# SMART: SELF-LEARNING META-STRATEGY AGENT FOR REASONING TASKS

Anonymous authors

Paper under double-blind review

#### ABSTRACT

Tasks requiring deductive reasoning, especially those involving multiple steps, often demand adaptive strategies such as intermediate generation of rationales or programs, as no single approach is universally optimal. While Language Models (LMs) can enhance their outputs through iterative self-refinement and strategy adjustments, they frequently fail to apply the most effective strategy in their first attempt. This inefficiency raises the question: Can LMs learn to select the optimal strategy in the first attempt, without a need for refinement? To address this challenge, we introduce SMART (Self-learning Meta-strategy Agent for Reasoning Tasks), a novel framework that enables LMs to autonomously learn and select the most effective strategies for various reasoning tasks. We model the strategy selection process as a Markov Decision Process and leverage reinforcement learningdriven continuous self-improvement to allow the model to find the suitable strategy to solve a given task. Unlike traditional self-refinement methods that rely on multiple inference passes or external feedback, SMART allows an LM to internalize the outcomes of its own reasoning processes and adjust its strategy accordingly, aiming for correct solutions on the first attempt. Our experiments across various reasoning datasets and with different model architectures demonstrate that SMART significantly enhances the ability of models to choose optimal strategies without external guidance (+15 points on the GSM8K dataset). By achieving higher accuracy with a single inference pass, SMART not only improves performance but also reduces computational costs for refinement-based strategies, paving the way for more efficient and intelligent reasoning in LMs.

031 032 033

034

035

004

010 011

012

013

014

015

016

017

018

019

021

023

024

025

026

027

028

029

#### 1 INTRODUCTION

When people first encounter complex reasoning tasks, such as solving mathematical problems, they
often make mistakes or approach them inefficiently (Brown et al., 2019). However, with experience,
humans tend to improve their performance by replacing ineffective or incorrect strategies with more
effective ones, using a mix of strategies tailored to the specific task (Adolph et al., 1998; Lemaire &
Callies, 2009; Torbeyns et al., 2009; Boncoddo et al., 2010, *inter alia*).

Language Models (LMs) similarly struggle with reasoning tasks, sometimes producing incoherent results (Madaan et al., 2023; Shridhar et al., 2023; Huang et al., 2024). A common remedy is to resample the output, a process known as *refinement*. This refinement may involve reusing the same reasoning approach (Madaan et al., 2023) or adopting an entirely new one (Shridhar et al., 2023). In addition, providing feedback on initial results has proven beneficial during resampling (Huang et al., 2024; Welleck et al., 2022; Shinn et al., 2024; Kim et al., 2024, *inter alia*). This raises a critical question: *Can LMs be taught to optimize their choice of reasoning strategy for specific tasks overtime on the first trial, much like humans do?*

To address this question, we propose a novel framework called SMART (Self-learning Meta-strategy
 Agent for Reasoning Tasks), which allows LMs to learn optimal strategy selection through a con tinuous self-learning approach. We model the task of identifying the optimal strategy as a *Markov Decision Process* (MDP) (Sutton et al., 1999; Puterman, 2014), where the agent (LM) starts with
 its pre-trained knowledge and iteratively improves its performance by learning from its own outputs
 and strategy choices. By integrating the LM's reasoning abilities with reinforcement learning-driven



073

074 Figure 1: **Our proposed methodology**: In the first step (initial sampling), an agent (LM) chooses a 075 strategy and solves the given task with it. If it is correct, the process ends successfully. If an incorrect 076 strategy is chosen, the agent iteratively refines its strategy, taking previous strategies into account. 077 The process stops when a correct strategy is chosen to solve a task, or when a stopping criterion 078 such as the number of attempts is reached. All correct strategies are used to further refine the model, and the process is repeated. During testing, we sample once from  $LM_t$  without refinement. 079

- 081

082 self-improvement, the agent can simulate different reasoning strategies, evaluate their effectiveness 083 based on past outcomes, and adjust its strategy choice accordingly.

084 Our approach differs from traditional methods by focusing on iterative reward-based learning, which 085 encourages the agent to produce the correct inference on the first attempt without resampling. This not only improves cost efficiency - only one sampling step is required during inference - but also 087 results in a more generalizable model capable of adapting its strategy selection based on the specific 880 task. We validate SMART on a variety of reasoning datasets and LM architectures and show that our method significantly improves the ability of LMs to select optimal strategies on the first try, outperforming baseline models that rely on traditional self-refinement techniques in both accuracy 090 and computational efficiency. On three mathematical datasets (GSM8K Cobbe et al. (2021), SVAMP 091 Patel et al. (2021b), ASDiv Miao et al. (2020)) over three LLM agents (Llama3 8B Dubey et al. 092 (2024), Gemma 7B Team et al. (2024), and Mistral 7B Jiang et al. (2023)), we demonstrate the effectiveness of our approach. On iterative refinement with SMART, we achieve gains of up to +15 094 points (a relative gain of +35%) on the GSM8K dataset without the need for refinement. In addition, 095 we improve refinement accuracy by +16 points over baselines.

- 096 097
- 2 METHODOLOGY
- 098 099

100 Let q be a problem best solved with multi-step reasoning. An example is presented in the form 101 of a mathematical word problem in Figure 1. An agent or language model (LM) can approach it 102 using various strategies, such as solving it step by step (Chain of Thought, CoT Wei et al. (2022)), 103 decomposing it into subproblems, and solving each one (Least to Most, L2M Zhou et al. (2023)), 104 or writing a program to solve it programmatically (Program of Thought, PoT Chen et al. (2023)), 105 among others. A common method is to prompt the LM with a specific strategy for solving the task. However, LMs are error-prone but can fix their answers when asked to do so, a process called 106 refinement. In the refinement process, LMs can either stick to the same strategy Madaan et al. 107 (2023) or switch to a more effective reasoning strategy Shridhar et al. (2023). Ideally, LMs could

learn to choose the best strategy and reasoning path on the first try, minimizing the need for costly refinement.

Our objective: The primary goal of our work is to enable language models (LMs) to autonomously learn and select the most effective strategies for various reasoning tasks on their first attempt, thereby improving both efficiency and accuracy. Unlike traditional self-refinement methods that require multiple inference passes or external feedback, our approach aims to internalize the learning process within the LM, allowing it to adjust its strategy selection based on past experience. This mirrors how humans learn to choose optimal strategies through experience when faced with complex tasks.

