Deliberate Reasoning for LLMs as Structure AWARE PLANNING WITH ACCURATE WORLD MODEL

Anonymous authors

Paper under double-blind review

ABSTRACT

Enhancing the reasoning capabilities of large language models (LLMs) remains a key challenge, especially for tasks that require complex, multi-step decisionmaking. Humans excel at these tasks by leveraging deliberate planning with an internal world model to simulate the potential outcomes of various actions. Inspired by this, we propose a novel multi-step reasoning framework for LLMs, referred to as Structure-aware Planning with Accurate World Model (SWAP). Unlike previous approaches that rely solely on Chain-of-Thought (CoT) reasoning in natural language, SWAP incorporates structural information to guide the reasoning process via a world model and provides a soft verification mechanism over the steps. Moreover, SWAP overcomes the challenge of accurate world state predictions in complex reasoning tasks by introducing a Generator-Discriminator architecture, which enables more reliable world modeling. Specifically, the generator predicts the next state, and the discriminator ensures alignment with the logical consistency required by the problem context. SWAP also encourages the policy model to explore a broad range of potential actions to prevent premature convergence. By resolving the bottlenecks of generation diversity for both actions and states using diversity-based modeling (DBM) and improving discrimination accuracy through contrastive ranking (CR), SWAP significantly enhances the reasoning performance of LLMs. We evaluate SWAP across diverse reasoning-intensive benchmarks including math reasoning, logical reasoning, and coding tasks. Extensive experiments demonstrate that SWAP achieves substantial improvements over the baselines and consistently outperforms existing methods.

031 032 033

034

004

010 011

012

013

014

015

016

017

018

019

021

023

024

025

026

027

028

029

1 INTRODUCTION

Large language models (LLMs) (OpenAI & et al., 2024; Dubey et al., 2024) have made remark-035 able progress in many fields. However, their ability to perform complex reasoning remains limited (Huang & Chang, 2023). Achieving human-level problem solving is viewed as the next milestone in 037 Artificial General Intelligence (AGI) (OpenAI, 2024). Unlike human cognition, the inference time of LLMs for reasoning tasks primarily depends on the number of input and output tokens rather than the complexity of the problem. For instance, while humans require multiple attempts, calculations, 040 and verification to solve difficult math problems, LLMs immediately begin generating responses 041 after reading the question. This indicates that they are not actually "thinking" but merely using intu-042 ition, *i.e.*, predicting the next token based on previous ones. In fact, there are two systems of thinking 043 in human mind (Kahneman, 2011): System 1 operates automatically and quickly, with little effort 044 and no sense of voluntary control; and System 2 allocates attention to the effortful mental activities 045 that demand it. In this paper, we aim to enhance the complex reasoning capabilities of LLMs, *i.e.*, turning thinking time into better outcome, with a planning-based approach that emulates System 2. 046

Recently, planning and decision-making frameworks (Yao et al., 2022) have been introduced into reasoning tasks for LLMs, where the model is required not only to propose actions but also to make adjustments based on feedback from the environment. However, in many real-world scenarios, environment feedback is either unavailable or difficult to scale. Inspired by human perception (Johnson-Laird, 1983; 2010), an internal world model is introduced to enable the model to simulate actions and their effects on the world state for deliberate planning (LeCun, 2022). Some recent approaches have demonstrated success in planning and reasoning tasks with a world model (Guan et al., 2023; Hao et al., 2023), which is implemented by prompting the same LLM with in-context demonstra-

tions. However, their system performance still falls short of expectations in complex tasks, since
constructing an accurate world model is inherently challenging. The predicted future state from an
inaccurate world model may lead to sub-optimal or even incorrect decisions. To address this limitation, we fine-tune the model using a Generator-Discriminator architecture. Furthermore, we achieve
substantial improvements on a diverse set of reasoning benchmarks by resolving the bottlenecks of
generation diversity and discrimination accuracy.

060 On the other hand, the Chain-of-Thought (CoT) approach (Wei et al., 2022), due to its high flexibility 061 and scalability, is widely adopted to enhance the reasoning capability of LLMs. However, it relies 062 purely on natural language, lacking an effective verification mechanism. To address this issue, 063 formal methods have been proposed, such as using first-order logic (Pan et al., 2023) or programs 064 (Chen et al., 2022). Nevertheless, these formal methods are often limited in their expressiveness for a variety of tasks (Yang et al., 2024). In this paper, we propose a semi-formal approach that introduces 065 structural information into the reasoning process, which provides a soft verification mechanism 066 for CoTs. These structures (Dalvi et al., 2021) describe how given premises are used to generate 067 intermediate conclusions that help validate the correctness of a particular answer. In our framework, 068 the multi-step reasoning process involves constructing a structure, *i.e.*, the policy model proposes 069 actions, and the world model predicts the next state and updates the structure. Specifically, new statements in the next state are introduced and linked to existing ones through entailment relations. 071 When the reasoning is complete, the system has built an entailment graph from the given premises 072 to the final answer, which itself serves as a justification of the reasoning process. 073

- 074 Specifically, our contributions mainly include:
 - We introduce **structure-aware planning**, which incorporates entailment graphs into multistep reasoning tasks. These graphs demonstrate how premises lead to intermediate conclusions and validate the correctness of the final answer, adding coherence and logical verification to the reasoning process.
 - Our framework, SWAP, augments the LLM with an **accurate world model**, which is implemented using a Generator-Discriminator architecture. In addition, we resolve the bottlenecks of generation diversity and discrimination accuracy with diversity-based modelling and contrastive process supervision, respectively.
 - Experiments on a diverse set of benchmarks, including math reasoning, logical reasoning and coding, show that SWAP is a general framework that achieves substantial improvements over recent popular reasoning and planning methods for LLMs.

2 RELATED WORK

075

076

077

078

079

081

082

084

085

087

088

Existing works that use advanced planning methods to enhance the multi-step problem-solving ca-090 pabilities of LLMs can be categorized into three types: re-ranking (Ni et al., 2023; Wang et al., 091 2023b; Li et al., 2023; Lei et al., 2024), iterative correction (Madaan et al., 2023; Shinn et al., 2023; 092 Yao et al., 2022; Chen et al., 2024a) and tree search (Chaffin et al., 2022; Gu et al., 2023; Hao et al., 2023; Yao et al., 2023; Zhou et al., 2023). Despite differences in their design, all these methods fun-094 damentally rely on a **discriminator** to evaluate the planning steps. Recent research (Huang et al., 095 2023; Chen et al., 2024b) has demonstrated that the discriminator plays a more crucial role than the 096 planning methods themselves. Consequently, using in-context learning to prompt the same LM as 097 both generator and discriminator may not sufficiently improve the model performance on complex 098 reasoning tasks.

099 To address this issue, prior research has explored various methodologies for designing the discrim-100 inator (or reward model). There are two primary types of reward models: Outcome Reward Model 101 (ORM) and Process Reward Model (PRM). The ORM evaluates the fully generated solution by 102 assigning a single scalar confidence score. Its training relies on outcome supervision by compar-103 ing generated answers with the ground truth. In contrast, the PRM (Lightman et al., 2023; Yuan 104 et al., 2024; Tian et al., 2024) provides stepwise rewards throughout the reasoning process, assign-105 ing a scalar confidence score to each intermediate steps. Empirical evidence shows that, compared with outcome supervision, process supervision ensures the correctness of each step, providing more 106 benefits to multi-step reasoning (Lightman et al., 2023). However, the training of PRM requires 107 process supervision, which is hard to obtain, e.g., collecting process annotation from humans is 108 inherently not scalable. Although recent research (Wang et al., 2023a; Luo et al., 2024) has increasingly explored automatic process annotations using tree search, training an effective PRM remains 110 challenging, as from a mathematical perspective, it assigns a **numerical value** within [0, 1] to each 111 state **independently**. To overcome this problem, we propose a novel strategy for **automatic ranking** 112 annotation, *i.e.*, given the current context and a set of candidate options, selecting the best option based on relative quality. Our ranking strategy offers significant advantages over traditional PRMs: 113 1) it emphasizes relative quality, making it more robust to noise; 2) it simplifies optimization and en-114 hances generalization. Notably, our high-quality automatic ranking annotation method is non-trivial 115 as it systemically incorporates three key factors: 1) structural information; 2) correctness; and 3) 116 semantical equivalence. 117

Furthermore, we notice that although some reasoning processes are inherently **non-linear**, existing methods mainly follow a linear problem-solving manner. Language models are expected to implicitly infer the non-linear structure from the linear representation of the reasoning process, which proves challenging for complex reasoning tasks (Ribeiro et al., 2023). To help the model, we integrate **structural information** into the reasoning process which explicitly represents the reasoning structure within the context. These structures provide the language model with additional **guidance** and **control**, enabling extra capabilities such as symbolic learning and verification.

125 126

127

133

134 135

136

137

138

139

140

141

142

143

144

3 PRELIMINARIES

128 3.1 TASK FORMULATION

When solving complex reasoning tasks that require multiple steps, LLMs must plan intelligently, anticipating future state and guiding their reasoning towards the desired outcome. We formulate this task as a Markov Decision Process (MDP) represented by (S, A, P, score) in which:

- State $s_t \in S$: Represents the current state, *i.e.*, all known or inferred information in the reasoning process. The initial state s_0 is extracted from the given context.
- Action $a_t \in \mathcal{A}$: Denotes a single action (produced by policy generator \mathcal{P}_{π_G}), *i.e.*, deriving new information or making inference based on current state, resulting in a state transition.
- **Transition probability** $\mathcal{P}(s_{t+1}|s_t, a_t)$: Describes the probability of transitioning to the next state s_{t+1} after taking action a_t in state s_t . We construct an enhanced world model (with generator \mathcal{P}_{wmG} and discriminator \mathcal{P}_{wmD}) to simulate the state change.
- Scoring function $score(a_t|s_t)$: Quantifies the quality of an action a_t given current state s_t . This function guides the reasoning process by prioritizing actions that are more likely to yield correct final answers. We adopt a ranking-based approach (with policy discriminator \mathcal{P}_{π_D}) instead of assigning explicit numerical scores (with PRM).

This MDP framework provides a foundation for applying planning methods to enhance the multistep reasoning capabilities of LLMs. Each reasoning step is viewed as a decision-making process, where the model generates the next action based on current state. By updating their parameters, the models gradually learn the optimal policy for each state, improving the overall performance of reasoning. Additionally, the policy must balance exploiting known optimal actions and exploring new action spaces, guided by the scoring function to help the model make the best choices.

150 151

152

3.2 STRUCTURED REASONING AS ENTAILMENT GRAPH CONSTRUCTION

The key innovation that distinguishes our approach from related work is conceptualizing the multi-153 step reasoning process $(s_0, a_0, s_1, \cdots, a_{T-1}, s_T)$ as entailment graph (Dalvi et al., 2021) con-154 struction (Figure 1), which outlines how the premises in s_0 lead to intermediate conclusions, ulti-155 mately validating the final answer in s_T . Formally, let $\mathcal{G} = (\mathcal{V}, \mathcal{E})$ represent the structure, where \mathcal{V} 156 is the set of nodes, with each node $v \in \mathcal{V}$ representing a statement, e.g., evidence, an assumption, 157 or a lemma/rule; \mathcal{E} is the set of directed (hyper) edges, where each (hyper) edge $e = (\mathcal{V}_{src}, \mathcal{V}_{tgt}) \in \mathcal{E}$ 158 represents an entailment relation from a source node set $V_{src} \subseteq V$ (the premises) to a target node set 159 $\mathcal{V}_{tgt} \subseteq \mathcal{V}$ (the conclusions). 160

Given s_0 , the world model generator \mathcal{P}_{wmG} first builds the initial graph \mathcal{G}_0 by extracting key statements and their relations. During the reasoning, \mathcal{P}_{wmG} incrementally grows the graph by adding new



Figure 1: SWAP performs multi-step reasoning through structure-aware planning in FOLIO (left) and MATH (right). At each step, given the current state, represented as a graph, and an action, the world model predicts the next state as an updated graph.

nodes and edges, ultimately forming \mathcal{G}_T , which includes the final answer. The generation process is 188 sampling-based, with the world model discriminator \mathcal{P}_{wmD} making decision at each step. Structural 189 verification is also introduced to ensure the quality of the graph. For simplicity, let us denote the 190 state with structural information as (s, \mathcal{G}) . Incorporating this structure provides two main benefits: 1) the policy model can make more informed decisions using the structural information; and 2) the world model can predict more accurate next state. 192

4 STRUCTURE-AWARE PLANNING WITH ACCURATE WORLD MODEL

4.1 FRAMEWORK

185

186 187

191

193 194

195 196

197

In this section, we present the Structure-aware Planning with Accurate World Model framework (SWAP) that enables LLMs to systematically construct and utilize an entailment graph for solving 199 a wide range of reasoning tasks. We use \mathcal{P}_{π_G} and \mathcal{P}_{π_D} to denote the policy generator & discrimi-200 nator, \mathcal{P}_{wmG} and \mathcal{P}_{wmD} to denote the world model generator & discriminator, and \mathcal{P}_{c} to denote the 201 controller based on pre-trained LMs. We consider $Q, A, G, H, s, \mathcal{G}, a$ as language sequences, *i.e.*, 202 $Q = (Q[1], \dots, Q[L])$ where each Q[l] is a token, so that $\mathcal{P}(Q) = \prod_{l=1}^{L} \mathcal{P}(Q[l]|Q[1..l-1])$. We 203 use (s, \mathcal{G}) to denote the state with structural information, and $c = (G, H, s, \mathcal{G})$ to denote the context 204 of goal G, plan H and state (s, \mathcal{G}) . 205

206 For notational convenience, we define the generation process as gen(model, input, N) where N 207 is the number of generations, and the discrimination process as dis (model, input, b) where b is the number of preserved candidates. To search potential plans and actions, we simulate the future 208 situations of the current state using the world model. Specifically, we use sim(c, t) to denote the 209 simulation starting from (s, \mathcal{G}) up to step t given the goal G and plan H from the context c. 210

211 Algorithm 1 outlines the **workflow**. Given a reasoning question Q, the world model generator 212 $\mathcal{P}_{wmG}(G, s_0, \mathcal{G}_0|Q)$ first generates the goal G and the initial state (s_0, \mathcal{G}_0) . The policy generator 213 then proposes a set of plans H by sampling N times from $\mathcal{P}_{\pi_G}(H|G, s_0, \mathcal{G}_0)$. The top b candidate plans are selected by the policy discriminator $\mathcal{P}_{\pi_{\mathrm{D}}}$ based on the simulation results (s_T, \mathcal{G}_T) under 214 each plan. Given the goal G, selected plan H and current state $(s_{t-1}, \mathcal{G}_{t-1})$, multi-step reasoning 215 at step t begins with the policy generator sampling N times from $\mathcal{P}_{\pi_G}(a_{t-1}|G, H, s_{t-1}, \mathcal{G}_{t-1})$ as

Algorithm 1 SWAP $(Q, \mathcal{P}_{\pi_{G}}, \mathcal{P}_{\pi_{D}}, \mathcal{P}_{wmG}, \mathcal{P}_{wmD}, \mathcal{P}_{c}, N, T, b)$	Algorithm 2 StatePredict $(\mathcal{C}, \mathcal{D}, \mathcal{P}_{wmG}, \mathcal{P}_{wmD}, \mathcal{P}_{c}, N)$
Require: Reasoning question Q, policy generator	Require: Context pool C , completed pool D ,
& discriminator $P_{\pi_{\rm G}}$ and $P_{\pi_{\rm D}}$, world model gen-	world model generator & discriminator P_{wmG}
erator & discriminator \mathcal{P}_{wmG} and \mathcal{P}_{wmD} , con-	and \mathcal{P}_{wmD} , controller \mathcal{P}_c , generation number
troller \mathcal{P}_{c} , generation number limit N, step limit	limit N.
T, breadth limit b .	parallel for $i = 1, \cdots, b$ do
$\mathcal{D} \leftarrow \{\}$	$(G, H, s, \mathcal{G}, a) = \mathcal{C}_i$
$G, s_0, \mathcal{G}_0 \leftarrow \texttt{gen}(\mathcal{P}_{wmG}, Q, 1)$	$\{(s'_i, \mathcal{G}'_i)\}_{i=1}^N \leftarrow \operatorname{gen}(\mathcal{P}_{wmG}, (s, \mathcal{G}, a),$
$\mathcal{C} \leftarrow \{(G, H, s_0, \mathcal{G}_0) \mid H \in \text{gen}(\mathcal{P}_{\pi_G},$	(N)
$(G, s_0, \mathcal{G}_0), N)\}$	$(s', \mathcal{G}') \leftarrow dis(\mathcal{P}_{wmD},$
$\mathcal{C} \leftarrow \operatorname{dis}(\mathcal{P}_{\pi_{\mathrm{D}}}, \{\operatorname{sim}(c, T) \mid c \in \mathcal{C}\}, b)$	$\{(s, \mathcal{G}, a, s'_i, \mathcal{G}'_i)\}_{i=1}^N, 1\}$
for $t = 1, \cdots, T$ do	$A \leftarrow \operatorname{qen}(\mathcal{P}_{c}, (G, s', \mathcal{G}'), 1)$
if $b = 0$ then break	if $A \neq None$ then
end if	$\mathcal{D}_{add}((G, s', G', A)) \triangleright \text{ collect the state}$
$\mathcal{C} \leftarrow \{ (G, H, s, \mathcal{G}, a) \mid (G, H, s, \mathcal{G}) \in \mathcal{C}, $	$C \operatorname{pon}(i)$ remove context C:
$a \in \operatorname{gen}(\mathcal{P}_{\pi_{\mathcal{C}}}, (G, H, s, \mathcal{G}), N)\}$	$b \neq b = 1$
$\mathcal{C} \leftarrow \operatorname{dis}(\mathcal{P}_{\tau}, \{\operatorname{sim}(c, t) \mid c \in \mathcal{C}\}, b)$	$1 \rightarrow 0$
StatePredict $(\mathcal{C}, \mathcal{D}, \mathcal{P}_{umc}, \mathcal{P}_{umc}, \mathcal{P}_{umc}, \mathcal{N})$	else (C, H, cl, cl) sur data content cl
end for	$C_i \leftarrow (G, H, S, G) \triangleright$ update context C_i
$A^* \leftarrow \operatorname{dis}(\mathcal{P} \mathcal{D} \ 1)$	
return A^*	ena lor

the next action pool. The policy discriminator $\mathcal{P}_{\pi_{D}}$ then evaluates and selects the top *b* candidate contexts (*G*, *H*, *s*_{*t*-1}, *G*_{*t*-1}, *a*_{*t*-1}) based on simulated states (*s*_{*t*}, *G*_{*t*}).

