SUPERVISED REINFORCEMENT LEARNING: FROM EX-PERT TRAJECTORIES TO STEP-WISE REASONING

Anonymous authors

000

001

002003004

010 011

012

013

014

015

016

017

018

019

021

023

025

026027028

029

031

034

039

040

041

042

043

044

045 046

047

052

Paper under double-blind review

ABSTRACT

Large Language Models (LLMs) often struggle with challenging, multi-step reasoning problems due to a fundamental learning gap – Reinforcement Learning with Verifiable Rewards (RLVR) suffers from sparse rewards when correct solutions are rarely sampled, while Supervised Fine-Tuning (SFT) tends to overfit to long demonstrations through rigid token mimicry. To bridge this gap, we introduce Supervised Reinforcement Learning (SRL), a framework that reformulates problem-solving as a sequence of logical actions. SRL trains the model to learn from each action, where the model first generates an internal reasoning monologue and then commits to an action. The model receives dense rewards based on the similarity between its actions and the expert's at each step, providing a richer signal than RLVR. More importantly, by only rewarding the action, SRL allows the model flexibility in its self-generated thought process, promoting stronger reasoning abilities than SFT. On challenging mathematical reasoning benchmarks, SRL significantly outperforms both methods. Furthermore, a curriculum that cold-start with SRL before refining with RLVR achieves the strongest results. SRL also generalizes effectively to agentic software engineering tasks, establishing it as a robust framework for various reasoning tasks.

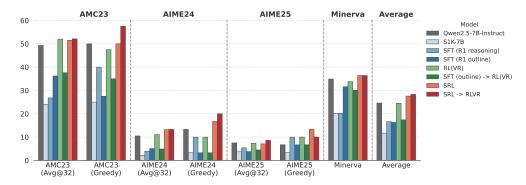


Figure 1: Performance of our method (SRL) against baselines on math reasoning benchmarks, with all models trained on the challenging s1k dataset (Muennighoff et al., 2025). Our key observations are: (1) Directly applying SFT on this dataset leads to performance degradation compared to the base model. (2) While RLVR can improve generalization over SFT, the gains are marginal. (3) Our proposed SRL method substantially outperforms these baselines, and the SRL \rightarrow RLVR pipeline achieves the highest performance, overcoming the challenges of training on difficult data.

1 Introduction

Large Language Models (LLMs) have shown impressive capabilities across a range of reasoning tasks, including solving math problems (Wang et al., 2025), generating code (Jiang et al., 2024), and agent planning (Xie et al., 2024; Li et al., 2025c). A significant recent advancement comes from leveraging reinforcement learning (RL) to enhance LLMs' complex reasoning abilities (Shao et al., 2024; Ahmadian et al., 2024; Lambert et al., 2024). By optimizing models with reward signals based on verifiable outcomes, such as the correctness of a final answer, RL offers a promising path to elicit beneficial problem-solving strategies such as self-reflection (Guo et al., 2025; Xie et al., 2025).

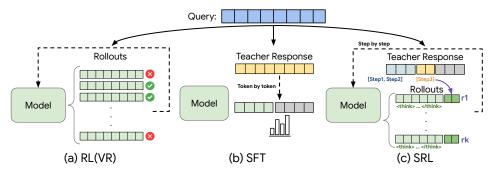


Figure 2: Illustration of SRL as compared to RL(VR) and SFT. (a) RL(VR) takes a query as input and performs k rollouts. The final answer correctness is used as the reward. (b) SFT uses both a query x and a complete teacher response y as input, training with a per-token loss to maximize the probability p(y|x). (c) SRL also uses a query and a teacher response. It breaks the response into step actions and, at each step, uses the previous steps as context. The model generates a next step action along with its step-wise inner thoughts, and the reward r_k is based on the similarity between the model's and the teacher's action.

The efficacy of these outcome-supervised RL methods relies on the policy model to discover a correct solution within the rollout budget (Brown et al., 2024). This learning approach falters on problems so difficult that the model's success rate is effectively zero (i.e., a pass@k rate of zero if rollout k times), a common scenario in tasks requiring complex, multi-step reasoning (Yue et al., 2025; Wang et al., 2024). On such problems, any incorrect intermediate step can derail the entire solution, leading to a sparse reward where positive feedback is too rare for learning to occur. Moreover, simply penalizing incorrect final outputs can be detrimental to model performance, making these difficult problems intractable for standard RL methods (Xiong et al., 2025).

An alternative approach is imitation learning, typically realized through Supervised Fine-Tuning (SFT) on expert demonstrations (Ross et al., 2011). While SFT can instill valuable reasoning behaviors, its next-token prediction objective promotes rigid, token-level mimicry, which hinders generalization. This issue is especially acute for less capable models learning from long, complex demonstrations (Chu et al., 2025a; Li et al., 2025b), an effect we show in Figure 1. This leaves a critical gap: both SFT and outcome-based RL fall short on difficult reasoning tasks, especially for base models susceptible to sparse rewards and brittle imitation.

To address this fundamental gap, we propose Supervised Reinforcement Learning (SRL), a novel framework that reformulates problem-solving as a sequential decision-making process. Instead of learning from a final outcome or mimicking a monolithic solution, SRL trains the model in a step-wise manner. We decompose expert demonstrations into a sequence of intermediate "actions". The learner is then instructed to generate its own "internal monologue" that articulates reasoning, followed by a logical action taken. At each of this step, the model is rewarded based on the similarity of its predicted action to the expert's action, which can be rapidly calculated to guarantee scalability.

SRL provides dense, granular feedback that mitigates the sparse reward problem of RLVR on complex reasoning problem. By only rewarding the action and not the monologue, SRL grants the model the flexibility to develop its own internal reasoning process, overcoming the rigid mimicry of SFT. SRL thus uniquely combines the benefits of both paradigms.

Our work makes the following contributions:

- We propose SRL, a novel framework designed to enable effective learning on complex reasoning tasks. By providing dense, step-wise rewards based on expert actions, SRL mitigates the sparse reward problem of RLVR while avoiding the rigid mimicry of SFT.
- We demonstrate the effectiveness and robustness of SRL through extensive experiments on challenging mathematical reasoning and agentic software engineering benchmarks. Our results show that SRL significantly outperforms strong baselines across both domains (5.1 & 5.3).
- Through detailed analysis, we confirm granular, step-wise guidance is vital to SRL's success and its impact on model behavior. We observe that SRL induces more flexible and sophisticated reasoning patterns, such as interleaved planning and self-verification, which improve solution quality without simply increasing output length (5.2).