- 117
- 118 119

#### 2.1 SMART: SELF-LEARNING META-STRATEGY AGENT FOR REASONING TASKS

We model the strategy selection process as a *Markov Decision Process* (MDP), where the LM acts as an agent that interacts with the environment (the reasoning tasks) by selecting strategies and observing the outcomes. Using reinforcement learning techniques, the LM can learn a policy that maximizes the expected reward, effectively learning to choose the optimal strategy for each task. This framework allows the LM to simulate different reasoning strategies, evaluate their effectiveness based on past outcomes, and adjust its strategy choice accordingly.

In the following sections, we formalize the problem formulation, define the agent's policy, and describe the learning objective. We then present our two-stage process with iterative refinement, which allows the agent to learn from its own reasoning processes and improve its strategy choice over time.

130 **MDP Setup:** We model the strategy selection framework as a Markov Decision Process (MDP) 131 given by the tuple  $\langle S, A, P, \mathcal{R}, \mu \rangle$ . Here, S represents the state space encapsulating all possible states of the environment, where each state  $s \in S$  is the current problem statement or the subsequent 132 LM response to it. The initial state distribution is given by  $\mu$ . The action space,  $\mathcal{A}$ , is the set of all 133 strategies available to the agent, where each action  $a_t \in \mathcal{A}$  corresponds to the choice of a particular 134 strategy at time t. The transition function  $\mathcal{P}: \mathcal{S} \times \mathcal{A} \to \mathcal{S}$  defines the probability of transitioning 135 to a new state  $s_{t+1}$  after applying strategy  $a_t$  in state  $s_t$ . Particularly for our case, the transition 136 function is non-deterministic as the next state is sampled by the agent (see Algorithm 1 later). The 137 reward function  $\mathcal{R}: \mathcal{S} \times \mathcal{A} \to \mathbb{R}$  assigns a scalar reward based on the correctness of the result after 138 applying the chosen strategy. 139

We start with the initial sampling step, where the LM chooses a strategy to solve the given task. In 140 other words, given a problem statement  $s_1 \sim \mu$ , the agent draws a strategy  $a_1$  to solve the problem 141 according to its policy (see below). Given the initial problem  $s_1$  and strategy  $a_1$ , the environment 142 transitions to the next state  $s_2 \sim \mathcal{P}(\cdot|s_1, a_1)$  with transition probability  $\mathcal{P}$  and receives a reward  $r_1$ 143 indicating the correctness of the response. If a correct strategy is chosen and the LM solves the task 144 using that strategy, the process terminates. Otherwise, the agent chooses another strategy  $a_2$  to solve 145 the problem, and the process repeats until the problem is solved. A realization of this stochastic process is a trajectory  $\tau = \left( (s_t, a_t, r_t)_{t=1}^{T-1}, s_T \right)$  up to a finite number of steps T. We denote a partial trajectory (history up to time t) as  $\tau_{1:t} = \left( (s_t, a_t, r_t)_{t=1}^{t-1}, s_t \right)$ . 146 147 148

**Agent's policy:** We start by defining a history  $h_t := \tau_{1:t}$  that includes the past actions up to time t - 1 and the observed states up to time t. Then, we model the agent using a non-Markovian stochastic policy  $\pi_{\theta}(a_t|h_t)$  parameterized by  $\theta$ , which models the probability of choosing the action  $a_t$  given the history  $h_t$ :

$$\pi_{\theta}(a_t|h_t) = \Pr(a_t = a \mid h_t; \theta)$$

**Objective function:** We want to optimize the following objective:

$$\theta^{\star} = \arg\max_{\theta} J(\pi_{\theta}) = \arg\max_{\theta} \mathbb{E}_{s_1, a_1 \sim \mu_{\text{test}}, \pi_{\theta}(a_1|s_1)}[r_1(s_1, a_1)] \tag{1}$$

159 160 161

154

156

157 158

where,  $\mu_{\text{test}}$  represents the state distribution over the test data.

162 **TWO-STAGE PROCESS WITH REFINEMENT** 163 164 STAGE 1 (INITIAL SAMPLING) 1. The process begins with a problem statement  $s_1 \sim \mu$  where  $\mu$  represents the initial state 166 distribution. 167 2. The agent samples an action  $a_1$  from the policy  $\pi_{\theta}(a_1|h_1)$ , where  $h_1 = s_1$  for the step 1. 3. The agent then generates an output based on strategy  $a_1$ 169 170 4. The agent receives reward  $r_1(h_1, a_1)$  based on it correctness as follows: 171  $r_1 = \begin{cases} 1, & \text{if the output is correct} \\ 0, & \text{otherwise} \end{cases}$ 172 173 174 5. Termination Check: If  $r_1 = 1$ , the process terminates successfully. 175 176 STAGE 2 (ITERATIVE REFINEMENT, IF  $r_1 = 0$ ) 177 178 For each time step t = 2 to T: 179 1. The agent observes history  $h_t$ . 181 2. Samples an action  $a_t \sim \pi_{\theta}(a_t | h_t)$ . 3. Generates an output based on strategy  $a_t$ . 183 4. Receives reward  $r_t(s_t, a_t)$ : 185  $r_t = \begin{cases} 1, & \text{if the output is correct} \\ 0, & \text{otherwise} \end{cases}$ 186 187 5. Termination Check: If  $r_t = 1$ , terminate the process. 188 189 6. Otherwise, transition to the next state  $s_{t+1}$ , which includes the reasoning choice from all 190 the previous steps. 191 **Trajectory and Reward Structure:** The resulting trajectory  $\tau$  with reward has the following struc-192 ture: 193  $\tau = ((s_1, a_1, r_1), (s_2, a_2, r_2), \dots, (s_{T'}, a_{T'}, r_{T'}), s_{T'+1})$ 194 where T' < T is the time step at which the process terminates (either by solving the problem or by reaching the maximum steps). We keep the number of trajectories one less than the number of 196 strategies in our work, since each time the model samples with a different strategy than the one 197 present in its history. For simplicity, we omit  $r_t = 0$  rewards in the trajectory above. The total reward for the trajectory is: 199  $R(\tau) = \sum_{t=1}^{T'} r_t$ 200 201 202 Given that  $r_t = 0$  for t < T' and  $r_{T'} = 1$ , if successful, total reward would be  $R(\tau) = 1$ . However, 203 in some cases, when a task is never solved, the total reward for that trajectory can be 0. 204 205 As typical in RL, the agent aims to learn a policy that maximizes the expected cumulative reward 206 (Sutton et al., 1999; Prajapat et al., 2024): 207  $\theta^{\star} = \arg\max_{\theta} J(\pi_{\theta}) = \arg\max_{\theta} \mathbb{E}_{\tau \sim f(\tau; \pi_{\theta})}[R(\tau)]$ (2)208 209 POLICY GRADIENT UPDATE 210 211 The gradient of the expected reward with respect to the policy parameters  $\theta$  is: 212 213 EL\_L 1 ٦

214  
215 
$$\nabla_{\theta} J(\pi_{\theta}) = \mathbb{E}_{\tau} \left[ \sum_{t=1}^{|\tau|-1} \nabla_{\theta} \log \pi_{\theta}(a_t | h_t) \cdot r_t \right]$$
(3)

The gradient is computed iteratively, and a step is only taken at time t only if all the previous attempts were incorrect.