241 Then the accurate state prediction (Algorithm 2) is performed in parallel for each selected context 242 $(G, H, s_{t-1}, \mathcal{G}_{t-1}, a_{t-1})$. Specifically, the world model generator predicts the next state (s_t, \mathcal{G}_t) 243 by sampling N times from $\mathcal{P}_{wmG}(s_t, \mathcal{G}_t | s_{t-1}, \mathcal{G}_{t-1}, a_{t-1})$. Then the world model discriminator 244 \mathcal{P}_{wmD} selects the top 1 candidate state. Based on the selected (s_t, \mathcal{G}_t) , the controller determines 245 whether to continue reasoning. If reasoning is complete, the controller $\mathcal{P}_{c}(A|G, s_{t}, \mathcal{G}_{t})$ generates the final answer A, stores $(G, s_t, \mathcal{G}_t, A)$ in the completed pool \mathcal{D} , and reduces b by 1. Otherwise, 246 $(G, H, s_t, \mathcal{G}_t)$ will be added to the context pool \mathcal{C} for the next step. The process continues until 247 the step limit T is reached or b becomes 0. Finally, the top answer A^* is selected by the policy 248 discriminator $\mathcal{P}_{\pi_{\mathrm{D}}}$ based on the completed states (with graphs) in \mathcal{D} . 249

250 251

216

4.2 SEEKING DIVERSITY IN ACTION GENERATION AND STATE PREDICTION

We identify two critical bottlenecks (generation diversity and discrimination accuracy) for the 253 Generator-Discriminator (G-D) architecture in SWAP. Improving generation diversity is essential 254 to allow the model to explore a broader solution space, increasing the chances of discovering the 255 global optimal solution. Thus, we propose a **Diversity-based Modelling** (DBM) approach (Figure 256 2). The key idea is to encourage the generator to produce steps that differ from existing ones, thereby 257 mitigating its inherent self-bias and promoting exploration. Compared to related work (Vijayakumar et al., 2016; Hu et al., 2023), DBM offers several advantages: 1) It builds on SFT, enabling 258 an end-to-end learning and scalable to large datasets; 2) It leverages the extensive world knowledge 259 embedded in pre-trained LMs. 260

²⁶¹ Diversity-based Modeling (DBM):

Given the current state (s_t, \mathcal{G}_t) , we use $\mathcal{P}_{\pi_G}^{ori}(a_t|G, H, s_t, \mathcal{G}_t)$ to denote the original distribution learned by supervised fine-tuning on the positive trajectories (that lead to correct final answers) during training. For *n*-th generation, we aim to introduce diversity by considering an additional distribution $\mathcal{P}_{\pi_G}^{sem}(a_t^n|a_t^{1..n-1})$, which represents steps that are semantically similar to those generated previously $a_t^{1..n-1}$. Specifically, the probability of *l*-th token $a_{t,l}^n$ in the *n*-th generation a_t^n is

268

269

$$\mathcal{P}_{\pi_{\rm G}}^{\rm sem}(a_{t,l}^{n}|a_{t}^{1..n-1},a_{t,1..l-1}^{n}) = \frac{1}{n-1}\sum_{j=1}^{n-1}\mathcal{P}_{\pi_{\rm G}}^{\rm sem}(a_{t,l}^{n}|a_{t}^{j},a_{t,1..l-1}^{n}),\tag{1}$$

where a_t^j denotes the *j*-th generation, and for notational simplicity, we **move the token index** to a subscript, so that $a_{t,1..l-1}^n$ denotes the preceding tokens of the *l*-th token $a_{t,l}^n$. We obtain $\mathcal{P}_{\pi_G}^{\text{sem}}(a_l'|a, a_{1..l-1}')$ by using supervised fine-tuning on the training data generated by GPT-40, where *a* and *a'* are pairs of actions that are semantically equivalent.

To encourage diversity, the generator adjusts the original distribution $\mathcal{P}_{\pi_{\rm G}}^{\rm ori}$ by reducing the probability mass assigned to steps that are semantically similar to previous generations, *i.e.*,

$$\mathcal{P}_{\pi_{G}}(a_{t,l}^{n}|G,H,s_{t},\mathcal{G}_{t},a_{t}^{1..n-1},a_{t,1..l-1}^{n}) = \operatorname{Norm}\left(\mathcal{P}_{\pi_{G}}^{\operatorname{ori}}(a_{t,l}^{n}|G,H,s_{t},\mathcal{G}_{t},a_{t,1..l-1}^{n}) - \gamma_{l}\mathcal{P}_{\pi_{G}}^{\operatorname{sem}}(a_{t,l}^{n}|a_{t}^{1..n-1},a_{t,1..l-1}^{n})\right)$$
(2)

where the decay factor, $\gamma_l = \gamma_0 \cdot \alpha^l$ with $\alpha \leq 1$, is introduced to emphasize diversity in early stages of generation while gradually reducing this effect. This ensures that the deduplication effect is stronger initially to explore different paths but weakens over time to avoid drifting too far from plausible solutions, thereby maintaining accuracy. Note that this discussion primarily focuses on action generation, while the process of plan generation using the policy generator, represented as $\mathcal{P}_{\pi_G}(H^n | G, s_0, \mathcal{G}_0, H^{1..n-1})$, follows a similar approach.

The normalization function

$$\operatorname{Norm}(\mathcal{P}) = \frac{\max(\mathcal{P}, 0)}{\mathbf{1}^{\top} \max(\mathcal{P}, 0)}$$
(3)

is applied to discard negative-valued tokens (that resemble previous generations or deviate from the
 intended progression of reasoning) and maintain a diverse and relevant generation. Other alternatives, such as Softmax, can distort the probability of irrelevant tokens by redistributing values across all tokens.

State Prediction Enhancement with Diversity:

By encouraging the generator to produce diverse predictions, we increase the likelihood of overcoming self-biases and discovering a more accurate future state. We then select the top 1 prediction from the diverse options generated. To achieve this, we apply a similar strategy to enhance diversity for state prediction $\mathcal{P}_{wmG}(s_t^n | s_{t-1}, \mathcal{G}_{t-1}, a_{t-1}, s_t^{1..n-1})$, that is,

$$\mathcal{P}_{wmG}(s_{t,l}^{n}|s_{t-1},\mathcal{G}_{t-1},a_{t-1},s_{t}^{1...n-1},s_{t,1..l-1}^{n}) = \operatorname{Norm}\left(\mathcal{P}_{wmG}^{\operatorname{ori}}(s_{t,l}^{n}|s_{t-1},\mathcal{G}_{t-1},a_{t-1},s_{t,1..l-1}^{n}) - \gamma_{l}\mathcal{P}_{wmG}^{\operatorname{sem}}(s_{t,l}^{n}|s_{t}^{1...n-1},s_{t,1..l-1}^{n})\right)$$
(4)

where s_t^j denotes the *j*-th generation, and $s_{t,1..l-1}^n$ is the preceding tokens of the *l*-th token $s_{t,l}^n$. Once the state s_t^n is generated, the corresponding graph \mathcal{G}_t^n is extracted from this state, allowing the model to maintain a consistent representation of entailment relationships as the reasoning progresses.

In addition to diversity-based modeling, we leverage a **dynamic context strategy** to further diversify the generation. This strategy involves randomly reframing the current state to create an alternative



Figure 2: Overview of the proposed Diversity-Based Modeling (DBM) method. The current context is processed by the language model, which is fine-tuned using Ori-LoRA and SemEquiv-LoRA. Previous generations are used to compute the semantic equivalence distribution, which is employed to adjust the original distribution to avoid repetition.

324 context for each step. For example, given the original state (s, \mathcal{G}) , we generate an alternative state 325 (s', \mathcal{G}') , where s' is sampled from the semantic equivalence distributions $\mathcal{P}_{wmG}^{sem}(s'|s)$. The corre-326 sponding graph \mathcal{G}' is then regenerated from s'. Our experiments show that this strategy significantly 327 contributes to generating diverse outputs, enhancing the model's robustness and performance on 328 reasoning tasks.

43 IMPROVING DISCRIMINATION ACCURACY IN REASONING

332 As highlighted in recent works (Huang et al., 2023; Chen et al., 2024b), discrimination accuracy is a critical aspect of various planning methods. However, training an effective PRM remains chal-333 lenging, as mathematically, it assigns a numerical value to each state independently. To address this 334 issue, our discriminator employs **Contrastive Ranking** (CR) to evaluate multiple candidate options 335 simultaneously. By focusing on relative comparisons, the model can effectively identify discrepan-336 cies between options, particularly erroneous parts, thereby simplifying the task. 337

338 **Contrastive Ranking (CR) for Enhanced Evaluation:**

339 To illustrate (Figure 3), given a positive trajectory $[(s_0, \mathcal{G}_0), a_0, \cdots, (s_T, \mathcal{G}_T)]$ that leads to the cor-340 rect final answer, we randomly select an intermediate step t, and finalize K subsequent reasoning 341 processes: $\{[a_t^j, \dots, (s_{T_i}^j, \mathcal{G}_{T_i}^j)]\}_{j=1}^K$, where T_j represents the length of the j-th trajectory. Among 342 these K trajectories, we identify the first erroneous steps in negative trajectories (which lead to incor-343 rect final answers) by determining which steps are semantically different from the positive trajectory 344 and then performing structural verification and N_{veri} completions for outcome verification, *i.e.*, if 345 none of the completions result in the correct final answer, we confirm these steps as erroneous. 346

Given the contrastive process annotations, we define the inputs and outputs of the discriminator 347 while incorporating **meta knowledge** $\mathcal{K}_{\text{meta}}$ to enhance model performance. Specifically, 348

$$E, a_t^{\text{best}} \sim \mathcal{P}_{\pi_{\text{D}}}\left(E, a_t^{\text{best}} \mid \mathcal{K}_{\text{meta}}, G, H, (s_t, \mathcal{G}_t), \{a_t^j, (s_{t+1}^j, \mathcal{G}_{t+1}^j)\}_{j=1}^K\right)$$
(5)

350 351 352

353

357

360

361 362

364

365

366 367

368 369

370

372

373

349

330

331

$$E, s_{t+1}^{\text{best}}, \mathcal{G}_{t+1}^{\text{best}} \sim \mathcal{P}_{\text{wmD}}\left(E, s_{t+1}^{\text{best}}, \mathcal{G}_{t+1}^{\text{best}} \mid \mathcal{K}_{\text{meta}}, (s_t, \mathcal{G}_t), a_t, \{(s_{t+1}^j, \mathcal{G}_{t+1}^j)\}_{j=1}^K\right)$$
(6)

where $(s_{t+1}^j, \mathcal{G}_{t+1}^j)$ is used for the selection of action a_t . We avoid using longer future trajecto-354 ries to prevent introducing new errors, which could interfere with action selection. For plan selec-355 tion $\mathcal{P}_{\pi_{\mathrm{D}}}(E, H^{\mathrm{best}} \mid \mathcal{K}_{\mathrm{meta}}, G, (s_0, \mathcal{G}_0), \{H^j, (s_{T_i}^j, \mathcal{G}_{T_i}^j)\}_{j=1}^K)$, we use the simulated completed states 356 $(s_{T_{z}}^{j}, \mathcal{G}_{T_{z}}^{j})$ for the selection of H. The discriminator generates an explanation E, highlighting differ-358 ences between the K future states before making a decision. We fine-tune the discriminator using 359 these explanations through bootstrapping from GPT-40. We use the superscript 'best' to denote the final selected option, *i.e.*, H^{best} , a_t^{best} and $(s_{t+1}^{\text{best}}, \mathcal{G}_{t+1}^{\text{best}})$, and the construction of meta knowledge \mathcal{K}_{meta} based on training data is provided in Appendix A.



374 Figure 3: Overview of our automatic ranking annotation. Starting from a selected step in the positive 375 trajectory, multiple future actions and states are generated to create candidate trajectories. Negative 376 trajectories, which lead to incorrect final answers, are analyzed to identify the first steps that are 377 semantically different from those in the positive trajectory. Structural verification and tree search for outcome are then employed to identify these potential erroneous steps.

382		Math Re	easoning	Logical Reasoning		Coding		
383	Model	GSM8K	MATH	FOLIO	FOLIO ReClor		MBPP	
384			LLaMA3-	8B-Instruct				
385	Zero-shot CoT	70.0 ± 2.0	27.6 ± 0.6	62.1 ± 1.8	57.8 ± 1.4	53.3 ± 0.6	51.8 ± 0.2	
386	Few-shot CoT (4-shot)	72.4 ± 1.8	23.6 ± 0.6	57.2 ± 1.4	52.1 ± 1.1	56.8 ± 0.2	53.6 ± 0.2	
387	SFT-CoT	71.3 ± 1.8	25.4 ± 0.4	66.0 ± 0.8	62.2 ± 0.8	51.6 ± 0.4	51.0 ± 0.3	
200	Self-consistency	74.1 ± 1.2	26.0 ± 0.4	66.2 ± 0.5	60.1 ± 0.6	-	-	
300		75.2 ± 1.1	28.8 ± 0.4	67.1 ± 0.8	60.6 ± 0.8	-	-	
389	RAP	76.0 ± 1.0	28.4 ± 0.3	67.5 ± 0.6	61.3 ± 0.6	-	-	
390	PRM (PRM800K*)	74.6 ± 0.8	28.8 ± 0.2	-	-	-	-	
391	PRM (Math-Shepherd)	76.2 ± 0.8	28.6 ± 0.3	-	-	-	-	
202	SWAP (w/o discriminator)	78.1 ± 1.0	37.3 ± 0.4	69.2 ± 0.8	69.1 ± 0.8	53.1 ± 0.8	53.4 ± 0.6	
392	SWAP	82.7 ± 0.6	42.3 ± 0.3	73.2 ± 0.5	74.1 ± 0.4	$\textbf{57.8} \pm 0.6$	58.6 ± 0.4	
393			Mistral-7	B -Instruct				
394	Zero-shot CoT	23.4 ± 1.8	12.0 ± 0.4	46.8 ± 1.5	38.8 ± 1.0	42.5 ± 0.5	38.8 ± 0.4	
395	Few-shot CoT (4-shot)	47.3 ± 1.6	12.7 ± 0.5	48.6 ± 1.6	36.2 ± 0.8	43.6 ± 0.4	44.8 ± 0.6	
396	SFT-CoT	48.0 ± 1.0	12.6 ± 0.3	52.0 ± 1.0	40.2 ± 0.6	43.8 ± 0.4	46.0 ± 0.4	
207	Self-consistency	52.1 ± 0.8	11.2 ± 0.2	51.2 ± 0.6	42.4 ± 0.4	-	-	
397	ТоТ	49.6 ± 1.2	12.3 ± 0.3	50.2 ± 1.2	40.8 ± 0.8	-	-	
398	RAP	56.1 ± 1.0	13.0 ± 0.2	52.1 ± 0.8	41.6 ± 0.6	-	-	
399	PRM (PRM800K*)	54.2 ± 0.8	14.2 ± 0.2	-	-	-	-	
400	PRM (Math-Shepherd)	55.4 ± 0.6	13.6 ± 0.2	-	-	-	-	
/01	SWAP (w/o discriminator)	54.0 ± 0.8	15.4 ± 0.3	54.0 ± 0.6	45.2 ± 0.6	45.0 ± 0.8	47.0 ± 0.4	
400	SWAP	$\textbf{60.4} \pm 0.6$	18.7 ± 0.2	$\textbf{58.0} \pm 0.3$	$\textbf{49.1} \pm 0.4$	$\textbf{48.4} \pm 0.6$	51.1 ± 0.3	

Table 1: Overall performance comparison across different benchmark datasets. The best performance for each task using the same base model is in bold. Note: We use the filtered PRM800K dataset (Sun et al., 2024) to evaluate performance on the full MATH test set.

