000
001SIKED:SELF-GUIDEDITERATIVEKNOWLEDGE002DISTILLATION FOR MATHEMATICAL REASONING

Anonymous authors

Paper under double-blind review

ABSTRACT

Large Language Models (LLMs) can transfer their reasoning skills to smaller models by teaching them to generate the intermediate reasoning process required to solve multistep reasoning tasks. While LLMs can accurately solve reasoning tasks through a variety of strategies, even without fine-tuning, smaller models are not expressive enough to fit the LLMs distribution on all strategies when distilled and tend to prioritize one strategy over the others. This reliance on one strategy poses a challenge for smaller models when attempting to solve reasoning tasks that may be difficult with their preferred strategy. To address this, we propose a distillation method SIKeD: Self-guided Iterative Knowledge Distillation¹, where the LLM teaches the smaller model to approach a task using different strategies and the smaller model uses its self-generated on-policy outputs to choose the most suitable strategy for the given task. The training continues in a *self-guided* iterative manner, where for each training iteration, a decision is made on how to combine the LLM data with the self-generated outputs. Unlike traditional distillation methods, SIKeD allows the smaller model to learn *which* strategy is suitable for a given task while continuously learning to solve a task using different strategies. Our experiments on various mathematical reasoning datasets show that SIKeD significantly outperforms traditional distillation techniques across smaller models of different sizes.

028 029

031

004

006

008 009

010 011

012

013

014

015

016

017

018

019

021

024

025

026

027

1 INTRODUCTION

Large language models (LLMs), with tens to hun-033 dreds of billions of parameters, generally outper-034 form smaller models (with billions of parameters or fewer) in a variety of reasoning tasks Touvron et al. (2023); Achiam et al. (2023). One notable 037 strength of large models is their ability to reason and perform multistep reasoning tasks, often considered an important aspect of intelligence Gómez-Veiga et al. (2018). However, the significant size 040 and computational demands of these large models 041 present several challenges. For example, LLaMA3 042 models Touvron et al. (2023) are trained using 043 clusters of 24,000 GPUs, limiting their accessibil-044 ity to many researchers and practitioners. 045

To bridge this gap, a key approach involves teaching smaller models to replicate the knowledge of a larger model, often referred to as *knowledge dis*-



Figure 1: Histogram of strategy choices for the LLM and the smaller model. LLM tends to select several reasoning strategies, but the smaller model is biased towards one strategy. The comparison was done on 1000 data points randomly sampled from the GSM8K train set.

tillation Hinton (2015). Typically, smaller models can be taught to replicate the multistep reasoning capabilities of larger models by incorporating a set of intermediate sequences (Kim & Rush, 2016; Shridhar et al., 2023). However, these intermediate steps can be derived from a number of different strategies, such as Chain of Thought (CoT) Wei et al. (2022), Subquestion Decomposition Shridhar et al. (2022); Zhou et al. (2023), and Program of Thoughts (PoT) Chen et al. (2023), among others. A

¹SIKeD is pronounced as "psyched"

054 viable solution is to distill these reasoning capabilities into smaller models either by distilling indi-055 vidual strategies (Magister et al., 2023; Shridhar et al., 2023; Hsieh et al., 2023) or by incorporating 056 multiple strategies simultaneously (Chenglin et al., 2023; Zhu et al., 2024). Although smaller models 057 have demonstrated impressive performance when distilled with a single strategy, they often struggle 058 to master multiple strategies equally well. An example is presented in Figure 1 where a larger model can use multiple strategies to generate data but upon distilling, a smaller model tends to favor one over the others. This is because reasoning through a variety of strategies tends to emerge as a result 060 of scaling language models, making it difficult for smaller models to replicate this behavior Lyu et al. 061 (2024).062

063 On the other hand, learning to solve a task using multiple strategies can help smaller models overcome 064 the limitations of relying on a single approach. However, a key challenge arises when, despite being trained on a fixed dataset containing various strategies, a distribution mismatch occurs between 065 the data generated by the LLM and the outputs produced by the smaller model during inference. 066 This mismatch can hinder the ability of the smaller model to generalize across different reasoning 067 strategies. This issue, often discussed in imitation learning (Pomerleau, 1991; Ross & Bagnell, 2010), 068 results in the student model consistently choosing one strategy, even when a different approach would 069 be more appropriate. As a result, the student generates outputs with strategy choices that are highly unlikely to match those produced by the teacher. 071

To address this challenge, we introduce our distillation methodology, SIKeD: Self-guided Iterative 072 Knowledge Distillation. The process begins with the LLM teaching the smaller model to approach 073 tasks using a variety of reasoning strategies, providing a strong foundation for the smaller model to 074 understand different problem-solving approaches. However, due to inherent biases and its limited 075 capacity, the smaller model may still struggle to match the LLM's distribution of strategy choices 076 effectively. To resolve this, we take inspiration from constructivist learning theory Narayan et al. 077 (2013), where the learner builds knowledge during the "assimilation phase" and refines their understanding during the "accommodation phase" to incorporate new insights. We propose generating 079 outputs using the smaller model in an on-policy setup and selecting the best strategies for the task. By mixing the LLM-generated data with self-generated outputs, we leverage the strengths of both 081 datasets. We iteratively fine-tune the smaller model allowing it to recognize strategies that it learned from the LLM but did not initially apply. With this approach, we align the smaller model with its 083 own learned knowledge rather than forcing its distribution to mirror that of the LLM's.

084 Our proposed method extends beyond traditional one-step distillation, as each iteration of SIKeD 085 leads to an updated policy that better grasps new information. We repeat multiple iterations of SIKeD based on the accuracy-cost tradeoff (does the improvement justify the cost of another iteration), 087 allowing for continuous refinement and improvement of the model's reasoning capabilities. We demonstrate the effectiveness of SIKeD on several mathematical reasoning tasks using models with 088 fewer than 7 billion parameters.² On four mathematical datasets-GSM8K Cobbe et al. (2021), 089 SVAMP Patel et al. (2021), ASDiv Miao et al. (2020), and MultiArith Roy & Roth (2016)—our ap-090 proach achieves improvements of up to +5 points over traditional distillation strategies. Additionally, 091 we show that multiple rounds of SIKeD allow the model to select the appropriate strategy for a given 092 problem, while traditional distillation using LLM's data tends to leave it biased. 093

2 PRELIMINARIES: LLM BASED DISTILLATION

Problem Setup We consider the standard setup of LLM-based distillation (also referred to as knowledge distillation), where the data is sampled from the larger model (LLM) with intermediate 098 reasoning and the smaller model is fine-tuned (distilled) over it Shridhar et al. (2023); Magister et al. (2023). Two auto-regressive sequence models are involved in the process: a larger model or the LLM 100 denoted as p_L and a smaller model to be distilled as p_{sm}^{θ} (with learnable parameters θ). In this work, 101 we consider a reasoning dataset \mathcal{D} consisting of a question q_i and a numerical answer a_i pairs for n 102 data points, i.e. $i \in \{1, ..., n\}$. Since our work focuses on improving the reasoning in smaller models 103 by teaching them to solve a variety of reasoning strategies, we consider three reasoning strategies in 104 this work: Chain-of-Thought (CoT), Least-to-Most (L2M), and Program-of-Thought (PoT). For a 105 specific reasoning strategy, denoted as $s \in S$, we generate the reasoning chain or rationale, denoted

106 107

094

²We acknowledge that "smaller model" is a relative term, and we consider models with fewer than 7 billion parameters to be smaller models.

as r_i leading to the final answer as: $r_i \sim p_L(. | pr_s, q_i)$, where, pr_s represents the strategy-specific prompt, and $s \in \{CoT, L2M, PoT\}$. Prompts used are provided in Appendix A.

110 111

112

123 124 125

126 127

128

129

130

2.1 LLM BASED DISTILLATION

113 We add the generated rationales to the the dataset \mathcal{D} to create an initial training dataset \mathcal{D}_{LLM} consisting 114 of a quadruple of $\{q_i, a_i, s, r_i\}$ for each data point. We perform a data filtering by extracting the final 115 answer \hat{a}_i from the generated rationale r_i and comparing it with the ground truth answer a_i . We 116 discard all samples that do not match, i.e., we keep samples where $\hat{a}_i = a_i$. This filtering process 117 eliminates incorrect rationales, ensuring that only high-quality data is used for distilling the smaller 118 models.

119 We start the distillation process by training the smaller model with the created dataset \mathcal{D}_{LLM} . The 120 question q_i is provided as input, and the smaller model p_{sm}^{θ} (with learnable parameters θ) is first 121 instructed to generate the strategy s, followed by the rationale r_i that leads to the final answer a_i . The 122 loss $\mathcal{L}_{L}(\theta)$ is defined as:

$$\mathcal{L}_{\mathbb{L}}(\theta) = -\mathbb{E}_{(q_i, s, r_i) \sim \mathcal{D}_{\text{LLM}}}\left[\log p_{\text{sm}}^{\theta}\left(s \mid q_i, I\right) + \sum_{t=1}^{M} \log p_{\text{sm}}^{\theta}\left(r_{i, t} \mid r_{i, < t}, s, q_i, I\right)\right], \quad (1)$$

where M represents the number of tokens decoded over time t in an autoregressive manner, and I is the instruction used during fine-tuning. Note that this is analogous to traditional knowledge distillation from LLMs except that we make a strategy choice before generating rationales.