**Implicit bias for Stage 1 to choose the right strategy in the first attempt:** As described in equation 1, our goal is to maximize the reward received as early as possible. To implicitly bias the model towards selecting the correct action in earlier steps, we adjust the dataset  $\mathcal{D}$  by replacing the sequences of unsuccessful actions with the final successful action taken at the initial state. Specifically, for trajectories where the problem is solved at time T' (T' > 1), we replace the samples ( $h_1, a_1, r_1$ ), ..., ( $h_{T'}, a_{T'}, r_{T'}$ ) with ( $h_1, a_{T'}, r_{T'}$ ). This encourages the model to learn to take the correct strategy  $a_{T'}$  for the problem statement  $s_1$  in the first step.

Since we only update the model based on correct outputs, the policy update can be viewed as maximizing the likelihood of the correct actions given the history:

$$\theta^{\star} = \arg \max_{\theta} \sum_{(h_i, a_i) \in D} \log \pi_{\theta}(a_i | h_i) \tag{4}$$

See 1 for a complete algorithm that illustrates our methodology in detail.

## ALGORITHM: SMART: SELF-LEARNING META-STRATEGY AGENT FOR REASONING TASKS

234 235 236

237

226

227 228

229 230

231 232 233

Algorithm 1 Training Procedure for SMART

**Require:** Initialized policy parameters  $\theta$ , learning rate  $\alpha$ , dataset of problems  $\mathcal{D}$ . 238 1: for each iteration e = 1, 2, ... do  $\triangleright$  For a fixed number of iterations or until convergence 239 2:  $\mathcal{D}_e \leftarrow \emptyset$ ▷ Initialize empty dataset 240 3: for each problem  $s_1 \in \mathcal{D}$  do 241 4: Stage 1: Initial Attempt 242 5: Observe initial state  $s_1$ 243 6: Sample action  $a_1 \sim \pi_{\theta}(a \mid s_1)$ 244 7: Generate output using strategy  $a_1$ Evaluate output to obtain reward  $r_1$ 8: 245 9: if  $r_1 = 1$  then ▷ Correct output 246 Collect sample  $(s_1, a_1, r_1)$ 10: 247  $\mathcal{D}_e \leftarrow \mathcal{D}_e \cup \{(s_1, a_1)\}$  $\triangleright$  Add sample to dataset 11: 248 Continue to next problem 12: 249 13: else 250 Stage 2: Iterative Refinement 14: 251 for each refinement iteration  $t = 2 \dots T$  do 15: 16: Observe state  $s_t$  and history  $h_t$ 252 Sample action  $a_t \sim \pi_{\theta}(a|h_t)$ . 17: 253 Generate output using strategy  $a_t$ 18: 254 Evaluate output to obtain reward  $r_t$ 19: 255 20: if  $r_t = 1$  then ▷ Correct output 256 21: Collect sample  $(s_1, a_1, \ldots, s_t, a_t, r_t)$ 257 22:  $\mathcal{D}_e \leftarrow \mathcal{D}_e \cup \{(s_1, a_1, \dots s_t, a_t)\}$  $\triangleright$  Add sample to dataset 258 Break loop and proceed to next problem 23: 259 24: end if end for 25: 260 26: end if 261 end for 27: 262 28: **Policy Update** 263 29: Update policy parameters  $\theta$  using collected data: 264  $\theta \leftarrow \theta + \alpha \cdot \sum_{(h_i, a_i) \in \mathcal{D}_e} \nabla_{\theta} \log \pi_{\theta}(a_i | h_i)$ 265 266 267 **Implicit Biasing** 30: 268  $\mathcal{D}_{e+1} \leftarrow \mathcal{D}_e \setminus \{(s_1, a_1, \dots, s_t, a_t)\} \bigcup \{(s_1, a_t)\}$ 31:  $\triangleright$  Update the dataset 269 32: end for

# 270 3 EXPERIMENTAL DETAILS

271 272

SMART operates within a *self-learn* framework, where we prompt a pre-trained model with 8-shot 273 examples to collect the initial training data. Prompts used for various strategies are provided in 274 Subsection 8.1. The 8-shot pre-trained model also serves as a baseline comparison for our method. 275 We use the MetaMath dataset Yu et al. (2024), a variant of the GSM8K Cobbe et al. (2021) training 276 set, which contains 110K reasoning problems. We evaluate our methodology on the GSM8K test set, which contains 1,319 samples. To show the generalization ability of our method, we also test on 278 two out-of-distribution datasets: the SVAMP dataset Patel et al. (2021a) with 1,000 samples where the questions are made more challenging by altering them in a non-trivial way, and the ASDiv 279 dataset Miao et al. (2020) with 2,300 samples, which consists of diverse mathematical problems 280 from elementary to middle school. 281

We ran all our experiments on models with 7-8 billion parameters, specifically Gemma 7B Team et al. (2024), Mistral 7B Jiang et al. (2023), Qwen2 7B Yang et al. (2024), and Llama3 8B Dubey et al. (2024). This is important for our *self-learn* setup, as the model needs a basic understanding of the task to begin the process. Smaller models often have difficulty starting the process, while very large models are too expensive to iterate over multiple times.

287 We employed three reasoning strategies: Chain of Thought (CoT) (Wei et al., 2022), Least to Most 288 (L2M) (Zhou et al., 2023), and Program of Thought (PoT) (Chen et al., 2023) in our work due to 289 their effectiveness on the multi-step reasoning tasks. We take the best out of these strategies as our baseline (underlined in Table 1). Since we used a different reasoning strategy during refinement 290 following Shridhar et al. (2023), the majority of our experiments are limited to two trajectories as the 291 third trajectory would make the remaining strategy the obvious choice. However, we also explore 292 our approach beyond 3 strategies later (in Section 5). We report the top-1 accuracy (maj@1) for all 293 experiments, as we want to test the accuracy of the output on the first try. We used a temperature of 294 0.7 to generate samples at each iteration. Since the models are already trained on the downstream 295 datasets during pre-training, we chose to train only part of the model and used LoRA (Hu et al., 296 2021) with rank 16, alpha 32, and a starting learning rate of 2e-4 with decay. We used the Unsloth 297 library (Unslothai, 2023) for training and we did the inference using the VLLM library (Kwon et al., 298 2023). Finally, we run SMART for multiple iterations until our accuracy gains no longer justify 299 the cost of training and sampling. In our experiments, 3 iterations worked well for most models, 300 although we trained Gemma 7B for 5 iterations before performance plateaued.

