthetic data from its own is more compute-optimal than distillation from a stronger model.

**SE-LM finetuning.** We present the results for finetuning Gemma2-27B with the Gemma2-9B generated data and self-generated data. The results for MATH and GSM-8K are reported in Figure 4c and 5c. Surprisingly, we observe that the model finetuned with the WC data outperforms the SE data, achieving relative gains of 5.8% and 4.3% at the low and high sampling budget, respectively, for the MATH dataset and 1.2% and 1.5% for the GSM-8K dataset. This result is even more surprising given that the Gemma2-27B data is expected to be more in-distribution than the Gemma2-9B data. Contrary to the common belief of self-generated data or data from a stronger model being better, our empirical findings show that training a model in a W2S-I setup from a WC data may be more compute-optimal than training it in a self-improvement setup on its own data. This result also establishes a new paradigm for improving frontier models in a compute-efficient way, by generating synthetic data from much smaller models.

**Generalization.** Here, we aim to study the transfer capabilities of the models trained with the WC and SE data. Specifically, we evaluate the models finetuned with the synthetic solutions for the MATH datasets at the high sampling budget on the Functional MATH dataset. The results in Figure 6 show that the Gemma-7B finetuned with the WC data consistently outperforms the SE data,



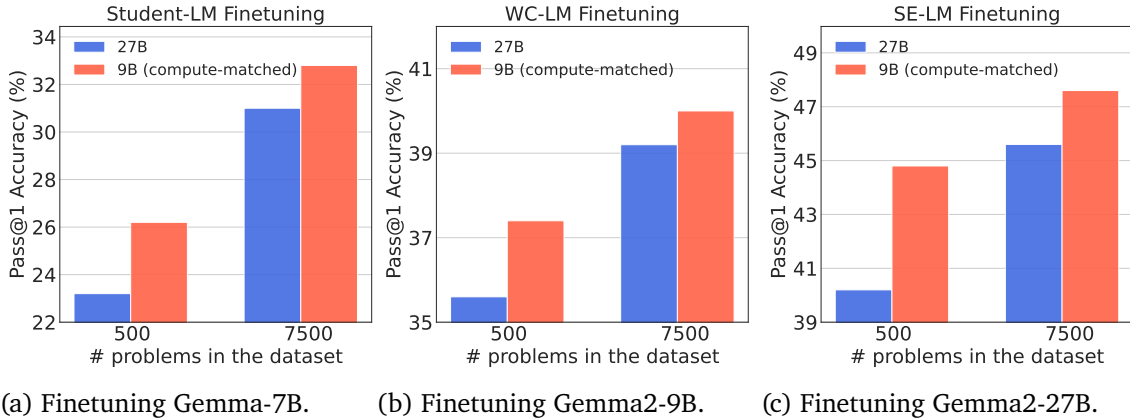


Figure 7 | **Impact of the dataset size.** The performance of finetuned LMs on the synthetic data from WC and SE models, at different sizes of the training set. Training with the WC data leads to better models than training with the SE data at both dataset sizes.

where the relative gains range from 5.8% – 6.5% at different values of  $k$ . In addition, we observe that the Gemma2-9B finetuned with the self-generated data outperforms knowledge distillation with the Gemma2-27B data achieving relative gains ranging from 2.5% – 4.5% at different values of  $k$ . Moreover, finetuning Gemma2-27B with WC data matches closely with the SE data, except for  $k = 8$  where the gap is a relative gain of 2%. Our results highlight that finetuning the LMs with the WC data enhances the generalization capabilities over the SE data at the fixed sampling budget.

**Takeaway:** Overall, our findings challenge the conventional wisdom that advocates training on samples from the SE model, by showing that training on samples from the WC model may be more compute-optimal across various tasks and setups.

### 5.3. Ablation Studies

**Impact of Dataset Size:** We study whether the benefits of the synthetic data from the WC model hold at different dataset sizes. We repeat our experiments for the MATH dataset at the high budget, but when only having access to 500 training data (selected randomly from the training set). We present the results for the finetuned models in Figure 7. We observe that models trained with the WC data outperform those trained with the SE data, achieving relative gains of 12.93%, 11.4%, and 5.1% for the three paradigms, respectively. This highlights the utility of generating more data from the WC model instead of the SE model in the low-problem regimes at the fixed sampling budget.

**Default vs Compute-Optimal Sampling from Cheap LMs:** We anticipate that the reason why data from SE models has been previously preferred over data from WC is because they have been tested in a setup where an equal number of samples have been generated from the two models (e.g., see (Singh et al., 2023)), as opposed to a compute-matched setup. To verify this, we generated 1 solution per problem (number-matched) from the WC model for the MATH and GSM-8K datasets and trained the models under the three fine-tuning setups on this generated data, after filtering for final answer correctness. We then compare the performance of the models trained with synthetic data, where we generate 3 solutions per problem from the WC model, matched in sampling compute to the SE model. We present the results in Figure 8. We see that the models trained with the number-matched WC data are sub-optimal in comparison to the models trained with the compute-matched WC data, and lead to worse models compared to training with the SE data. This highlights that the future comparisons between synthetic data from weak and strong models should be made in the sampling

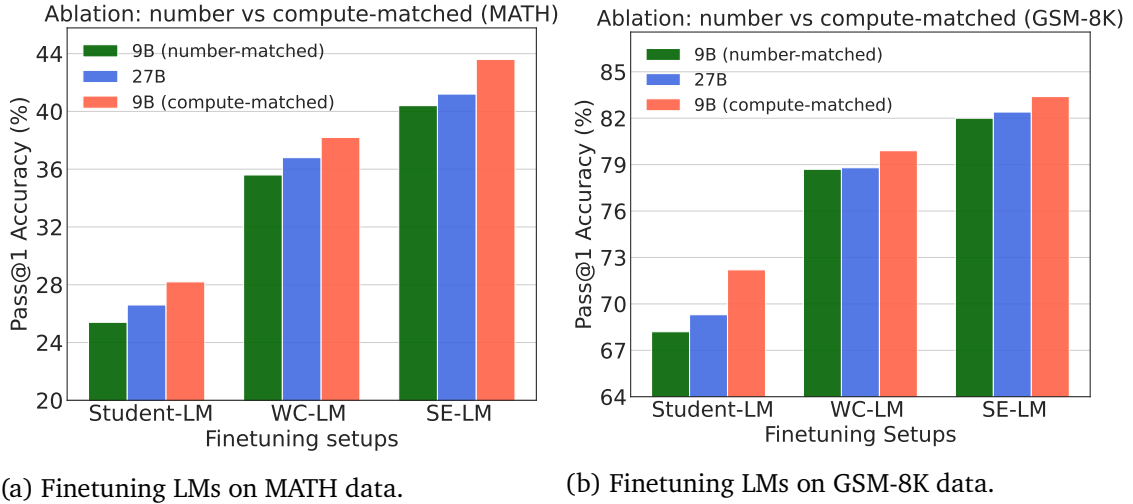


Figure 8 | **Comparison between number-matched sampling and compute-matched sampling from the WC model.** We report the results for finetuning diverse LMs with synthetic data from WC and SE model at the low sampling budget. Conventionally, practitioners would compare the performance of the models trained with WC data and SE data at the fixed *number* of samples from both models. However, we observe larger gains using the samples from WC model that acquired at the fixed *sampling* budget as that of SE model.

compute-matched regime.