2 RELATED WORK

2.1 SFT (DISTILLATION) FOR LLM REASONING.

Distilling reasoning into smaller models via SFT on teacher-generated long Chain-of-Thought (CoT) rationales has proven highly effective for transferring complex problem-solving skills (Li et al., 2023; Huang et al., 2024; Min et al., 2024; Yeo et al., 2025), as exemplified by the small models distilled from DeepSeek R1 (Guo et al., 2025). Research indicates this process is surprisingly data-efficient, with small, high-quality datasets often being sufficient (Ye et al., 2025; Muennighoff et al., 2025). Given the success, research has focused on the underlying factor for effective SFT distillation (Chen et al., 2025a). Some emphasized the logical structure of the reasoning trace rather than its semantic correctness (Luo et al., 2025; Stechly et al., 2025), as models can learn from demonstrations with factual errors (Li et al., 2025a). Moreover, significant challenges remain in the student-teacher gap where the student fails to learn from overly complex data (Li et al., 2025b), and the risk of teacher hacking, where the student overfits to a teacher's specific flaws (Tiapkin et al., 2025). Ultimately, distillation from a teacher model imposes a performance ceiling on the student (Huang et al., 2024).

2.2 RL FOR LLM REASONING.

The development of DeepSeek-R1 (Guo et al., 2025) showed the effectiveness of rule-based RL for enhancing the reasoning capabilities of LLMs. This approach utilizes a scalable reward system based on final answer correctness, exemplified by the Group Relative Policy Optimization (GRPO) algorithm (Shao et al., 2024) and parallel algorithms (Ahmadian et al., 2024; Lambert et al., 2024; Xie et al., 2025). Building on this foundation, subsequent research has introduced numerous algorithmic refinements. For example, Dr. GRPO (Liu et al., 2025) mitigates bias by removing variance normalization, while DAPO (Yu et al., 2025) introduces a token-level loss and relaxes the policy update constraint by increasing the clipping threshold. Other notable advancements include modifications to clipping methods, normalization techniques, the KL divergence loss, and dynamic sampling strategies (Chu et al., 2025b; Zhang & Zuo, 2025; Zhang et al., 2025; Chen et al., 2025b). Despite these algorithmic variations, these approaches primarily rely on the final outcome's reward signal. A critical challenge arises when the rollouts fail to identify a correct solution trajectory, particularly for difficult queries. DAPO (Yu et al., 2025), for instance, addresses this by filtering out instructions that do not yield any successful rollouts.

3 PRELIMINARIES

A Large Language Model (LLM) is formally defined by a probability distribution p_{θ} over sequences of tokens, parameterized by a set of model weights θ . Given an input prompt, represented as a token sequence $\mathbf{x} = [x_1, \dots, x_n]$, the model generates a response sequence $\mathbf{y} = [y_1, \dots, y_m]$. The response is produced autoregressively, where the generation of the token y_j at any step j is conditioned on the initial prompt \mathbf{x} and all preceding tokens in the generated sequence, (y_1, \dots, y_{j-1}) . The joint probability of the entire response sequence \mathbf{y} given the prompt \mathbf{x} is thus factorized as: $p_{\theta}(\mathbf{y}|\mathbf{x}) = \prod_{j=1}^m p_{\theta}(y_j|\mathbf{x},y_1,\dots,y_{j-1})$.

Supervised Fine-Tuning (SFT). SFT is typically employed to specialize LLM for downstream applications or domains. It is also commonly used to establish a cold start for subsequent RL training phases that requires certain reply format/pattern, such as RL for reasoning (Deng et al., 2025) or tool use (Feng et al., 2025). Specifically, the process utilizes a dataset $\mathcal{D} = \{(\mathbf{x}^{(i)}, \mathbf{y}^{(i)})\}_{i=1}^N$, where each $\mathbf{x}^{(i)}$ is an input prompt and $\mathbf{y}^{(i)}$ is the corresponding desired model output. The primary objective is to update the parameters θ to maximize the conditional probability of generating the target response $\mathbf{y}^{(i)}$ given the input prompt $\mathbf{x}^{(i)}$. This goal is formally achieved by minimizing the negative log-likelihood loss function: $\mathcal{L}_{\text{SFT}}(\theta) = -\sum_{i=1}^N \log p_{\theta}(\mathbf{y}^{(i)}|\mathbf{x}^{(i)})$ over the entire dataset. By minimizing this loss, the model learns to produce responses that are closely aligned with the exact words demonstrated in the labeled training examples.

Reinforcement Learning (RL). Recent literature on improving model reasoning capability has focused on RL with verifiable reward (RLVR), where the policy model receives reward signals purely based on the final answer correctness. Building on this principle, Group Relative Policy Optimization

(GRPO) (Shao et al., 2024) involves sampling a group of G response trajectories, $\{\mathbf{o}_i\}_{i=1}^G$, from the previous policy model, θ_{old} , for each input query \mathbf{x} . The objective function for GRPO is:

$$\mathbb{E}\left[\frac{1}{G}\sum_{i=1}^{G}\frac{1}{|\mathbf{o}_{i}|}\sum_{t=1}^{|\mathbf{o}_{i}|}\min\left(\frac{p_{\theta}(o_{i,t}\mid\mathbf{x},\mathbf{o}_{i,< t})}{p_{\theta_{old}}(o_{i,t}\mid\mathbf{x},\mathbf{o}_{i,< t})}\hat{A}_{i,t},\operatorname{clip}\left(\frac{p_{\theta}(o_{i,t}\mid\mathbf{x},\mathbf{o}_{i,< t})}{p_{\theta_{old}}(o_{i,t}\mid\mathbf{x},\mathbf{o}_{i,< t})},1-\epsilon,1+\epsilon\right)\hat{A}_{i,t}\right)\right]$$

$$-\beta\,\mathbb{D}_{\mathrm{KL}}\left[p_{\theta}\parallel p_{\mathrm{ref}}\right].\tag{1}$$

The hyperparameter $\epsilon>0$ defines the clipping range for the policy update ratio, and the coefficient $\beta>0$ modulates the influence of the KL-divergence penalty against the policy update. The term θ_{old} refers to the policy from the previous iteration. The advantage function, $\hat{A}_{i,t}=(\tilde{r}_i-\text{mean}(\tilde{r}))/\text{std}(\tilde{r})$, is defined as the group-level normalized reward.

A key challenge for these RL algorithms emerges when input queries are either too easy or too hard, resulting in uniform correctness within policy rollouts $\{\mathbf{o}_i\}_{i=1}^G$. In such cases, the advantage estimate $\hat{A}_{i,t}$ vanishes, yielding an uninformative policy gradient and preventing model updates. A common strategy to mitigate this is to dynamically sample the batches, filtering out samples and re-sampling until the data sample satisfies $0 < |\{\mathbf{o}_i|\text{is_correct}(\mathbf{o}_i)\}| < G$ (Yu et al., 2025).