131 **Limitations** Training solely on LLM-generated data \mathcal{D}_{LLM} can lead to a distribution mismatch 132 between the training data and the smaller model's own output distribution. Specifically, the larger 133 model due to its larger capacity, may produce correct reasoning across multiple strategies that the 134 smaller model can find difficult to replicate directly Agarwal et al. (2024). A comparison of the 135 strategy selected by the LLM and the smaller model on 1K samples is presented in Figure 1. The 136 smaller model performs poorly when generating outputs on its own, as the training data distribution 137 $P_{\text{train}}(x)$ is different from the model's output distribution $P_{\text{sm}}^{\theta}(x)$ as $P_{\text{train}}^{(1)}(x) = P_{\text{LLM}}(x)$, where x represents the samples (q_i, s, r_i) , and $P_{\text{LLM}}(x)$ is the distribution of data generated by the LLM p_L . 138 139

Proposed Solution To mitigate the distributional shift in strategy choice between the LLM and the smaller model, we propose to incorporate the smaller model's own correct outputs into the training data. This *self-guided* training with data mixing aligns the training data distribution more closely with the smaller model's output distribution, making learning more effective. A visualization of the data mixing approach is presented in Figure 2 that demonstrates that data mixing reduces the distribution shift, bringing the LLM and the smaller model's output distribution closer. This allows the smaller model to choose the right strategy for a given task, much like the LLM.

147 148

149

3 SIKED: SELF-GUIDED ITERATIVE KNOWLEDGE DISTILLATION

150 We propose SIKeD, an *iterative* training where smaller models can take advantage of their own 151 generations to refine their strategy choices for a given task. In a nutshell, we generate data from the 152 smaller model, filter out the correct samples based on whether the generated solutions are correct, and 153 mix this data with the LLM-generated data to adjust its strategy preferences. The smaller distilled model is used to iteratively generate data in an on-policy setting where it updates itself by leveraging 154 both the LLM data and its own generations. This iterative process allows the smaller model to improve its reasoning abilities and strategy selection over time by leveraging the LLM's knowledge 156 and its own prior learning. The following paragraphs discuss the steps involved in our proposed 157 iterative distillation methodology and the training objective. 158

159

Data generation For each question q_i and its associated reasoning strategy s, we first generate Krationales using the current smaller model p_{sm}^{θ} as: $r_i^{(k)} \sim p_{sm}^{\theta}(\cdot \mid s, q_i, I)$, for k = 1, ..., K. Note that we generate multiple samples K as the likelihood of a correct answer being present in one of the



Figure 2: Alignment of the smaller model's strategy distribution with the LLM over iterations. Each subplot represents an iteration in the training process, showing the probability distributions over reasoning strategies: PoT, L2M, and CoT. The blue bars depict the LLM's distribution P_L , while the orange bars represent the smaller model's distribution P_{SM} , which is biased towards CoT. The green 176 bars show the training data distribution $P_{\text{train}}^{(t)}$, a mixture of P_L and P_{SM} weighted by the mixing rate α . As α decreases over iterations (from 0.90 to 0.20), $P_{\text{train}}^{(t)}$ shifts from being similar to the LLM's distribution towards the smaller model's distribution. The KL divergence between the training data and the smaller model distributions decreases accordingly, indicating increased similarity.

rationales increases significantly with additional generations for smaller models (Jain & Shridhar, 2023; Wang et al., 2023).

Data Filtering Next, we extract the predicted answer $\hat{a}_i^{(k)}$ from each rationale $r_i^{(k)}$ and compare it with the ground truth a_i . We collect the correct samples, where $\hat{a}_i^{(k)} = a_i$, into a new dataset $\mathcal{D}_{\text{self}} = \left\{ (q_i, s, r_i^{(k)}) \mid \hat{a}_i^{(k)} = a_i \right\}.$

Data mixing We combine the LLM-generated dataset \mathcal{D}_{LLM} with the self-generated dataset \mathcal{D}_{self} 191 to form the mixed dataset $\mathcal{D}_{\text{mix}} = \mathcal{D}_{\text{LLM}} \cup \mathcal{D}_{\text{self}}$. 192

193 Note that, we do not always use all the data from LLM in \mathcal{D}_{mix} , and study two variations: All when all LLM data is used in \mathcal{D}_{mix} , and *Sparse* when only queries that have no correct generations 194 in \mathcal{D}_{self} are taken from \mathcal{D}_{LLM} . Sparse uses less generated data from the LLM, resulting in more 195 computationally efficient training. 196

197 The corresponding training data distribution changes to a mixture of the LLM data distribution and the model's own output distribution: 199

$$P_{\text{train}}^{(2)}(x) = \alpha P_{\text{LLM}}(x) + (1-\alpha) P_{\text{sm}}^{\theta}(x),$$

where $\alpha = \frac{|\mathcal{D}_{\text{LLM}}|}{|\mathcal{D}_{\text{LLM}}| + |\mathcal{D}_{\text{self}}|}$ serves as a normalized mixing rate between the two datasets.

Training objective By including \mathcal{D}_{self} in the training data, we reduce the divergence between $P_{\text{train}}^{(2)}(x)$ and the model's own output distribution $P_{\text{sm}}^{\theta}(x)$, thus minimizing the distribution shift and improving training effectiveness of choosing the right strategy for a given task.

We continue training the smaller model on \mathcal{D}_{mix} using the following loss function:

$$\mathcal{L}_{\min}(\theta) = -\mathbb{E}_{(q_i, s, r_i) \sim \mathcal{D}_{\min}} \left[\log p_{\text{sm}}^{\theta}\left(s \mid q_i, I\right) + \sum_{t=1}^{M} \log p_{\text{sm}}^{\theta}\left(r_{i,t} \mid r_{i, < t}, s, q_i, I\right) \right]$$
(2)

The expected loss over the training data is: 214

173

174

175

177

178 179

181 182

183

185

186

187 188

189 190

200 201 202

203 204

205

206

207 208

209 210 211

212 213

 \mathcal{L}

$$\min(\theta) = -\mathbb{E}_{x \sim P_{\text{train}}^{(2)}(x)} \left[\log p_{\text{sm}}^{\theta}(x) \right]$$

216 where $x = (q_i, s, r_i)$, and $p_{em}^{\theta}(x)$ denotes the probability assigned by the model to the sample x. 217

Analogous to minimizing the KL divergence Mixing the data is analogous to minimizing the Kullback-Leibler (KL) divergence Kullback & Leibler (1951) between the training data distribution $P_{\text{train}}^{(2)}(x)$ and the model's output distribution $P_{\text{sm}}^{\theta}(x)$:

$$D_{\mathrm{KL}}(P_{\mathrm{train}}^{(2)}(x) \parallel P_{\mathrm{sm}}^{\theta}(x)) = \sum_{x} P_{\mathrm{train}}^{(2)}(x) \log \frac{P_{\mathrm{train}}^{(2)}(x)}{P_{\mathrm{sm}}^{\theta}(x)}.$$

As we include more self-generated data, $(1 - \alpha)$ increases, and $P_{\text{train}}^{(2)}(x)$ becomes closer to $P_{\text{sm}}^{\theta}(x)$. This reduces the KL divergence and aligns the training data distribution with the model output distribution, leading to more effective learning. Figure 2 demonstrates that as the value of α decreases, the distribution of the training data strategy choices aligns with the smaller model, with a reduction in their KL value over iterations. This allows the smaller model to better capture the strategy distribution of the larger model.

3.1 ITERATIVE SELF-TRAINING OF SIKED

234 We repeat the data generation, filtering, mixing, and training steps iteratively. In each iteration t, the 235 smaller model potentially generates new correct rationales that are added to the training data. The training data distribution at iteration t becomes: 236

$$P_{\text{train}}^{(t)}(x) = \alpha^{(t)} P_{\text{LLM}}(x) + (1 - \alpha^{(t)}) P_{\text{sm}}^{\theta^{(t-1)}}(x),$$

where $\theta^{(t-1)}$ are the model parameters from the previous iteration, and $\alpha^{(t)}$ is updated based on 240 the sizes of \mathcal{D}_{LLM} and $\mathcal{D}_{\text{self}}^{(t)}$ at iteration t. Note that the generated samples from the smaller model 242 automatically govern the value of $\alpha^{(t)}$ based on the size of $\mathcal{D}_{self}^{(t)}$.