During inference, given the discriminator, we apply a **voting** strategy to decide the top *b* candidates as mentioned in Algorithm 1, 2. To further enhance robustness, we reframe the candidate options and reorder them to have **multiple comparisons** within the same group. In addition, we further enhance the discrimination accuracy with **structural verification** on the graphs $\{\mathcal{G}^j\}_{j=1}^K$. Details of these strategies are given in Appendix A.

407 408 409

410

403

404

405

406

381

5 EXPERIMENTS

411 412 5.1 EXPERIMENTAL SETUP

413 We conduct experiments on various types of reasoning tasks. Dataset statistics and examples are 414 provided in Appendix B. For each dataset, we use different types of models (GPT-40 (OpenAI & 415 et al., 2024), DeepSeek-V2 (DeepSeek-AI & et al., 2024), LLaMA3 (Dubey et al., 2024)) to generate 416 multiple trajectories for the training and validation sets. We label the trajectories as positive or 417 negative based on their final answers. To improve the model stability, we augment training questions 418 using GPT-40. Given the positive and negative trajectories of the same question, we automatically generate contrastive process annotations (Figure 3) using DeepSeek-V2. Additionally, to address the 419 class imbalance in contrastive ranking data, we apply pre-processing and post-processing techniques 420 (see Appendix D for details). With the complete training data, SWAP is fine-tuned from LLaMA3-421 8B-Instruct using LoRA (Hu et al., 2021). The parameter settings are as follows: For DBM, $\gamma_0 = 0.7$ 422 and $\alpha = 0.95$. For CR, $N_{\text{veri}} = 3$; we choose $K = \{2, 3\}$ for discriminator training, and during 423 inference, multiple options are divided into groups of size 2 or 3; for meta knowledge, we use 424 $\mathcal{M} = 5$. To ensure the effectiveness of training, we also employ specialized strategies such as 425 curriculum learning and self-improving training (details in Appendix D). During evaluation, we 426 compare our SWAP against popular strategies, CoT, Self-consistency (SC) (Wang et al., 2023b), ToT 427 (Yao et al., 2023), and RAP (Hao et al., 2023)) as well as SFT on CoTs and verification with PRMs 428 (Lightman et al., 2023; Wang et al., 2023a), using different base models (LLaMA3-8B-Instruct and Mistral-7B-Instruct (Jiang et al., 2023)). The number of candidate solutions for self-consistency and 429 PRMs is set to 8. For SWAP, ToT, and RAP (utilizing MCTS), the generation number and step limits 430 are set to 5 and 10, respectively. The number of rollouts (breadth limit) is set as 8. More details 431 about data generation, model training and evaluation are provided in Appendix C, D.

432 5.2 MAIN RESULTS

434 Overall performance is shown in Table 1, with fine-grained results and examples provided in Ap-435 pendix E and Appendix G, respectively. We summarize the key findings as follows:

436 SWAP consistently achieves the best or comparable performance among different methods. 437 One-pass CoT and verification methods, such as self-consistency and PRMs, do not involve search-438 ing through intermediate steps during the reasoning process. In contrast, our framework empowers 439 the model to reason more like humans' conscious planning, which significantly improves perfor-440 mance on multi-step reasoning tasks. This planning ability becomes especially crucial in more chal-441 lenging tasks, where deliberate reasoning is necessary to avoid intermediate errors. For instance, 442 our framework with LLaMA3-8B-Instruct achieves a 14.7% better accuracy compared to CoT (53% relative improvement) on the more difficult MATH dataset, and a 10.3% improvement on GSM8k. 443

444 Structure-aware planning and an accurate world model further enhance the effectiveness of 445 planning in LLMs. Methods such as ToT and RAP, which also incorporate planning or search-based 446 strategies, do not match our approach in performance. They lack the deeper structural understanding 447 and precise state modeling that our framework provides. SWAP explicitly introduces the structure 448 that describes the relationship between key statements, which facilitates both action generation and 449 state prediction. In addition, Diversity-based Modeling (DBM) enables the generator to explore a broader solution space, increasing the likelihood of finding optimal steps. Contrastive Ranking 450 (CR), on the other hand, significantly improves the accuracy of the discriminator by focusing on 451 relative comparisons between candidate solutions. This combination of enhanced exploration and 452 more precise discrimination is key to the substantial performance improvements observed in our 453 experiments, especially on challenging datasets like MATH. 454

455

457

475

476

477

478

479

480

481

482

5.3 ANALYSIS

We investigate the effect of search tree width and depth on overall accuracy, providing insights for both parameter selection and dynamic evaluation across various tasks.

460 The benefit of increasing search tree width, *i.e.*, the number of search attempts per step, becomes marginal after a certain point. For planning-based approaches, the width of the search 461 tree directly influences the thoroughness of exploring the solution space at each step. We analyze 462 the effect of search width on accuracy in SWAP (Figure 4). As shown, there is a consistent upward 463 trend across all datasets. However, the benefits diminish beyond a certain point, e.g., after 5-7 search 464 attempts in FOLIO and GSM8K, since most of the promising options have already been explored. 465 We found that using a search width of 5 offers the best trade-off between computational cost and per-466 formance. We also observed some variability between datasets. GSM8K and MATH show a sharper 467 initial increase in accuracy with fewer search attempts, while FOLIO and HumanEval exhibit a 468 more gradual improvement. This discrepancy likely arises from the variations in task complexity 469 and dataset size. 470

Model performance improves gradually with increasing search tree depth, *i.e.*, the number of searched steps in the trajectory. Another important factor is the search tree depth, which refers to the number of searched steps. We analyze how accuracy changes with search depth in SWAP (Figure 5). For each value of N_{search} , we search and optimize the first N_{search} steps and allow the model to



483 484 Figu





Figure 5: Effect of increasing search tree depth on overall accuracy for different benchmark datasets in SWAP. More searched steps lead to improved accuracy and reduced variance.

Table 2: Ablation studies. The complete framework achieves the highest performance across all tasks, demonstrating that each component contributes positively to overall accuracy.

	Math Re	asoning	Logical I	Reasoning	Coding		
Method	GSM8K	MATH	FOLIO	ReClor	HumanEval	MBPP	
SWAP (Ours)	82.7	42.3	73.2	74.1	57.8	58.6	
w/o structure info	81.2	40.4	72.5	71.8	56.1	57.5	
w/o DBM	79.0	38.5	70.0	71.2	55.0	55.9	
w/o meta knowledge	81.3	40.8	72.4	73.0	56.2	57.1	
w/ PRM instead of CR	81.0	39.1	71.6	72.0	55.9	56.8	
w/o discriminator	78.1	37.3	69.2	69.1	53.1	53.4	

complete the remaining trajectory directly. As seen, accuracy steadily increases with the number of searched steps across all datasets, indicating that our planning brings benefits. Notably, the benefits of planning depend on the difficulty of each stage, as more challenging steps yield greater accuracy improvements after searching. Toward the end of the trajectory, the accuracy curve begins to flatten, and its variance is reduced as the trajectory converges to the optimal one.

5.4 ABLATION STUDY

We analyze the impact of the key components proposed in this paper (Table 2). The complete framework achieves the highest performance across all tasks, demonstrating that each component contributes positively to overall accuracy. Notably, the discriminator has the most significant im-pact by effectively selecting optimal actions and state predictions. The incorporation of structural information is also crucial, particularly for complex reasoning tasks like math and logical reason-ing. DBM enhances generation diversity by promoting the exploration of diverse solution paths, while CR outperforms PRM in multi-step reasoning, as selecting the optimal solution by compar-ing different options is more reliable than scoring each option independently. Finally, incorporating meta knowledge further improves discrimination accuracy. These improvements are consistently observed across different task types.

6 CONCLUSION

In this paper, we introduce SWAP, a novel framework for enhancing the multi-step reasoning capabilities of LLMs through structure-aware planning with an accurate world model. Our approach consistently outperforms existing methods in extensive experiments, demonstrating significant improvements on reasoning-heavy benchmarks, including math, logical reasoning, and coding tasks. In this work, we primarily adopt a re-ranking strategy, as it provides a good balance between computational cost and model performance. For future research, exploring reinforcement learning (RL) methods to enable dynamic interaction with the world model could further optimize LLMs for longterm rewards. Additionally, teaching the model to recognize and correct its own mistakes represents another promising direction, potentially leading to even more robust reasoning capabilities.

540 REFERENCES

- Jacob Austin, Augustus Odena, Maxwell Nye, Maarten Bosma, Henryk Michalewski, David Dohan,
 Ellen Jiang, Carrie Cai, Michael Terry, Quoc Le, et al. Program synthesis with large language
 models. arXiv preprint arXiv:2108.07732, 2021.
- Antoine Chaffin, Vincent Claveau, and Ewa Kijak. PPL-MCTS: Constrained textual generation through discriminator-guided MCTS decoding. In Marine Carpuat, Marie-Catherine de Marneffe, and Ivan Vladimir Meza Ruiz (eds.), Proceedings of the 2022 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, pp. 2953–2967, Seattle, United States, July 2022. Association for Computational Linguistics. doi: 10.18653/v1/2022.naacl-main.215. URL https://aclanthology.org/2022. naacl-main.215.
- 552 Mark Chen, Jerry Tworek, Heewoo Jun, Qiming Yuan, Henrique Ponde de Oliveira Pinto, Jared 553 Kaplan, Harri Edwards, Yuri Burda, Nicholas Joseph, Greg Brockman, Alex Ray, Raul Puri, 554 Gretchen Krueger, Michael Petrov, Heidy Khlaaf, Girish Sastry, Pamela Mishkin, Brooke Chan, Scott Gray, Nick Ryder, Mikhail Pavlov, Alethea Power, Lukasz Kaiser, Mohammad Bavarian, 555 Clemens Winter, Philippe Tillet, Felipe Petroski Such, Dave Cummings, Matthias Plappert, Fo-556 tios Chantzis, Elizabeth Barnes, Ariel Herbert-Voss, William Hebgen Guss, Alex Nichol, Alex Paino, Nikolas Tezak, Jie Tang, Igor Babuschkin, Suchir Balaji, Shantanu Jain, William Saunders, 558 Christopher Hesse, Andrew N. Carr, Jan Leike, Josh Achiam, Vedant Misra, Evan Morikawa, Alec 559 Radford, Matthew Knight, Miles Brundage, Mira Murati, Katie Mayer, Peter Welinder, Bob Mc-Grew, Dario Amodei, Sam McCandlish, Ilya Sutskever, and Wojciech Zaremba. Evaluating large 561 language models trained on code, 2021. 562
- Wenhu Chen, Xueguang Ma, Xinyi Wang, and William W Cohen. Program of thoughts prompt ing: Disentangling computation from reasoning for numerical reasoning tasks. *arXiv preprint arXiv:2211.12588*, 2022.
- 566 Xinyun Chen, Maxwell Lin, Nathanael Schärli, and Denny Zhou. Teaching large language models
 567 to self-debug. In *The Twelfth International Conference on Learning Representations*, 2024a. URL
 568 https://openreview.net/forum?id=KuPixIqPiq.
- Ziru Chen, Michael White, Raymond Mooney, Ali Payani, Yu Su, and Huan Sun. When is tree search useful for llm planning? it depends on the discriminator. *arXiv preprint arXiv:2402.10890*, 2024b.
- Karl Cobbe, Vineet Kosaraju, Mohammad Bavarian, Mark Chen, Heewoo Jun, Lukasz Kaiser,
 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*,
 2021.
- 577
 578
 578
 578
 579
 579
 580
 580
 580
 581
 582
 582
 583
 584
 584
 584
 584
 584
 585
 586
 586
 586
 586
 586
 586
 586
 586
 586
 586
 586
 586
 586
 586
 586
 586
 586
 586
 586
 586
 586
 586
 586
 586
 586
 586
 586
 586
 586
 586
 586
 586
 586
 586
 586
 586
 586
 586
 586
 586
 586
 586
 586
 586
 586
 586
 586
 586
 586
 586
 586
 586
 586
 586
 586
 586
 586
 586
 586
 586
 586
 586
 586
 586
 586
 586
 586
 586
 586
 586
 586
 586
 586
 586
 586
 586
 586
 586
 586
 586
 586
 586
 586
 586
 586
 586
 586
 586
 586
 586
 586
 586
 586
 586
 586
 586
 586
 586
 586
 586
 586
 586
 586
 586
 586
 586
 586
 586
 586
 586
 586
 586
 586
 586
 586
 586
 586
 586
 586
 586
 586
 586
 586
 586
 586
 586
 586
 586
 586
 586
 586
 586
 586
 586
 586
 586
 586
 586
- DeepSeek-AI and et al. Deepseek-v2: A strong, economical, and efficient mixture-of-experts language model, 2024. URL https://arxiv.org/abs/2405.04434.
- Abhimanyu Dubey, Abhinav Jauhri, and et al. The llama 3 herd of models, 2024. URL https: //arxiv.org/abs/2407.21783.
- Yu Gu, Xiang Deng, and Yu Su. Don't generate, discriminate: A proposal for grounding language
 models to real-world environments. In *Proceedings of the 61st Annual Meeting of the Association* for Computational Linguistics (Volume 1: Long Papers), pp. 4928–4949, 2023.
- Lin Guan, Karthik Valmeekam, Sarath Sreedharan, and Subbarao Kambhampati. Leveraging pretrained large language models to construct and utilize world models for model-based task planning. In A. Oh, T. Naumann, A. Globerson, K. Saenko, M. Hardt, and S. Levine (eds.), Advances in Neural Information Processing Systems, volume 36, pp. 79081–79094. Curran Associates, Inc., 2023. URL https://proceedings.neurips.cc/paper_files/paper/2023/ file/f9f54762cbb4fe4dbffdd4f792c31221-Paper-Conference.pdf.