4 METHODOLOGY

4.1 THE CHALLENGE OF HARD REASONING PROBLEM

RL with verifiable reward is a prominent technique for enhancing the reasoning capabilities of LLMs. The strategy is to close the gap between a model's potential performance across multiple attempts (pass@k) (Brown et al., 2024; Yue et al., 2025). However, this paradigm falters on problems where the model's pass@k rate is already near zero. For this set of difficult problems, which we term \mathcal{D}_{hard} , positive reward signals are too sparse for RLVR to be effective (Xiong et al., 2025). Moreover, simply penalizing incorrect outputs can be detrimental to model performance (Yu et al., 2025; Xiong et al., 2025), creating a significant challenge for improving model reasoning.

Formally, we define $\mathcal{D}_{\text{hard}} = \{\mathbf{x}^{(i)}, a^{(i)}\}_{i=1}^{N}$ as the set of problems (\mathbf{x}, a) where policy model's success rate is low with k samples: $\frac{1}{k} \sum_{j=1}^{k} \mathbb{I}\left(\text{ExtractAnswer}(\mathbf{y}^{(j)}) == a\right) \leq \epsilon$, where each solution attempt $\mathbf{y}^{(j)}$ is sampled from the policy $p_{\theta}(\cdot|\mathbf{x})$ and $\epsilon > 0$ is a small constant.

Due to the scarcity of successful trajectories, standard RL with verifiable reward struggles on \mathcal{D}_{hard} . Such data is further difficult to be learned by SFT, due to its limited amount and complexity in teacher reasoning trajectories (Li et al., 2025b).

4.2 SUPERVISED REINFORCEMENT LEARNING (SRL)

To address the challenge of learning from \mathcal{D}_{hard} , we introduce Supervised Reinforcement Learning (SRL), a framework that decomposes complex problem-solving as a sequential decision-making process, and thus can be easily learned on how to properly operate step-wise. Instead of generating a monolithic solution, the model learns to take actions similar to the expert while producing their own inner reasoning process in a step-by-step manner. The whole framework is illustrated in Figure 3. 1

Action-based problem formulation. Given an expert solution trajectory \mathbf{y} that leads to a correct final answer, we decompose \mathbf{y} into a sequence of tuples: $\mathbf{y} = \{\mathbf{y}_{\text{step}_n}\}_{n=1}^N$. Each steps represents a *logical action*: the concrete action to be operated. This formulation is domain-agnostic; for instance, an action in mathematical reasoning could be an algebraic manipulation, while for a software agent, it could be a command executed in a code repository.

Step-wise training data construction. To create training data for SRL, we leverage a powerful teacher model, θ_{expert} to generate solution trajectories. From a single complete solution with N steps, we construct N-1 partial trajectories. For each step $k \in \{1,\ldots,N-1\}$, we create a new input prompt $\mathbf{x}_{\text{step}_k} = [\mathbf{x},\mathbf{y}_{\text{step}_k},\ldots,\mathbf{y}_{\text{step}_{k-1}}]$, where the model's task is to predict the subsequent

¹Empirically, we found that providing the subsquent step title (e.g., "2. **Coprime Pairs**" in Figure 3) as the additional context for the learner to predict the rest of the step content can further boost the performance.

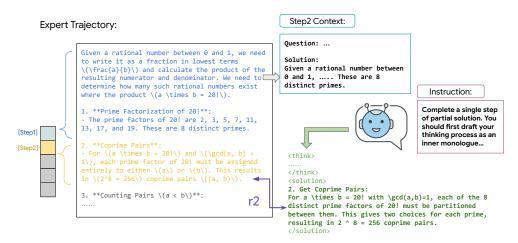


Figure 3: Given a solution trajectory, we take each summarized step as an action to be learned and take the partial solution before the step as the context of our newly created data. The model is then prompted to generate its thinking process followed by the action for the current step. A reward (r_2 in the figure) is then calculated based on the similarity between the model's and the expert's action.

step, $\mathbf{y}_{\text{step}_k}$. This process transforms one expert solution into a rich set of training instances that teach the model to proceed correctly from various intermediate states.

Learning with a sequence similarity reward with own inner monologue. Given a partial context $\mathbf{x}_{\text{step}_k}$ containing the problem and a partial solution, the policy model p_{θ} is prompted to generate the subsequent action step with their own inner monologue $\mathbf{y}'_{\text{think}}$, which is encapsulated by "<think>" tags. We then provide a dense reward based on the quality of the generated logical action $\mathbf{y}'_{\text{step}_k}$. The prediction can be formally specified as: $\mathbf{y}' \sim p_{\theta}(\cdot|\mathbf{x}_{\text{step}_k}) = [\mathbf{y}'_{\text{think}}, \mathbf{y}'_{\text{step}_k}]$.

To guide training, we consider the reward function that measures the similarity between the generated action: $R(\mathbf{y}'_{\text{step}_k}, \mathbf{y}_{\text{step}_k}) = \frac{2M}{T}$, where

- T (Total elements): This is the total number of elements in both sequences combined. It is calculated as the sum of the lengths of the two sequences: $T = |S_1| + |S_2|$.
- M (Matched elements): The total count of elements found in all non-overlapping matching blocks between the two sequences. The algorithm first finds the longest contiguous matching subsequence and then recursively searches for more matches in the segments to the left and right of that block. If we represent the set of all such matching blocks as a list of tuples (i, j, n), where n is the length of the matching block, then M is the sum of all lengths n: $M = \sum_{(i,j,n) \in \text{MatchingBlocks}} n$.

Combining these definitions, we can calculate the similarity ratio $R \in [0,1]$ as:

$$R = \frac{2\sum_{(i,j,n) \in \mathsf{MatchingBlocks}} n}{|S_1| + |S_2|}$$

In practice, we use Python's difflib. SequenceMatcher for this comparison, and assign a negative reward if the generated output \mathbf{y}' fails to follow the required format. Hence, the final reward used is:

$$r(\mathbf{y}_{\mathsf{step}_k}', \mathbf{y}_{\mathsf{step}_k}) = \begin{cases} R(\mathbf{y}_{\mathsf{step}_k}', \mathbf{y}_{\mathsf{step}_k}) & \text{if } \mathbf{y}' \text{ follows format,} \\ -1 & \text{otherwise.} \end{cases}$$

The policy p_{θ} is then optimized using this reward signal with the GRPO objective function defined in Equation 1. Notably, our reward is computed only on the logical action, not the internal monologue. This grants the model flexibility to develop its own internal reasoning style while ensuring its external actions align with the expert's strategy. This design provides dense, step-level feedback and allows for rapid reward calculation, making the SRL framework both effective and scalable.