## 4 Results

304 SMART significantly improves results on in-distribution dataset: We compared SMART with 305 baselines on the GSM8K dataset, which we also consider to be an in-distribution dataset since the 306 train and test sets have the same distribution. Table 1 shows that SMART outperformed the baseline 307 in its first iteration on the GSM8K dataset, achieving a gain of +6 points for both the Gemma 7B and 308 Mistral 7B models ( $40.4 \rightarrow 46.5$  and  $56.9 \rightarrow 63.8$ , respectively). Although Qwen2's performance is 309 already very strong on the GSM8K dataset, we still observed a gain of +2.6 points (81.9  $\rightarrow$  84.5) in 310 the first iteration. After a few more iterations, we saw a total gain of +15 points for Gemma 7B (40.4  $\rightarrow$  55.4), +11 points for Mistral 7B (56.9  $\rightarrow$  67.9), and +4 points for Qwen2 7B (81.9  $\rightarrow$  85.4). 311

312 **SMART serves as a great refinement strategy:** Since SMART involves iterative refinement to find 313 the optimal action given the trajectory, we also compare against two refinement baselines: refine-314 ment with the same strategy (Madaan et al., 2023) and refinement with a strategy change (Shridhar 315 et al., 2023). We used the Oracle Verifier, which identifies incorrect samples and refines them either 316 with the same strategy or by choosing a different one. <sup>1</sup> Table 1 compares the refinement accuracy with SMART and our proposed methodology shows significant improvements over the baselines. 317 Gemma 7B gains over +16 points (48.9  $\rightarrow$  67.5) compared to the best refinement baseline, Mistral 318 7B gains +8 points (66.5  $\rightarrow$  78.0), and Qwen2 7B gains +1.5 points (86.9  $\rightarrow$  91.9). 319

320 SMART generalizes well to out-of-distribution dataset: We test our trained checkpoints using
 321 our proposed approach SMART on two out-of-distribution datasets: ASDiv and SVAMP. These are

301 302

303

<sup>1</sup>Note that the Oracle Verifier is used to set the upper bound for all the methods compared and cannot be

<sup>322</sup> 323

<sup>&</sup>lt;sup>1</sup>Note that the Oracle Verifier is used to set the upper bound for all the methods compared and cannot be used during inference. It is only to compare the capabilities of different approaches.

#### 324

Table 1: Test Accuracy (maj@1) comparison between different baselines using three strategies CoT, L2M, PoT with our approach SMART on the GSM8K dataset. Baselines used 8-shot in-context examples to generate the output. Additionally, we report refinement accuracy obtained through the Oracle verifier for both the baselines and our approach. Refinements for the baselines can follow the *same* strategy as in Madaan et al. (2023) and a *different* strategy than the initial one as in Shridhar et al. (2023). Results are presented for three models: Gemma 7B, Mistral 7B, and Qwen2 7B. The best results among the baseline are <u>underlined</u> while the best overall results are in **bold**.

Model	Mathad	Test Accuracy (%)	Refinement	
WIGUEI	Wethou	Test Accuracy (70)	Strategy	Accuracy (%)
	Baseline (Using 8-shot examples)			
	Chain of Thought (CoT)	40.0	same	44.6
	chain of Thought (CoT)	10.0	different	49.4
	Least to Most (L2M)	34.9	same	43.4
	()		different	48.7
	Program of Thought (PoT)	40.4	same	46.1
C			different	<u>51.3</u>
Gemma /B	CMADT (Proposed Approach)			
	Iteration 1	16.5	CMADT	64.6
	Iteration 2	40.5	SMARI	64.0
	Final Iteration - Iteration 5	<b>55 6</b> († 115 2)	SMART	67 5 († 116 2)
	That herailon - herailon 5	<b>33.0</b> (  +13.2)	SMART	07.3 (  +10.2)
	<b>Baseline</b> (Using 8 shot argmplas)			
	Basenne (Using 0-snot examples)		same	59.0
	Chain of Thought (CoT)	50.6	different	67.5
			same	56.6
	Least to Most (L2M)	52.4	different	67.2
		56.0	same	61.3
	Program of Thought (PoT)	<u>56.9</u>	different	70.1
Mistral 7B				
	SMART (Proposed Approach)			
	Iteration 1	63.8	SMART	74.1
	Iteration 2	66.3	SMART	76.4
	Final Iteration - Iteration 3	<b>67.9</b> († +11.0)	SMART	<b>78.0</b> († +7.9)
	Baseline (Using 8-shot examples)			
	Chain of Thought (CoT)	81.9	same	90.4
	chain of Thought (CoT)	01.9	different	89.3
	Least to Most (L2M)	80.0	same	84.9
			different	88.7
	Program of Thought (PoT)	76.9	same	86.0
Ouron 2 7D			different	88.6
Qwen2 /B	CMADT (Droposed Approach)			
	Iteration 1	84.5	CMADT	01.4
	Iteration 2	04.J 85 1 (4. 25)	SMAKI	91.4 01 0 (4.1.5
	Final Iteration - Iteration 3	<b>03.4</b> († +3.5) <b>85 2</b>	SMARI	01 2
	rinul heration - heration 5	03.2	SMARI	91.2

366

367 368

369 referred to as out-of-distribution because the model was not trained on these datasets but only eval-370 uated on them. Table 2 compares the results of SMART with the baselines (created using 8-shot in 371 context examples) across the three reasoning strategies: CoT, L2M, and PoT. Among the three mod-372 els, Gemma 7B showed the most improvement, with a gain of +2.6 points on the SVAMP dataset 373  $(65.6 \rightarrow 68.2)$  and +2.9 points on the ASDiv dataset  $(68.0 \rightarrow 70.9)$ . Mistral 7B, on the other hand, 374 gained +1 point on the ASDiv dataset (71.3  $\rightarrow$  72.3) but performed worse on the SVAMP dataset. We suspect that the test dataset for SVAMP is different from GSM8K, and since the model was not 375 optimized for SVAMP-style questions, this could explain the drop in performance. Finally, for the 376 Qwen2 7B model, we observed a modest gain of +1 point across both datasets (89.3  $\rightarrow$  90.3 on 377 SVAMP and  $89.5 \rightarrow 90.5$  on ASDiv).