This iterative process continues until the model's performance converges or a predefined stopping 244 criterion is met. Over multiple iterations, the model's own output distribution $P_{sm}^{\theta^{(t)}}(x)$ gradually 245 improves, and the training data distribution becomes increasingly aligned with it. We present an 246 end-to-end training methodology in Algorithm 1. 247

248 249

250

218

219

220

221

222

224 225

226

227

228

229

230

231 232

233

237 238 239

241

243

EXPERIMENTAL DETAILS 4

251 **Dataset** Our work demonstrates the effectiveness of selecting an appropriate strategy for a given 252 task. We consider multi-step mathematical reasoning datasets in our work, as various strategies 253 can solve the task fairly well. We trained SIKeD on the GSM8K training set Cobbe et al. (2021), 254 which includes 7,473 samples, and tested it on the corresponding test set of 1,319 samples. To assess the domain transferability of our distillation method, we also evaluated it on three additional 255 mathematical datasets: SVAMP Patel et al. (2021) with 1,000 samples, ASDiv Miao et al. (2020) with 256 2,300 test samples, and MultiArith Roy & Roth (2016) with 180 samples. As the GSM8K training 257 set was used to train the smaller model, we classify it as an *in-distribution* dataset. In contrast, no 258 training data from SVAMP, ASDiv, or MultiArith was used, as they were exclusively employed for 259 testing purposes and thus considered out-of-distribution.

260

261 **Implementation Details** We used the Llama3 70B model Dubey et al. (2024) as the large language 262 model (LLM) to generate the rationales. We performed distillation on different smaller models 263 ranging from 0.5B to 7B parameters, including Qwen2 0.5B Bai et al. (2023), Qwen2 1.5B Bai et al. 264 (2023), SmolLM 1.7B Hugging Face (2023), Gemma 2B Team et al. (2024), and Gemma 7B Team 265 et al. (2024). All smaller models were fine-tuned using LoRA Hu et al. (2021) with a rank of 16, 266 and alpha of 32. We used a learning rate of 3e-4 for Qwen models with a cyclic scheduler, while we 267 set 2e-4 as the learning rate for other models and used a linear scheduler. We train all models for 3 epochs. We implemented all our experiments using the Unsloth FastLanguageModel Unslothai 268 (2023) and used the VLLM library Kwon et al. (2023) for inference. We set the temperature t = 0 for 269 data generation from the LLM while t = 0.7 was used for generating samples from the smaller model

270 Algorithm 1: SIKeD: Self-guided Iterative Knowledge Distillation 271 **Input:** \mathcal{D} : Reasoning dataset with questions $\{q_i\}_{i=1}^N$ and answers $\{a_i\}_{i=1}^N$, \mathcal{D}_{LLM} : Reasoning dataset 272 generated using the LLM with questions $\{q_i\}$, answers $\{a_i\}$, strategy $\{s\}$, rationales $\{r_i\}$, S: Set 273 of reasoning strategies, I: instruction, p_L : LLM for rationale generation, $p_{sm}^{\theta^{(0)}}$: Smaller model 274 with initial parameters $\theta^{(0)}$, K: Number of samples per question and strategy, T: Maximum 275 number of iterations, Variation: All or Sparse, 276 // LLM-Based Distillation 277 Train $p_{sm}^{\theta^{(0)}}$ on \mathcal{D}_{LLM} by minimizing $\mathcal{L}_{L}(\theta^{(0)})$ (equation 1) 278 // SIKeD: Self-quided Iterative Knowledge Distillation 279 **for** *iteration* t = 1 *to* T **do** Initialize dataset $\mathcal{D}_{\text{sel}}^{(t)}$ $- \leftarrow \emptyset$ 281 for each question $q_i \in \mathcal{D}$ do for each strategy $s \in S$ do for k = 1 to K do Generate rationale $r_i^{(k)}$ using $p_{sm}^{\theta^{(t-1)}}$: $r_i^{(k)} \sim p_{sm}^{\theta^{(t-1)}}$ ($\cdot \mid s, q_i, I$) 284 Extract answer $\hat{a}_i^{(k)}$ from $r_i^{(k)}$ if $\hat{a}_i^{(k)} = a_i$ then Add $(q_i, s, r_i^{(k)})$ to $\mathcal{D}_{\text{self}}^{(t)}$ 287 end end 289 end 290 end 291 if Variation is All then Combine datasets: $\mathcal{D}_{mix}^{(t)} = \mathcal{D}_{LLM} \cup \mathcal{D}_{self}^{(t)}$ 292 293 else Identify questions with no correct self-generated rationales: $\mathcal{I} = \{i \mid \text{no correct } r_i^{(k)} \text{ in } \mathcal{D}_{self}^{(t)}\}$ Include corresponding LLM data: $\mathcal{D}_{\text{LLM}}^{(t)} = \{(q_i, s, r_i) \in \mathcal{D}_{\text{LLM}} \mid i \in \mathcal{I}\}$ Combine datasets: $\mathcal{D}_{\text{mix}}^{(t)} = \mathcal{D}_{\text{LLM}}^{(t)} \cup \mathcal{D}_{\text{self}}^{(t)}$ 295 296 end 297 Update $\alpha^{(t)} = \frac{|\mathcal{D}_{\text{LLM}}^{(t)}|}{|\mathcal{D}_{\text{LLM}}^{(t)}| + |\mathcal{D}_{\text{self}}^{(t)}|}$ 298 299 Retrain $p_{sm}^{\theta^{(t)}}$ on $\mathcal{D}_{mix}^{(t)}$ by minimizing $\mathcal{L}_{mix}^{(t)}(\theta^{(t)})$ (equation 2) 300 end 301 **Output:** Updated smaller model $p_{sm}^{\theta^{(T)}}$ 302 303

at each iteration. We set the number of generated samples or K to 10. We report Top-1 accuracy (maj@1).

Our proposed approach was compared against a set of single-strategy distillation methods. In this work, we employed three reasoning strategies: Chain-of-Thought (CoT) Wei et al. (2022), Program-of-Thought (PoT) Chen et al. (2023), and Least-to-Most (L2M) Zhou et al. (2023).

310 311 312

304

305

306 307

308

5 RESULTS AND DISCUSSION

313 **LLM Based Distillation** We start by distilling smaller models using the reasoning dataset generated 314 using the LLM in two variations: using data from a single strategy (CoT, PoT, or L2M), and a 315 combination of all three strategies (referred to as "Combined"). Table 1 compares the accuracies 316 of the approaches across four mathematical datasets. The "Combined" approach benefited smaller 317 models, yielding slight improvements for the Qwen 0.5B, Qwen 1.5B, and SmolLM 1.7B models. 318 However, it showed little to no improvement, and sometimes even worse performance, for the larger 319 Gemma 2B and 7B models. This indicates that simply merging the distillation data for each strategy 320 is not sufficient for effective multi-strategy distillation.

321

Consistent improvement across in-distribution dataset Compared to the traditional LLM-based distillation approaches, we observe consistent improvements with SIKeD across all models, ranging from 0.5B to 7B parameters as shown in Table 1. On the in-distribution GSM8K dataset, both Gemma

Table 1: Top-1 (maj@1) accuracy comparison across different sized models (Gemma 2B and 7B, SmolLM 1.7B, Qwen 0.5B and 1.5B) on four mathematical datasets: GSM8K, ASDiv, MultiArith, and SVAMP. "Combined" refers to the scenario where data from all three reasoning strategies are merged and then used for distillation. We report the two variants of SIKeD: "Sparse" and "All".

			Datas	et	
Model	Method	In Distribution	0	ut of Distributi	on
		GSM8K	ASDiv	MultiArith	SVAMP
	CoT	67.40	68.76	98.33	66.80
	L2M	69.29	64.69	96.11	64.80
	PoT	71.34	67.85	98.89	75.00
Gemma 7B	Combined	70.74	<u>69.11</u>	<u>99.44</u>	69.40
	SIKeD				
	Sparse	73.84 († +2.5)	70.59 (↑ +1.5)	99.44 (-)	72.90 (↓ -2.1)
	All	72.55 (++1.2)	70.33 (++1.2)	100.0 (↑ +0.6)	75.00 (↑ +0.0)
	CoT	36.54	54.01	87.22	41.90
	L2M	36.92	43.47	81.67	31.60
	PoT	44.05	58.13	<u>90.56</u>	56.80
Gemma 2B	Combined	44.05	57.96	84.44	56.20
	SIKeD				
	Sparse	47.23 († +3.2)	59.05 (++0.8)	91.11 († +0.6)	58.60 (↑ +1.8)
	All	46.02 (++2.0)	59.39 (↑ +1.3)	91.11 († +0.6)	57.50 (++0.7)
	СоТ	16.38	30.37	58.89	22.60
	L2M	18.73	22.13	53.89	17.90
	PoT	23.73	43.77	61.11	34.50
SmolLM 1.7B	Combined	24.56	<u>46.77</u>	<u>67.22</u>	<u>35.90</u>
	SIKeD				
	Sparse	27.98 († +3.4)	47.20 (↑ +0.4)	72.22 († +5.0)	37.80 (↑ +1.9)
	All	26.54 (+2.0)	45.47 (↓-1.3)	70.56 (+3.2)	36.30 (++0.4)
	СоТ	55.57	68.76	99.44	66.30
	L2M	54.59	63.69	96.67	62.30
	PoT	64.22	66.94	95.56	<u>74.30</u>
Qwen 1.5B	Combined	<u>64.44</u>	<u>67.64</u>	98.89	73.20
	SIKeD				
	Sparse	64.97 († +0.5)	68.98 († +1.3)	99.44 (-)	75.40 (↑ +1.1)
	All	64.14 (↓ -0.3)	67.72 († +0.1)	99.44 (-)	74.50 (↑+0.2)
	CoT	36.47	54.66	83.89	43.00
	L2M	33.59	49.76	76.67	44.60
	PoT	41.62	56.83	92.22	<u>51.40</u>
Qwen 0.5B	Combined	42.38	<u>57.79</u>	90.56	<u>51.40</u>
	SIKeD	12.1.1			-1 -0
	Sparse	43.14 († +0.8)	58.44 († +0.7)	93.33 († +1.1)	51.70 († +0.3)
	All	44.28 († +1.9)	58.05 († +0.3)	95.00 († +2.8)	51.70 († +0.3)