Dynamic sampling for SRL. As our reward signal $r \in [0,1] \cup \{-1\}$ is dense, we generalize the dynamic sampling strategy previously designed for outcome accuracy and implement it to filter

samples with less meaningful updates. Specifically, a sample should be filtered out if its rollouts yield rewards with near-zero variance, providing a weak advantage and thus weak learning signal. With the sequence similarity reward in SRL, we retain a sample if the standard deviation of the reward scores of its rollouts exceeds a threshold $\epsilon>0$:

$$\sqrt{\frac{\sum_{i=1}^{G} (r(\mathbf{o}_i, \mathbf{y}) - \bar{r})^2}{G}} > \epsilon$$

where G is the number of generated rollouts, $r(\mathbf{o}_i, \mathbf{y})$ is the sequence similarity reward for the i-th rollout \mathbf{o}_i given the expert trajectory \mathbf{y} , and \bar{r} is the mean reward for the sample. To maintain a consistent batch size of B, we continuously sample and filter until the batch is filled.

5 EXPERIMENTS

5.1 MAIN RESULTS: MATH REASONING

Setup. We finetue *Qwen2.5-7B-Instruct* (Yang et al., 2024) on the *s1K-1.1* dataset (Muennighoff et al., 2025). This dataset contains 1,000 diverse and challenging questions, each accompanied by a detailed reasoning trace and a final solution generated by DeepSeek R1. The solutions from DeepSeek R1 are formatted with structured, numbered steps (e.g., "1. Title of Step 1"). We leverage this structure to generate intermediate training targets by parsing these solutions and treating each complete step as a ground-truth continuation. Any data points that do not adhere to this format are excluded. We hold out 60 questions from the dataset to form our validation set.

Baselines. We benchmark our proposed methods against several baselines, all initialized from the *Qwen2.5-7B-Instruct* model. These baselines include: (i) SFT on either the complete reasoning traces (R1 reasoning) or the final solutions from the s1K-1.1 dataset (R1 outline); (ii) s1K-7B, the official distilled model released by the dataset's authors; and (iii) RLVR, which we implement using the GRPO algorithm. To ensure fair comparison, we implement additional dynamic sampling as in Yu et al. (2025), which removes samples with all correct or incorrect rollouts. We evaluate RLVR in two distinct settings: applied directly to the base model and applied after an initial SFT phase. Our proposed method, SRL, is likewise evaluated both as a standalone technique and in a sequential configuration where it precedes RLVR (SRL then RLVR). All models are trained for up to 30 epochs, and for each method, we select the checkpoint with the best performance on the validation set.

Evaluation. We evaluate all models on the following four competition-level mathematical reasoning benchmarks: AMC23², AIME24³, AIME25⁴ and Minerva Math (Lewkowycz et al., 2022). Our evaluation protocol for all benchmarks strictly follows the setup established by Qwen2.5-Math⁵ and report the accuracy of greedy sampling. In addition, for AMC23, AIME24 and AIME25, we report the average@32 score with a temperature of 1.0 for all baselines to ensure a more robust evaluation.

Performance. The performance results of our models are summarized in Table 1. Consistent with the officially released S1K-7B model, our model trained with SFT on the same dataset exhibited a notable performance degradation. In contrast, methods based on RL maintained generalization on the evaluation benchmarks. Specifically, while RLVR maintained the performance, SRL provided a substantial boost of 3.0% on average. Furthermore, applying RLVR after SRL training yielded a 3.7% increase on average, leveraging only 1k training data.

5.2 ANALYSIS: MATH REASONING

Effect of dynamic sampling in SRL. In Table 2, we analyze the impact of the dynamic sampling component in SRL, based on thresholding the standard deviation of sequence similarity rewards within rollouts. For both models, we train until the training reward converges and select checkpoint based on validation scores. Our results are consistent the findings of DAPO (Yu et al., 2025), which stated that removing samples that provide a zero learning signal is critical in the effectiveness of the RL training loop, showing non-trivial d performance improvement.

²https://huggingface.co/datasets/AI-MO/aimo-validation-amc

³https://huggingface.co/datasets/AI-MO/aimo-validation-aime

⁴https://huggingface.co/datasets/math-ai/aime25

⁵https://github.com/QwenLM/Qwen2.5-Math

Table 1: Evaluation results across competition-level math benchmarks. We take Qwen2.5-7B-Instruct as the base model and report the performance of different training schemes (SFT, RLVR via GRPO, and SRL) using the same set of training data. The **bold** numbers indicate the best results among the open-source models and the <u>underscored</u> numbers represent the second-best results.

Model	AMC23		AIME24		AIME25		Minerva Math	Average
	Avg@32	Greedy	Avg@32	Greedy	Avg@32	Greedy		Tiverage
			Base	Model				
Qwen2.5-7B-Instruct	49.3	50.0	10.5	13.3	<u>7.5</u>	6.7	<u>34.9</u>	24.6
			Training	with SFT				
S1K-7B SFT (R1 reasoning) SFT (R1 outline)	24.1 26.8 36.2	25.0 40.0 27.5	2.2 3.9 5.1	3.3 10.0 3.3	3.7 5.4 3.8	3.3 10.0 6.7	20.2 20.2 31.6	11.7 16.6 16.3
			Training w	vith RL(VR)			
$\begin{array}{c} RL(VR) \\ SFT \ (outline) \rightarrow RL(VR) \end{array}$	<u>52.0</u> 37.6	47.5 35.0	11.1 4.9	10.0 3.3	7.4 4.5	10.0 6.7	33.8 30.1	24.5 17.4
			Training	with SRL				
$\begin{array}{c} \text{SRL} \\ \text{SRL} \rightarrow \text{RLVR} \end{array}$	51.5 52.1	50.0 57.5	13.2 13.3	16.7 20.0	7.1 8.6	13.3 10.0	36.4 36.4	27.6 28.3

Table 2: The effect of dynamic filtering on SRL. Filtering out samples with less meaningful updates provided nontrivial performance improvement. DS stands for dynamic sampling.

Model AMC		C23 AIM		E24 AIM		E25	Minerva Math	Average
1110401	Avg@32	Greedy	Avg@32	Greedy	Avg@32	Greedy	1/ 11101 / W 1/ 1W01	11,01,00
SRL w/out DS	48.5	52.5	11.1	13.3	6.8	6.7	33.8	24.7
SRL w/ DS	51.5	50.0	13.2	16.7	7.1	13.3	36.4	27.6

Disentangling the impact of guidance granularity in SRL. To isolate the effects of guidance granularity from the benefits of sequence similarity rewards, we evaluated our multi-step method against two simplified, one-step baselines:

- Final-answer reward: This baseline utilizes RLVR trained with GRPO, rewarding the model solely based on the correctness of the final answer.
- Holistic sequence similarity reward: The model generates a complete solution in a single step.
 This entire solution is then evaluated for similarity against the complete ground-truth trajectory.