#### 378

Table 2: Test accuracy (maj@1) comparison between different baselines using three strategies CoT, L2M, PoT with our proposed approach SMART on the out-of-domain datasets of SVAMP and ASDiv. Baselines used 8-shot in-context examples to generate the output. Results are presented for three models: Gemma 7B, Mistral 7B, and Qwen2 7B. The best baseline results are <u>underlined</u>, while the best overall results are in **bold**.

Model	Method	SVAMP Acc (in %)	ASDiv Acc (in %)
	Baseline (Using 8-shot examples)		
	Chain of Thought (CoT)	<u>65.6</u>	<u>68.0</u>
	Least to Most (L2M)	63.4	62.7
Gemma 7B	Program of Thought (PoT)	51.6	54.2
	SMART (Proposed Approach)		
	Final Iteration	<b>68.2</b> († +2.6)	<b>70.9</b> († +2.9)
	Baseline (Using 8-shot examples)		
	Chain of Thought (CoT)	65.8	71.3
	Least to Most (L2M)	67.9	67.3
Mistral 7B	Program of Thought (PoT)	<u>71.7</u>	70.1
	SMART (Proposed Approach)		
	Final Iteration	70.9 (1 -0.8)	<b>72.3</b> (↑ +1.0)
	Baseline (Using 8-shot examples)		
	Chain of Thought (CoT)	<u>89.3</u>	<u>89.5</u>
	Least to Most (L2M)	87.5	84.2
Qwen2 7B	Program of Thought (PoT)	82.8	80.8
	SMART (Proposed Approach)		
	Final Iteration	<b>90.3</b> († +1.0)	<b>90.5</b> († +1.0)

#### 5 DISCUSSION

406 407

396 397

How the strategy distribution changes over iterations: With SMART, the goal is to select the 408 desired strategy at the first attempt, i.e., over iterations the LM should learn to select the appropriate 409 strategy for a given task. Figure 2 shows the changes in strategy distribution over iterations for the 410 Gemma 7B model on the GSM8K dataset. As indicated by the baseline in Table 1, PoT emerges as 411 the best strategy for the Gemma 7B model, followed by CoT, with L2M being the least effective. 412 A similar pattern is observed in Figure 2, where the model increasingly favors PoT (indicated by 413 the upward trend in the red line) while decreasing its preference for the other two strategies. Cor-414 respondingly, the accuracy for PoT improves the most, followed by CoT and L2M, demonstrating 415 that over iterations the model learns to select the optimal strategy for a given task. A qualitative 416 example where over iterations the model learned to choose the right strategy in its first attempt is provided in Figure 4. Initially, the model chose the wrong strategy but was fixed during refinement, 417 and over iterations, the model picked the correct strategy in its initial sampling, without the need for 418 refinement. 419

420 Extending the strategies beyond three strategies in SMART: We want to test if it is possible to 421 go beyond the three strategies explored in this paper by extending our approach to other strategies. 422 Although we could not find a strategy with performance comparable to those explored in this paper, 423 we introduced a <Unsolvable> tag for questions that the model could not answer during either sampling or refinement attempts and tried to use it as a fourth strategy. However, extending this strat-424 egy did not yield promising results due to the limited number of samples in the <Unsolvable> 425 category, as the model was able to solve the task using one of the strategies during refinement. Sim-426 ilarly, we experimented with a <Answer Only> strategy, where the model directly predicts the 427 answer to very simple questions without any intermediate reasoning. As with the <Unsolvable> 428 strategy, the skewed data distribution led to poorer performance. 429

Is the strategy selection effective?: We compare the strategy selection made by SMART during inference with a fixed strategy (CoT) on the GSM8K dataset using the Gemma 7B model, as shown in Table 3. Across all five iterations, the results highlight the importance of selecting the appropriate



Figure 2: Strategy distribution change over iterations for Gemma 7B model on GSM8K dataset.

Table 3: Comparison of fixed strategy (CoT) vs SMART based strategy during inference for GSM8K dataset using Gemma 7B model.

Strategy	<b>Iteration 1</b>	Iteration 2	Iteration 3	<b>Iteration 4</b>	<b>Iteration 5</b>
Fixed (CoT)	45.5	47.5	48.5	50.6	49.7
SMART	46.5	50.6	52.9	55.0	55.6

strategy, with the SMART approach outperforming the fixed strategy (CoT) by up to 6 points when the optimal strategy is selected. While the fixed strategy shows improvements over iterations, illustrating the *self-learning* effect over multiple generations, our work demonstrates that performance can be further improved by selecting the most appropriate strategy for each task.

# 460 SMART with better starting samples:

Since we start with the model-generated 462 samples, for a weaker model the starting 463 samples can be improved if those samples 464 come from the stronger model. We in-465 vestigated whether the SMART based self-466 *learning* approach could be extended to 467 a setup where the initial data points are 468 collected from a stronger model to initi-469 ate training for a weaker model. This is 470 particularly beneficial for weaker models that cannot independently initiate the self-471 learning process using SMART due to their 472 limited capabilities on a given task. Ini-473 tially, we collected data using 8-shot in-474 context examples generated by the Llama3 475 8B model (Dubey et al., 2024). This 476 approach yielded an average accuracy of 477



Figure 3: Figure showing a comparison of the effects of different starting data points for the Gemma 7B model. SMART is compared against two baselines: the pre-trained Gemma 7B model and the fine-tuned Gemma 7B model on Llama3 8B data.

about 80% across all three strategies, significantly higher than Gemma 7B's initial accuracy of about 478 40%. Figure 3 illustrates the comparison between different iterations of SMART and the pre-trained 479 Gemma 7B baseline at 40.4%, as well as the fine-tuned Gemma 7B baseline at 68.7%. Using SMART 480 resulted in an additional improvement of +3.1 points (from  $68.7\% \rightarrow 71.8\%$ ).