594 Simeng Han, Hailey Schoelkopf, Yilun Zhao, Zhenting Qi, Martin Riddell, Luke Benson, Lucy 595 Sun, Ekaterina Zubova, Yujie Qiao, Matthew Burtell, David Peng, Jonathan Fan, Yixin Liu, Brian 596 Wong, Malcolm Sailor, Ansong Ni, Linyong Nan, Jungo Kasai, Tao Yu, Rui Zhang, Shafiq Joty, 597 Alexander R. Fabbri, Wojciech Kryscinski, Xi Victoria Lin, Caiming Xiong, and Dragomir Radev. 598 Folio: Natural language reasoning with first-order logic. arXiv preprint arXiv:2209.00840, 2022. URL https://arxiv.org/abs/2209.00840. 600 Shibo Hao, Yi Gu, Haodi Ma, Joshua Jiahua Hong, Zhen Wang, Daisy Zhe Wang, and Zhiting Hu. 601 Reasoning with language model is planning with world model. arXiv preprint arXiv:2305.14992, 602 2023. 603 Dan Hendrycks, Collin Burns, Saurav Kadavath, Akul Arora, Steven Basart, Eric Tang, Dawn Song, 604 and Jacob Steinhardt. Measuring mathematical problem solving with the math dataset. NeurIPS, 605 2021. 606 607 Edward J Hu, Yelong Shen, Phillip Wallis, Zeyuan Allen-Zhu, Yuanzhi Li, Shean Wang, Lu Wang, 608 and Weizhu Chen. Lora: Low-rank adaptation of large language models. arXiv preprint 609 arXiv:2106.09685, 2021. 610 Edward J Hu, Moksh Jain, Eric Elmoznino, Younesse Kaddar, Guillaume Lajoie, Yoshua Bengio, 611 and Nikolay Malkin. Amortizing intractable inference in large language models. arXiv preprint 612 arXiv:2310.04363, 2023. 613 614 Jie Huang and Kevin Chen-Chuan Chang. Towards reasoning in large language models: A survey. 615 In 61st Annual Meeting of the Association for Computational Linguistics, ACL 2023, pp. 1049– 616 1065. Association for Computational Linguistics (ACL), 2023. 617 Jie Huang, Xinyun Chen, Swaroop Mishra, Huaixiu Steven Zheng, Adams Wei Yu, Xinying Song, 618 and Denny Zhou. Large language models cannot self-correct reasoning yet. arXiv preprint 619 arXiv:2310.01798, 2023. 620 621 Albert Q Jiang, Alexandre Sablayrolles, Arthur Mensch, Chris Bamford, Devendra Singh Chaplot, 622 Diego de las Casas, Florian Bressand, Gianna Lengyel, Guillaume Lample, Lucile Saulnier, et al. Mistral 7b. arXiv preprint arXiv:2310.06825, 2023. 623 624 Philip N Johnson-Laird. Mental models and human reasoning. Proceedings of the National Academy 625 of Sciences, 107(43):18243-18250, 2010. 626 Philip Nicholas Johnson-Laird. Mental models: Towards a cognitive science of language, inference, 627 and consciousness. Number 6. Harvard University Press, 1983. 628 629 Daniel Kahneman. Thinking, fast and slow. Farrar, Straus and Giroux, 2011. 630 631 Vladimir Karpukhin, Barlas Oguz, Sewon Min, Patrick Lewis, Ledell Wu, Sergey Edunov, Danqi Chen, and Wen-tau Yih. Dense passage retrieval for open-domain question answering. In Bonnie 632 Webber, Trevor Cohn, Yulan He, and Yang Liu (eds.), Proceedings of the 2020 Conference on 633 Empirical Methods in Natural Language Processing (EMNLP), pp. 6769-6781, Online, Novem-634 ber 2020. Association for Computational Linguistics. doi: 10.18653/v1/2020.emnlp-main.550. 635 URL https://aclanthology.org/2020.emnlp-main.550. 636 637 Yann LeCun. A path towards autonomous machine intelligence version 0.9. 2, 2022-06-27. Open 638 *Review*, 62(1):1–62, 2022. 639 Bin Lei, Yi Zhang, Shan Zuo, Ali Payani, and Caiwen Ding. Macm: Utilizing a multi-agent 640 system for condition mining in solving complex mathematical problems, 2024. URL https: 641 //arxiv.org/abs/2404.04735. 642 643 Yifei Li, Zeqi Lin, Shizhuo Zhang, Qiang Fu, Bei Chen, Jian-Guang Lou, and Weizhu Chen. Mak-644 ing language models better reasoners with step-aware verifier. In Anna Rogers, Jordan Boyd-Graber, and Naoaki Okazaki (eds.), Proceedings of the 61st Annual Meeting of the Association 645 for Computational Linguistics (Volume 1: Long Papers), pp. 5315–5333, Toronto, Canada, July 646 2023. Association for Computational Linguistics. doi: 10.18653/v1/2023.acl-long.291. URL 647

https://aclanthology.org/2023.acl-long.291.

651

685

686

687

688

- Hunter Lightman, Vineet Kosaraju, Yura Burda, Harri Edwards, Bowen Baker, Teddy Lee, Jan Leike, John Schulman, Ilya Sutskever, and Karl Cobbe. Let's verify step by step. *arXiv preprint arXiv:2305.20050*, 2023.
- Liangchen Luo, Yinxiao Liu, Rosanne Liu, Samrat Phatale, Harsh Lara, Yunxuan Li, Lei Shu, Yun
 Zhu, Lei Meng, Jiao Sun, et al. Improve mathematical reasoning in language models by automated
 process supervision. *arXiv preprint arXiv:2406.06592*, 2024.
- Aman Madaan, Niket Tandon, Prakhar Gupta, Skyler Hallinan, Luyu Gao, Sarah Wiegreffe, Uri
 Alon, Nouha Dziri, Shrimai Prabhumoye, Yiming Yang, Shashank Gupta, Bodhisattwa Prasad
 Majumder, Katherine Hermann, Sean Welleck, Amir Yazdanbakhsh, and Peter Clark. Self-refine:
 Iterative refinement with self-feedback. In *Thirty-seventh Conference on Neural Information Processing Systems*, 2023. URL https://openreview.net/forum?id=S37hOerQLB.
- Ansong Ni, Srini Iyer, Dragomir Radev, Veselin Stoyanov, Wen-Tau Yih, Sida Wang, and Xi Victoria Lin. LEVER: Learning to verify language-to-code generation with execution. In Andreas Krause, Emma Brunskill, Kyunghyun Cho, Barbara Engelhardt, Sivan Sabato, and Jonathan Scarlett (eds.), *Proceedings of the 40th International Conference on Machine Learning*, volume 202 of *Proceedings of Machine Learning Research*, pp. 26106–26128. PMLR, 23–29 Jul 2023. URL https://proceedings.mlr.press/v202/ni23b.html.
- 666 667 OpenAI. Openai ol system card, 2024. URL https://cdn.openai.com/ ol-system-card.pdf.
- OpenAI and et al. Gpt-4 technical report, 2024. URL https://arxiv.org/abs/2303.
 08774.
- Liangming Pan, Alon Albalak, Xinyi Wang, and William Yang Wang. Logic-lm: Empowering large language models with symbolic solvers for faithful logical reasoning. *arXiv preprint arXiv:2305.12295*, 2023.
- Danilo Ribeiro, Shen Wang, Xiaofei Ma, Henry Zhu, Rui Dong, Deguang Kong, Juliette Burger, Anjelica Ramos, William Wang, Zhiheng Huang, et al. Street: A multi-task structured reasoning and explanation benchmark. *arXiv preprint arXiv:2302.06729*, 2023.
- Noah Shinn, Federico Cassano, Ashwin Gopinath, Karthik R Narasimhan, and Shunyu Yao. Re flexion: language agents with verbal reinforcement learning. In *Thirty-seventh Conference on Neural Information Processing Systems*, 2023. URL https://openreview.net/forum?
 id=vAElhFcKW6.
- Zhiqing Sun, Longhui Yu, Yikang Shen, Weiyang Liu, Yiming Yang, Sean Welleck, and Chuang
 Gan. Easy-to-hard generalization: Scalable alignment beyond human supervision. *arXiv preprint arXiv:2403.09472*, 2024.
 - Ye Tian, Baolin Peng, Linfeng Song, Lifeng Jin, Dian Yu, Haitao Mi, and Dong Yu. Toward self-improvement of llms via imagination, searching, and criticizing. *arXiv preprint arXiv:2404.12253*, 2024.
- Ashwin K Vijayakumar, Michael Cogswell, Ramprasath R Selvaraju, Qing Sun, Stefan Lee, David
 Crandall, and Dhruv Batra. Diverse beam search: Decoding diverse solutions from neural se quence models. arXiv preprint arXiv:1610.02424, 2016.
- Peiyi Wang, Lei Li, Zhihong Shao, RX Xu, Damai Dai, Yifei Li, Deli Chen, Y Wu, and Zhifang
 Sui. Math-shepherd: A label-free step-by-step verifier for llms in mathematical reasoning. *arXiv preprint arXiv:2312.08935*, 2023a.
- Kuezhi Wang, Jason Wei, Dale Schuurmans, Quoc V Le, Ed H. Chi, Sharan Narang, Aakanksha
 Chowdhery, and Denny Zhou. Self-consistency improves chain of thought reasoning in language
 models. In *The Eleventh International Conference on Learning Representations*, 2023b. URL
 https://openreview.net/forum?id=1PL1NIMMrw.
- Jason Wei, Xuezhi Wang, Dale Schuurmans, Maarten Bosma, Fei Xia, Ed Chi, Quoc V Le, Denny Zhou, et al. Chain-of-thought prompting elicits reasoning in large language models. *Advances in neural information processing systems*, 35:24824–24837, 2022.

702 703 704	Yuan Yang, Siheng Xiong, Ali Payani, Ehsan Shareghi, and Faramarz Fekri. Can llms reason in the wild with programs? <i>arXiv preprint arXiv:2406.13764</i> , 2024.
704 705 706	Shunyu Yao, Jeffrey Zhao, Dian Yu, Nan Du, Izhak Shafran, Karthik Narasimhan, and Yuan Cao. React: Synergizing reasoning and acting in language models. <i>arXiv preprint arXiv:2210.03629</i> ,
707	2022.
708	Shunyu Yao, Dian Yu, Jeffrey Zhao, Izhak Shafran, Thomas L Griffiths, Yuan Cao, and Karthik
709	Narasimhan. Tree of thoughts: deliberate problem solving with large language models. In Pro-
710	ceedings of the 37th International Conference on Neural Information Processing Systems, pp.
711	11809–11822, 2023.
712	Weihao Yu, Zihang Jiang, Yanfei Dong, and Jiashi Feng. Reclor: A reading comprehension dataset
714	requiring logical reasoning. In International Conference on Learning Representations (ICLR),
715	April 2020.
716	Lifan Yuan, Gangu Cui, Hanbin Wang, Ning Ding, Xingyao Wang, Jia Deng, Boji Shan, Huimin
717 718	Chen, Ruobing Xie, Yankai Lin, et al. Advancing llm reasoning generalists with preference trees. <i>arXiv preprint arXiv:2404.02078</i> , 2024.
719	$\mathbf{X} = \mathbf{X}$
720	agent tree search unifies reasoning acting and planning in language models 2023
721	agent nee searen annes reasoning acting and planning in language models, 2023.
722	
723	
724	
725	
726	
728	
729	
730	
731	
732	
733	
734	
735	
736	
737	
730	
740	
741	
742	
743	
744	
745	
746	
747	
748	
749	
750	
752	
753	
754	
755	

A Additional Methodological Details

The **meta knowledge** \mathcal{K}_{meta} , which helps verify answers and identify errors, is derived from training questions. Formally, $\mathcal{K}_{meta} = \operatorname{concat}_m(\mathcal{K}_m)$, where \mathcal{K}_m represents stored knowledge from the *m*th training sample, and we select the top \mathcal{M} samples based on the cosine similarity between the training query embedding q_m and the test query embedding q.

We fine-tune our discriminator using the contrastive process annotations, which helps it accurately identify subtle differences between trajectories and improve its ability to distinguish between valid and invalid reasoning steps. To further enhance robustness during inference, we randomly **group the candidate options** and reorder them, then apply a **voting** strategy to determine the final ranking to decide the top *b* candidates as mentioned in Algorithm 1, 2. This approach ensures that the model is not biased by specific sequences and provides a more reliable assessment of the best candidate.

In addition, we introduce structural verification for generated entailment graphs G to further enhance discrimination. Key steps involves: 1) Syntax Verification: Validates the format of nodes and edges. 2) Node Dependency Analysis: Examines the dependencies between nodes (assumptions, lemmas, facts, or inferences derived from prior nodes). 3) Cycle Detection: Ensures acyclic structures to maintain logical consistency. 4) Redundancy Check: Detects redundant or disconnected nodes. All of them are implemented according to standard graph algorithms.

774 775

B DATASET OVERVIEW

776 777

In this section, we present statistics and examples for all benchmark datasets used in our study. We 778 consider GSM8K (Cobbe et al., 2021), MATH (Hendrycks et al., 2021) for math reasoning, FOLIO 779 (Han et al., 2022), ReClor (Yu et al., 2020) for logical reasoning, and HumanEval (Chen et al., 780 2021), MBPP (Austin et al., 2021) for coding. For GSM8K, there are 7,473 training samples and 781 1,319 test samples. For MATH, there are 7,500 training samples and 5,000 test samples. For FOLIO, 782 the training and validation sets consist of 1,001 and 203 samples, respectively. For ReClor, we use 783 4,638 training samples, 500 validation samples (used as test set, as the original test set answers are 784 not publicly available), and 1,000 test samples. HumanEval contains 164 test samples, and since 785 it lacks a training set, we use the entire MBPP dataset (after format transformation) for training. MBPP consists of 374 training samples, 90 validation samples, and 500 test samples. 786

787 788

Example - GSM8K

789	Problem: Weng earns \$12 an hour for babysitting. Yesterday, she just did 50 minutes of babysitting. How
790	much did she earn?
791	Solution:
792	Weng earns $12/60 = \$ < < 12/60 = 0.2 > 0.2$ per minute
793	Working 50 minutes, she earned 0.2 x 50 = $\$ < 0.2 \times 50 = 10 > 10$.
794	#### 10
795	
796	Problem: Janet hires six employees. Four of them are warehouse workers who make \$15/hour, and the other
797	everyone works 25 days a month and 8 hours a day, how much does lanet owe total for their wages and taxes.
798	for one month?
799	
800	Solution:
801	First figure out how many hours each worker works per month by multiplying the number of days they work
802	by the number of hours a day they work: 25 days $*$ 8 hours/day = $<<25*8=200>>200$ hours
803	Then calculate how much one warehouse worker makes per month by multiplying their hourly rate by the number of hours they work: 200 hours $\frac{15}{500} = \frac{5}{200} + \frac{15}{2000} = \frac{15}$
804	Then multiply that number by 4 to find out how much all the warehouse workers make: $3000/$ worker * 4
805	workers = $\$ < 3000 * 4 = 12000 > 12,000$
806	Now multiply the hours each manager works (also 200) by their hourly wage to find out how much one
807	manager makes per month: 200 hours $\$20$ /hour = $\$<200 \times 20=4000 >>4,000$
808	Now multiply one manager's wages by the number of managers (2) to find their total wage amount: $$4,000$
000	manager * 2 managers = $\phi < 4000 * 2=000 > 0,000$
003	<<8000+12000=20000>>20,000

ax+3,

x-5

810 Now multiply the total wage bill by 10% to find how much the FICA taxes are: 20,000 * .1 =811 <<20000*.1=2000>>2,000 812 Now add the total wage bill to the total tax amount to find the grand total: 2,000 + 20,000 =813 <<2000+20000=22000>>22,000 #### 22000 814

if x > 2,

if $-2 \le x \le 2$,

Example - MATH

Problem: Let f(x) =

2x - bif x < -2. Find a + b if the piecewise function is continuous (which means that its graph can be drawn without lifting your pencil from the paper).

Solution:

825 826 827

828 829

830

831

832 833

834

835

836

837

815 816

817

818

819

820

821

822 823

824

- For the piecewise function to be continuous, the cases must "meet" at 2 and -2. For example, ax + 3 and x-5 must be equal when x=2. This implies a(2)+3=2-5, which we solve to get $2a = -6 \Rightarrow a = -3$. Similarly, x - 5 and 2x - b must be equal when x = -2. Substituting, we get -2-5=2(-2)-b, which implies b=3. So a+b=-3+3=0
- **Problem:** Square ABCD has its center at (8, -8) and has an area of 4 square units. The top side of the square is horizontal. The square is then dilated with the dilation center at (0,0) and a scale factor of 2. What are the coordinates of the vertex of the image of square ABCD that is farthest from the origin? Give your answer as an ordered pair.

Solution:

With the center of dilation at the origin and a scale factor of 2, all the coordinates of square ABCD are twice the coordinates of its preimage. The preimage has an area of 4 square units, so its side length is 2 units. Since the center of the preimage is at (8, -8), the four vertices of the preimage are at (7, -9), (7, -7), (9, -7) and (9, -9). The point (9, -9) is the farthest from the origin on the preimage, so the point farthest from the origin on the image of square ABCD is (18, -18)

838 839 840

841

850

853

854

855

Example - FOLIO

842 Problem:

843 Premises: All customers in James' family who subscribe to AMC A-List are eligible to watch three movies every week 844

- without any additional fees. 845
- Some of the customers in James' family go to the cinema every week. 846
- Customers in James' family subscribe to AMC A-List or HBO service. 847
- Customers in James' family who prefer TV series will not watch TV series in cinemas.
- 848 All customers in James' family who subscribe to HBO services prefer TV series to movies.
- Lily is in James' family; she watches TV series in cinemas. 849

Conclusion:

851 Lily goes to cinemas every week or watches 3 movies every week without any additional fees. 852

Solution: True

Problem:

Premises:

- 856 If a legislator is found guilty of stealing government funds, they will be suspended from office.
- 857 Tiffany T. Alston was a legislator in Maryland's House of Delegates from 2011 to 2013.
- 858 Tiffany T. Alston was found guilty of stealing government funds in 2012.

859 Conclusion: 860

Tiffany T. Alston went to prison for stealing government funds.