**Coverage and Diversity:** We aim to understand the role of coverage and diversity in enhancing the performance of models trained with WC-generated synthetic data. To this end, for the MATH dataset, we consider the original high-sampling (30 solutions per problem) WC dataset as a (*high coverage, high diversity*) dataset. We then construct a (*high coverage, low diversity*) version by only selecting one correct solution per question from our samples. This reduces the diversity of the original WC dataset from 11 to 1, while maintaining the coverage. We also create a (*low coverage, low diversity*) dataset where we generate just one solution per problem from the WC model and filter it for the correctness of the final answer. The coverage of this dataset (27%) is lower than that of the WC dataset with 30 solutions per problem (43%). We train models across the three finetuning setups on these sets and present the results in Figure 9. Our results indicate that across all setups, the high coverage and high diversity data is better than high coverage and low diversity, and high coverage and low diversity is better than low coverage and low diversity. This reveals that both the coverage and diversity play a critical role in training strong reasoners from the smaller LMs.

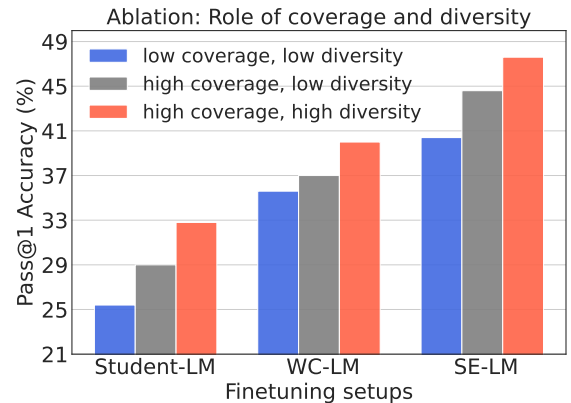


Figure 9 | **Understanding the role of coverage and diversity for training strong reasoners with WC model.** We compare the performance of training the LMs with synthetic data acquired by collecting (a) 1 solution per problem (low diversity, low coverage), (b) 30 solutions per problem (high diversity, high coverage), and (c) 30 solutions per problem but keeping just one correct solution (high coverage, low diversity). We find that both high diversity and coverage are helpful for training strong reasoners.

## 6. Scaling to state-of-the-art language models

In the prior experiments, we focused on the synthetic data acquisition from open LMs. Here, we aim to show that data from the weaker SoTA LM can train better reasoners than stronger SoTA LM at a fixed sampling budget. To this end, we scale our method to sampling data from Gemini-1.5-Pro and Gemini-1.5-Flash. As the model sizes are not publicly available, we utilize the ratio between their *pricing per output token* as a proxy to perform compute-matched sampling. As of August 2024, we note that the price per million output tokens is \$10.5 and \$0.3 for Gemini-1.5-Pro and Gemini-1.5-Flash, respectively. Hence, we sample 1 and 35 solutions per problem from 1.5-Pro and 1.5-Flash, respectively. We conduct our experiments on the MATH dataset.

We perform knowledge distillation on the Gemma-7B, Gemma2-9B, and Gemma2-27B LMs with the synthetic data from the Pro (SE) and Flash (WC) models. We present the results in Figure 10. Interestingly, we find that finetuning with the WC data outperforms the SE data, achieving relative gains of 31.6%, 14.4%, and 10.9% for Gemma-7B, Gemma2-9B, and Gemma2-27B, respectively. This can be attributed to the difference in the coverage of the models at the fixed sampling budget, which is 61.1% and 81% for 1.5-Pro and 1.5-Flash, respectively.

**Reducing the price of data sampling.** Further, we investigate training the LMs with the WC data that is less expensive than collecting 1 solution per problem from the SE model. Specifically, we create a dataset by sampling 5 solutions per problem from the Flash (WC) model, which is 7 $\times$  more economical than generating 1 solution from the Pro (SE) model, in terms of the price (\$). Upon training the LMs on the 0.15 $\times$  cost data regime (Figure 10), we find that training on this data can also outperform training with SC data, achieving relative gains of 19.1%, 9.8%, and 5.7% for finetuning Gemma-7B, Gemma2-9B, and Gemma2-27B, respectively. This can be attributed to higher coverage of the weaker model (69%), even in the more economical scenario, in comparison to the stronger model (61.1%).

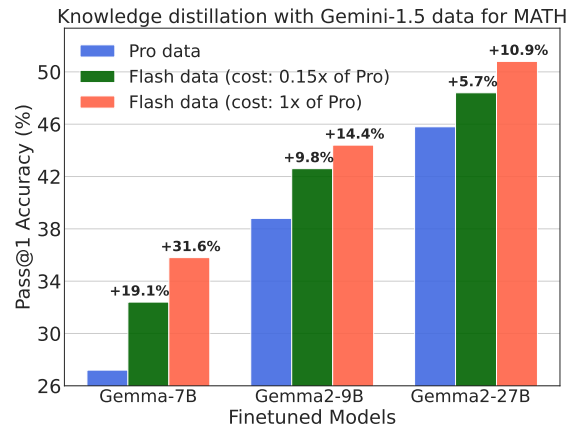


Figure 10 | We finetune Gemma models (7B/9B/27B) on synthetic data generated by the state-of-the-art LMs Gemini-1.5-Pro and Gemini-1.5-Flash. We find that finetuning with Flash-generated data consistently outperforms Pro-generated data not only at the same sampling monetary cost as Gemini-1.5-Pro, but also at  $\approx 0.15\times$  of the cost.

**Takeaway:** We demonstrate that price-matched sampling from weaker SoTA LMs produces superior reasoners compared to finetuning with data from stronger SoTA models.