2B and 7B show significant gains of +3.2 points and +2.5 points respectively ($44.05 \rightarrow 47.23$ and $71.34 \rightarrow 73.84$, respectively). Similarly, SmolLM showed the largest improvement of +3.4 points ($24.56 \rightarrow 27.98$). In contrast, the smaller Qwen models see gains of +0.5 points for the larger variant (1.5B) and +1.9 points for the smaller variant (0.5B).

SIKeD performs well on out-of-distribution datasets For the out-of-distribution datasets, there is a steady improvement on the ASDiv dataset, with Gemma 7B gaining +1.5 points (69.11 \rightarrow 70.59), +1.3 points for Gemma 2B (58.13 \rightarrow 59.39), +0.4 points for SmolLM, +1.3 points for Qwen 1.5B $(67.64 \rightarrow 68.98)$, and +0.7 points for Qwen 0.5B (57.79 \rightarrow 58.44). A similar trend is seen for the MultiArith dataset, where SmolLM shows the largest gain of +5 points. It is followed by Qwen 0.5B with +2.8 points, while other models outperform the baseline. In particular, Gemma 7B achieves a perfect score of 100. The results are similar for the SVAMP dataset, with Qwen 0.5B, Qwen 1.5B, SmolLM 1.7B, and Gemma 2B gaining +0.3, +1.1, +1.9, and +1.8 points, respectively. On the other hand, Gemma 7B maintains its baseline score of 75. Upon further analysis, we find that the SVAMP dataset tends to favor the PoT strategy, which outperforms other strategies by up to 10 points for



Figure 3: Accuracy comparison between single distillation strategies of CoT, PoT, and L2M with SIKeD biased training using the same strategy using the Gemma 7B model.



Figure 4: Iterative accuracy comparison for the Gemma 2B model across all datasets. The process is stopped when the gains diminish or when it is no longer cost effective to continue.

Gemma 7B. Although our approach often selects PoT, it does not always do so, leading to results similar to the baseline score.

Biasing SIKeD in favor of our strategy of choice For some tasks, one strategy might be better than the others due to its performance, lower cost, or better suitability for some use cases (for example, PoT is significantly better for SVAMP compared to other strategies). In such cases, it would be beneficial to bias the selection towards that strategy³. This can be done by taking only the sample from our choice of biased strategy when more than one strategy is correct from the model-generated samples. For example, if for a given data point, a smaller model samples both CoT and PoT correctly, and our biased strategy choice is PoT, we will ignore the CoT output and take only the PoT. Figure 3 compares the individual distillation strategy with the biased SIKeD. Using Gemma 7B as a smaller model across all datasets, SIKeD outperforms individual distillation strategies by a margin of 2-4 points, highlighting the effectiveness of SIKeD over other distillation approaches.

How many iterations to run for SIKeD With each iteration of SIKeD, the model learns to solve
 a task using different strategies and adjusts its strategy choice for a given task. This allows for
 continuous training of SIKeD. Figure 4 illustrates the accuracy improvements across iterations for
 the Gemma 2B model on various datasets. The iterative training is stopped when accuracy shows
 only marginal improvements or declines. Three iterations have consistently proven to be the optimal
 balance across different models and datasets in our experiments.

How the strategy distribution changes over iterations Figure 5 illustrates the strategy distribution across different iterations for the GSM8K dataset using the SmolLM 1.7B model. Iteration 0 represents the baseline "combined" training from Table 1, and as expected, the smaller model is initially biased towards one strategy (PoT in this case). Iterations 1, 2, and 3 show the model's progression using SIKeD, where it learns to diversify and select the suitable strategy for the given problem. Notably, while PoT remains the dominant strategy, the model improves its usage of the

³Note that this is different from the already biased selection of the smaller model, as our biased strategy may not be the default biased choice of the smaller model.



Figure 5: Strategy distribution over iterations for GSM8K dataset using SmolLM 1.7B model.

other two strategies—CoT and L2M—which were absent in the biased baseline. This diversification of strategies results in an overall gain of +3 points over the baseline.

449 **Training from the last checkpoint vs training from pre-trained checkpoint** In our work, we iteratively train from the last checkpoint using on-policy training as we expect continuous improve-450 ments in the model performance with a newly learned strategy. However, an alternative approach uses 451 off-policy training (training the pre-trained model at each iteration) to achieve strong performance 452 Gulcehre et al. (2023). We compared on-policy training (our proposed approach) with off-policy 453 training (as in Gulcehre et al. (2023)), utilizing both LLM-generated and self-generated data, and 454 observed a notable decrease in the overall accuracy with off-policy training. Note that we used all of 455 the LLM data at each iteration for off-policy training as the training is done on the pre-trained model. 456 On the GSM8K dataset, our on-policy approach outperformed off-policy training by +6 points (45.19 457 vs 38.90) using the Gemma 2B model. A similar trend was observed on out-of-distribution datasets, 458 where SIKeD outperformed with a gain of +4 points on both the SVAMP and ASDiv datasets (55.4 459 vs 51.3 for SVAMP and 58.96 vs 55.44 for ASDiv) and a gain of +2 points on the MultiArith dataset 460 (91.11 vs 88.33).

461

443

444 445 446

447

448

Is data mixing better than LLM-based distil-462 lation or self-distillation Mixing data from 463 the smaller model with the LLM helps in bridg-464 ing the distributional gap between the LLM and 465 the smaller model. We explore the role of the 466 data mixing rate alpha (α) to validate our hy-467 pothesis regarding data mixing. When α is set 468 to 1, only LLM-generated data is used, while at 469 the other extreme, $\alpha = 0$ means only data gener-470 ated by the smaller model is used. As shown in 471 Figure 6, an α value between 0 and 1 performs 472 better than either of the extreme cases.



473 474

Figure 6, an α value between 0 and 1 performs better than either of the extreme cases. Figure 6: Accuracy comparison with different data mix controlled by α . $\alpha = 1$ means all data is from the LLM only, while $\alpha = 0$ means only smaller model data is used.

smaller models learn to choose the right strategy over iterations The increase in the performance of smaller models can be attributed to the change in the distribution of reasoning strategies
as the model becomes more aligned with its strategy choices. Figure 7 shows that a model can solve
a given problem in the third iteration by switching the reasoning strategy which it initially struggled
with. Furthermore, Figure 11 shows an example of a model correcting its reasoning chain over time
due to its improved reasoning capabilities.

481

6 RELATED WORK

482 483

Knowledge Distillation for Reasoning Tasks Knowledge distillation (Buciluă et al., 2006; Hinton, 2015) is a widely-used technique for transferring knowledge from a large language model (LLM) to a smaller model. Previous research has focused on transferring intermediate reasoning steps to

486	Question There are tryice as me	where a girle of Dr. Wartz's school. If there	are 60 girls and 5 students to avery tageher how
487	many teachers are there?	iny boys as gins at Dr. wenz's school. If there	are of girls and 5 students to every teacher, now
488			Iteration2
489	Baseline	Iteration1	girls = 60 boys = 2 * girls
490	girls = 60 boys = 2 * girls	girls = 60 boys = 2 * girls	total_students = girls + boys
491	teachers = $girls/5$	teachers = $girls/5$	teachers_per_student = 5 total teachers = total students * teach-
492	answer = teachers	answer = teachers	ers_per_student
493	Final Answer: 12	Final Answer: 12 A	answer = total_teachers
101			Final Answer: 900 🗡
	Iteration3		
495	Sub-question 1: How many boy	ys are there?	
496	Answer to Sub-question 1: 2 *	60 = << 2 * 60 = 120 >> 120	
107	Sub-question 2: How many stu	dents are there in total?	
497	Answer to Sub-question 2: 60	+120 = << 60 + 120 = 180 >> 180	
498	Sub-question 3: How many tea	chers are there?	
499	Answer to Sub-question 3: 180 Final Answer: 36 ✓	0/5 = << 180/5 = 36 >> 36	

Figure 7: A qualitative example showing how Gemma 2B shifts its strategy selection over iterations to solve a given task.