481 482

445

446 447

448

457

458

459

- 6 RELATED WORK
- 483 484

**Refinement in LLMs:** Refinement refers to the process of improving the initial output of large language models (LLMs) through iterative adjustments. This refinement can be achieved by following

9

486 487 Question Brian's friend Bobby has 5 fewer than 3 times as many video games as Brian does. If Brian has 20 video games but lost 5 right before the comparison was made, how many does Bobby have? 488 Iteration 1 Stage 1 (Sampling) 489 Program of Thought:  $brian_games = 20$ 490  $games_lost = 5$ 491 bobby\_games = 3 \* brian\_games - games\_lost - 5 492 answer = bobby\_games X 493 Iteration 1 Stage 2 (Refinement) 494 Chain of Thought: Brian has 20 video games and lost 5 games, so he has 20 - 5 = 15 games left. 495 Bobby has 5 fewer than 3 times as many video games as Brian does, so he has 3 \* 15 - 5 = 40 video games. Final Answer: 40 🗸 496 Iteration 5 Stage 1 (Sampling) 497 Chain of Thought: Brian has 20 - 5 = 15 video games. 498 Bobby has 5 fewer than 3 times as many video games as Brian does, so Bobby has 3 \* 15 - 5 = 40 video games. 499 Final Answer: 40 🗸 500 501

Figure 4: Qualitative example demonstrating that Gemma 7B model learnt to refinement strategy in its initial sampling stage, removing the need for refinement.

504

505 the same method used initially (Madaan et al., 2023), by incorporating feedback while using the 506 same approach (Welleck et al., 2022; Shinn et al., 2024; Kim et al., 2024; Huang et al., 2024), or by 507 using an alternative method (Shridhar et al., 2023; 2024). However, recent studies have shown that naively applying self-correction can sometimes degrade performance (Huang et al., 2024; Qu et al., 508 2024; Tyen et al., 2024), highlighting the need for more effective strategies. Supervised fine-tuning 509 with feedback from larger models (Ye et al., 2023; Qu et al., 2024), or using an ensemble of mod-510 els Havrilla et al. (2024), has produced notable results. Nevertheless, relying on larger or multiple 511 models for feedback presents challenges. In contrast, our approach takes advantage of *self-learning*, 512 eliminating the need for external model feedback. 513

Self-Training in LLMs: Self-training is a semi-supervised learning method in which the model's own predictions are used as additional data to improve its performance (Scudder, 1965; Yarowsky, 1995). This technique has been applied to various NLP tasks, such as machine translation (He et al., 2020; Sun et al., 2021; Gulcehre et al., 2023). In contrast, we use self-learning principles to generate new data and continually update the policy in an on-policy fashion. This approach can be viewed as an on-policy counterpart to Self Imitation Learning (Oh et al., 2018), where the policy learns from prospective successful trajectories in its initial stages, rather than imitating past successful behavior.

521 522

523

524

525

526

527

528

529

### 7 CONCLUSION

We present SMART: Self-learning Meta-strategy Agent for Reasoning Tasks, a solution to the challenges LMs face in selecting strategies for complex reasoning tasks. By modeling the strategy selection process as a Markov decision process and leveraging reinforcement learning, SMART enables LMs to autonomously learn and apply the most effective reasoning strategies on the first trial, thereby reducing the reliance on iterative self-refinement. Our proposed approach not only improves the accuracy of LMs, as demonstrated by significant performance gains across multiple datasets but also enhances computational efficiency by minimizing the need for multiple inference passes.

- 530 531
- 532
- 534
- 535
- 536
- 537
- 538
- 539

#### 540 REFERENCES 541

550

559

566

571

577

585

586

- 542 Karen E Adolph, Beatrix Vereijken, and Mark A Denny. Learning to crawl. *Child development*, 69 (5):1299–1312, 1998. 543
- 544 Rebecca Boncoddo, James A Dixon, and Elizabeth Kelley. The emergence of a novel representation from action: Evidence from preschoolers. *Developmental science*, 13(2):370–377, 2010. 546
- 547 Sarah A Brown, David Menendez, and Martha W Alibali. Strategy adoption depends on character-548 istics of the instruction, learner, and strategy. Cognitive Research: Principles and Implications, 4 549 (1):9, 2019.
- Wenhu Chen, Xueguang Ma, Xinyi Wang, and William W. Cohen. Program of thoughts prompt-551 ing: Disentangling computation from reasoning for numerical reasoning tasks. Transactions on 552 Machine Learning Research, 2023. ISSN 2835-8856. URL https://openreview.net/ 553 forum?id=YfZ4ZPt8zd. 554
- 555 Karl Cobbe, Vineet Kosaraju, Mohammad Bavarian, Mark Chen, Heewoo Jun, Lukasz Kaiser, 556 Matthias Plappert, Jerry Tworek, Jacob Hilton, Reiichiro Nakano, Christopher Hesse, and John Schulman. Training verifiers to solve math word problems. arXiv preprint arXiv:2110.14168, 558 2021.
- Abhimanyu Dubey, Abhinav Jauhri, Abhinav Pandey, Abhishek Kadian, Ahmad Al-Dahle, Aiesha 560 Letman, Akhil Mathur, Alan Schelten, Amy Yang, Angela Fan, et al. The llama 3 herd of models. 561 *arXiv preprint arXiv:2407.21783*, 2024. 562
- 563 Caglar Gulcehre, Tom Le Paine, Srivatsan Srinivasan, Ksenia Konyushkova, Lotte Weerts, Abhishek 564 Sharma, Aditya Siddhant, Alex Ahern, Miaosen Wang, Chenjie Gu, et al. Reinforced self-training 565 (rest) for language modeling. arXiv preprint arXiv:2308.08998, 2023.
- Alex Havrilla, Sharath Raparthy, Christoforus Nalmpantis, Jane Dwivedi-Yu, Maksym Zhuravin-567 skyi, Eric Hambro, and Roberta Railneau. Glore: When, where, and how to improve llm reasoning 568 via global and local refinements. arXiv preprint arXiv:2402.10963, 2024. 569
- 570 Junxian He, Jiatao Gu, Jiajun Shen, and Marc'Aurelio Ranzato. Revisiting self-training for neural sequence generation. In International Conference on Learning Representations, 2020. URL 572 https://openreview.net/forum?id=SJqdnAVKDH. 573
- 574 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. arXiv preprint 575 arXiv:2106.09685, 2021. 576
- Jie Huang, Xinyun Chen, Swaroop Mishra, Huaixiu Steven Zheng, Adams Wei Yu, Xinying Song, 578 and Denny Zhou. Large language models cannot self-correct reasoning yet. In The Twelfth 579 International Conference on Learning Representations, 2024. URL https://openreview. 580 net/forum?id=IkmD3fKBPQ. 581
- 582 AQ Jiang, A Sablayrolles, A Mensch, C Bamford, DS Chaplot, D de las Casas, F Bressand, 583 G Lengyel, G Lample, L Saulnier, et al. Mistral 7b (2023). arXiv preprint arXiv:2310.06825, 2023. 584
  - Geunwoo Kim, Pierre Baldi, and Stephen McAleer. Language models can solve computer tasks. Advances in Neural Information Processing Systems, 36, 2024.
- 588 Woosuk Kwon, Zhuohan Li, Siyuan Zhuang, Ying Sheng, Lianmin Zheng, Cody Hao Yu, Joseph E. 589 Gonzalez, Hao Zhang, and Ion Stoica. Efficient memory management for large language model 590 serving with pagedattention. In Proceedings of the ACM SIGOPS 29th Symposium on Operating Systems Principles, 2023.