862 Solution: Uncertain

863

861

864	Example DeClar
865	Example - ReClor
866	Problem: Paula will visit the deptist tomorrow morning only if Bill goes golfing in the morning. Bill will not as calling
867	Paula will visit the definist tomorrow morning only it bill goes goining in the morning. Bill will not go goining unless Damien agrees to go golfing too. However, Damien has decided not to go golfing. Therefore, Paula
868	will not be visiting the dentist tomorrow morning.
869	
870	0. If Marge goes to the bank today, Lauren will not cash her check tomorrow. Marge will not wash her car
871	unless it is sunny. However, it is sunny, so Marge will wash her car and go shopping with Lauren.
872	1. Kevin will wash his car tomorrow only if Brittany has to go visit her grandmother. Unless Aunt Susan has
873	errands. Kevin will not wash his car tomorrow.
874	2. Renee will do her homework tonight if there is nothing good on television and if her neighbors do not have
875	a party. Although, there is something good on television; her neighbors are also having a party. Ttherefore,
876	Renee will attend the party.
877	3. Maddie will plan a picnic only if one of her friends, Lisa or Kenny, will come. Kenny will not come to the
878	picnic, but Lisa will. I therefore, Maddle will plan a picnic.
879	The pattern of reasoning displayed above most closely parallels which of the following?
880	
881	Solution: 1
882	
883	Problem:
884	the need for stricter safety standards for the oil industry. Since the industry refuses to take action it is the
885	national government that must regulate industry safety standards. In particular, the government has to at least
886	require oil companies to put double hulls on their tankers and to assume financial responsibility for accidents.
887	Industry representative: The industry alone should be responsible for devising safety standards because of its
888	expertise in handling oil and its understanding of the cost entailed. Implementing the double-hull proposal is
880	not currently feasible because it creates new safety issues. Furthermore, the cost would be burdensome to the industry and consumers
200	industry and consumers.
201	0. Government safety regulations are developed in a process of negotiation with industry leaders and
091	independent experts.
092	1. Environmental concerns outweigh all financial considerations when developing safety standards.
093	2. The requirement of two hulls on oil tankers, although initially costly, will save money over time by
094	3 The only effective sources of increased stringency in safety standards for oil tankers are action by the
090	industry itself or national government regulation.
896	
097	Which one of the following is an assumption on which the argument of the environmentalist depends?
898	
899	Solution: 3
900	
901	Example - HumanEval
902	Problem
903	from typing import List
904	
905	
906	def separate_paren_groups(paren_string: str) -> List[str]:
907	input to this function is a string containing multiple groups of nested parentheses. Your goal is to separate those group into separate strings and return the list of those.
908	Separate trose group into separate surings and return the first of those. Separate groups are balanced (each open brace is properly closed) and not nested within each other
909	Ignore any spaces in the input string.
910	>>> separate_paren_groups('() (()) (()())')
911	['0', '(0)', '(00)']
912	
913	Solution:
914	result = []
915	current_string = []
916	$current_depth = 0$
917	
	for c in paren_string:

918	if a V'
919	n c = -(.)
920	current_string.append(c)
921	elif $c = ')'$:
922	$current_depth = 1$
923	current_string.append(c)
924	
925	if current_depth == 0: result append(" join(current string))
926	current_string_clear()
927	
028	return result
020	
020	Droblom
031	def count up to(n):
022	"Implement a function that takes an non-negative integer and returns an array of the first n
932	integers that are prime numbers and less than n.
933	for example:
934	$count_up_to(5) => [2,3]$
935	$count_up_to(11) => [2,3,3,7]$
936	$count_up_{10}(0) => [1]$ $count_up_{10}(20) => [2 3 5 7 11 13 17 19]$
937	count up to(1) => []
938	$count_up_to(18) => [2,3,5,7,11,13,17]$
939	m
940	
941	Solution:
942	primes = [] for i in range(2, n):
943	is prime = True
944	for j in range(2, i):
945	if i $\%$ j == 0:
946	is_prime = False
947	break
948	II IS_prime:
949	return primes
950	retuin printes
951	
952	Example MBDD
953	
954	Problem: Write a function to find the minimum cost path to reach (m, n) from $(0, 0)$ for the given cost matrix cost
955	and a position (m, n) in cost[]].
956	
957	assert min_cost([[1, 2, 3], [4, 8, 2], [1, 5, 3]], 2, 2) == 8

assert min_cost([[2, 3, 4], [5, 9, 3], [2, 6, 4]], 2, 2) == 12 assert min_cost([[3, 4, 5], [6, 10, 4], [3, 7, 5]], 2, 2) == 16

Solution:

958

959 960

```
R = 3
961
          C = 3
962
          def min_cost(cost, m, n):
963
               tc = [[0 for x in range(C)] for x in range(R)]
964
               tc[0][0] = cost[0][0]
965
               for i in range(1, m+1):
                    tc[i][0] = tc[i-1][0] + cost[i][0]
966
               for j in range(1, n+1):
967
                    tc[0][j] = tc[0][j-1] + cost[0][j]
968
               for i in range(1, m+1):
969
                    for j in range(1, n+1):
970
                        tc[i][j] = min(tc[i-1][j-1], tc[i-1][j], tc[i][j-1]) + cost[i][j]
971
               return tc[m][n]
```

Problem:

972 973

977

978

979 980

981

982

983

984

985

986

987

988

974 Write a function to count the longest repeating subsequences such that the two subsequences don't have same string characters at same positions.
976 (1.1.1) (1.

assert find_longest_repeating_subseq("AABEBCDD") == 3
assert find_longest_repeating_subseq("aabb") == 2
assert find_longest_repeating_subseq("aab") == 1

Solution:

```
def find_longest_repeating_subseq(str):

n = len(str)

dp = [[0 \text{ for } k \text{ in } range(n+1)] \text{ for } l \text{ in } range(n+1)]

for i in range(1, n+1):

for j in range(1, n+1):

if (str[i-1] == str[j-1] \text{ and } i != j):

dp[i][j] = 1 + dp[i-1][j-1]

else:

dp[i][j] = max(dp[i][j-1], dp[i-1][j])

return dp[n][n]
```

989 990 991

992 993

994

995

996

C PROMPTS FOR DATA GENERATION

In this section, we present all the prompts used in our data generation process. These prompts include those for plan generation, action generation, state generation, final answer generation, semantic equivalence data generation, semantic equivalence evaluation, meta-knowledge generation, and contrastive process supervision for plan, action, and state generation.

997 Prompt - Plan Generation 998 Based on the goal, and the initial state (including the graph), propose a plan. Do not solve the problem; just 999 outline the steps for proceeding. 1000 1001 Example: 1002 ### Input: 1003 "Problem": "Solve for a: $\sqrt{4 + \sqrt{16 + 16a}} + \sqrt{1 + \sqrt{1 + a}} = 6$." 1004 "Goal": "Solve a." 1005 "Initial state": "We know that $\sqrt{4 + \sqrt{16 + 16a}} + \sqrt{1 + \sqrt{1 + a}} = 6$." 1006 "Initial graph": {"Statement": {"s1": " $\sqrt{4 + \sqrt{16 + 16a}} + \sqrt{1 + \sqrt{1 + a}} = 6$ "}, "Entailment": {"s1": "Given 1007 condition"}} 1008 ### Output: 1009 "Plan": "To solve a, we begin by simplifying $\sqrt{4 + \sqrt{16 + 16a}}$. This simplification may also help us simplify 1010 the left side of the equation further." 1011 1012 1013 Prompt - Action Generation 1014 Based on the goal, the plan, and the history of actions and states (including graphs), propose the next action. 1015 Only specify the action itself; do not provide the outcome. 1016 1017 Example 1: 1018 ### Input: "Problem": "Solve for a: $\sqrt{4 + \sqrt{16 + 16a}} + \sqrt{1 + \sqrt{1 + a}} = 6$." 1019 1020 "Goal": "Solve a." "Initial state": "We know that $\sqrt{4 + \sqrt{16 + 16a}} + \sqrt{1 + \sqrt{1 + a}} = 6$." 1021 1022 "Initial graph": {"Statement": {"s1": " $\sqrt{4} + \sqrt{16 + 16a} + \sqrt{1 + \sqrt{1 + a}} = 6$ "}, "Entailment": {"s1": "Given 1023 condition"}} 1024 "Plan": "To solve a, we begin by simplifying $\sqrt{4} + \sqrt{16 + 16a}$. This simplification may also help us simplify 1025 the left side of the equation further."