504

514

500 501

502

486

smaller models, either step-by-step (Shridhar et al., 2023; Magister et al., 2023; Hsieh et al., 2023) 505 or by combining multiple strategies simultaneously Chenglin et al. (2023); Zhu et al. (2024); Li 506 et al. (2024). These approaches can be viewed as aggregating diverse data sources for distillation, 507 similar to the LLM data approach in our work. In contrast, Hahn & Choi (2019) and Xu et al. 508 (2020) focus on *self-distillation*, where a model improves its performance without external data or 509 knowledge. Specifically, Hahn & Choi (2019) leverages word embeddings, while Xu et al. (2020) 510 uses temporal model checkpoints as a proxy for ground truth. However, both approaches rely solely 511 on data generated by the smaller model and exclude LLM data. Our method strikes a balance between 512 these two extremes by using LLM data to learn multiple strategies and self-generated data to optimize 513 for the right strategy choice.

515 **Self-learning** Previous studies, such as (He et al., 2019; Sun et al., 2020; Gulcehre et al., 2023), 516 have shown the effectiveness of the *self-training paradigm* in tasks such as machine translation. 517 While ReST Gulcehre et al. (2023) uses off-policy training, we find on-policy training more suitable for our case both in terms of data efficieny and performance. On-policy training also allows a better 518 choice of learning strategies, since the model can use its most recent learning. Agarwal et al. (2024) 519 introduces Generalized Knowledge Distillation (GKD), an on-policy training method that aligns the 520 distributions of large language models (LLMs) and smaller models by incorporating output sequences 521 sampled from the student during training. However, the task was limited to the distribution alignment 522 and not to aligning the strategy choices in a multi-strategy distillation. Simply applying GKD would 523 not address this issue, as it would force the smaller model to learn all strategies, which is impractical 524 given its limited capacity.

Finally, we compare our distillation strategies with LLM-based distillation using both individual strategies Shridhar et al. (2023); Magister et al. (2023); Hsieh et al. (2023) and a combination of several strategies at once Chenglin et al. (2023); Zhu et al. (2024).

529 530

531

525

7 CONCLUSION

We propose SIKeD: Self-guided Iterative Knowledge Distillation, that addresses the challenge of distilling multistep reasoning skills from large language models (LLMs) to smaller models. Unlike traditional methods, which often leave smaller models biased towards a single strategy, SIKeD uses iterative *self-guided* training, combining LLM and self-generated data to improve overall reasoning in smaller models. We demonstrate our approach across various mathematical datasets and demonstrate that SIKeD improves the ability of smaller models to handle complex reasoning tasks, achieving significant performance gains.

540 REFERENCES 541

549

550

551

556

563

570

574

575

576

- Josh Achiam, Steven Adler, Sandhini Agarwal, Lama Ahmad, Ilge Akkaya, Florencia Leoni Aleman, 542 Diogo Almeida, Janko Altenschmidt, Sam Altman, Shyamal Anadkat, et al. Gpt-4 technical report. 543 *arXiv preprint arXiv:2303.08774*, 2023. 544
- Rishabh Agarwal, Nino Vieillard, Yongchao Zhou, Piotr Stanczyk, Sabela Ramos Garea, Matthieu 546 Geist, and Olivier Bachem. On-policy distillation of language models: Learning from self-547 generated mistakes. In The Twelfth International Conference on Learning Representations, 2024. 548 URL https://openreview.net/forum?id=3zKtaqxLhW.
 - Jinze Bai, Shuai Bai, Yunfei Chu, Zeyu Cui, Kai Dang, Xiaodong Deng, Yang Fan, Wenbin Ge, Yu Han, Fei Huang, et al. Qwen technical report. arXiv preprint arXiv:2309.16609, 2023.
- 552 Cristian Buciluă, Rich Caruana, and Alexandru Niculescu-Mizil. Model compression. In Proceedings 553 of the 12th ACM SIGKDD international conference on Knowledge discovery and data mining, pp. 554 535-541, Philadelphia PA USA, August 2006. ACM. ISBN 978-1-59593-339-3. doi: 10.1145/ 555 1150402.1150464. URL https://dl.acm.org/doi/10.1145/1150402.1150464.
- Wenhu Chen, Xueguang Ma, Xinyi Wang, and William W. Cohen. Program of thoughts prompting: Disentangling computation from reasoning for numerical reasoning tasks. Transactions on Machine 558 Learning Research, 2023. ISSN 2835-8856. URL https://openreview.net/forum? 559 id=YfZ4ZPt8zd. 560
- 561 Li Chenglin, Chen Qianglong, Wang Caiyu, and Zhang Yin. Mixed distillation helps smaller language 562 model better reasoning. arXiv preprint arXiv:2312.10730, 2023.
- Karl Cobbe, Vineet Kosaraju, Mohammad Bavarian, Mark Chen, Heewoo Jun, Lukasz Kaiser, 564 Matthias Plappert, Jerry Tworek, Jacob Hilton, Reiichiro Nakano, et al. Training verifiers to solve 565 math word problems. arXiv preprint arXiv:2110.14168, 2021. 566
- 567 Abhimanyu Dubey, Abhinav Jauhri, Abhinav Pandey, Abhishek Kadian, Ahmad Al-Dahle, Aiesha 568 Letman, Akhil Mathur, Alan Schelten, Amy Yang, Angela Fan, et al. The llama 3 herd of models. 569 arXiv preprint arXiv:2407.21783, 2024.
- Caglar Gulcehre, Tom Le Paine, Srivatsan Srinivasan, Ksenia Konyushkova, Lotte Weerts, Abhishek 571 Sharma, Aditya Siddhant, Alex Ahern, Miaosen Wang, Chenjie Gu, et al. Reinforced self-training 572 (rest) for language modeling. arXiv preprint arXiv:2308.08998, 2023. 573
- Isabel Gómez-Veiga, José O. Vila Chaves, Gonzalo Duque, and Juan A. García Madruga. A new look to a classic issue: Reasoning and academic achievement at secondary school. Frontiers in Psychology, 9, 2018. ISSN 1664-1078. 10.3389/ doi: fpsyg.2018.00400. URL https://www.frontiersin.org/journals/psychology/ articles/10.3389/fpsyg.2018.00400. 578
- 579 Sangchul Hahn and Heeyoul Choi. Self-Knowledge Distillation in Natural Language Processing, 580 2019. URL https://arxiv.org/abs/1908.01851. Version Number: 1. 581
- 582 Junxian He, Jiatao Gu, Jiajun Shen, and Marc'Aurelio Ranzato. Revisiting Self-Training for Neu-583 ral Sequence Generation, 2019. URL https://arxiv.org/abs/1909.13788. Version Number: 3. 584
- 585 Geoffrey Hinton. Distilling the knowledge in a neural network. arXiv preprint arXiv:1503.02531, 586 2015.
- 588 Cheng-Yu Hsieh, Chun-Liang Li, Chih-kuan Yeh, Hootan Nakhost, Yasuhisa Fujii, Alex Ratner, 589 Ranjay Krishna, Chen-Yu Lee, and Tomas Pfister. Distilling step-by-step! outperforming larger 590 language models with less training data and smaller model sizes. In Anna Rogers, Jordan Boyd-Graber, and Naoaki Okazaki (eds.), Findings of the Association for Computational Linguistics: ACL 2023, pp. 8003–8017, Toronto, Canada, July 2023. Association for Computational Linguis-592 tics. doi: 10.18653/v1/2023.findings-acl.507. URL https://aclanthology.org/2023. findings-acl.507.