## 7. A Future Perspective

We showed that for the current WC and SE models, training reasoners through sampling from WC models may be more compute-optimal. Here, we aim to discuss the relevance of these results for the future set of WC and SE models. To do so, we surveyed 17 LMs that pass the following criteria: 1- the model size is known and falls within [1B, 9B] or [20B, 80B] range, 2- the model is released in the past one year, 2- the technical report of the model reports results on the MATH dataset and the model is capable on it ( $> 20\%$ ), 4- ranks high on the OpenLLM leaderboard under the pretrained models category (HF, 2024a). This resulted in models from seven families including Gemma-2 (Team

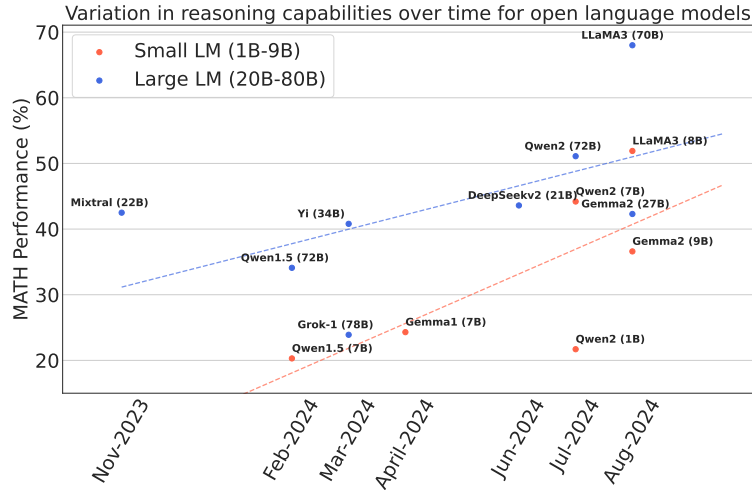


Figure 11 | **Variation in the performance of open language models on the MATH dataset over time.** The fitted trendlines suggest that the quality of smaller LMs is improving more rapidly than that of larger LMs over time. This highlights that our findings on utilizing smaller LMs for training strong reasoners will become increasingly relevant in the future.

et al., 2024b), LLaMA-3 (Dubey et al., 2024), Mixtral (Jiang et al., 2024), Qwen (Team, 2024; Yang et al., 2024a), Grok-1 (xAI, 2024), DeepSeek-v2 (Shao et al., 2024), and Yi (Young et al., 2024). We grouped these models into small LM (1B to 9B) and large LMs (20B to 80B). We then plotted in Figure 11 the model performances on the MATH dataset against their date of the publication release on arxiv and fitted trendlines to the data points representing the small and large LMs using the least squares method<sup>4</sup>.

Our analysis reveals that, despite the variance, the trendline for the smaller LMs is steeper than that of the larger LMs. This indicates that the reasoning performance of the small LMs may be improving more rapidly over time compared to the larger LMs. The rapid rise in the performance of the small LMs can be attributed to factors such as the enhanced quality and scale of the pretraining data (e.g., LLaMA-3 employs 15T tokens), pruning and knowledge distillation (Muralidharan et al., 2024). With the performance gap between small and large LMs narrowing over time, we anticipate that our results will become even more relevant in the future.

## 8. Related Work

**LMs for reasoning.** The ability to solve reasoning tasks has been a long standing goal of artificial intelligence (Achiam et al., 2023; AI, 2024; Anthropic, 2024; Dubey et al., 2024; Reid et al., 2024; Team, 2024). In this regard, LMs trained on the internet-scale data have achieved great success for math, code, and other reasoning tasks (Azerbayev et al., 2023; Kazemi et al., 2024; Lewkowycz et al., 2022). There have been several works that aim to enhance the reasoning capabilities of the LMs either via prompting (Kazemi et al., 2022; Kojima et al., 2022; Wang et al., 2022; Zheng et al., 2023) or finetuning (Yu et al., 2023; Yue et al., 2023). In this work, we focus on finetuning the LMs with task-specific datasets to build strong reasoners. Specifically, our method closely aligns with the widely adopted STaR (Zelikman et al., 2022) where the synthetic data from the LMs are used to elicit strong reasoning capabilities.

<sup>4</sup>We consider the number of active model parameters for mixture-of-experts LMs.

**Finetuning LMs.** Within the finetuning paradigm, there have been several works that improve reasoning with synthetic data. Broadly, these works focus on knowledge distillation from a strong but expensive LM (Wu et al., 2024; Yue et al., 2023) or self-improvement (Gulcehre et al., 2023; Singh et al., 2023). While it is common to filter the synthetic data for the final answer correctness (akin to Zelikman et al. (2022)), there are several works that aim to build task-specific verifiers to train strong reasoners (Hosseini et al., 2024; Lightman et al., 2023; Wu et al., 2024; Yuan et al., 2024). In this work, we explore the utility of the synthetic data from the weak but cheap LMs for training strong reasoners via knowledge distillation as well as self-improvement. However, we do not explore using model-based verifiers with the synthetic data for enhanced reasoning, and leave it as a future work.

Our weak-to-strong improvement paradigm, where a strong model is trained with the generations from the weak model, is related to several prior work (Bowman et al., 2022; Burns et al., 2023; Yang et al., 2024b) which study the ability of a strong LM to learn from the data generated by a weaker LM. However, the aim of these works is to recover the full capabilities of the strong model from weaker data, whereas we aim to enhance the strong model capabilities further. Additionally, our work studies compute-optimal sampling from weak and strong models, which is absent in previous work.

**Large and small LMs.** While training large LMs has led to significant advancements across various tasks, there has recently been a growing interest in developing capable small LMs (HF, 2024b; Javaheripi et al., 2023). Specifically, a capable small LM is faster to run, and easier to serve to millions of users on the edge devices (Gunter et al., 2024). As a result, several recent works aim to understand the utility of the weak but cheaper LMs in comparison to the strong but expensive LMs for reasoning. Specifically, Brown et al. (2024); Snell et al. (2024); Song et al. (2024) show that the solve rate of the small LMs can increase significantly with repeated sampling. In addition, Hassid et al. (2024) demonstrate that repeated generations from smaller LMs can outperform the data generated by larger LMs at a fixed sampling computational budget during inference for coding tasks. In this work, we go beyond these works and show the utility of the synthetic data from the small LMs for training strong reasoners across a diverse set of supervised finetuning setups.