As shown in Table 3, the results underscore the value of fine-grained guidance for sequence similarity reward. While the holistic sequence-similarity reward improved performance on average, providing granular, step-by-step guidance led to markedly superior performance across benchmarks.

Table 3: Model Performance of different reward functions and density. For sequence similarity reward, we implement it with the entire expert output as an one-step supervision. The model benefits from our multi-step decomposition on the small set of challenging training data.

Model	AMC23		AIME24		AIME25		Minerva Math	Average
	Avg@32	Greedy	Avg@32	Greedy	Avg@32	Greedy	Willier va Wiath	Tiverage
			Final A	Answer Rev	ward			
One-step (RLVR)	52.0	47.5	11.1	10.0	7.4	10.0	33.8	24.5
Sequence Similarity Reward								
One-step Multi-step (SRL)	51.8 51.5	52.5 50.0	11.7 13.2	13.3 16.7	6.9 7.1	10.0 13.3	34.9 36.4	25.9 27.6

Interleaved reasoning behaviors. Beyond achieving superior performance, models trained with SRL demonstrate distinctive and flexible reasoning patterns. Conventional models typically generate a single, monolithic block of reasoning at the beginning of a solution. In contrast, our models, particularly those fine-tuned with RLVR, dynamically interleave reasoning steps with the solution-

generation process. This behavior manifests in three primary ways: (1) upfront planning by outlining all subsequent steps in a comprehensive initial plan; (2) on-the-fly adjustments by inserting multiple, distinct reasoning blocks throughout its solution to make iterative plans and adjustments; and (3) reflective verification: the model pauses to reflect on and verify its answer before delivering the final output. The examples (Example 1 below and 2 in the Appendix), drawn from a model trained with our $SRL \rightarrow RL$ approach, illustrate these emergent reasoning patterns.

Example 1: Interleaved thinking and self-verification

Question: What is the least positive integer value of x for which the inequality 3x > 2x + 1 is true? <think>

To solve the inequality 3x > 2x + 1, we can start by isolating x. Let's subtract 2x from both sides of the inequality.

</think>

Subtract 2x from both sides:

$$3x - 2x > 2x + 1 - 2x$$

This simplifies to:

<think>

We can verify this by substituting x=2 into the original inequality to ensure it holds true. </think>

The inequality x > 1 tells us that x must be greater than 1. The smallest positive integer greater than 1 is 2. Therefore, the least positive integer value of x that satisfies the inequality is 2. Substitute x = 2 into the original inequality:

$$3(2) > 2(2) + 1$$

This simplifies to:

This is true, so x = 2 is indeed a solution.

Therefore, the least positive integer value of x for which the inequality 3x > 2x + 1 is true is 2

Reasoning length. We further investigate whether the performance gains from SRL can be simply attributed to an increase in reasoning length. In Figure 4, we plot the distribution of reasoning lengths (i.e., the number of words in the solution) for both the base model and our SRL-finetuned model. We observe no significant difference between the two distributions. This finding suggests that the performance improvement stems from enhanced planning and higher-quality reasoning, rather than merely an increased token budget at inference time.

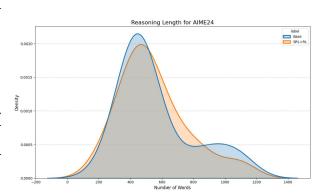


Figure 4: Reasoning length distribution for base model and model trained with SRL.

5.3 EXTENSION: SOFTWARE ENGINEERING AGENTIC REASONING

Task. We extend our SRL framework to the domain of software engineering, training agents to resolve real-world programming issues. These tasks are commonly evaluated on benchmarks like SWE-Bench (Jimenez et al., 2023), which requires agents to perform complex, multi-turn interactions with large codebases and reason deeply about code functionality.

However, unlike in math domain, the direct application of online RL to software engineering is impeded by significant practical challenges. These include handling long context windows, high-latency environment feedback, and slow patch verification (Wei et al., 2025; Golubev et al., 2025). Consequently, these obstacles have hindered the development of stable and scalable end-to-end RL

Table 4: Performance of SRL on SWE-Bench-Verified. Results in the table are using greedy decoding.

	Oracle File Edit	End-to-End
Qwen2.5-Coder-Instruct (Base)	5.8	3.2
SWE-Gym-7B	8.4	4.2
SRL (ours)	14.8	8.6

methods, leading to a prevailing approach of collecting expert agent trajectories and distilling them into a policy via SFT (Pan et al., 2024; Yang et al., 2025).

Setup. We apply SRL to further fine-tune *Qwen2.5-Coder-7B-Instruct* (Hui et al., 2024), a model already specialized for coding tasks. We use a dataset from Yang et al. (2025), which consists of 5,000 expert agent trajectories. These trajectories were generated by *claude-3-7-sonnet-20250219* (Anthropic, 2025) and subsequently verified to ensure they produce correct code patches.

Each trajectory is composed of multiple steps defined by the agent's interactions with the coding environment. As the example below illustrates, a single step contains natural language reasoning followed by an executable action:

In line with our SRL formulation (Section 4.2), we define the "action" as the environment-consumable command (e.g., the bash call). Following this decomposition, we process the full trajectories to create 134k step-wise training instances. For validation, we hold out 30 full trajectories, from which we curate a validation set of 650 step-wise instances.

Evaluation. We evaluate our model's patch generation performance by measuring its resolve rate (%) under two distinct configurations, following Wei et al. (2025): (1) Oracle file editing evaluation: The model is provided with the oracle code files to repair. This configuration isolates and measures the model's core patch generation capability; (2) End-to-end evaluation: This setting uses the Agentlessmini agent scaffold (Wei et al., 2025) to first identify the file(s) to modify and subsequently generate the patch. It tests the model's fault localization and code repair abilities in conjunction.

We compare our SRL-trained model against two crucial baselines: the original base model (*Qwen2.5-Coder-Instruct*) and SWE-Gym-7B (Pan et al., 2024). Since SWE-Gym-7B is an SFT-based model finetuned from the same base model, this provides a direct, fair comparison between SFT and our SRL training methodology. As shown in Table 4, SRL substantially outperforms both baselines. In the oracle setting, SRL achieves a 14.8% resolve rate, representing a 74% relative improvement over the strong SWE-Gym-7B baseline. The performance gain is consistent when evaluating in the challenging end-to-end setting, where SRL can obtain twice the performance.