594 595 596	Edward J Hu, Yelong Shen, Phillip Wallis, Zeyuan Allen-Zhu, Yuanzhi Li, Shean Wang, Lu Wang, and Weizhu Chen. Lora: Low-rank adaptation of large language models. <i>arXiv preprint arXiv:2106.09685</i> , 2021.
597 598 599	Hugging Face. smol-llm: Train a small llm from scratch. https://huggingface.co/blog/ smollm, 2023. Accessed: 2024-09-23.
600 601 602	Kushal Jain and Kumar Shridhar. First step advantage: Importance of starting right in multi-step reasoning. <i>arXiv preprint arXiv:2311.07945</i> , 2023.
603 604 605 606 607	Yoon Kim and Alexander M. Rush. Sequence-level knowledge distillation. In Jian Su, Kevin Duh, and Xavier Carreras (eds.), <i>Proceedings of the 2016 Conference on Empirical Methods in Natural Language Processing</i> , pp. 1317–1327, Austin, Texas, November 2016. Association for Computational Linguistics. doi: 10.18653/v1/D16-1139. URL https://aclanthology.org/D16-1139.
608 609	Solomon Kullback and Richard A Leibler. On information and sufficiency. <i>The annals of mathematical statistics</i> , 22(1):79–86, 1951.
611 612 613 614	Woosuk Kwon, Zhuohan Li, Siyuan Zhuang, Ying Sheng, Lianmin Zheng, Cody Hao Yu, Joseph E. Gonzalez, Hao Zhang, and Ion Stoica. Efficient memory management for large language model serving with pagedattention. In <i>Proceedings of the ACM SIGOPS 29th Symposium on Operating Systems Principles</i> , 2023.
615 616 617	Chenglin Li, Qianglong Chen, Liangyue Li, Caiyu Wang, Yicheng Li, Zulong Chen, and Yin Zhang. Mixed Distillation Helps Smaller Language Model Better Reasoning, February 2024. URL http://arxiv.org/abs/2312.10730. arXiv:2312.10730 [cs].
618 619 620 621	Qing Lyu, Kumar Shridhar, Chaitanya Malaviya, Li Zhang, Yanai Elazar, Niket Tandon, Marianna Apidianaki, Mrinmaya Sachan, and Chris Callison-Burch. Calibrating large language models with sample consistency. <i>arXiv preprint arXiv:2402.13904</i> , 2024.
622 623 624 625 626 627	Lucie Charlotte Magister, Jonathan Mallinson, Jakub Adamek, Eric Malmi, and Aliaksei Sev- eryn. Teaching small language models to reason. In Anna Rogers, Jordan Boyd-Graber, and Naoaki Okazaki (eds.), <i>Proceedings of the 61st Annual Meeting of the Association for</i> <i>Computational Linguistics (Volume 2: Short Papers)</i> , pp. 1773–1781, Toronto, Canada, July 2023. Association for Computational Linguistics. doi: 10.18653/v1/2023.acl-short.151. URL https://aclanthology.org/2023.acl-short.151.
628 629 630 631 632	Shen-yun Miao, Chao-Chun Liang, and Keh-Yih Su. A diverse corpus for evaluating and developing English math word problem solvers. In Dan Jurafsky, Joyce Chai, Natalie Schluter, and Joel Tetreault (eds.), <i>Proceedings of the 58th Annual Meeting of the Association for Computational</i> <i>Linguistics</i> , pp. 975–984, Online, July 2020. Association for Computational Linguistics. doi: 10. 18653/v1/2020.acl-main.92. URL https://aclanthology.org/2020.acl-main.92.
633 634 635	Ratna Narayan, Cynthia Rodriguez, Juan Araujo, Ali Shaqlaih, and Glenda Moss. Constructivism—constructivist learning theory. <i>IAP Information Age Publishing</i> , 2013.
636 637 638 639 640	Arkil Patel, Satwik Bhattamishra, and Navin Goyal. Are NLP models really able to solve simple math word problems? In <i>Proceedings of the 2021 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies</i> , pp. 2080–2094, Online, June 2021. Association for Computational Linguistics. doi: 10.18653/v1/2021.naacl-main. 168. URL https://aclanthology.org/2021.naacl-main.168.
641 642 643	Dean A Pomerleau. Efficient training of artificial neural networks for autonomous navigation. <i>Neural computation</i> , 3(1):88–97, 1991.
644 645 646 647	Stephane Ross and Drew Bagnell. Efficient reductions for imitation learning. In Yee Whye Teh and Mike Titterington (eds.), <i>Proceedings of the Thirteenth International Conference on Artificial Intelligence and Statistics</i> , volume 9 of <i>Proceedings of Machine Learning Research</i> , pp. 661–668, Chia Laguna Resort, Sardinia, Italy, 13–15 May 2010. PMLR. URL https://proceedings.mlr.press/v9/ross10a.html.

648 Subhro Roy and Dan Roth. Solving general arithmetic word problems. arXiv preprint 649 arXiv:1608.01413, 2016. 650

Kumar Shridhar, Jakub Macina, Mennatallah El-Assady, Tanmay Sinha, Manu Kapur, and Mrinmaya 651 Sachan. Automatic generation of socratic subquestions for teaching math word problems. In 652 Yoav Goldberg, Zornitsa Kozareva, and Yue Zhang (eds.), Proceedings of the 2022 Conference on 653 Empirical Methods in Natural Language Processing, pp. 4136–4149, Abu Dhabi, United Arab 654 Emirates, December 2022. Association for Computational Linguistics. doi: 10.18653/v1/2022. 655 emnlp-main.277. URL https://aclanthology.org/2022.emnlp-main.277. 656

657 Kumar Shridhar, Alessandro Stolfo, and Mrinmaya Sachan. Distilling reasoning capabilities into smaller language models. In Anna Rogers, Jordan Boyd-Graber, and Naoaki Okazaki (eds.), 658 Findings of the Association for Computational Linguistics: ACL 2023, pp. 7059–7073, Toronto, 659 Canada, July 2023. Association for Computational Linguistics. doi: 10.18653/v1/2023.findings-acl. 660 441. URL https://aclanthology.org/2023.findings-acl.441. 661

- 662 Haipeng Sun, Rui Wang, Kehai Chen, Masao Utiyama, Eiichiro Sumita, and Tiejun Zhao. Self-663 Training for Unsupervised Neural Machine Translation in Unbalanced Training Data Scenarios, 664 2020. URL https://arxiv.org/abs/2004.04507. Version Number: 2.
- Gemma Team, Thomas Mesnard, Cassidy Hardin, Robert Dadashi, Surya Bhupatiraju, Shreya Pathak, 666 Laurent Sifre, Morgane Rivière, Mihir Sanjay Kale, Juliette Love, et al. Gemma: Open models 667 based on gemini research and technology. arXiv preprint arXiv:2403.08295, 2024. 668
- 669 Hugo Touvron, Thibaut Lavril, Gautier Izacard, Xavier Martinet, Marie-Anne Lachaux, Timothée 670 Lacroix, Baptiste Rozière, Naman Goyal, Eric Hambro, Faisal Azhar, et al. Llama: Open and efficient foundation language models. arXiv preprint arXiv:2302.13971, 2023. 671
- 672 Unslothai. Unsloth. https://github.com/unslothai/unsloth, 2023. URL https: 673 //github.com/unslothai/unsloth. GitHub repository. 674
- Xuezhi Wang, Jason Wei, Dale Schuurmans, Quoc V Le, Ed H. Chi, Sharan Narang, Aakanksha 675 Chowdhery, and Denny Zhou. Self-consistency improves chain of thought reasoning in language 676 models. In The Eleventh International Conference on Learning Representations, 2023. URL 677 https://openreview.net/forum?id=1PL1NIMMrw. 678
- 679 Jason Wei, Xuezhi Wang, Dale Schuurmans, Maarten Bosma, Fei Xia, Ed Chi, Quoc V Le, Denny 680 Zhou, et al. Chain-of-thought prompting elicits reasoning in large language models. Advances in 681 neural information processing systems, 35:24824–24837, 2022.
- Yige Xu, Xipeng Qiu, Ligao Zhou, and Xuanjing Huang. Improving BERT Fine-Tuning via Self-683 Ensemble and Self-Distillation, 2020. URL https://arxiv.org/abs/2002.10345. Version Number: 1. 685
- 686 Denny Zhou, Nathanael Schärli, Le Hou, Jason Wei, Nathan Scales, Xuezhi Wang, Dale Schuurmans, 687 Claire Cui, Olivier Bousquet, Quoc V Le, and Ed H. Chi. Least-to-most prompting enables complex 688 reasoning in large language models. In The Eleventh International Conference on Learning Representations, 2023. URL https://openreview.net/forum?id=WZH7099tgfM. 689
 - Xunyu Zhu, Jian Li, Yong Liu, Can Ma, and Weiping Wang. Improving small language models' mathematical reasoning via mix thoughts distillation. arXiv preprint arXiv:2401.11864, 2024.
- 693

690

691

692

682

684

- 696
- 697
- 699
- 700

702 A APPENDIX 703

B tł	below is an instruction that describes a task, paired with an input that provides further context. Write a response nat appropriately completes the request.
Iı	nstruction: Solve the given math problem step by step. Put your final answer after 'Final answer:'.
lı H	nput: John buys 3 dress shirts. They sell for \$20 each. He also has to pay 10% tax on everything. Iow much did he pay in total?
R S	Response: The shirts cost $3*$20=$<<3*20=60>>60$ before tax The tax cost $60*.1=$<<60*.1=6$ 6 o in total they paid $60+6=$<<60+6=66$ 66 Final Answer: $66<\cos>$
	[7 more examples randomly sampled from the training set]
Iı a	nput: Thomas is training at the gym to prepare for a competition. He trained for 5 hours every day for month (30 days). If he continues to train for the next 12 days, how many hours will he spend on training in total?
R	esponse: Model generated response
	Figure 8: 8-shot Chain of Thought demonstration.
B tł	Below is an instruction that describes a task, paired with an input that provides further context. Write a response nat appropriately completes the request.
Iı P	nstruction: Solve the given math problem by decomposing it into smaller, manageable sub-questions. Put your final answer after 'Final answer: '.
lı H	nput: John buys 3 dress shirts. They sell for \$20 each. He also has to pay 10% tax on everything. Iow much did he pay in total?
R A S A S A F	tesponse: Sub-question 1: How much did John spend on shirts? Answer to Sub-question 1: $3*$20=$<<3*20=60>>60$ ub-question 2: How much did John spend on tax? Answer to Sub-question 2: $$60*.1=$_{i,i}60*.1=6_{i,i}6$ ub-question 3: How much did John spend in total? Answer to Sub-question 3: $$60+$6=$<<60+6=66>>66$ Tinal Answer: 66
	[7 more examples randomly sampled from the training set]
Iı a	nput: Thomas is training at the gym to prepare for a competition. He trained for 5 hours every day for month (30 days). If he continues to train for the next 12 days, how many hours will he spend on training in total?
R	Response: Model generated response
	Figure 9: 8-shot Least to Most demonstration.