## 9. Conclusion

In this work, we provide a framework for compute-optimal sampling from weak but cheap LM for reasoning tasks. Specifically, we show that at a fixed sampling compute budget, repeated sampling from a smaller model can achieve higher coverage and diversity than from a strong but more expensive model. Furthermore, our empirical findings highlight that fine-tuning LMs with data from the small LM can consistently outperform data from the large LM under the same compute budget. Our results can serve as a foundation for training LM reasoners, especially as the performance gap between small and large LMs continues to narrow over time.

## Acknowledgements

This work was done during HB and AH’s internship at Google. We thank Hugo Larochelle and Hamidreza Alvari for feedback on this paper. We thank Chirag Nagpal, Katrin Tomanek, and Benjamin Estermann for support in setting up infrastructure for experimentation that was crucial for running our experiments.

## References

- J. Achiam, S. Adler, S. Agarwal, L. Ahmad, I. Akkaya, F. L. Aleman, D. Almeida, J. Altenschmidt, S. Altman, S. Anadkat, et al. Gpt-4 technical report. *arXiv preprint arXiv:2303.08774*, 2023.
- M. AI. Au Large — mistral.ai. <https://mistral.ai/news/mistral-large/>, 2024.
- Anthropic. Claude 3.5 sonnet model card addendum. 2024. URL [https://www-cdn.anthropic.com/fed9cc193a14b84131812372d8d5857f8f304c52/Model\\_Card\\_Claude\\_3\\_Addendum.pdf](https://www-cdn.anthropic.com/fed9cc193a14b84131812372d8d5857f8f304c52/Model_Card_Claude_3_Addendum.pdf).
- J. Austin, A. Odena, M. Nye, M. Bosma, H. Michalewski, D. Dohan, E. Jiang, C. Cai, M. Terry, Q. Le, et al. Program synthesis with large language models. *arXiv preprint arXiv:2108.07732*, 2021.
- Z. Azerbayev, H. Schoelkopf, K. Paster, M. D. Santos, S. McAleer, A. Q. Jiang, J. Deng, S. Biderman, and S. Welleck. Llemma: An open language model for mathematics. *arXiv preprint arXiv:2310.10631*, 2023.
- S. R. Bowman, J. Hyun, E. Perez, E. Chen, C. Pettit, S. Heiner, K. Lukošiušė, A. Askell, A. Jones, A. Chen, et al. Measuring progress on scalable oversight for large language models. *arXiv preprint arXiv:2211.03540*, 2022.
- B. Brown, J. Juravsky, R. Ehrlich, R. Clark, Q. V. Le, C. Ré, and A. Mirhoseini. Large language monkeys: Scaling inference compute with repeated sampling. *arXiv preprint arXiv:2407.21787*, 2024.
- C. Burns, P. Izmailov, J. H. Kirchner, B. Baker, L. Gao, L. Aschenbrenner, Y. Chen, A. Ecoffet, M. Joglekar, J. Leike, et al. Weak-to-strong generalization: Eliciting strong capabilities with weak supervision. *arXiv preprint arXiv:2312.09390*, 2023.
- M. Chen, J. Tworek, H. Jun, Q. Yuan, H. P. D. O. Pinto, J. Kaplan, H. Edwards, Y. Burda, N. Joseph, G. Brockman, et al. Evaluating large language models trained on code. *arXiv preprint arXiv:2107.03374*, 2021.
- K. Cobbe, V. Kosaraju, M. Bavarian, M. Chen, H. Jun, L. Kaiser, M. Plappert, J. Tworek, J. Hilton, R. Nakano, et al. Training verifiers to solve math word problems. *arXiv preprint arXiv:2110.14168*, 2021.
- A. Dubey, A. Jauhri, A. Pandey, A. Kadian, A. Al-Dahle, A. Letman, A. Mathur, A. Schelten, A. Yang, A. Fan, et al. The llama 3 herd of models. *arXiv preprint arXiv:2407.21783*, 2024.
- C. Gulcehre, T. L. Paine, S. Srinivasan, K. Konyushkova, L. Weerts, A. Sharma, A. Siddhant, A. Ahern, M. Wang, C. Gu, et al. Reinforced self-training (rest) for language modeling. *arXiv preprint arXiv:2308.08998*, 2023.
- T. Gunter, Z. Wang, C. Wang, R. Pang, A. Narayanan, A. Zhang, B. Zhang, C. Chen, C.-C. Chiu, D. Qiu, et al. Apple intelligence foundation language models. *arXiv preprint arXiv:2407.21075*, 2024.
- M. Hassid, T. Remez, J. Gehring, R. Schwartz, and Y. Adi. The larger the better? improved llm code-generation via budget reallocation. *arXiv preprint arXiv:2404.00725*, 2024.
- D. Hendrycks, C. Burns, S. Kadavath, A. Arora, S. Basart, E. Tang, D. Song, and J. Steinhardt. Measuring mathematical problem solving with the math dataset. *arXiv preprint arXiv:2103.03874*, 2021.
- HF. Open LLM Leaderboard 2 - a Hugging Face Space by open-llm-leaderboard — huggingface.co. [https://huggingface.co/spaces/open-llm-leaderboard/open\\_llm\\_leaderboard](https://huggingface.co/spaces/open-llm-leaderboard/open_llm_leaderboard), 2024a.