6 Conclusion

In conclusion, we introduce Supervised Reinforcement Learning (SRL), a method for teaching complex reasoning to LLMs using expert demonstrations, especially where conventional RL or SFT struggles. By breaking down problems and using a step-wise sequence similarity reward, SRL bridges imitation and reinforcement learning. Empirically, SRL significantly surpasses baselines on math and software engineering tasks and enables a powerful curriculum learning strategy when paired with RLVR. SRL is a robust technique for training agents to solve challenging multi-step problems, paving the way for training more capable and versatile AI agents.

REFERENCES

- Arash Ahmadian, Chris Cremer, Matthias Gallé, Marzieh Fadaee, Julia Kreutzer, Olivier Pietquin, Ahmet Üstün, and Sara Hooker. Back to basics: Revisiting reinforce style optimization for learning from human feedback in llms. *arXiv* preprint arXiv:2402.14740, 2024.
- Anthropic. Introducing claude 3.7 sonnet, 2025., 2025. URL https://www.anthropic.com/news/claude-3-7-sonnet.
- Bradley Brown, Jordan Juravsky, Ryan Ehrlich, Ronald Clark, Quoc V Le, Christopher Ré, and Azalia Mirhoseini. Large language monkeys: Scaling inference compute with repeated sampling. *arXiv preprint arXiv:2407.21787*, 2024.
- Xinghao Chen, Zhijing Sun, Wenjin Guo, Miaoran Zhang, Yanjun Chen, Yirong Sun, Hui Su, Yijie Pan, Dietrich Klakow, Wenjie Li, et al. Unveiling the key factors for distilling chain-of-thought reasoning. *arXiv preprint arXiv:2502.18001*, 2025a.
- Yi Chen, Yuying Ge, Rui Wang, Yixiao Ge, Junhao Cheng, Ying Shan, and Xihui Liu. Grpocare: Consistency-aware reinforcement learning for multimodal reasoning. *arXiv preprint arXiv:2506.16141*, 2025b.
- Tianzhe Chu, Yuexiang Zhai, Jihan Yang, Shengbang Tong, Saining Xie, Dale Schuurmans, Quoc V Le, Sergey Levine, and Yi Ma. Sft memorizes, rl generalizes: A comparative study of foundation model post-training. *arXiv preprint arXiv:2501.17161*, 2025a.
- Xiangxiang Chu, Hailang Huang, Xiao Zhang, Fei Wei, and Yong Wang. Gpg: A simple and strong reinforcement learning baseline for model reasoning. *arXiv preprint arXiv:2504.02546*, 2025b.
- Yihe Deng, Hritik Bansal, Fan Yin, Nanyun Peng, Wei Wang, and Kai-Wei Chang. Openvlthinker: An early exploration to complex vision-language reasoning via iterative self-improvement. *arXiv* preprint arXiv:2503.17352, 2025.
- Jiazhan Feng, Shijue Huang, Xingwei Qu, Ge Zhang, Yujia Qin, Baoquan Zhong, Chengquan Jiang, Jinxin Chi, and Wanjun Zhong. Retool: Reinforcement learning for strategic tool use in llms. *arXiv preprint arXiv:2504.11536*, 2025.
- Alexander Golubev, Maria Trofimova, Sergei Polezhaev, Ibragim Badertdinov, Maksim Nekrashevich, Anton Shevtsov, Simon Karasik, Sergey Abramov, Andrei Andriushchenko, Filipp Fisin, Sergei Skvortsov, and Boris Yangel. Training long-context, multi-turn software engineering agents with reinforcement learning. *arXiv* preprint arXiv:2508.03501, 2025.
- Daya Guo, Dejian Yang, Haowei Zhang, Junxiao Song, Ruoyu Zhang, Runxin Xu, Qihao Zhu, Shirong Ma, Peiyi Wang, Xiao Bi, et al. Deepseek-r1: Incentivizing reasoning capability in llms via reinforcement learning. *arXiv* preprint arXiv:2501.12948, 2025.
- Zhen Huang, Haoyang Zou, Xuefeng Li, Yixiu Liu, Yuxiang Zheng, Ethan Chern, Shijie Xia, Yiwei Qin, Weizhe Yuan, and Pengfei Liu. O1 replication journey–part 2: Surpassing o1-preview through simple distillation, big progress or bitter lesson? *arXiv preprint arXiv:2411.16489*, 2024.
- Binyuan Hui, Jian Yang, Zeyu Cui, Jiaxi Yang, Dayiheng Liu, Lei Zhang, Tianyu Liu, Jiajun Zhang, Bowen Yu, Kai Dang, et al. Qwen2. 5-coder technical report. *arXiv preprint arXiv:2409.12186*, 2024.
- Juyong Jiang, Fan Wang, Jiasi Shen, Sungju Kim, and Sunghun Kim. A survey on large language models for code generation. *arXiv preprint arXiv:2406.00515*, 2024.
- Carlos E Jimenez, John Yang, Alexander Wettig, Shunyu Yao, Kexin Pei, Ofir Press, and Karthik Narasimhan. Swe-bench: Can language models resolve real-world github issues? *arXiv preprint arXiv:2310.06770*, 2023.
- Nathan Lambert, Jacob Morrison, Valentina Pyatkin, Shengyi Huang, Hamish Ivison, Faeze Brahman, Lester James V Miranda, Alisa Liu, Nouha Dziri, Shane Lyu, et al. Tulu 3: Pushing frontiers in open language model post-training. *arXiv preprint arXiv:2411.15124*, 2024.