Below	is an instruction that describes a task, paired with an input that provides further context. Write a response
.пат ар	propriately completes the request.
Instruc	ction: Solve the given math problem by writing a python program. Store your result as a va
able na	amed 'answer'.
Input:	John buys 3 dress shirts. They sell for \$20 each. He also has to pay 10% tax on everythi
How n	nuch did he pay in total?
Resno	nse: total shirts $= 3$
cost_o	f_{one} -shirt = 20
total_c	ost_shirts = total_shirts * cost_of_one_shirt
tax_rat	e = 0.1
tax_am	nount = tax_rate * total_cost_shirts
iotal_c	$ost = total_cost_shirts + tax_amount$
answei	$r = total_cost < eos>$
	[7 more examples randomly sampled from the training set]
input:	Thomas is training at the gym to prepare for a competition. He trained for 5 hours every day the (20 days). If he continues to training for the part 12 days hours many hours will be grand on training in tot
1 mont	in (50 days). If he continues to train for the next 12 days, now many hours will be spend on training in to
Respo [,]	nse: Model generated response
	Figure 10: 8-shot Program of Thought demonstration.
Questi	ion In a dance class of 20 students, 20% enrolled in contemporary dance, 25% of the remaining enrol
Questi n jazz	ion In a dance class of 20 students, 20% enrolled in contemporary dance, 25% of the remaining enrol dance, and the rest enrolled in hip-hop dance. What percentage of the entire students enrolled in hip-
Questi in jazz lance?	ion In a dance class of 20 students, 20% enrolled in contemporary dance, 25% of the remaining enrol dance, and the rest enrolled in hip-hop dance. What percentage of the entire students enrolled in hip-ł?
Questi in jazz dance? Baseli	ion In a dance class of 20 students, 20% enrolled in contemporary dance, 25% of the remaining enrol dance, and the rest enrolled in hip-hop dance. What percentage of the entire students enrolled in hip-l? ne
Questi in jazz dance? Baseli total_st	ion In a dance class of 20 students, 20% enrolled in contemporary dance, 25% of the remaining enrol dance, and the rest enrolled in hip-hop dance. What percentage of the entire students enrolled in hip-l? ne tudents = 20 tudents = 20
Questi in jazz dance? Baseli total_si percen	ion In a dance class of 20 students, 20% enrolled in contemporary dance, 25% of the remaining enrol c dance, and the rest enrolled in hip-hop dance. What percentage of the entire students enrolled in hip-l ? ne tudents = 20 tage_enrolled_contemporary = 0.2 ts enrolled_contemporary = total students = percentage enrolled_contemporary
Questi in jazz dance? Baseli total_st percen studen	ion In a dance class of 20 students, 20% enrolled in contemporary dance, 25% of the remaining enrol c dance, and the rest enrolled in hip-hop dance. What percentage of the entire students enrolled in hip-l? ne tudents = 20 tage_enrolled_contemporary = 0.2 ts_enrolled_contemporary = total_students * percentage_enrolled_contemporary bing_students = students = students enrolled_contemporary
Questi in jazz dance? Baseli total_si percen studen remain percen	ion In a dance class of 20 students, 20% enrolled in contemporary dance, 25% of the remaining enrol dance, and the rest enrolled in hip-hop dance. What percentage of the entire students enrolled in hip-l? ne tudents = 20 tage_enrolled_contemporary = 0.2 ts_enrolled_contemporary = total_students * percentage_enrolled_contemporary hing_students = total_students - students_enrolled_contemporary tage_enrolled_jazz = 0.25
Questi in jazz dance? Baseli total_st percen studen remain percen studen	ion In a dance class of 20 students, 20% enrolled in contemporary dance, 25% of the remaining enrol dance, and the rest enrolled in hip-hop dance. What percentage of the entire students enrolled in hip-l? ne tudents = 20 ttage_enrolled_contemporary = 0.2 ts_enrolled_contemporary = total_students * percentage_enrolled_contemporary hing_students = total_students - students_enrolled_contemporary ttage_enrolled_jazz = 0.25 ts_enrolled_jazz = remaining_students * percentage_enrolled_jazz
Questi in jazz dance Baselii total_si percen studen remain percen studen remain	ion In a dance class of 20 students, 20% enrolled in contemporary dance, 25% of the remaining enrol dance, and the rest enrolled in hip-hop dance. What percentage of the entire students enrolled in hip-l? ne tudents = 20 tage_enrolled_contemporary = 0.2 ts_enrolled_contemporary = total_students * percentage_enrolled_contemporary hing_students = total_students - students_enrolled_contemporary tage_enrolled_jazz = 0.25 ts_enrolled_jazz = remaining_students * percentage_enrolled_jazz hing_students = remaining_students - students_enrolled_jazz
Questi in jazz dance Baseli total_st percen studen remain percen studen remain percen	ion In a dance class of 20 students, 20% enrolled in contemporary dance, 25% of the remaining enrol dance, and the rest enrolled in hip-hop dance. What percentage of the entire students enrolled in hip-l? ne tudents = 20 tage_enrolled_contemporary = 0.2 ts_enrolled_contemporary = total_students * percentage_enrolled_contemporary hing_students = total_students - students_enrolled_contemporary tage_enrolled_jazz = 0.25 ts_enrolled_jazz = remaining_students * percentage_enrolled_jazz hing_students = remaining_students - students_enrolled_jazz hing_students = remaining_students - students_enrolled_jazz tage_enrolled_hip_hop = 1 - percentage_enrolled_contemporary - percentage_enrolled_jazz
Questi in jazz dance ⁵ Baseli total_si percen studen remain percen answer	ion In a dance class of 20 students, 20% enrolled in contemporary dance, 25% of the remaining enrol dance, and the rest enrolled in hip-hop dance. What percentage of the entire students enrolled in hip-l ne tudents = 20 ttage_enrolled_contemporary = 0.2 ts_enrolled_contemporary = total_students * percentage_enrolled_contemporary ning_students = total_students - students_enrolled_contemporary ttage_enrolled_jazz = 0.25 ts_enrolled_jazz = remaining_students * percentage_enrolled_jazz ning_students = remaining_students - students_enrolled_jazz tage_enrolled_hip_hop = 1 - percentage_enrolled_contemporary - percentage_enrolled_jazz r = percentage_enrolled_hip_hop
Questi in jazz dance Baseli total_si percen studen remain percen answer Final A	ion In a dance class of 20 students, 20% enrolled in contemporary dance, 25% of the remaining enrol dance, and the rest enrolled in hip-hop dance. What percentage of the entire students enrolled in hip-l ne tudents = 20 ttage_enrolled_contemporary = 0.2 ts_enrolled_contemporary = total_students * percentage_enrolled_contemporary hing_students = total_students - students_enrolled_contemporary ttage_enrolled_jazz = 0.25 ts_enrolled_jazz = remaining_students * percentage_enrolled_jazz hing_students = remaining_students - students_enrolled_jazz ttage_enrolled_hip_hop = 1 - percentage_enrolled_contemporary - percentage_enrolled_jazz r = percentage_enrolled_hip_hop Answer: 0.55 ×
Questi in jazz dance Baseli total_s percen studen remain percen answer Final <i>A</i> [terati	ion In a dance class of 20 students, 20% enrolled in contemporary dance, 25% of the remaining enrol of dance, and the rest enrolled in hip-hop dance. What percentage of the entire students enrolled in hip-le ne tudents = 20 ttage_enrolled_contemporary = 0.2 tts_enrolled_contemporary = total_students * percentage_enrolled_contemporary hing_students = total_students - students_enrolled_contemporary ttage_enrolled_jazz = 0.25 tts_enrolled_jazz = remaining_students * percentage_enrolled_jazz hing_students = remaining_students - students_enrolled_jazz ttage_enrolled_hip_hop = 1 - percentage_enrolled_contemporary - percentage_enrolled_jazz r = percentage_enrolled_hip_hop Answer: 0.55 ×
Questi in jazz dance Baseli total_s percen studen remain percen answer Final <i>A</i> [terati total_s	ion In a dance class of 20 students, 20% enrolled in contemporary dance, 25% of the remaining enrol dance, and the rest enrolled in hip-hop dance. What percentage of the entire students enrolled in hip-le ne tudents = 20 ttage_enrolled_contemporary = 0.2 ts_enrolled_contemporary = total_students * percentage_enrolled_contemporary ting_students = total_students - students_enrolled_contemporary ttage_enrolled_jazz = 0.25 ts_enrolled_jazz = remaining_students * percentage_enrolled_jazz hing_students = remaining_students - students_enrolled_jazz tage_enrolled_hip_hop = 1 - percentage_enrolled_contemporary - percentage_enrolled_jazz r = percentage_enrolled_hip_hop Answer: 0.55 × ion1 tudents = 20 tudents = 20