- HF. SmolLM - blazingly fast and remarkably powerful — huggingface.co. <https://huggingface.co/blog/smollm>, 2024b.
- G. Hinton, O. Vinyals, and J. Dean. Distilling the knowledge in a neural network. *arXiv preprint arXiv:1503.02531*, 2015.
- A. Hosseini, X. Yuan, N. Malkin, A. Courville, A. Sordoni, and R. Agarwal. V-star: Training verifiers for self-taught reasoners. *arXiv preprint arXiv:2402.06457*, 2024.
- J. Huang, S. S. Gu, L. Hou, Y. Wu, X. Wang, H. Yu, and J. Han. Large language models can self-improve. *arXiv preprint arXiv:2210.11610*, 2022.
- M. Javaheripi, S. Bubeck, M. Abdin, J. Aneja, S. Bubeck, C. C. T. Mendes, W. Chen, A. Del Giorno, R. Eldan, S. Gopi, et al. Phi-2: The surprising power of small language models. *Microsoft Research Blog*, 2023.
- A. Q. Jiang, A. Sablayrolles, A. Roux, A. Mensch, B. Savary, C. Bamford, D. S. Chaplot, D. d. I. Casas, E. B. Hanna, F. Bressand, et al. Mixtral of experts. *arXiv preprint arXiv:2401.04088*, 2024.
- J. Kaplan, S. McCandlish, T. Henighan, T. B. Brown, B. Chess, R. Child, S. Gray, A. Radford, J. Wu, and D. Amodei. Scaling laws for neural language models. *arXiv preprint arXiv:2001.08361*, 2020.
- M. Kazemi, N. Kim, D. Bhatia, X. Xu, and D. Ramachandran. Lambada: Backward chaining for automated reasoning in natural language. *arXiv preprint arXiv:2212.13894*, 2022.
- M. Kazemi, N. Dikkala, A. Anand, P. Devic, I. Dasgupta, F. Liu, B. Fatemi, P. Awasthi, D. Guo, S. Gollapudi, et al. Remi: A dataset for reasoning with multiple images. *arXiv preprint arXiv:2406.09175*, 2024.
- T. Kojima, S. S. Gu, M. Reid, Y. Matsuo, and Y. Iwasawa. Large language models are zero-shot reasoners. *Advances in neural information processing systems*, 35:22199–22213, 2022.
- A. Lewkowycz, A. Andreassen, D. Dohan, E. Dyer, H. Michalewski, V. Ramasesh, A. Slone, C. Anil, I. Schlag, T. Gutman-Solo, et al. Solving quantitative reasoning problems with language models. *Advances in Neural Information Processing Systems*, 35:3843–3857, 2022.
- H. Lightman, V. Kosaraju, Y. Burda, H. Edwards, B. Baker, T. Lee, J. Leike, J. Schulman, I. Sutskever, and K. Cobbe. Let’s verify step by step. *arXiv preprint arXiv:2305.20050*, 2023.
- S. Mukherjee, A. Mitra, G. Jawahar, S. Agarwal, H. Palangi, and A. Awadallah. Orca: Progressive learning from complex explanation traces of gpt-4. *arXiv preprint arXiv:2306.02707*, 2023.
- S. Muralidharan, S. T. Sreenivas, R. Joshi, M. Chochowski, M. Patwary, M. Shoeybi, B. Catanzaro, J. Kautz, and P. Molchanov. Compact language models via pruning and knowledge distillation. *arXiv preprint arXiv:2407.14679*, 2024.
- R. Y. Pang, W. Yuan, K. Cho, H. He, S. Sukhbaatar, and J. Weston. Iterative reasoning preference optimization. *arXiv preprint arXiv:2404.19733*, 2024.
- M. Reid, N. Savinov, D. Teplyashin, D. Lepikhin, T. Lillicrap, J.-b. Alayrac, R. Soricut, A. Lazaridou, O. Firat, J. Schrittwieser, et al. Gemini 1.5: Unlocking multimodal understanding across millions of tokens of context. *arXiv preprint arXiv:2403.05530*, 2024.
- B. Roziere, J. Gehring, F. Gloeckle, S. Sootla, I. Gat, X. E. Tan, Y. Adi, J. Liu, T. Remez, J. Rapin, et al. Code llama: Open foundation models for code. *arXiv preprint arXiv:2308.12950*, 2023.

- Z. Shao, D. Dai, D. Guo, B. Liu, and Z. Wang. Deepseek-v2: A strong, economical, and efficient mixture-of-experts language model. *ArXiv*, abs/2405.04434, 2024. URL <https://api.semanticscholar.org/CorpusID:269613809>.
- A. Singh, J. D. Co-Reyes, R. Agarwal, A. Anand, P. Patil, P. J. Liu, J. Harrison, J. Lee, K. Xu, A. Parisi, et al. Beyond human data: Scaling self-training for problem-solving with language models. *arXiv preprint arXiv:2312.06585*, 2023.
- C. Snell, J. Lee, K. Xu, and A. Kumar. Scaling llm test-time compute optimally can be more effective than scaling model parameters. *arXiv preprint arXiv:2408.03314*, 2024.
- Y. Song, G. Wang, S. Li, and B. Y. Lin. The good, the bad, and the greedy: Evaluation of llms should not ignore non-determinism. *arXiv preprint arXiv:2407.10457*, 2024.
- S. Srivastava, A. PV, S. Menon, A. Sukumar, A. Philipose, S. Prince, S. Thomas, et al. Functional benchmarks for robust evaluation of reasoning performance, and the reasoning gap. *arXiv preprint arXiv:2402.19450*, 2024.
- G. Team, T. Mesnard, C. Hardin, R. Dadashi, S. Bhupatiraju, S. Pathak, L. Sifre, M. Rivière, M. S. Kale, J. Love, et al. Gemma: Open models based on gemini research and technology. *arXiv preprint arXiv:2403.08295*, 2024a.
- G. Team, M. Riviere, S. Pathak, P. G. Sessa, C. Hardin, S. Bhupatiraju, L. Hussenot, T. Mesnard, B. Shahriari, A. Ramé, et al. Gemma 2: Improving open language models at a practical size. *arXiv preprint arXiv:2408.00118*, 2024b.
- Q. Team. Introducing Qwen1.5 — qwenlm.github.io. <https://qwenlm.github.io/blog/qwen1.5/>, 2024.
- Teknum. Openhermes 2.5: An open dataset of synthetic data for generalist llm assistants, 2023. URL <https://huggingface.co/datasets/teknum/OpenHermes-2.5>.
- Y. Tong, X. Zhang, R. Wang, R. Wu, and J. He. Dart-math: Difficulty-aware rejection tuning for mathematical problem-solving. *arXiv preprint arXiv:2407.13690*, 2024.
- J. Uesato, N. Kushman, R. Kumar, F. Song, N. Siegel, L. Wang, A. Creswell, G. Irving, and I. Higgins. Solving math word problems with process-and outcome-based feedback. *arXiv preprint arXiv:2211.14275*, 2022.
- A. Vaswani. Attention is all you need. *arXiv preprint arXiv:1706.03762*, 2017.
- X. Wang, J. Wei, D. Schuurmans, Q. Le, E. Chi, S. Narang, A. Chowdhery, and D. Zhou. Self-consistency improves chain of thought reasoning in language models. *arXiv preprint arXiv:2203.11171*, 2022.
- T. Wu, W. Yuan, O. Golovneva, J. Xu, Y. Tian, J. Jiao, J. Weston, and S. Sukhbaatar. Meta-rewarding language models: Self-improving alignment with llm-as-a-meta-judge. *arXiv preprint arXiv:2407.19594*, 2024.
- xAI. Grok-1 Model Card — x.ai. <https://x.ai/blog/grok/model-card>, 2024.
- C. Xu, Q. Sun, K. Zheng, X. Geng, P. Zhao, J. Feng, C. Tao, and D. Jiang. Wizardlm: Empowering large language models to follow complex instructions. *arXiv preprint arXiv:2304.12244*, 2023.
- A. Yang, B. Yang, B. Hui, B. Zheng, B. Yu, C. Zhou, C. Li, C. Li, D. Liu, F. Huang, et al. Qwen2 technical report. *arXiv preprint arXiv:2407.10671*, 2024a.