- Aitor Lewkowycz, Anders Andreassen, David Dohan, Ethan Dyer, Henryk Michalewski, Vinay Ramasesh, Ambrose Slone, Cem Anil, Imanol Schlag, Theo Gutman-Solo, et al. Solving quantitative reasoning problems with language models. *Advances in neural information processing systems*, 35:3843–3857, 2022.
 - Dacheng Li, Shiyi Cao, Tyler Griggs, Shu Liu, Xiangxi Mo, Eric Tang, Sumanth Hegde, Kourosh Hakhamaneshi, Shishir G Patil, Matei Zaharia, et al. Llms can easily learn to reason from demonstrations structure, not content, is what matters! *arXiv preprint arXiv:2502.07374*, 2025a.
 - Liunian Harold Li, Jack Hessel, Youngjae Yu, Xiang Ren, Kai-Wei Chang, and Yejin Choi. Symbolic chain-of-thought distillation: Small models can also" think" step-by-step. *arXiv preprint arXiv:2306.14050*, 2023.
 - Yuetai Li, Xiang Yue, Zhangchen Xu, Fengqing Jiang, Luyao Niu, Bill Yuchen Lin, Bhaskar Ramasubramanian, and Radha Poovendran. Small models struggle to learn from strong reasoners. arXiv preprint arXiv:2502.12143, 2025b.
 - Zhiwei Li, Yong Hu, and Wenqing Wang. Encouraging good processes without the need for good answers: Reinforcement learning for llm agent planning. *arXiv preprint arXiv:2508.19598*, 2025c.
 - Zichen Liu, Changyu Chen, Wenjun Li, Penghui Qi, Tianyu Pang, Chao Du, Wee Sun Lee, and Min Lin. Understanding r1-zero-like training: A critical perspective. *arXiv preprint arXiv:2503.20783*, 2025.
 - Yijia Luo, Yulin Song, Xingyao Zhang, Jiaheng Liu, Weixun Wang, GengRu Chen, Wenbo Su, and Bo Zheng. Deconstructing long chain-of-thought: A structured reasoning optimization framework for long cot distillation. *arXiv* preprint arXiv:2503.16385, 2025.
 - Yingqian Min, Zhipeng Chen, Jinhao Jiang, Jie Chen, Jia Deng, Yiwen Hu, Yiru Tang, Jiapeng Wang, Xiaoxue Cheng, Huatong Song, et al. Imitate, explore, and self-improve: A reproduction report on slow-thinking reasoning systems. *arXiv preprint arXiv:2412.09413*, 2024.
 - Niklas Muennighoff, Zitong Yang, Weijia Shi, Xiang Lisa Li, Li Fei-Fei, Hannaneh Hajishirzi, Luke Zettlemoyer, Percy Liang, Emmanuel Candès, and Tatsunori Hashimoto. s1: Simple test-time scaling, 2025. URL https://arxiv.org/abs/2501.19393.
 - Jiayi Pan, Xingyao Wang, Graham Neubig, Navdeep Jaitly, Heng Ji, Alane Suhr, and Yizhe Zhang. Training software engineering agents and verifiers with swe-gym. *arXiv preprint arXiv:2412.21139*, 2024.
 - Stéphane Ross, Geoffrey Gordon, and Drew Bagnell. A reduction of imitation learning and structured prediction to no-regret online learning. In *Proceedings of the fourteenth international conference on artificial intelligence and statistics*, pp. 627–635. JMLR Workshop and Conference Proceedings, 2011.
 - Zhihong Shao, Peiyi Wang, Qihao Zhu, Runxin Xu, Junxiao Song, Xiao Bi, Haowei Zhang, Mingchuan Zhang, YK Li, Yang Wu, et al. Deepseekmath: Pushing the limits of mathematical reasoning in open language models. *arXiv* preprint arXiv:2402.03300, 2024.
 - Kaya Stechly, Karthik Valmeekam, Atharva Gundawar, Vardhan Palod, and Subbarao Kambhampati. Beyond semantics: The unreasonable effectiveness of reasonless intermediate tokens. *arXiv* preprint arXiv:2505.13775, 2025.
 - Daniil Tiapkin, Daniele Calandriello, Johan Ferret, Sarah Perrin, Nino Vieillard, Alexandre Ramé, and Mathieu Blondel. On teacher hacking in language model distillation. *arXiv preprint arXiv:2502.02671*, 2025.
 - Chaojie Wang, Yanchen Deng, Zhiyi Lyu, Liang Zeng, Jujie He, Shuicheng Yan, and Bo An. Q*: Improving multi-step reasoning for llms with deliberative planning. *arXiv preprint arXiv:2406.14283*, 2024.
 - Peng-Yuan Wang, Tian-Shuo Liu, Chenyang Wang, Yi-Di Wang, Shu Yan, Cheng-Xing Jia, Xu-Hui Liu, Xin-Wei Chen, Jia-Cheng Xu, Ziniu Li, and Yang Yu. A survey on large language models for mathematical reasoning. *arXiv* preprint arXiv:2506.08446, 2025.

- Yuxiang Wei, Olivier Duchenne, Jade Copet, Quentin Carbonneaux, Lingming Zhang, Daniel Fried, Gabriel Synnaeve, Rishabh Singh, and Sida I. Wang. Swe-rl: Advancing llm reasoning via reinforcement learning on open software evolution. *arXiv preprint arXiv:2502.18449*, 2025.
- Jian Xie, Kai Zhang, Jiangjie Chen, Tinghui Zhu, Renze Lou, Yuandong Tian, Yanghua Xiao, and Yu Su. Travelplanner: A benchmark for real-world planning with language agents. In *Forty-first International Conference on Machine Learning*, 2024.
- Tian Xie, Zitian Gao, Qingnan Ren, Haoming Luo, Yuqian Hong, Bryan Dai, Joey Zhou, Kai Qiu, Zhirong Wu, and Chong Luo. Logic-rl: Unleashing llm reasoning with rule-based reinforcement learning. *arXiv preprint arXiv:2502.14768*, 2025.
- Wei Xiong, Jiarui Yao, Yuhui Xu, Bo Pang, Lei Wang, Doyen Sahoo, Junnan Li, Nan Jiang, Tong Zhang, Caiming Xiong, and Hanze Dong. A minimalist approach to llm reasoning: from rejection sampling to reinforce. *arXiv preprint arXiv:2504.11343*, 2025.
- An Yang, Baosong Yang, Binyuan Hui, Bo Zheng, Bowen Yu, Chang Zhou, Chengpeng Li, Chengyuan Li, Dayiheng Liu, Fei Huang, Guanting Dong, Haoran Wei, Huan Lin, Jialong Tang, Jialin Wang, Jian Yang, Jianhong Tu, Jianwei Zhang, Jianxin Ma, Jin Xu, Jingren Zhou, Jinze Bai, Jinzheng He, Junyang Lin, Kai Dang, Keming Lu, Keqin Chen, Kexin Yang, Mei Li, Mingfeng Xue, Na Ni, Pei Zhang, Peng Wang, Ru Peng, Rui Men, Ruize Gao, Runji Lin, Shijie Wang, Shuai Bai, Sinan Tan, Tianhang Zhu, Tianhao Li, Tianyu Liu, Wenbin Ge, Xiaodong Deng, Xiaohuan Zhou, Xingzhang Ren, Xinyu Zhang, Xipin Wei, Xuancheng Ren, Yang Fan, Yang Yao, Yichang Zhang, Yu Wan, Yunfei Chu, Yuqiong Liu, Zeyu Cui, Zhenru Zhang, and Zhihao Fan. Qwen2 technical report. arXiv preprint arXiv:2407.10671, 2024.
- John Yang, Kilian Leret, Carlos E Jimenez, Alexander Wettig, Kabir Khandpur, Yanzhe Zhang, Binyuan Hui, Ofir Press, Ludwig Schmidt, and Diyi Yang. Swe-smith: Scaling data for software engineering agents. *arXiv* preprint arXiv:2504.21798, 2025.
- Yixin Ye, Zhen Huang, Yang Xiao, Ethan Chern, Shijie Xia, and Pengfei Liu. Limo: Less is more for reasoning. *arXiv preprint arXiv:2502.03387*, 2025.
- Edward Yeo, Yuxuan Tong, Morry Niu, Graham Neubig, and Xiang Yue. Demystifying long chain-of-thought reasoning in llms. *arXiv preprint arXiv:2502.03373*, 2025.
- Qiying Yu, Zheng Zhang, Ruofei Zhu, Yufeng Yuan, Xiaochen Zuo, Yu Yue, Weinan Dai, Tiantian Fan, Gaohong Liu, Lingjun Liu, et al. Dapo: An open-source llm reinforcement learning system at scale. *arXiv preprint arXiv:2503.14476*, 2025.
- Yang Yue, Zhiqi Chen, Rui Lu, Andrew Zhao, Zhaokai Wang, Shiji Song, and Gao Huang. Does reinforcement learning really incentivize reasoning capacity in llms beyond the base model? *arXiv* preprint arXiv:2504.13837, 2025.
- Jixiao Zhang and Chunsheng Zuo. Grpo-lead: A difficulty-aware reinforcement learning approach for concise mathematical reasoning in language models. *arXiv preprint arXiv:2504.09696*, 2025.
- Xiaojiang Zhang, Jinghui Wang, Zifei Cheng, Wenhao Zhuang, Zheng Lin, Minglei Zhang, Shaojie Wang, Yinghan Cui, Chao Wang, Junyi Peng, et al. Srpo: A cross-domain implementation of large-scale reinforcement learning on llm. *arXiv preprint arXiv:2504.14286*, 2025.