Questi in jazz dance ² Baseli total_s percen studen remain studen remain studen studen remain stude	<pre>ion In a dance class of 20 students, 20% enrolled in contemporary dance, 25% of the remaining enrol dance, and the rest enrolled in hip-hop dance. What percentage of the entire students enrolled in hip-le? ne tudents = 20 ttage_enrolled_contemporary = 0.2 ts_enrolled_contemporary = total_students * percentage_enrolled_contemporary hing_students = total_students - students_enrolled_contemporary ttage_enrolled_jazz = 0.25 ts_enrolled_jazz = remaining_students * percentage_enrolled_jazz hing_students = remaining_students - students_enrolled_jazz r = percentage_enrolled_hip_hop = 1 - percentage_enrolled_contemporary - percentage_enrolled_jazz r = percentage_enrolled_hip_hop Answer: 0.55 X ion 1 tudents = 20 nporary_students = total_students * 0.2 iven students = total_students * 0.2</pre>
Questi in jazz dance ⁵ Baseli total_s percen studen remain percen studen remain percen final <i>A</i> Iterati conten remain	<pre>ion In a dance class of 20 students, 20% enrolled in contemporary dance, 25% of the remaining enrol dance, and the rest enrolled in hip-hop dance. What percentage of the entire students enrolled in hip-le? ne tudents = 20 ttage_enrolled_contemporary = 0.2 ts_enrolled_contemporary = total_students * percentage_enrolled_contemporary hing_students = total_students - students_enrolled_contemporary ttage_enrolled_jazz = 0.25 ts_enrolled_jazz = remaining_students * percentage_enrolled_jazz hing_students = remaining_students - students_enrolled_jazz tage_enrolled_hip_hop = 1 - percentage_enrolled_jazz r = percentage_enrolled_hip_hop Answer: 0.55 X ion1 tudents = 20 nporary_students = total_students * 0.2 hing_students = total_students - contemporary_students weaters = remaining_students - contemporary_students weaters = remaining_students - contemporary_students weaters = remaining_students - contemporary_students weaters = remaining_students - contemporary_students</pre>
Questi in jazz dance ⁵ Baseli total_s ⁷ percen studen remain percen studen remain percen studen remain parcen final <i>A</i> Iterati total_si conten remain pazz_st	<pre>ion In a dance class of 20 students, 20% enrolled in contemporary dance, 25% of the remaining enrol dance, and the rest enrolled in hip-hop dance. What percentage of the entire students enrolled in hip-l ne tudents = 20 ttage_enrolled_contemporary = 0.2 tts_enrolled_contemporary = total_students * percentage_enrolled_contemporary hing_students = total_students - students_enrolled_contemporary ttage_enrolled_jazz = 0.25 tts_enrolled_jazz = remaining_students * percentage_enrolled_jazz ttage_enrolled_hip_hop = 1 - percentage_enrolled_jazz ttage_enrolled_hip_hop = 1 - percentage_enrolled_contemporary - percentage_enrolled_jazz r = percentage_enrolled_hip_hop Answer: 0.55 × ion1 tudents = 20 nporary_students = total_students * 0.2 ing_students = total_students - contemporary_students udents = remaining_students * 0.25 ing_students = remaining_students * 0.25 ing_students = remaining_students - 0.25 ing_students = remaining_students * 0.25 ing_students = remaining_students * 0.25 ing_students = remaining_students * 0.25 ing_students = remaining_students * 0.25 ing_students = remaining_students = 0.25 ing_students = 0.25 ing_studen</pre>
Questi in jazz dance Baseli total_s percen studen remain percen answer Final <i>A</i> Iterati total_st conten iazz_st	<pre>ion In a dance class of 20 students, 20% enrolled in contemporary dance, 25% of the remaining enrol dance, and the rest enrolled in hip-hop dance. What percentage of the entire students enrolled in hip-l ? ne tudents = 20 ttage_enrolled_contemporary = 0.2 tts_enrolled_contemporary = total_students * percentage_enrolled_contemporary hing_students = total_students - students_enrolled_contemporary ttage_enrolled_jazz = 0.25 ts_enrolled_jazz = remaining_students * percentage_enrolled_jazz ttage_enrolled_hip_hop = 1 - percentage_enrolled_jazz ttage_enrolled_hip_hop = 1 - percentage_enrolled_contemporary - percentage_enrolled_jazz r = percentage_enrolled_hip_hop Answer: 0.55 X ion1 tudents = 20 nporary_students = total_students * 0.2 ing_students = total_students * 0.2 ing_students = total_students - contemporary_students udents = remaining_students - jazz_students tage_hip_hop = hip_hop_students / total_students = 100</pre>
Questi in jazz dance Baseli total_s percen studen remain percen answei Final A Iterati total_si conten emain jazz_st uip_ho percen answei	ion In a dance class of 20 students, 20% enrolled in contemporary dance, 25% of the remaining enrol a dance, and the rest enrolled in hip-hop dance. What percentage of the entire students enrolled in hip-l ne tudents = 20 ttage_enrolled_contemporary = 0.2 ts_enrolled_contemporary = total_students * percentage_enrolled_contemporary ttage_enrolled_jazz = 0.25 ts_enrolled_jazz = remaining_students * percentage_enrolled_jazz hing_students = remaining_students * percentage_enrolled_jazz ttage_enrolled_hip_hop = 1 - percentage_enrolled_contemporary - percentage_enrolled_jazz r = percentage_enrolled_hip_hop Answer: 0.55 × ion 1 tudents = total_students * 0.2 hing_students = total_students * 0.2 hing_students = total_students * 0.2 hip_students = remaining_students - jazz_students udents = remaining_students * 0.25 inp_students = remaining_students
Questi in jazz dance Baseli total_s percen studen remain percen answei Final A Iterati total_si conten remain jazz_st nip_ho percen answei Final A	ion In a dance class of 20 students, 20% enrolled in contemporary dance, 25% of the remaining enrol a dance, and the rest enrolled in hip-hop dance. What percentage of the entire students enrolled in hip-le ne tudents = 20 ttage_enrolled_contemporary = 0.2 ts_enrolled_contemporary = total_students * percentage_enrolled_contemporary ining_students = total_students - students_enrolled_contemporary ttage_enrolled_jazz = 0.25 ts_enrolled_jazz = remaining_students * percentage_enrolled_jazz ining_students = remaining_students - students_enrolled_jazz ining_students = remaining_students - students_enrolled_jazz ining_students = remaining_students - students_enrolled_jazz tage_enrolled_hip_hop = 1 - percentage_enrolled_contemporary - percentage_enrolled_jazz r = percentage_enrolled_hip_hop Answer: 0.55 × ion 1 tudents = 20 nporary_students = total_students * 0.2 ing_students = total_students * 0.2 ing_students = total_students * 0.25 p_students = total_students - contemporary_students udents = remaining_students * 0.25 inp_students = remaining_students * 0.25 inp_students = total_students - jazz_students tage_hip_hop = hip_hop_students/total_students * 100 r = percentage_hip_hop Answer: 60 ✓
Questi in jazz dance Baseli total_s percen studen remain percen answer Final A Iterati total_si conten remain jazz_st nip_ho percen answer Final A	ion In a dance class of 20 students, 20% enrolled in contemporary dance, 25% of the remaining enrol i dance, and the rest enrolled in hip-hop dance. What percentage of the entire students enrolled in hip-hop tudents = 20 tage_enrolled_contemporary = 0.2 ts_enrolled_contemporary = total_students * percentage_enrolled_contemporary ining_students = total_students - students_enrolled_contemporary tage_enrolled_jazz = 0.25 ts_enrolled_jazz = remaining_students * percentage_enrolled_jazz tage_enrolled_hip_hop = 1 - percentage_enrolled_jazz tage_enrolled_hip_hop = 1 - percentage_enrolled_contemporary - percentage_enrolled_jazz r = percentage_enrolled_hip_hop Answer: 0.55 ★ ion1 tudents = 20 nporary_students = total_students * 0.2 ting_students = total_students * 0.2 ting_students = total_students * 0.2 p_students = remaining_students * 0.25 p_students = remaining_students - jazz_students tage_hip_hop = hip_hop_students/total_students * 100 r = percentage_hip_hop

Figure 11: Over iteration, Gemma 2B model learnt to fix its mistake using the same strategy