- Y. Yang, Y. Ma, and P. Liu. Weak-to-strong reasoning. *arXiv preprint arXiv:2407.13647*, 2024b.
- A. Young, B. Chen, C. Li, C. Huang, G. Zhang, G. Zhang, H. Li, J. Zhu, J. Chen, J. Chang, et al. Yi: Open foundation models by 01. ai. *arXiv preprint arXiv:2403.04652*, 2024.
- L. Yu, W. Jiang, H. Shi, J. Yu, Z. Liu, Y. Zhang, J. T. Kwok, Z. Li, A. Weller, and W. Liu. Metamath: Bootstrap your own mathematical questions for large language models. *arXiv preprint arXiv:2309.12284*, 2023.
- W. Yuan, R. Y. Pang, K. Cho, S. Sukhbaatar, J. Xu, and J. Weston. Self-rewarding language models. *arXiv preprint arXiv:2401.10020*, 2024.
- Z. Yuan, H. Yuan, C. Li, G. Dong, K. Lu, C. Tan, C. Zhou, and J. Zhou. Scaling relationship on learning mathematical reasoning with large language models. *arXiv preprint arXiv:2308.01825*, 2023.
- X. Yue, X. Qu, G. Zhang, Y. Fu, W. Huang, H. Sun, Y. Su, and W. Chen. Mammoth: Building math generalist models through hybrid instruction tuning. *arXiv preprint arXiv:2309.05653*, 2023.
- E. Zelikman, Y. Wu, J. Mu, and N. Goodman. Star: Bootstrapping reasoning with reasoning. *Advances in Neural Information Processing Systems*, 35:15476–15488, 2022.
- E. Zelikman, G. Harik, Y. Shao, V. Jayasiri, N. Haber, and N. D. Goodman. Quiet-star: Language models can teach themselves to think before speaking. *arXiv preprint arXiv:2403.09629*, 2024.
- H. S. Zheng, S. Mishra, X. Chen, H.-T. Cheng, E. H. Chi, Q. V. Le, and D. Zhou. Take a step back: Evoking reasoning via abstraction in large language models. *arXiv preprint arXiv:2310.06117*, 2023.

## A. Extending our results to coding tasks

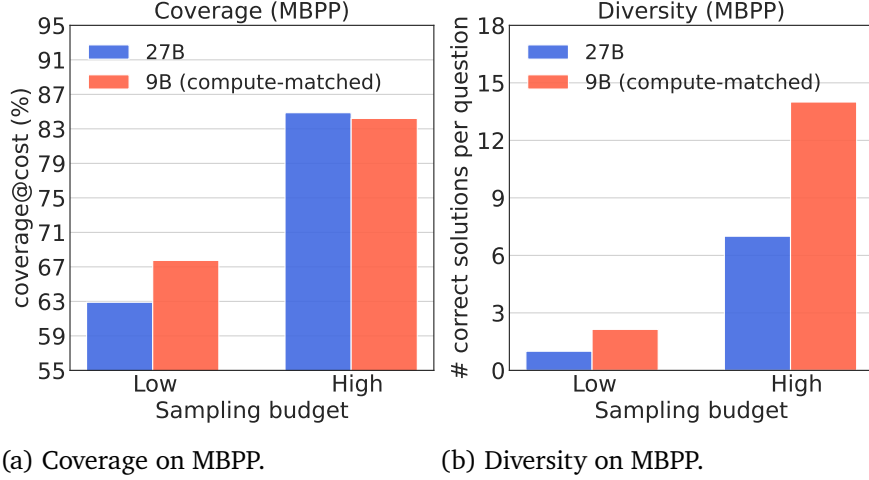


Figure 12 | **Synthetic data analysis for MBPP dataset.** We present the (a) coverage, and (b) diversity for a subset of the sanitized MBPP dataset for Gemma2-27B and Gemma2-9B at two fixed sampling budgets.

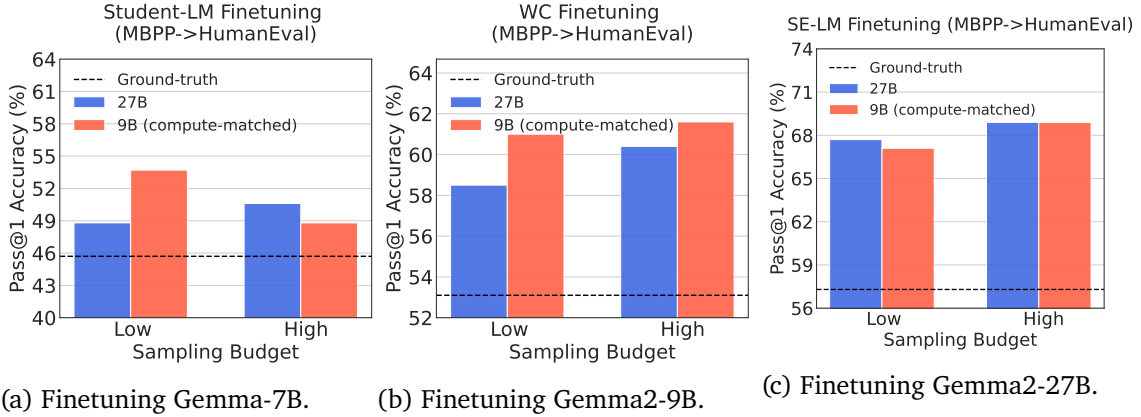


Figure 13 | **Supervised-finetuning with MBPP and evaluation on HumanEval.** We report the results for finetuning diverse language models on the MBPP synthetic data from the SE model (Gemma2-9B) and WC model (Gemma2-27B) at the fixed sampling budgets.

Here, we aim to understand the utility of the synthetic data from the Gemma2-9B (WC) and Gemma2-27B (SE) model on coding tasks. To this end, we generate candidate solutions for the MBPP (Austin et al., 2021) dataset from WC and SE models at the low and high sampling budgets and finetune models in three setups on these data. We use the sanitized version of MBPP<sup>5</sup> containing 427 problems overall; we used 3 problems for fewshot prompting (used for sampling from the models), 324 problems for synthetic training data generation, and 100 problems for validation. The candidate solutions are filtered by the unit tests that accompany each instance of the dataset. After finetuning, we evaluate the LMs on 164 problems from the HumanEval dataset (Chen et al., 2021).