592

Patrick Lemaire and Sophie Callies. Children's strategies in complex arithmetic. Journal of Experimental Child Psychology, 103(1):49-65, 2009.

- 594 Aman Madaan, Niket Tandon, Prakhar Gupta, Skyler Hallinan, Luyu Gao, Sarah Wiegreffe, Uri 595 Alon, Nouha Dziri, Shrimai Prabhumoye, Yiming Yang, Shashank Gupta, Bodhisattwa Prasad 596 Majumder, Katherine Hermann, Sean Welleck, Amir Yazdanbakhsh, and Peter Clark. Self-refine: 597 Iterative refinement with self-feedback. In Thirty-seventh Conference on Neural Information Pro-598 cessing Systems, 2023. URL https://openreview.net/forum?id=S37hOerQLB. Shen-yun Miao, Chao-Chun Liang, and Keh-Yih Su. A diverse corpus for evaluating and developing 600 English math word problem solvers. In Proceedings of the 58th Annual Meeting of the Associ-601 ation for Computational Linguistics, pp. 975–984, Online, July 2020. Association for Compu-602 tational Linguistics. doi: 10.18653/v1/2020.acl-main.92. URL https://aclanthology. 603 org/2020.acl-main.92. 604 Junhyuk Oh, Yijie Guo, Satinder Singh, and Honglak Lee. Self-imitation learning. In Jennifer 605 Dy and Andreas Krause (eds.), Proceedings of the 35th International Conference on Machine 606 Learning, volume 80 of Proceedings of Machine Learning Research, pp. 3878–3887. PMLR, 607 10-15 Jul 2018. URL https://proceedings.mlr.press/v80/oh18b.html. 608 609 Arkil Patel, Satwik Bhattamishra, and Navin Goyal. Are NLP models really able to solve simple 610 math word problems? In Proceedings of the 2021 Conference of the North American Chapter 611 of the Association for Computational Linguistics: Human Language Technologies, pp. 2080-612 2094, Online, June 2021a. Association for Computational Linguistics. doi: 10.18653/v1/2021. 613 naacl-main.168. URL https://aclanthology.org/2021.naacl-main.168. 614 Arkil Patel, Satwik Bhattamishra, and Navin Goyal. Are NLP models really able to solve simple 615 math word problems? In Proceedings of the 2021 Conference of the North American Chapter 616 of the Association for Computational Linguistics: Human Language Technologies, pp. 2080-617 2094, Online, June 2021b. Association for Computational Linguistics. doi: 10.18653/v1/2021. 618 naacl-main.168. URL https://aclanthology.org/2021.naacl-main.168. 619 Manish Prajapat, Mojmir Mutny, Melanie N. Zeilinger, and Andreas Krause. Submodular reinforce-620 ment learning. In The Twelfth International Conference on Learning Representations, ICLR 2024, 621 Vienna, Austria, May 7-11, 2024. OpenReview.net, 2024. URL https://openreview. 622 net/forum?id=loYSzjSaAK. 623 624 Martin L Puterman. Markov decision processes: discrete stochastic dynamic programming. John 625 Wiley & Sons, 2014. 626 Yuxiao Qu, Tianjun Zhang, Naman Garg, and Aviral Kumar. Recursive introspection: Teaching 627 language model agents how to self-improve. arXiv preprint arXiv:2407.18219, 2024. 628 629 Henry Scudder. Probability of error of some adaptive pattern-recognition machines. IEEE Transac-630 tions on Information Theory, 11(3):363-371, 1965. 631 Noah Shinn, Federico Cassano, Ashwin Gopinath, Karthik Narasimhan, and Shunyu Yao. Reflexion: 632 Language agents with verbal reinforcement learning. Advances in Neural Information Processing 633 Systems, 36, 2024. 634 635 Kumar Shridhar, Harsh Jhamtani, Hao Fang, Benjamin Van Durme, Jason Eisner, and Patrick Xia. 636 Screws: A modular framework for reasoning with revisions. arXiv preprint arXiv:2309.13075, 637 2023. 638 Kumar Shridhar, Koustuv Sinha, Andrew Cohen, Tianlu Wang, Ping Yu, Ramakanth Pasunuru, 639 Mrinmaya Sachan, Jason Weston, and Asli Celikyilmaz. The ART of LLM refinement: Ask, 640 refine, and trust. In Proceedings of the 2024 Conference of the North American Chapter of the 641 Association for Computational Linguistics: Human Language Technologies (Volume 1: Long Pa-642 pers), pp. 5872–5883, Mexico City, Mexico, June 2024. Association for Computational Linguis-643
- 643 tics. doi: 10.18653/v1/2024.naacl-long.327. URL https://aclanthology.org/2024.
   644 naacl-long.327.
- Haipeng Sun, Rui Wang, Kehai Chen, Masao Utiyama, Eiichiro Sumita, and Tiejun Zhao. Self training for unsupervised neural machine translation in unbalanced training data scenarios. In
   Kristina Toutanova, Anna Rumshisky, Luke Zettlemoyer, Dilek Hakkani-Tur, Iz Beltagy, Steven