	### Output:
1027	"Action 1": "Simplify $\sqrt{4 + \sqrt{16 + 16a}}$."
1028	
1029	
1030	Example 2: ### Input:
1031	""" input. "Problem": "Solve for a: $\sqrt{4 + \sqrt{16 + 16a}} + \sqrt{1 + \sqrt{1 + a}} = 6$ "
1032	"Goal": "Solve to a ."
1034	"Initial state": "We know that $\sqrt{4 + \sqrt{16 + 16a}} + \sqrt{1 + \sqrt{1 + a}} = 6$."
1035	"Initial graph": {"Statement": {"s1": " $\sqrt{4 + \sqrt{16 + 16a}} + \sqrt{1 + \sqrt{1 + a}} = 6$ "}, "Entailment": {"s1": "Given
1036	condition"}}
1037 1038	"Plan": "To solve a, we begin by simplifying $\sqrt{4 + \sqrt{16 + 16a}}$. This simplification may also help us simplify the left side of the equation further."
1039	"Action 1": "Simplify $\sqrt{4 + \sqrt{16 + 16a}}$."
1040	"State 1": "We start by simplifying the expression under the square root. $n\$ Tirst, observe that: n
1041	$16 + 16a = 16(1+a) \setminus \text{nSubstitute this back:} \sqrt{4 + \sqrt{16(1+a)}} \setminus \text{nNow, simplify } \sqrt{16(1+a)} \setminus \text{n}$
1042	$\sqrt{16(1+a)} = 4\sqrt{1+a} \ln s$ the expression becomes: $\sqrt{4+4\sqrt{1+a}} \ln s$, now the 4 from
1043	inside the square root: $\ln\sqrt{4(1+\sqrt{1+a})} \ln$ the square root of 4, we get 2: $\ln 2\sqrt{1+\sqrt{1+a}} \ln S_0$
1044	the simplified form of the expression is: $\frac{1}{\sqrt{1+\sqrt{1+a}}}$
1046	"Graph 1": {"Statement": {"s?": " $\sqrt{4} + \sqrt{16 + 16a} = 2\sqrt{1 + \sqrt{1 + a}}$ "} "Entailment": {"s?": "Identity
1047	equation"}
1048	
1049	$\frac{\#\#\#}{2}$ Output:
1050	"Action 2": "Combine like terms $\sqrt{1 + \sqrt{1 + a}}$."
1051	
1052	Prompt - State Generation
1053	Based on the goal, the plan, and the history of actions and states (including graphs), predict the next state
1054	based on the goal, the plan, and the instory of actions and states (including graphs), predict the next state.
1054 1055	Additionally, extract the graph that illustrates key statements and their entailment relations.
1054 1055 1056	Additionally, extract the graph that illustrates key statements and their entailment relations.
1054 1055 1056 1057	Additionally, extract the graph that illustrates key statements and their entailment relations. Example 1: #### Input:
1054 1055 1056 1057 1058 1059	Additionally, extract the graph that illustrates key statements and their entailment relations. Example 1: ### Input: "Problem": "Solve for $a: \sqrt{4 + \sqrt{16 + 16a}} + \sqrt{1 + \sqrt{1 + a}} = 6$."
1054 1055 1056 1057 1058 1059 1060	Additionally, extract the graph that illustrates key statements and their entailment relations. Example 1: ### Input: "Problem": "Solve for $a: \sqrt{4 + \sqrt{16 + 16a}} + \sqrt{1 + \sqrt{1 + a}} = 6$." "Goal": "Solve a ."
1054 1055 1056 1057 1058 1059 1060 1061 1062	Additionally, extract the graph that illustrates key statements and their entailment relations. Example 1: ### Input: "Problem": "Solve for $a: \sqrt{4 + \sqrt{16 + 16a}} + \sqrt{1 + \sqrt{1 + a}} = 6$." "Goal": "Solve a ." "Initial state": "We know that $\sqrt{4 + \sqrt{16 + 16a}} + \sqrt{1 + \sqrt{1 + a}} = 6$." "Initial graph": {"Statement": {"s1": " $\sqrt{4 + \sqrt{16 + 16a}} + \sqrt{1 + \sqrt{1 + a}} = 6$ "}, "Entailment": {"s1": "Given condition"}}
1054 1055 1056 1057 1058 1059 1060 1061 1062 1063	Additionally, extract the graph that illustrates key statements and their entailment relations. Example 1: ### Input: "Problem": "Solve for $a: \sqrt{4 + \sqrt{16 + 16a}} + \sqrt{1 + \sqrt{1 + a}} = 6$." "Goal": "Solve $a.$ " "Initial state": "We know that $\sqrt{4 + \sqrt{16 + 16a}} + \sqrt{1 + \sqrt{1 + a}} = 6$." "Initial graph": {"Statement": {"s1": " $\sqrt{4 + \sqrt{16 + 16a}} + \sqrt{1 + \sqrt{1 + a}} = 6$ ", "Entailment": {"s1": "Given condition"}} "Plan": "To solve a , we begin by simplifying $\sqrt{4 + \sqrt{16 + 16a}}$. This simplification may also help us simplify the laft side of the acoustion further "
1054 1055 1056 1057 1058 1059 1060 1061 1062 1063 1064	Additionally, extract the graph that illustrates key statements and states (including graphs), predict the lickt state. Example 1: ### Input: "Problem": "Solve for $a: \sqrt{4 + \sqrt{16 + 16a}} + \sqrt{1 + \sqrt{1 + a}} = 6$." "Goal": "Solve a ." "Initial state": "We know that $\sqrt{4 + \sqrt{16 + 16a}} + \sqrt{1 + \sqrt{1 + a}} = 6$." "Initial graph": {"Statement": {"s1": " $\sqrt{4 + \sqrt{16 + 16a}} + \sqrt{1 + \sqrt{1 + a}} = 6$ "}, "Entailment": {"s1": "Given condition"}} "Plan": "To solve a , we begin by simplifying $\sqrt{4 + \sqrt{16 + 16a}}$. This simplification may also help us simplify the left side of the equation further."
1054 1055 1056 1057 1058 1059 1060 1061 1062 1063 1064 1065 1066	Additionally, extract the graph that illustrates key statements and their entailment relations. Example 1: ### Input: "Problem": "Solve for $a: \sqrt{4 + \sqrt{16 + 16a}} + \sqrt{1 + \sqrt{1 + a}} = 6$." "Goal": "Solve $a.$ " "Initial state": "We know that $\sqrt{4 + \sqrt{16 + 16a}} + \sqrt{1 + \sqrt{1 + a}} = 6$." "Initial graph": {"Statement": {"s1": " $\sqrt{4 + \sqrt{16 + 16a}} + \sqrt{1 + \sqrt{1 + a}} = 6$ "}, "Entailment": {"s1": "Given condition"}} "Plan": "To solve a , we begin by simplifying $\sqrt{4 + \sqrt{16 + 16a}}$. This simplification may also help us simplify the left side of the equation further." "Action 1": "Simplify $\sqrt{4 + \sqrt{16 + 16a}}$."
1054 1055 1057 1058 1059 1060 1061 1062 1063 1064 1065 1066 1067	Additionally, extract the graph that illustrates key statements and their entailment relations. Example 1: ### Input: "Problem": "Solve for $a: \sqrt{4 + \sqrt{16 + 16a}} + \sqrt{1 + \sqrt{1 + a}} = 6$." "Goal": "Solve $a.$ " "Initial state": "We know that $\sqrt{4 + \sqrt{16 + 16a}} + \sqrt{1 + \sqrt{1 + a}} = 6$." "Initial graph": {"Statement": {"s1": " $\sqrt{4 + \sqrt{16 + 16a}} + \sqrt{1 + \sqrt{1 + a}} = 6$ "}, "Entailment": {"s1": "Given condition"}} "Plan": "To solve a , we begin by simplifying $\sqrt{4 + \sqrt{16 + 16a}}$. This simplification may also help us simplify the left side of the equation further." "Action 1": "Simplify $\sqrt{4 + \sqrt{16 + 16a}}$."
1054 1055 1056 1057 1058 1059 1060 1061 1062 1063 1064 1065 1066 1067 1068	Additionally, extract the graph that illustrates key statements and their entailment relations. Example 1: ### Input: "Problem": "Solve for $a: \sqrt{4 + \sqrt{16 + 16a}} + \sqrt{1 + \sqrt{1 + a}} = 6$." "Goal": "Solve $a.$ " "Initial state": "We know that $\sqrt{4 + \sqrt{16 + 16a}} + \sqrt{1 + \sqrt{1 + a}} = 6$." "Initial graph": {"Statement": {"s1": " $\sqrt{4 + \sqrt{16 + 16a}} + \sqrt{1 + \sqrt{1 + a}} = 6$ "}, "Entailment": {"s1": "Given condition"}} "Plan": "To solve a , we begin by simplifying $\sqrt{4 + \sqrt{16 + 16a}}$. This simplification may also help us simplify the left side of the equation further." "Action 1": "Simplify $\sqrt{4 + \sqrt{16 + 16a}}$." ### Output: "State 1": "We start by simplifying the expression under the square root.\n\nFirst, observe that:\n
1054 1055 1057 1058 1059 1060 1061 1062 1063 1064 1065 1066 1067 1068 1069	Additionally, extract the graph that illustrates key statements and their entailment relations. Example 1: ### Input: "Problem": "Solve for a: $\sqrt{4 + \sqrt{16 + 16a}} + \sqrt{1 + \sqrt{1 + a}} = 6$." "Goal": "Solve a." "Initial state": "We know that $\sqrt{4 + \sqrt{16 + 16a}} + \sqrt{1 + \sqrt{1 + a}} = 6$." "Initial graph": {"Statement": {"s1": " $\sqrt{4 + \sqrt{16 + 16a}} + \sqrt{1 + \sqrt{1 + a}} = 6$ "}, "Entailment": {"s1": "Given condition"}} "Plan": "To solve a, we begin by simplifying $\sqrt{4 + \sqrt{16 + 16a}}$. This simplification may also help us simplify the left side of the equation further." "Action 1": "Simplify $\sqrt{4 + \sqrt{16 + 16a}}$." ### Output: "State 1": "We start by simplifying the expression under the square root.\n\nFirst, observe that:\n $16 + 16a = 16(1 + a) \setminus n$ Substitute this back: $\ln \sqrt{4 + \sqrt{16(1 + a)}} \setminus n$ Now, simplify $\sqrt{16(1 + a)}: \setminus n$
1054 1055 1057 1058 1059 1060 1061 1062 1063 1064 1065 1066 1067 1068 1069 1070	Additionally, extract the graph that illustrates key statements and states (including graphs), predict the next state. Additionally, extract the graph that illustrates key statements and their entailment relations. Example 1: ### Input: "Problem": "Solve for a: $\sqrt{4 + \sqrt{16 + 16a}} + \sqrt{1 + \sqrt{1 + a}} = 6$." "Goal": "Solve a." "Initial state": "We know that $\sqrt{4 + \sqrt{16 + 16a}} + \sqrt{1 + \sqrt{1 + a}} = 6$." "Initial graph": {"Statement": {"s1": " $\sqrt{4 + \sqrt{16 + 16a}} + \sqrt{1 + \sqrt{1 + a}} = 6$ "}, "Entailment": {"s1": "Given condition"}} "Plan": "To solve a, we begin by simplifying $\sqrt{4 + \sqrt{16 + 16a}}$. This simplification may also help us simplify the left side of the equation further." "Action 1": "Simplify $\sqrt{4 + \sqrt{16 + 16a}}$." ### Output: "State 1": "We start by simplifying the expression under the square root.\n\nFirst, observe that:\n $16 + 16a = 16(1 + a)$ \nSubstitute this back:\n $\sqrt{4 + \sqrt{16(1 + a)}}$ \n\nNow, simplify $\sqrt{16(1 + a)}$:\n $\sqrt{16(1 + a)} = 4\sqrt{1 + a}$ \n\nSo the expression becomes:\n $\sqrt{4 + 4\sqrt{1 + a}}$ \n\nNext, factor out the 4 from
1054 1055 1056 1057 1058 1059 1060 1061 1062 1063 1064 1065 1066 1067 1068 1069 1070 1071	Additionally, extract the graph that illustrates key statements and states (including graphs), product the next state. Additionally, extract the graph that illustrates key statements and their entailment relations. Example 1: ### Input: "Problem": "Solve for a : $\sqrt{4 + \sqrt{16 + 16a}} + \sqrt{1 + \sqrt{1 + a}} = 6$." "Goal": "Solve a ." "Initial state": "We know that $\sqrt{4 + \sqrt{16 + 16a}} + \sqrt{1 + \sqrt{1 + a}} = 6$." "Initial graph": {"Statement": {"s1": " $\sqrt{4 + \sqrt{16 + 16a}} + \sqrt{1 + \sqrt{1 + a}} = 6$ "}, "Entailment": {"s1": "Given condition"}} "Plan": "To solve a , we begin by simplifying $\sqrt{4 + \sqrt{16 + 16a}}$. This simplification may also help us simplify the left side of the equation further." "Action 1": "Simplify $\sqrt{4 + \sqrt{16 + 16a}}$." #### Output: "State 1": "We start by simplifying the expression under the square root.\n\nFirst, observe that:\n $16 + 16a = 16(1 + a)$ \nSubstitute this back:\n $\sqrt{4 + \sqrt{16(1 + a)}}$ \n\nNow, simplify $\sqrt{16(1 + a)}$:\n $\sqrt{16(1 + a)} = 4\sqrt{1 + a}$ \n\nSo the expression becomes:\n $\sqrt{4 + 4\sqrt{1 + a}}$ \n\nNext, factor out the 4 from inside the square root:\n $\sqrt{4(1 + \sqrt{1 + a})}$ \n\nTaking the square root of 4, we get 2:\n $2\sqrt{1 + \sqrt{1 + a}}$ \nSo,
1054 1055 1057 1058 1059 1060 1061 1062 1063 1064 1065 1066 1067 1068 1069 1070 1071 1072	Additionally, extract the graph that illustrates key statements and their entailment relations. Example 1: ### Input: "Problem": "Solve for $a: \sqrt{4 + \sqrt{16 + 16a}} + \sqrt{1 + \sqrt{1 + a}} = 6$." "Goal": "Solve $a.$ " "Initial state": "We know that $\sqrt{4 + \sqrt{16 + 16a}} + \sqrt{1 + \sqrt{1 + a}} = 6$." "Initial graph": {"Statement": {"s1": " $\sqrt{4 + \sqrt{16 + 16a}} + \sqrt{1 + \sqrt{1 + a}} = 6$ "}, "Entailment": {"s1": "Given condition"}} "Plan": "To solve a , we begin by simplifying $\sqrt{4 + \sqrt{16 + 16a}}$. This simplification may also help us simplify the left side of the equation further." "Action 1": "Simplify $\sqrt{4 + \sqrt{16 + 16a}}$." ### Output: "State 1": "We start by simplifying the expression under the square root.\n\nFirst, observe that:\n $16 + 16a = 16(1 + a)$ \nSubstitute this back:\n $\sqrt{4 + \sqrt{16(1 + a)}}$ \n\nNow, simplify $\sqrt{16(1 + a)}$:\n $\sqrt{16(1 + a)} = 4\sqrt{1 + a}$ \n\nSo the expression becomes:\n $\sqrt{4 + 4\sqrt{1 + a}}$ \n\nNext, factor out the 4 from inside the square root:\n $\sqrt{4(1 + \sqrt{1 + a})}$ \n\nTaking the square root of 4, we get 2:\n2 $\sqrt{1 + \sqrt{1 + a}}$ \nSo, the simplified form of the expression is:\n2 $\sqrt{1 + \sqrt{1 + a}}$."
1054 1055 1057 1058 1059 1060 1061 1062 1063 1064 1065 1066 1067 1068 1069 1070 1071 1072 1073 1074	Discontince goar, the print, and the history of actions and states (including graphs), predict the next state. Additionally, extract the graph that illustrates key statements and their entailment relations. Example 1: ### Input: "Problem": "Solve for $a: \sqrt{4 + \sqrt{16 + 16a}} + \sqrt{1 + \sqrt{1 + a}} = 6$." "Initial state": "We know that $\sqrt{4 + \sqrt{16 + 16a}} + \sqrt{1 + \sqrt{1 + a}} = 6$." "Initial graph": {"Statement": {"s1": " $\sqrt{4 + \sqrt{16 + 16a}} + \sqrt{1 + \sqrt{1 + a}} = 6$ "}, "Entailment": {"s1": "Given condition"}} "Plan": "To solve a , we begin by simplifying $\sqrt{4 + \sqrt{16 + 16a}}$. This simplification may also help us simplify the left side of the equation further." "Action 1": "Simplify $\sqrt{4 + \sqrt{16 + 16a}}$." ### Output: "State 1": "We start by simplifying the expression under the square root.\n\nFirst, observe that:\n $16 + 16a = 16(1 + a)$ \nSubstitute this back:\n $\sqrt{4 + \sqrt{16(1 + a)}}$ \n\nNow, simplify $\sqrt{16(1 + a)}$:\n $\sqrt{16(1 + a)} = 4\sqrt{1 + a}$ \n\nSo the expression becomes:\n $\sqrt{4 + 4\sqrt{1 + a}}$ \n\nNext, factor out the 4 from inside the square root:\n $\sqrt{4(1 + \sqrt{1 + a})}$ \n\nTaking the square root of 4, we get 2:\n $2\sqrt{1 + \sqrt{1 + a}}$ \nSo, the simplified form of the expression is:\n $2\sqrt{1 + \sqrt{1 + a}}$." "Graph 1": {"Statement": {"s2": " $\sqrt{4 + \sqrt{16 + 16a}} = 2\sqrt{1 + \sqrt{1 + a}}$, "Entailment": {"s2": "Identity
1054 1055 1057 1058 1059 1060 1061 1062 1063 1064 1065 1066 1067 1068 1069 1070 1071 1072 1073 1074 1075	based on the goal, the plan, and the instory of actions and states (increding graphs), predict the next state. Additionally, extract the graph that illustrates key statements and their entailment relations. Example 1: ### Input: "Problem": "Solve for $a: \sqrt{4 + \sqrt{16 + 16a}} + \sqrt{1 + \sqrt{1 + a}} = 6$." "Goal": "Solve a ." "Initial state": "We know that $\sqrt{4 + \sqrt{16 + 16a}} + \sqrt{1 + \sqrt{1 + a}} = 6$." "Initial state": "We know that $\sqrt{4 + \sqrt{16 + 16a}} + \sqrt{1 + \sqrt{1 + a}} = 6$." "Initial graph": {"Statement": {"s1": " $\sqrt{4 + \sqrt{16 + 16a}} + \sqrt{1 + \sqrt{1 + a}} = 6$." "Initial state": "We know that $\sqrt{4 + \sqrt{16 + 16a}} + \sqrt{1 + \sqrt{1 + a}} = 6$." "Initial graph": {"Statement": {"s1": " $\sqrt{4 + \sqrt{16 + 16a}} + \sqrt{1 + \sqrt{1 + a}} = 6$." "Initial state": "We know that $\sqrt{4 + \sqrt{16 + 16a}} + \sqrt{1 + \sqrt{16 + 16a}}$. This simplification may also help us simplify the left side of the equation further." "Action 1": "Simplify $\sqrt{4 + \sqrt{16 + 16a}}$." ### Output: "State 1": "We start by simplifying the expression under the square root.\n\nFirst, observe that:\n $16 + 16a = 16(1 + a)$ \nSubstitute this back:\n $\sqrt{4 + \sqrt{16(1 + a)}}$ \n\nNow, simplify $\sqrt{16(1 + a)}$:\n $ \sqrt{16(1 + a)} = 4\sqrt{1 + a}$ \nNo the expression becomes:\n $\sqrt{4 + 4\sqrt{1 + a}}$ \n\nNext, factor out the 4 from inside the square root:\n $\sqrt{4(1 + \sqrt{1 + a})}$ \n\nTaking the square root of 4, we get 2:\n2 $\sqrt{1 + \sqrt{1 + a}}$ \nSo, the simplified form of the expression is:\n2 $\sqrt{1 + \sqrt{1 + a}}$." "Graph 1": {"Statement": {"s2": " $\sqrt{4 + \sqrt{16 + 16a}} = 2\sqrt{1 + \sqrt{1 + a}}$, "Entailment": {"s2": "Identity equation"}}
1054 1055 1057 1058 1059 1060 1061 1062 1063 1064 1065 1066 1067 1068 1069 1070 1071 1072 1073 1074 1075 1076	based on the goal, the plan, and the instory of actions and states (including graphs), predict the first state. Additionally, extract the graph that illustrates key statements and their entailment relations. Example 1: ### Input: "Problem": "Solve for a: $\sqrt{4 + \sqrt{16 + 16a}} + \sqrt{1 + \sqrt{1 + a}} = 6$." "Goal": "Solve a." "Initial state": "We know that $\sqrt{4 + \sqrt{16 + 16a}} + \sqrt{1 + \sqrt{1 + a}} = 6$." "Initial graph": {"Statement": {"s1": " $\sqrt{4 + \sqrt{16 + 16a}} + \sqrt{1 + \sqrt{1 + a}} = 6$ "}, "Entailment": {"s1": "Given condition"}} "Plan": "To solve a, we begin by simplifying $\sqrt{4 + \sqrt{16 + 16a}}$. This simplification may also help us simplify the left side of the equation further." "Action 1": "Simplify $\sqrt{4 + \sqrt{16 + 16a}}$." ### Output: "State 1": "We start by simplifying the expression under the square root.\n\nFirst, observe that:\n $16 + 16a = 16(1 + a)$ \nSubstitute this back:\n $\sqrt{4 + \sqrt{16(1 + a)}}$ \n\nNow, simplify $\sqrt{16(1 + a)}$:\n $\sqrt{16(1 + a)} = 4\sqrt{1 + a}$ \n\nSo the expression becomes:\n $\sqrt{4 + 4\sqrt{1 + a}}$ \n\nNext, factor out the 4 from inside the square root:\n $\sqrt{4(1 + \sqrt{1 + a})}$ \n\nTaking the square root of 4, we get 2:\n2 $\sqrt{1 + \sqrt{1 + a}}$ \nSo, the simplified form of the expression is:\n2 $\sqrt{1 + \sqrt{1 + a}}$." "Graph 1": {"Statement": {"s2": " $\sqrt{4 + \sqrt{16 + 16a}} = 2\sqrt{1 + \sqrt{1 + a}}$, "Entailment": {"s2": "Identity equation"}}
1054 1055 1057 1058 1059 1060 1061 1062 1063 1064 1065 1066 1067 1068 1069 1070 1071 1072 1073 1074 1075 1076 1077	based on the goal, the plan, and the history of actions and states (including graphs), predict the next state. Additionally, extract the graph that illustrates key statements and their entailment relations. Example 1: ### Input: "Problem": "Solve for $a: \sqrt{4 + \sqrt{16 + 16a}} + \sqrt{1 + \sqrt{1 + a}} = 6$." "Initial state": "We know that $\sqrt{4 + \sqrt{16 + 16a}} + \sqrt{1 + \sqrt{1 + a}} = 6$." "Initial graph": {"Statement": {"s1": " $\sqrt{4 + \sqrt{16 + 16a}} + \sqrt{1 + \sqrt{1 + a}} = 6$ "}, "Entailment": {"s1": "Given condition"}} "Plan": "To solve a , we begin by simplifying $\sqrt{4 + \sqrt{16 + 16a}}$. This simplification may also help us simplify the left side of the equation further." "Action 1": "Simplify $\sqrt{4 + \sqrt{16 + 16a}}$." ### Output: "State 1": "We start by simplifying the expression under the square root.\n\nFirst, observe that:\n $16 + 16a = 16(1 + a)$ \nSubstitute this back:\n $\sqrt{4 + \sqrt{16(1 + a)}}$ \n\nNow, simplify $\sqrt{16(1 + a)}$:\n $\sqrt{16(1 + a)} = 4\sqrt{1 + a}$ \n\nSo the expression becomes:\n $\sqrt{4 + 4\sqrt{1 + a}}$ \n\nNext, factor out the 4 from inside the square root:\n $\sqrt{4(1 + \sqrt{1 + a})}$ \n\nTaking the square root of 4, we get 2:\n2 $\sqrt{1 + \sqrt{1 + a}}$ \nSo, the simplified form of the expression is:\n2 $\sqrt{1 + \sqrt{1 + a}}$." "Graph 1": {"Statement": {"s2": " $\sqrt{4 + \sqrt{16 + 16a}} = 2\sqrt{1 + \sqrt{1 + a}}$, "Entailment": {"s2": "Identity equation"}} Example 2:
1054 1055 1057 1058 1059 1060 1061 1062 1063 1064 1065 1066 1067 1068 1069 1070 1071 1072 1073 1074 1075 1076 1077 1078	Placed on the goal, the plan, and the instory of actions and states (including graphs), pleater the next state: Additionally, extract the graph that illustrates key statements and their entailment relations. Example 1: ### Input: "Problem": "Solve for a : $\sqrt{4 + \sqrt{16 + 16a}} + \sqrt{1 + \sqrt{1 + a}} = 6$." "Initial state": "We know that $\sqrt{4 + \sqrt{16 + 16a}} + \sqrt{1 + \sqrt{1 + a}} = 6$." "Initial graph": {"Statement": {"s1": " $\sqrt{4 + \sqrt{16 + 16a}} + \sqrt{1 + \sqrt{1 + a}} = 6$ ", "Entailment": {"s1": "Given condition"}} "Plan": "To solve a , we begin by simplifying $\sqrt{4 + \sqrt{16 + 16a}}$. This simplification may also help us simplify the left side of the equation further." "Action 1": "Simplify $\sqrt{4 + \sqrt{16 + 16a}}$." ### Output: "State 1": "We start by simplifying the expression under the square root.\n\nFirst, observe that:\n $16 + 16a = 16(1 + a) \setminus n$ Substitute this back: $\ln \sqrt{4 + \sqrt{16(1 + a)}} \setminus n \setminus n$ Now, simplify $\sqrt{16(1 + a)} \ge n$ $\sqrt{16(1 + a)} = 4\sqrt{1 + a} \setminus n \setminus n$ So the expression becomes: $\ln \sqrt{4 + 4\sqrt{1 + a}} \setminus n \setminus n$ Next, factor out the 4 from inside the square root: $\ln \sqrt{4(1 + \sqrt{1 + a})} \setminus n \setminus n$ Taking the square root of 4, we get $2: \ln 2\sqrt{1 + \sqrt{1 + a}} \setminus n$ So, the simplified form of the expression is: $\ln 2\sqrt{1 + \sqrt{1 + a}}$." "Graph 1": {"Statement": {"s2": " $\sqrt{4 + \sqrt{16 + 16a}} = 2\sqrt{1 + \sqrt{1 + a}}$, "Entailment": {"s2": "Identity equation"}} Example 2: ### Input:

1080 "Initial state": "We know that $\sqrt{4 + \sqrt{16 + 16a}} + \sqrt{1 + \sqrt{1 + a}} = 6$." 1081 "Initial graph": {"Statement": {"s1": " $\sqrt{4 + \sqrt{16 + 16a}} + \sqrt{1 + \sqrt{1 + a}} = 6$ "}, "Entailment": {"s1": "Given 1082 condition"}} 1083 "Plan": "To solve a, we begin by simplifying $\sqrt{4 + \sqrt{16 + 16a}}$. This simplification may also help us simplify 1084 the left side of the equation further." "Action 1": "Simplify $\sqrt{4 + \sqrt{16 + 16a}}$." 1086 "State 1": "We start by simplifying the expression under the square root.\n\nFirst, observe that:\n 16 + 16a = 16(1+a)\nSubstitute this back: $\sqrt{4} + \sqrt{16(1+a)}$ \n\nNow, simplify $\sqrt{16(1+a)}$:\n 1088 $\sqrt{16(1+a)} = 4\sqrt{1+a} \ln \theta$ expression becomes: $\sqrt{4+4\sqrt{1+a}} \ln \theta$ 1089 1090 | inside the square root: $\ln\sqrt{4(1+\sqrt{1+a})}$ \n\nTaking the square root of 4, we get 2: $\ln2\sqrt{1+\sqrt{1+a}}$ \nSo, 1091 the simplified form of the expression is: $\ln 2\sqrt{1 + \sqrt{1 + a}}$." 1092 "Graph 1": {"Statement": {"s2": " $\sqrt{4 + \sqrt{16 + 16a}} = 2\sqrt{1 + \sqrt{1 + a}}$ "}, "Entailment": {"s2": "Identity 1093 equation"} 1094 "Action 2": "Combine like terms $\sqrt{1 + \sqrt{1 + a}}$." 1095 1096 ### Output: 1097 "State 2": "We can obtain $\sqrt{4 + \sqrt{16 + 16a}} + \sqrt{1 + \sqrt{1 + a}} = 2\sqrt{1 + \sqrt{1 + a}} + \sqrt{1 + \sqrt{1 + a}}$ $= 3\sqrt{1+\sqrt{1+a}}$." 1099 "Graph 2": {"Statement": {"s3": " $\sqrt{4 + \sqrt{16 + 16a}} + \sqrt{1 + \sqrt{1 + a}} = 3\sqrt{1 + \sqrt{1 + a}}$ "}, "Entailment": {" 1100 s3": ["s1", "s2"]}} 1101 1102 1103 Prompt - Final Answer Generation 1104 Based on the goal and the current state (including the graph), determine if the goal has been achieved. If it has , generate the final answer; otherwise, return "Not yet". 1105 1106 Example 1: 1107 ### Input: 1108 "Problem": "Solve for $a: \sqrt{4 + \sqrt{16 + 16a}} + \sqrt{1 + \sqrt{1 + a}} = 6$." 1109 "Goal": "Solve a." 1110 "Initial state": "We know that $\sqrt{4 + \sqrt{16 + 16a}} + \sqrt{1 + \sqrt{1 + a}} = 6$." 1111 "Initial graph": {"Statement": {"s1": " $\sqrt{4 + \sqrt{16 + 16a}} + \sqrt{1 + \sqrt{1 + a}} = 6$ "}, "Entailment": {"s1": "Given 1112 condition"}} 1113 "Plan": "To solve a, we begin by simplifying $\sqrt{4+\sqrt{16+16a}}$. This simplification may also help us simplify 1114 the left side of the equation further." 1115 "Action 1": "Simplify $\sqrt{4 + \sqrt{16 + 16a}}$." 1116 "State 1": "We start by simplifying the expression under the square root.\n\nFirst, observe that:\n 1117 16 + 16a = 16(1+a)\nSubstitute this back: $\sqrt{4} + \sqrt{16(1+a)}$ \n\nNow, simplify $\sqrt{16(1+a)}$: 1118 $\sqrt{16(1+a)} = 4\sqrt{1+a} \ln \sin \theta$ the expression becomes: $\sqrt{4+4\sqrt{1+a}} \ln \sin \theta$ 1119 1120 inside the square root: $\ln\sqrt{4(1+\sqrt{1+a})} \ln \alpha$ in Taking the square root of 4, we get 2: $\ln 2\sqrt{1+\sqrt{1+a}} \ln S_0$, 1121 the simplified form of the expression is: $\ln 2\sqrt{1 + \sqrt{1 + a}}$." 1122 "Graph 1": {"Statement": {"s2": " $\sqrt{4 + \sqrt{16 + 16a}} = 2\sqrt{1 + \sqrt{1 + a}}$ "}, "Entailment": {"s2": "Identity 1123 equation"} 1124 "Action 2": "Combine like terms $\sqrt{1 + \sqrt{1 + a}}$." 1125 "State 2": "We can obtain $\sqrt{4 + \sqrt{16 + 16a}} + \sqrt{1 + \sqrt{1 + a}} = 2\sqrt{1 + \sqrt{1 + a}} + \sqrt{1 + \sqrt{1 + a}}$ 1126 $= 3\sqrt{1+\sqrt{1+a}}$." 1127 "Graph 2": {"Statement": {"s3": " $\sqrt{4 + \sqrt{16 + 16a}} + \sqrt{1 + \sqrt{1 + a}} = 3\sqrt{1 + \sqrt{1 + a}}$ "}, "Entailment": {" 1128 s3": ["s1", "s2"]}} 1129 1130 ### Output: 1131 "Not yet" 1132 1133 Example 2:

1134	### Input:
1135	"Deablewells "Solve for $\alpha_1 = \sqrt{A + \sqrt{12 + 12 - 2}} + \sqrt{1 + \sqrt{1 + 2}} + C$ "
1136	From the solution of a : $\sqrt{4} + \sqrt{10} + 10a + \sqrt{1} + \sqrt{1} + a = 0$.
1137	
1138	"Initial state": "We know that $\sqrt{4} + \sqrt{16 + 16a} + \sqrt{1} + \sqrt{1 + a} = 6$."
1120	"Initial graph": {"Statement": {"s1": " $\sqrt{4} + \sqrt{16 + 16a} + \sqrt{1 + \sqrt{1 + a}} = 6$ "}, "Entailment": {"s1": "Given
1133	condition"}}
1140	"Plan": "To solve a, we begin by simplifying $\sqrt{4 + \sqrt{16 + 16a}}$. This simplification may also help us simplify
1141	the left side of the equation further."
1142	"Action 1": "Simplify $\sqrt{4 + \sqrt{16 + 16a}}$ "
1143	"State 1": "We start by simplifying the expression under the square root λ First observe that λ
1144	
1145	$16 + 16a = 16(1+a) \ \text{Nouse}$ back: $\sqrt{4} + \sqrt{16(1+a)} \ \text{Now}$, simplify $\sqrt{16(1+a)} \ \text{Now}$
1146	$\sqrt{16(1+a)} = 4\sqrt{1+a} \ln \sin \theta$ the expression becomes: $\sqrt{4+4\sqrt{1+a}} \ln \theta$ (n/nNext, factor out the 4 from
1147	$ $ inside the square root: $\ln\sqrt{4(1+\sqrt{1+a})} \ln \alpha$, $\ln \alpha$, we get 2: $\ln 2\sqrt{1+\sqrt{1+a}} \ln \alpha$,
1148	the simplified form of the expression is: $\ln 2\sqrt{1 + \sqrt{1 + a}}$."
1150	"Graph 1": {"Statement": {"s2": " $\sqrt{4 + \sqrt{16 + 16a}} = 2\sqrt{1 + \sqrt{1 + a}}$ "}, "Entailment": {"s2": "Identity equation"}
1151	"Action 2": "Combine like terms $\sqrt{1+\sqrt{1+a}}$ "
1152	
1153	State 2": "We can obtain $\sqrt{4} + \sqrt{16} + 16a + \sqrt{1} + \sqrt{1} + a = 2\sqrt{1} + \sqrt{1} + a + \sqrt{1} + \sqrt{1} + a$
1154	$ - 0 \vee 1 + \vee 1 + u. $
1155	"Graph 2": {"Statement": {" $\sqrt{4} + \sqrt{16} + 16a + \sqrt{1} + \sqrt{1} + a = 3\sqrt{1} + \sqrt{1} + a^{*}$ }, "Entailment": {"
1156	[\$5 : [\$1 , "\$2] } "A science 20" [\$5 - 10" - 10"
1157	
1158	"State 3": "Isolate the square root term by dividing both sides by $3:\ln\sqrt{1 + \sqrt{1 + a}} = 2\ln\operatorname{Square}$ both
1159	$ $ sides: $(\sqrt{1+\sqrt{1+a}})^2 = 2^2 \ln 1 + \sqrt{1+a} = 4 \ln \ln 1$ solate the inner square root: $\ln \sqrt{1+a} = 4 - 1 \ln 1$
1160	$\sqrt{1+a} = 3 \ln s$ guare both sides again: $\ln(\sqrt{1+a})^2 = 3^2 \ln 1 + a = 9 \ln s$ over for $a: \ln a = 9 - 1 \ln a$
1161	$a = 8 \ln n$ the solution is $a = 8$."
1162	"Graph 3": {"Statement": {"s4": " $\sqrt{1 + \sqrt{1 + a}} = 2$ ", "s5": " $a = 8$ "}, "Entailment": {"s4": ["s1", "s3"], "s5":
1163	["s4"]}}
1164	
1165	### Output:
1166	"Final answer": "8"
1167	
1168	
1169	Prompt - Semantic Equivalence Data Generation
1170	Rewrite the given sentence into two or three different versions. Ensure that each version is distinct in wording
1171	and structure. Provide both a thought process and a final answer. In the thought process, include as many
11/1	details as possible, ensuring that no steps are omitted.
11/2	Example 1
1173	Example 1.
1174	
1175	Sentence : $\overline{XY} = \overline{XY} = \overline{2}$.
1176	### Output:
1177	$ \pi\pi\pi$ output.
1178	Thought : There are various ways to rewrite the given sentence: 1. $\frac{1}{XY} = \frac{1}{2}$ and $\frac{1}{XY} = \frac{1}{2}$. $AZ = ZY$
1179	$ = \frac{1}{2}A r$. These paraphrased versions capture the same relationship as the original equation.",
1180	["Answer": [" $\frac{XZ}{XY} = \frac{1}{2}$ and $\frac{ZI}{XY} = \frac{1}{2}$.", " $XZ = ZY = \frac{1}{2}XY$."]
1100	
101	Example 2:
1182	### input: "Santanaa", "125 a is a thread digit number "
1183	semence : 120 <i>a</i> is a infee-digit number.
4404	
1184	- ### Output:
1184 1185	### Output: "Thought": "To rewrite the sentence '125g is a three-digit number' we need to express it in different forms
1184 1185 1186	### Output: "Thought": "To rewrite the sentence ' $125a$ is a three-digit number', we need to express it in different forms that convey the same meaning. The given sentence implies that $125a$ is within the range of three-digit

1188 satisfies the condition $100 \le 125a \le 999'$. Each of these reframed versions captures the same relationship as 1189 the original sentence." 1190 "Answer": ["100 < 125a < 999.", "The value of 125a lies within the interval [100, 999]", "125a is an integer that satisfies the condition $100 \le 125a \le 999$."] 1191 1192 1193 Prompt - Semantic Equivalence Evaluation 1194 1195 Compare the provided candidate options, considering both their current attributes and potential future 1196 outcomes (if applicable). Determine whether they are semantically equivalent, and respond with either "same" or "different". 1197 1198 Example: 1199 ### Input: 1200 "Problem": "Given the polynomials p(x) = 2x + 3 and $q(x) = x^2 - x + 4$, find the polynomial resulting 1201 from multiplying p(x) and q(x) and express it in standard polynomial form. What is the coefficient of x^2 in 1202 the resulting polynomial?" 1203 "Goal": "Find the coefficient of x^2 in the resulting polynomial from multiplying p(x) and q(x).", 1204 "Initial state": "We have the polynomials p(x) = 2x + 3 and $q(x) = x^2 - x + 4$.", "Initial graph": {"Statement": {"s1": "p(x) = 2x + 3", "s2": " $q(x) = x^2 - x + 4$ "}, "Entailment": {"s1": " 1205 Given condition", "s2": "Given condition"}}, 1206 "Plan": "First, multiply the polynomials p(x) and q(x). Then, identify the coefficient of x^2 in the resulting 1207 polynomial.", 1208 "Action 1": "Multiply the polynomials p(x) and q(x).", 1209 "State 1": "To multiply p(x) = 2x + 3 and $q(x) = x^2 - x + 4$, distribute each term of p(x) to each term of 1210 $q(x): (2x+3)(x^2-x+4) = 2x(x^2) + 2x(-x) + 2x(4) + 3(x^2) + 3(-x) + 3(4) = 2x^3 - 2x^2 + 8x + 3(x^2) + 3(x^2)$ 1211 $3x^2 - 3x + 12.$ ", 1212 "Graph 1": {"Statement": {"s3": " $p(x)q(x) = 2x^3 - 2x^2 + 8x + 3x^2 - 3x + 12$ "}, "Entailment": {"s3": ["s11213 ', "s2"]}}, "Action 2": "Combine like terms in the resulting polynomial.", 1214 "State 2": { 1215 "Option 1": $2x^{3} + (2x^{2} - 2x^{2}) + 8x + 3x^{2} + (-3x) + 12 = 2x^{3} + 3x^{2} + 5x + 12$ ", 1216 "Option 2": $2x^{3} + (-2x^{2} + 3x^{2}) + (8x - 3x) + 12 = 2x^{3} + x^{2} + 5x + 12$ ", 1217 1218 "Future": { 1219 "Future 1": {"Graph 2": {"Statement": {"s4": " $p(x)q(x) = 2x^3 + 3x^2 + 5x + 12$ "}, "Entailment": {"s4": 1220 ["s3"]}}, 1221 "Future 2": {"Graph 2": {"Statement": {"s4": " $p(x)q(x) = 2x^3 + x^2 + 5x + 12$ "}, "Entailment": {"s4": 1222 ["s3"]}} 1223 ł 1224 ### Output: 1225 "Different" 1226 1227 1228 Prompt - Meta Knowledge Generation 1229 Provide the general knowledge needed to solve the problem. Do not provide a specific solution plan. Instead, 1230 offer knowledge that can help verify answers and identify errors. 1231 1232 Example 1: 1233 ### Input: "Problem": "Consider the quadratic equation $ax^2 + bx + c = 0$ where a, b, and c are integers. For this 1234 equation to have integer solutions, the discriminant $b^2 - 4ac$ must be a perfect square. Suppose the quadratic 1235 equation $2x^2 + kx + 3 = 0$ has integer solutions. What is the value of k?" 1236 1237 ### Output: 1238 "General knowledge": "To avoid errors when determining the value of variables in algebraic equations, follow 1239 these general guidelines: 1. Understand conditions for solutions: Know the conditions required for specific 1240 solution types and how they relate to equation parameters. 2. Discriminant analysis: Use the discriminant to 1241 determine solution types and ensure it meets conditions for the desired type. 3. Factorization techniques: Use

factorization to simplify equations and find relationships between parameters. 4. Systematic approach to

1242
 1243
 1244
 1244
 1244
 1244
 1244
 1244
 1244
 1244
 1244
 1244
 1244
 1244
 1244
 1244
 1244
 1244
 1244
 1244
 1244
 1244
 1244
 1244
 1244
 1244
 1244
 1244
 1244
 1244
 1244
 1244
 1244
 1244
 1244
 1244
 1244
 1244
 1244
 1244
 1244
 1244
 1244
 1244
 1244
 1244
 1244
 1244
 1244
 1244
 1244
 1244
 1244
 1244
 1244
 1244
 1244
 1244
 1244
 1244
 1244
 1244
 1244
 1244
 1244
 1244
 1244
 1244
 1244
 1244
 1244
 1244
 1244
 1244
 1244
 1244
 1244
 1244
 1244
 1244
 1244
 1244
 1244
 1244
 1244
 1244
 1244
 1244
 1244
 1244
 1244
 1244
 1244
 1244
 1244
 1244
 1244
 1244
 1244
 1244
 1244
 1244
 1244
 1244
 1244
 1244
 1244
 1244
 1244
 1244
 1244
 1244
 1244
 <li

1245 Example 2:

1246 ### Input:

1247 | "Problem": "Given the polynomials p(x) = 2x + 3 and $q(x) = x^2 - x + 4$, find the polynomial resulting 1248 | from multiplying p(x) and q(x) and express it in standard polynomial form. What is the coefficient of x^2 in 1249 | the resulting polynomial?"