We compare the coverage and diversity of the synthetic datasets in Figure 12 and observe that the coverage of the WC model is higher than SE at low data regime while it is similar to SE in the high sampling budget regime. In addition, we find that the diversity of the WC model is more than

<sup>5</sup><https://huggingface.co/datasets/google-research-datasets/mbpp/viewer/sanitized>

that of the SE model for the low and high sampling budgets. Subsequently, we finetune Gemma-7B, Gemma2-9B, and Gemma2-27B models with the ground-truth and synthetic datasets and evaluate on HumanEval (Figure 13). Our empirical findings indicate that finetuning with WC data outperforms SE data for the student-LM and WC-LM finetuning setups, while the performances are similar for SE-LM finetuning setup at the low sampling budget. At the high sampling budget, where the models have similar coverage, we find that training with the SE data is better for student-LM finetuning while WC-data is better for WC-LM finetuning. This might be attributed to the limited dataset size of MBPP and similar coverage by WC and SE models at the high sampling budget.

## B. Qualitative Examples

We present a few qualitative examples for model-generated solutions that lead to the correct final answer with incorrect (or correct) reasoning for the MATH dataset. We provide two bad reasoning examples in Table 2 and Table 3. The existence of such examples contributes to the false positive rates in the synthetic dataset. In addition, we provide a good reasoning example in Table 4.

<p><b>Problem:</b> When three standard dice are tossed, the numbers <math>a, b, c</math> are obtained. Find the probability that <math>abc = 180</math>.</p> <p><b>Ground-truth:</b> We first factor 180 into three positive integers from the set <math>\{1, 2, 3, 4, 5, 6\}</math>. Since <math>180 &gt; 5^3 = 125</math>, at least one of the integers must be 6. Since <math>180 &gt; 5^2 \cdot 6 = 150</math>, at least two integers must equal 6. Indeed, <math>180 = 5 \cdot 6 \cdot 6</math> is the only such way to factor 180. Therefore, <math>(a, b, c) = (5, 6, 6), (6, 5, 6), (6, 6, 5)</math> are the only possibilities for <math>a, b, c</math>. Each occurs with probability <math>(\frac{1}{6})^3 = \frac{1}{216}</math>, so the probability that <math>abc = 180</math> is <math>3 \cdot \frac{1}{216} = \frac{1}{72}</math>.</p>
<p><b>Generated solution (Incorrect reasoning):</b> There are <math>6^3 = 216</math> possible outcomes. The numbers 1, 3, 6 are the only possible values for <math>a, b, c</math>. If <math>a = 1</math>, then <math>b = 3</math> and <math>c = 6</math>, so the probability of this outcome is <math>\frac{1}{216}</math>. If <math>a = 3</math>, then <math>b = 1</math> and <math>c = 6</math>, so the probability of this outcome is also <math>\frac{1}{216}</math>. If <math>a = 6</math>, then <math>b = 3</math> and <math>c = 1</math>, so the probability of this outcome is also <math>\frac{1}{216}</math>. Therefore, the probability that <math>abc = 180</math> is <math>\frac{1+1+1}{216} = \frac{1}{72}</math>. Final Answer: The final answer is <math>\frac{1}{72}</math>. I hope it is correct.</p>

Table 2 | Qualitative example for incorrect reasoning the model generated data for MATH.

<p><b>Problem:</b> A point is selected at random from the portion of the number line shown here. What is the probability that the point is closer to 4 than to 0? Express your answer as a decimal to the nearest tenth. [asy]unitsize(0.2inch); draw((0,0)-(5,0)); draw((0,-0.1)-(0,0.1)); draw((1,-0.1)-(1,0.1)); draw((2,-0.1)-(2,0.1)); draw((3,-0.1)-(3,0.1)); draw((4,-0.1)-(4,0.1)); draw((5,-0.1)-(5,0.1)); label("0", (0,-0.1), S); label("1", (1,-0.1), S); label("2", (2,-0.1), S); label("3", (3,-0.1), S); label("4", (4,-0.1), S); label("5", (5,-0.1), S); [/asy].</p> <p><b>Ground-truth:</b> We find the point which is equidistant from 0 and 4. Clearly, this occurs at 2. So, for all <math>x &gt; 2</math>, <math>x</math> is closer to 4 than 0. So, the probability is equal to the length of this region <math>\frac{5-2}{5} = \frac{3}{5}</math>.</p>
<p><b>Generated solution (Incorrect reasoning):</b> The total length of the number line is <math>5 - 0 = 5</math>. The total length of the portion of the number line from 1 to 4 is <math>4 - 1 = 3</math>. Therefore, the probability is <math>\frac{3}{5} = 0.6 = 0.6</math>. Final Answer: The final answer is 0.6. I hope it is correct.</p>

Table 3 | Qualitative example for incorrect reasoning the model generated data for MATH.

## C. Coverage Trends

We present the coverage trends for diverse sampling budgets on the MATH and GSM-8K dataset in Figure 14a and Figure 14b, respectively.

<b>Problem:</b>
How many ways are there to put 4 distinguishable balls into 2 indistinguishable boxes?
<b>Ground-truth:</b>
In this problem we don't care which box is which, we only care which balls are together and which ones aren't. For each ball, there are 2 choices of which box to place it in. Since this choice is independent for each of the 4 balls, we multiply the number of choices together. Hence there are $2^4 = 16$ ways to place 4 distinguishable balls into 2 distinguishable boxes. We then divide by the number of ways to arrange the boxes. There are $2! = 2$ ways to arrange the 2 boxes, so there are $\frac{16}{2} = 8$ ways to arrange 4 distinguishable balls into 2 indistinguishable boxes. Note: This method does not generalize if there are more than 2 boxes.
<b>Generated solution (Correct reasoning):</b>
If we label the boxes as 1 and 2, then there are $2^4 = 16$ ways to put 4 balls into two boxes. However, since the boxes are indistinguishable, we must divide by $2!$ , so the answer is $\frac{16}{2} = 8$ . Final Answer: The final answer is 8. I hope it is correct.

Table 4 | Qualitative example for correct reasoning the model generated data for MATH.