A ILLUSTRATION OF SRL ON SWE TASKS.

In Figure 5, we illustrate how we approach the SWE tasks with SRL. We take two consecutive action-observation pairs from the expert trajectories in the given SFT data as context. We prompt the LLM to first think in monologues and then conclude with its action. Sequence similarity score is thus computed between model action and expert action in the trajectory.

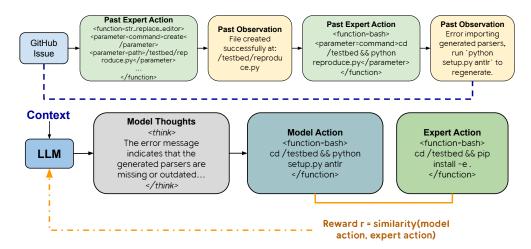


Figure 5: Illustration of applying SRL to SWE tasks. We take two rounds of the past expert actions and corresponding observations in context and prompt the LLM to think before reaches the next action. The model action is compared with the expert action to compute the sequence similarity reward.

B EXPERIMENT DETAILS

We use VeRL as the training framework for GRPO and implement SRL based on it. In Table 5, 6 and 7, we detail the hyperparameters that we used for SFT, GRPO and SRL. The batch size for GRPO is set to be smaller (128) than SRL (512) as the filter rate is high based on final answer correctness, for most of the questions are very difficult. We maintain a smaller batch size to avoid repeated sampling and achieve better performance. We further lay out the system prompts we used for SRL. Experiments were conducted on GPU clusters to the similar level of NVIDIA H100 80GB GPU.

Table 5: Supervised fine-tuning hyperparameters.

Data type	bf16
Learning rate	5e-6
Global batch size	64
Scheduler	Cosine
Warmup ratio	0.3
Num train epochs	3

System Prompt for SRL

You are a helpful assistant for solving mathematical problems. A user will provide a math problem, which may include a partial solution. Your task is to continue the solution by providing the very next logical step.

A user will ask you to solve a task. You should first draft your thinking process (inner monologue). Then, generate the solution. Your response format must follow the template below: <think> Your thoughts or/and draft, like working through an exercise on scratch paper. Be as casual and as long as you want until you are confident to generate a correct solution.
 /think> Provide only the single, next step to continue the solution. Do not solve the entire problem.

Table 6: GRPO hyperparameters.

Batch size	128
Max grad norm	1.0
Data type	bf16
Learning rate	5e-7
Warmup ratio	0.0
Rollout temperature	1.0
Rollout number	8
KL loss coeff	0

Table 7: SRL hyperparameters.

512
1.0
bf16
5e-7
0.0
1.0
8
0

C ADDITIONAL OUTPUT EXAMPLES

Example 2: Planning

Question: Josh decides to try flipping a house. He buys a house for \$80,000 and then puts in \$50,000 in repairs. This increased the value of the house by 150%. How much profit did he make? <think>

To determine Josh's profit, we need to follow these steps:

- 1. Calculate the new value of the house after the repairs.
- 2. Subtract the initial cost of the house and the repair costs from the new value to find the profit.

First, let's calculate the new value of the house after the repairs. The value of the house increased by 150% after the repairs.

The initial cost of the house is \$80,000, and the repairs cost \$50,000. The total cost before the value increase is:

$$80,000 + 50,000 = 130,000$$

The value of the house increased by 150% of its initial value of \$80,000. So, the increase in value is:

$$150\% \times 80,000 = 1.5 \times 80,000 = 120,000$$

Therefore, the new value of the house is:

$$80,000 + 120,000 = 200,000$$

Now, we subtract the total cost (initial cost of the house plus the repair costs) from the new value of the house to find the profit:

$$200,000 - 130,000 = 70,000$$

So, Josh's profit is:

Example 3: Self-verification

Question: Which one of these constellations is not located along the Milky Way in the sky? Answer Choices: (A) Perseus (B) Cygnus (C) Scorpius (D) Leo

First, let's recall which constellations are along the Milky Way. Perseus, Cygnus, and Scorpius are all well-known constellations that lie along the Milky Way. However, Leo is not one of them; it is located in the northern sky and does not align with the Milky Way. Therefore, the correct answer is (D) Leo.

<think>

 To double-check, I can recall the general locations of these constellations:

- Perseus is in the northern sky and is part of the Perseus Arm of the Milky Way.
- Cygnus is also in the northern sky and is part of the Orion Arm of the Milky Way.
- Scorpius is in the southern sky and is part of the Sagittarius Arm of the Milky Way.
- Leo is in the northern sky and is not aligned with the Milky Way.
- </think>

The answer is (D) Leo.

D LLM USAGE

Throughout the preparation of this manuscript, we utilized a LLM as an assistant with refinement of writing. Its role was strictly limited to improving grammar, rephrasing sentences for clarity, and correcting spelling. All intellectual contributions, ideas, and conclusions presented in this paper are entirely the work of the authors.