648 649	Bethard, Ryan Cotterell, Tanmoy Chakraborty, and Yichao Zhou (eds.), <i>Proceedings of the 2021</i>
650	Conference of the North American Chapter of the Association for Computational Linguistics: Hu-
651	Linguistics doi: 10.18652/v1/2021.paged.main 211. URL https://calenthele.main.
652	2021 papel-main 311
653	
654	Richard S Sutton, Andrew G Barto, et al. Reinforcement learning. Journal of Cognitive Neuro-
655	science, 11(1):126–134, 1999.
656	Commo Toom Thomas Magnard Cassidy Hardin Dobart Dadachi Surva Dhunatiraju Shraya
657	Pathak Laurent Sifre Morgane Rivière Mihir Saniay Kale Juliette Love et al. Gemma: Open
658	models based on gemini research and technology. arXiv preprint arXiv:2403.08295, 2024.
659	Joke Torbeyns, Bert De Smedt, Pol Ghesquière, and Lieven Verschaffel. Acquisition and use of
660	shortcut strategies by traditionally schooled children. <i>Educational Studies in Mathematics</i> , 71(1):
661	1–17, 2009.
662	
663	Gladys Tyen, Hassan Mansoor, Victor Carbune, Peter Chen, and Tony Mak. LLMs cannot find
664	reasoning errors, but can correct them given the error location. In Lun-wei Ku, Andre Mar- ting, and Viyak Srikyman (ads), Findings of the Association for Computational Linguistics
665	ACL 2024 pp 13804–13008 Bangkok Thailand and virtual meeting August 2024 Association
666	ation for Computational Linguistics doi: 10.18653/v1/2024 findings-acl.826 LIRL https:
667	//aclanthology.org/2024.findings-acl.826.
668	, ,
669	Unslothai. Unsloth. https://github.com/unslothai/unsloth, 2023. URL https:
670	//github.com/unslothai/unsloth. GitHub repository.
671	Jason Wei, Xuezhi Wang, Dale Schuurmans, Maarten Bosma, Fei Xia, Ed Chi, Ouoc V Le, Denny
672	Zhou, et al. Chain-of-thought prompting elicits reasoning in large language models. Advances in
673	neural information processing systems, 35:24824–24837, 2022.
674	
675 676	Sean Welleck, Ximing Lu, Peter West, Faeze Brahman, Tianxiao Shen, Daniel Khashabi, and Yejin Choi. Generating sequences by learning to self-correct. <i>arXiv preprint arXiv:2211.00053</i> , 2022.
677	
678	An Yang, Baosong Yang, Binyuan Hui, Bo Zheng, Bowen Yu, Chang Zhou, Chengpeng Li,
679	chengyuan Li, Dayineng Liu, Fei Huang, et al. Qwen2 technical report. arXiv preprint
680	<i>urxiv.2407.10071, 202</i> 4.
681	David Yarowsky. Unsupervised word sense disambiguation rivaling supervised methods. In 33rd
682	Annual Meeting of the Association for Computational Linguistics, pp. 189–196, Cambridge, Mas-
683	sachusetts, USA, June 1995. Association for Computational Linguistics. doi: 10.3115/981658.
684	981684. URL https://aclanthology.org/P95-1026.
685	Seonghyeon Ye, Yongrae Jo, Doyoung Kim, Sungdong Kim, Hyeonbin Hwang and Minioon Seo
686	Selfee: Iterative self-revising llm empowered by self-feedback generation. <i>Blog post</i> , 2023.
687	
688	Longhui Yu, Weisen Jiang, Han Shi, Jincheng YU, Zhengying Liu, Yu Zhang, James Kwok, Zhen-
689	guo Li, Adrian Weller, and Weiyang Liu. Metamath: Bootstrap your own mathematical questions
690	for large language models. In <i>The Twelfth International Conference on Learning Representations</i> , 2024 JIPI https://openrousi.com/formum2id=NONOhemDDt
691	2024. OKL HUUPS://OPEHIEVIEW.HEU/IOLUM:IQ=N8NUNGNDKU.
692	Denny Zhou, Nathanael Schärli, Le Hou, Jason Wei, Nathan Scales, Xuezhi Wang, Dale Schu-
693	urmans, Claire Cui, Olivier Bousquet, Quoc V Le, and Ed H. Chi. Least-to-most prompting
694	enables complex reasoning in large language models. In The Eleventh International Confer-
695	ence on Learning Representations, 2023. URL https://openreview.net/forum?id=
696	WZH7099tgfM.
697	
698	
699	
700	

# 702 8 APPENDIX

#### 8.1 PROMPTS

Below is an instruction that describes a task, paired with an input that provides further context. Write a response that appropriately completes the request. Instruction: Solve the given math problem step by step. Put your final answer after 'Final answer:'. Input: John buys 3 dress shirts. They sell for \$20 each. He also has to pay 10% tax on everything. How much did he pay in total? Response: The shirts cost 3\*20= <<3\*20= 60 >>60 before tax The tax cost 60\*.1= <<60\*.1= 6So in total they paid \$60+\$6=\$<<60+6=66 66 Final Answer: 66<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 for a month (30 days). If he continues to train for the next 12 days, how many hours will he spend on training in total? Response: Model generated response .. Figure 5: 8-shot Chain of Thought demonstration. Below is an instruction that describes a task, paired with an input that provides further context. Write a response that appropriately completes the request. Instruction: Solve the given math problem by decomposing it into smaller, manageable sub-questions. Put your final answer after 'Final answer: '. Input: John buys 3 dress shirts. They sell for \$20 each. He also has to pay 10% tax on everything. How much did he pay in total? Response: Sub-question 1: How much did John spend on shirts? Answer to Sub-question 1: 3\*\$20=\$<<3\*20=60>>60 Sub-question 2: How much did John spend on tax? Answer to Sub-question 2: \$60\*.1=\$;;60\*.1=6;;6 Sub-question 3: How much did John spend in total? Answer to Sub-question 3: \$60+\$6=\$<<60+6=66>>66 Final Answer: 66<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 for a month (30 days). If he continues to train for the next 12 days, how many hours will he spend on training in total? Response: Model generated response ... Figure 6: 8-shot Least to Most demonstration. 

756	
757	
758	
759	
760	
761	
762	
762	
703	
764	
765	
766	
767	
768	
769	
770	
771	
772	
773 774	Below is an instruction that describes a task, paired with an input that provides further context. Write a response that appropriately completes the request.
775	
776	Instruction: Solve the given math problem by writing a python program. Store your result as a variable
777	named answer.
778	Input: John buys 3 dress shirts. They sell for \$20 each. He also has to pay 10% tax on everything.
779	How much did he pay in total?
780	
700	Response: $total_shirts = 3$
701	$cost_of_one\_shirt = 20$
782	total_cost_shirts = total_shirts * cost_of_one_shirt
783	$lax_rale = 0.1$
784	$tax_amount = tax_rate + total_cost_smits$
785	answer = total_cost $<$ eos $>$
786	
787	[7 more examples randomly sampled from the training set]
788 789 700	Input: Thomas is training at the gym to prepare for a competition. He trained for 5 hours every day for a month (30 days). If he continues to train for the next 12 days, how many hours will he spend on training in total?
790 791	Response: Model generated response
792	
793	Figure 7: 8-shot Program of Thought demonstration
794	rigute // o shot rogram of rhought demonstration.
795	
796	
797	
708	
700	
000	
000	
001	
802	
803	
804	
805	
806	
807	
808	
809	