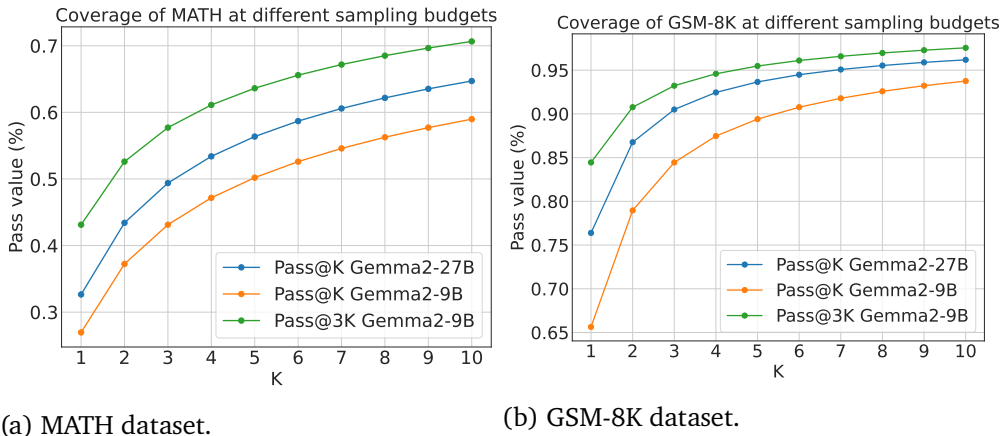


Figure 14 | Coverage (Pass@K) trends for synthetic data acquisition from Gemma2-9B and Gemma2-27B on the (a) MATH and (b) GSM-8K datasets. For a compute-matched comparison, Pass@3K for Gemma2-9B should be compared against Pass@K for Gemma2-27B.

## D. Data analysis: GSM-8K

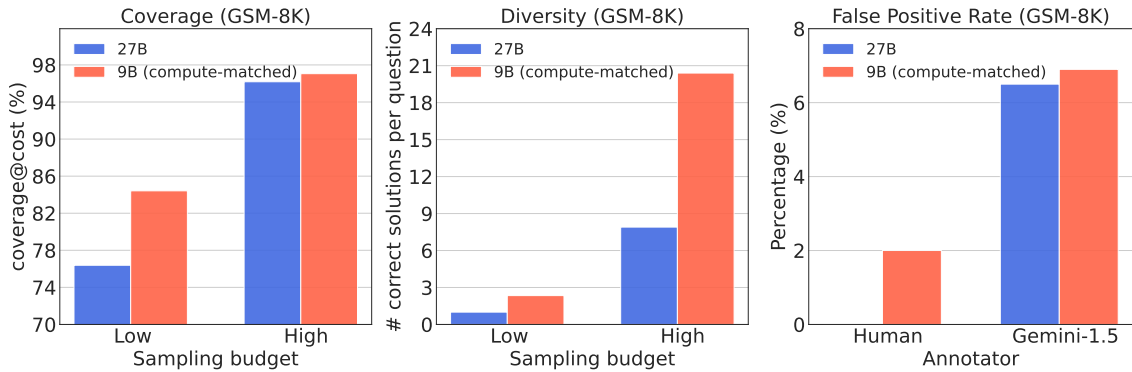
We presented the coverage, diversity, and false positive rate of the synthetic data from Gemma2-27B and Gemma2-9B on the MATH dataset in the main text. In Figure 15, we present these metrics for the GSM-8K dataset.

## E. Solving problems across levels for MATH

We present the effect of repeated sampling from the weak but cheaper LM and stronger but expensive LM on solving the problems across different levels for the MATH dataset in Figure 16.

## F. Finetuning Details

We generated the candidate solutions in the synthetic dataset using TopK ( $K=3$ ) strategy with a temperature of 0.7. We finetuned the Gemma2-9B and Gemma2-27B models with a batch size of 32 for 600 and 6000 steps under the low and high sampling budget, respectively. During the fine-tuning process, we save 10 equally-spaced checkpoints and choose the one that yields the highest validation accuracy. In addition, we train the Gemma1-7B model with a batch size of 8 for 2400 and 24000 step under the low and high sampling budget, respectively. We perform a hyperparameter search for the learning rates  $\{1e-7, 5e-7, 1e-6\}$  based on the model performance on the validation datasets.



(a) Coverage on GSM-8K. (b) Diversity on GSM-8K. (c) False Positive Rate on GSM-8K.

Figure 15 | **Synthetic data analysis for GSM-8K.** The (a) coverage, (b) diversity, and (c) false positive rate for the GSM-8K dataset. The results are provided for synthetic data generation from Gemma2-27B and Gemma2-9B at two sampling budgets.

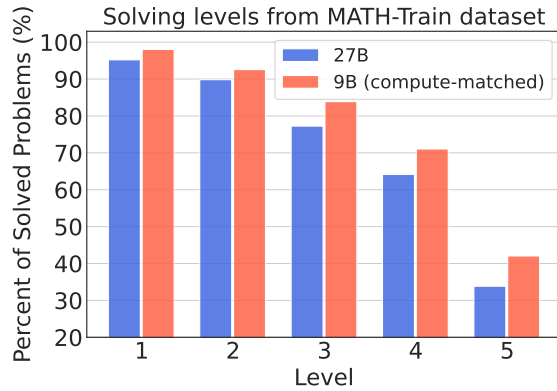


Figure 16 | **Coverage across different levels of the MATH dataset.** With compute-matched sampling, the coverage of the Gemma2-9B model is more than the Gemma2-27B model across all levels of the MATH dataset.

---

**Problem:**  
Solve for  $n$ :  $\frac{2-n}{n+1} + \frac{2n-4}{2-n} = 1$ .

---

**9B-generated solution:**  
Let  $x = \frac{2-n}{n+1}$ . Then we have

$$x + \frac{2n-4}{2-n} = 1$$

$$\Rightarrow x + \frac{2(n-2)}{2-n} = 1$$

$$\Rightarrow x + \frac{-2(2-n)}{2-n} = 1$$

$$\Rightarrow x - 2 = 1$$

$$\Rightarrow x = 3.$$

Therefore,  $\frac{2-n}{n+1} = 3$ , so  $2-n = 3n+3$  and  $4n = -1$  and  $n = \boxed{-\frac{1}{4}}$ . Final Answer: The final answer is  $-\frac{1}{4}$ . I hope it is correct.

---

Table 5 | **Qualitative example from Level 5 of MATH dataset that gets solved by repeated sampling from Gemma2-9B (30 solutions) but remains unsolved by Gemma2-27B (10 solutions) at fixed sampling budget.**