1250 ### Output:

"General knowledge": "To avoid mistakes when solving polynomial multiplication problems, follow these
general guidelines: 1. Understand Polynomial Terms and Their Degrees: Recognize the individual terms of a
polynomial and their respective degrees. 2. Distributive Property in Polynomial Multiplication: Apply the
distributive property correctly by multiplying each term of the first polynomial by each term of the second
polynomial. 3. Combine Like Terms: After distributing, combine the like terms, which are terms with the
same degree. Be systematic in organizing terms to ensure all like terms are combined correctly. 4. Pay
Attention to Signs: Be careful with positive and negative signs during multiplication and when combining like
terms. Ensure that the signs of the terms are handled correctly during the distribution process."

1258 1259

1260

Prompt - Contrastive Process Supervision for Plan Generation

Compare the provided candidate options, considering both their current attributes and potential future
 outcomes (if applicable). Pay attention to the meta knowledge. First, present a detailed comparison, showing
 every step without skipping any. Then, provide a conclusion, selecting only one answer.

1264 Example:

1265 ### Input:

"Meta knowledge": "To avoid errors when determining the value of variables in algebraic equations, follow
these general guidelines: 1. Understand conditions for solutions: Know the conditions required for specific
solution types and how they relate to equation parameters. 2. Discriminant analysis: Use the discriminant to
determine solution types and ensure it meets conditions for the desired type. 3. Factorization techniques: Use
factorization to simplify equations and find relationships between parameters. 4. Systematic approach to
finding solutions: Solve related equations systematically and consider all potential factor pairs. 5. Verification:
Verify solutions by substituting back to confirm correctness and catch errors.",

1272 | "Problem": "Consider the quadratic equation $ax^2 + bx + c = 0$ where a, b, and c are integers. For this 1273 | equation to have integer solutions, the discriminant $b^2 - 4ac$ must be a perfect square. Suppose the quadratic 1274 | equation $2x^2 + kx + 3 = 0$ has integer solutions. What is the value of k?",

1275 "Goal": "Find the value of k.",

"Initial state": "We know that $2x^2 + kx + 3 = 0$ has integer solutions.",

1277 |"Initial graph": {"Statement": {"s1": " $2x^2 + kx + 3 = 0$ has integer solutions."}, "Entailment": {"s1": "Given condition"}}, "Plan": {

1279 "Option 1": "We can use the fact that the discriminant of a quadratic equation must be a perfect square to 1280 write an equation in terms of k and then solve for k.",

1281 "Option 2": "The discriminant $b^2 - 4ac$ must be a perfect square. We will set up the discriminant and 1282 solve for k under the condition.",

"Option 3": "We know that the discriminant $b^2 - 4ac$ must be a perfect square. Let's try to find the value of k for which the discriminant is a perfect square."

| }, | "Future": {

1283

1284

1285

1286 "Future 1": {"Action 1": "Find the discriminant.", "State 1": "The discriminant is $b^2 - 4ac$, so in this case, 1287 it is $k^2 - 24$.", "Graph 1": {"Statement": {"s2": "The discriminant is $k^2 - 24$."}, "Entailment": {"s2": ["s1" 1288 "]}}, "Action 2": "Since the discriminant must be a perfect square, we can set $k^2 - 24 = x^2$ for some integer 1289 x.", "State 2": "Now, we need to find k by solving the equation $k^2 - 24 = x^2$. We can do this by completing 1290 the square: $k^2 - 24 = x^2$. Add 24 to both sides: $k^2 = x^2 + 24$. Take the square root of both sides: 1291 $k = \pm \sqrt{x^2 + 24}$.", "Graph 2": {"Statement": {"s3": " $k = \pm \sqrt{x^2 + 24}$ "}, "Entailment": {"s3": ["s1", "s2"]}}, " 1292 Action 3": "Since k is an integer, $x^2 + 24$ must also be a perfect square. Let $x^2 + 24 = y^2$ for some integer y. Then, $k = \pm y$.", "State 3": "The only values of y that work are $y = \pm 4$, so $k = \pm 4$. The final answer is 1293 k = 4 or k = -4.", "Graph 3": {"Statement": {"s4": " $k = \pm 4$ "}, "Entailment": {"s4": ["s3"]}}}, 1294 "Future 2": {"Action 1": "Calculate the discriminant $b^2 - 4ac$ for the quadratic equation.", "State 1": "The 1295

discriminant Δ for the equation is given by: $\Delta = k^2 - 4 \cdot 2 \cdot 3 = k^2 - 24$.", "Graph 1": {"Statement": {"s2":

1296	$ ^{\Lambda} - k^2 - 24 $ "Entailment": $\{ _{2}2 _{2} $
1297	$\Delta = \kappa - 24$. Entaiment . { S2 . [S1]}}, Action 2 . Set the discriminant $\kappa - 24$ to be a perfect
1298	square. State 2. For $k = 24$ to be a perfect square, let $k = 24 = n$ for some integer n . This gives us the equation: $k^2 = n^2 = 24$ $(k = n)(k + n) = 24$ " "Graph 2": J"Statement": J"s 2": " $(k = n)(k + n) = 24$
1299	"In equation: $k = n = 24$ ($k = n$)($k + n$) = 24. , Graph 2. (Statement : $\{s, s\}$: $(k = n)(k + n) = 24$ "}. "Entailment": $\{s, s\}$: $[s_1]$. " s_2 "]}. "Action 3": "Solve for k by finding pairs $(k - n)$ and $(k + n)$ that
1300	multiply to 24.", "State 3": "We need to find integer pairs $(k - n)$ and $(k + n)$ such that their product is 24.
1301	Possible pairs are: $(1, 24), (2, 12), (3, 8), (4, 6), (-1, -24), (-2, -12), (-3, -8), (-4, -6)$. Solving for k
1302	in each case, we find: $k - n = 2, k + n = 12, 2k = 14, k = 7$.", "Graph 3": {"Statement": {"s4": " $k = 7$ "}, "
1303	Entailment": {"s4": ["s3"]}}},
1304	"Future 3": {"Action 1": "Find the discriminant of the equation.", "State 1": "The discriminant is $b^2 - 4ac$
1305	$=k^2 - 4(2)(3) = k^2 - 24$.", "Graph 1": {"Statement": {"s2": "The discriminant is $k^2 - 24$."}, "Entailment": {"s2": ["a1"]]
1306	{ s2. [s1]}, Action 2. Make the discriminant a perfect square, state 2. For the discriminant $k^2 = 24$ to be a perfect square, k^2 must be a perfect square plue 24. Let's call $k^2 = m^2$ for some integer me
1307	$k^2 - 24$ to be a perfect square, k must be a perfect square plus 24. Let s can $k^2 = m$ for some integer m .
1308	Then, $m = 24$ must be a perfect square. This means that m must be 25, 49, 64, 81, or 100, since these are the perfect squares closest to 24 " "Graph 2"; ("Statement"; ("s ² ; "k ² must be 25, 40, 64, 81, or 100,") "
1309	In period squares closest to 24. , Graph 2. { Statement : { S5. k must be 25, 49, 64, 61, 61, 100. }, Entailment": ["a1" "scull]] "Action 2": "Eind the corresponding values of k " "State 2": "If $m^2 = 25$
1310	Entaiment . { s5 . [s1 , s2]}}, Action 5 . Find the corresponding values of k. , state 5 . If $m^2 = 23$, then $m = 5$ and $k = 5$. If $m^2 = 40$, then $m = 7$ and $k = 7$. If $m^2 = 64$, then $m = 8$ and $k = 8$. If $m^2 = 81$, then
1311	$m = 0$ and $k = 0$. If $m^2 = 100$, then $m = 10$ and $k = 10$." "Graph 2": ["Statement": ["ad": " $k = 5.7.8.0.10$]
1312	$m = 9$, and $\kappa = 9$. If $m = 100$, then $m = 10$, and $\kappa = 10$. Chapter 5 . { Statement . { 54 . $\kappa = 5, 7, 8, 9, 10$ ["} "Entailment": {"s4". [s4". [s4".]s4".
1313	}
1314	5
1315	### Output:
1216	"Comparison": "All options mention that the discriminant must be a perfect square. Based on this observation,
1217	they will solve for k. All candidates are valid. However, the futures are different. In Future 1, $k^2 = x^2 + 24$
1010	leads to $k = 4$ or $k = -4$, which is incorrect. Future 2 is correct. In Future 3, it mentions $m^2 - 24$ must be a
1010	perfect square, then m^2 must be 25, 49, 64, 81, or 100, which is incorrect.",
1319	Conclusion : Option 2
1020	
1222	Prompt Contractive Process Supervision for Action Generation
1322	1 Tompt - Contrastive 1 Toess Supervision for Action Ocheration
1222	
1323	Compare the provided candidate options, considering both their current attributes and potential future outcomes (if applicable). Pay attention to the meta knowledge First present a detailed comparison, showing
1323 1324 1225	Compare the provided candidate options, considering both their current attributes and potential future outcomes (if applicable). Pay attention to the meta knowledge. First, present a detailed comparison, showing every step without skipping any. Then, provide a conclusion, selecting only one answer.
1323 1324 1325	Compare the provided candidate options, considering both their current attributes and potential future outcomes (if applicable). Pay attention to the meta knowledge. First, present a detailed comparison, showing every step without skipping any. Then, provide a conclusion, selecting only one answer.
1323 1324 1325 1326	Compare the provided candidate options, considering both their current attributes and potential future outcomes (if applicable). Pay attention to the meta knowledge. First, present a detailed comparison, showing every step without skipping any. Then, provide a conclusion, selecting only one answer. Example:
1323 1324 1325 1326 1327	Compare the provided candidate options, considering both their current attributes and potential future outcomes (if applicable). Pay attention to the meta knowledge. First, present a detailed comparison, showing every step without skipping any. Then, provide a conclusion, selecting only one answer. Example: #### Input:
1323 1324 1325 1326 1327 1328 1329	Compare the provided candidate options, considering both their current attributes and potential future outcomes (if applicable). Pay attention to the meta knowledge. First, present a detailed comparison, showing every step without skipping any. Then, provide a conclusion, selecting only one answer. Example: ### Input: "Meta knowledge": "To avoid errors when determining the value of variables in algebraic equations, follow these general guidelines: 1. Understand conditions for solutions: Know the conditions required for specific
1323 1324 1325 1326 1327 1328 1329 1320	Compare the provided candidate options, considering both their current attributes and potential future outcomes (if applicable). Pay attention to the meta knowledge. First, present a detailed comparison, showing every step without skipping any. Then, provide a conclusion, selecting only one answer. Example: ### Input: "Meta knowledge": "To avoid errors when determining the value of variables in algebraic equations, follow these general guidelines: 1. Understand conditions for solutions: Know the conditions required for specific solution types and how they relate to equation parameters, 2. Discriminant analysis: Use the discriminant to
1323 1324 1325 1326 1327 1328 1329 1330	Compare the provided candidate options, considering both their current attributes and potential future outcomes (if applicable). Pay attention to the meta knowledge. First, present a detailed comparison, showing every step without skipping any. Then, provide a conclusion, selecting only one answer. Example: ### Input: "Meta knowledge": "To avoid errors when determining the value of variables in algebraic equations, follow these general guidelines: 1. Understand conditions for solutions: Know the conditions required for specific solution types and how they relate to equation parameters. 2. Discriminant analysis: Use the discriminant to determine solution types and ensure it meets conditions for the desired type. 3. Factorization techniques: Use
1323 1324 1325 1326 1327 1328 1329 1330 1331	Compare the provided candidate options, considering both their current attributes and potential future outcomes (if applicable). Pay attention to the meta knowledge. First, present a detailed comparison, showing every step without skipping any. Then, provide a conclusion, selecting only one answer. Example: #### Input: "Meta knowledge": "To avoid errors when determining the value of variables in algebraic equations, follow these general guidelines: 1. Understand conditions for solutions: Know the conditions required for specific solution types and how they relate to equation parameters. 2. Discriminant analysis: Use the discriminant to determine solution types and ensure it meets conditions for the desired type. 3. Factorization techniques: Use factorization to simplify equations and find relationships between parameters. 4. Systematic approach to
1323 1324 1325 1326 1327 1328 1329 1330 1331 1332	Compare the provided candidate options, considering both their current attributes and potential future outcomes (if applicable). Pay attention to the meta knowledge. First, present a detailed comparison, showing every step without skipping any. Then, provide a conclusion, selecting only one answer. Example: #### Input: "Meta knowledge": "To avoid errors when determining the value of variables in algebraic equations, follow these general guidelines: 1. Understand conditions for solutions: Know the conditions required for specific solution types and how they relate to equation parameters. 2. Discriminant analysis: Use the discriminant to determine solution types and ensure it meets conditions for the desired type. 3. Factorization techniques: Use factorization to simplify equations and find relationships between parameters. 4. Systematic approach to finding solutions: Solve related equations systematically and consider all potential factor pairs. 5. Verification:
1323 1324 1325 1326 1327 1328 1329 1330 1331 1332 1333	Compare the provided candidate options, considering both their current attributes and potential future outcomes (if applicable). Pay attention to the meta knowledge. First, present a detailed comparison, showing every step without skipping any. Then, provide a conclusion, selecting only one answer. Example: #### Input: "Meta knowledge": "To avoid errors when determining the value of variables in algebraic equations, follow these general guidelines: 1. Understand conditions for solutions: Know the conditions required for specific solution types and how they relate to equation parameters. 2. Discriminant analysis: Use the discriminant to determine solution types and ensure it meets conditions for the desired type. 3. Factorization techniques: Use factorization to simplify equations and find relationships between parameters. 4. Systematic approach to finding solutions: Solve related equations systematically and consider all potential factor pairs. 5. Verification: Verify solutions by substituting back to confirm correctness and catch errors.",
1323 1324 1325 1326 1327 1328 1329 1330 1331 1332 1333 1334 1325	Compare the provided candidate options, considering both their current attributes and potential future outcomes (if applicable). Pay attention to the meta knowledge. First, present a detailed comparison, showing every step without skipping any. Then, provide a conclusion, selecting only one answer. Example: #### Input: "Meta knowledge": "To avoid errors when determining the value of variables in algebraic equations, follow these general guidelines: 1. Understand conditions for solutions: Know the conditions required for specific solution types and how they relate to equation parameters. 2. Discriminant analysis: Use the discriminant to determine solution types and ensure it meets conditions for the desired type. 3. Factorization techniques: Use factorization to simplify equations and find relationships between parameters. 4. Systematic approach to finding solutions: Solve related equations systematically and consider all potential factor pairs. 5. Verification: Verify solutions by substituting back to confirm correctness and catch errors.", "Problem": "Consider the quadratic equation $ax^2 + bx + c = 0$ where $a, b,$ and c are integers. For this equation to have integer colutions the discriminant b^2 .
1323 1324 1325 1326 1327 1328 1329 1330 1331 1332 1333 1334 1335	Compare the provided candidate options, considering both their current attributes and potential future outcomes (if applicable). Pay attention to the meta knowledge. First, present a detailed comparison, showing every step without skipping any. Then, provide a conclusion, selecting only one answer. Example: ### Input: "Meta knowledge": "To avoid errors when determining the value of variables in algebraic equations, follow these general guidelines: 1. Understand conditions for solutions: Know the conditions required for specific solution types and how they relate to equation parameters. 2. Discriminant analysis: Use the discriminant to determine solution types and ensure it meets conditions for the desired type. 3. Factorization techniques: Use factorization to simplify equations and find relationships between parameters. 4. Systematic approach to finding solutions: Solve related equations systematically and consider all potential factor pairs. 5. Verification: Verify solutions by substituting back to confirm correctness and catch errors.", "Problem": "Consider the quadratic equation $ax^2 + bx + c = 0$ where $a, b,$ and c are integers. For this equation to have integer solutions, the discriminant $b^2 - 4ac$ must be a perfect square. Suppose the quadratic equation $2a^2 + bx + c = 0$ where a, b^2 .
1323 1324 1325 1326 1327 1328 1329 1330 1331 1332 1333 1334 1335 1336	Compare the provided candidate options, considering both their current attributes and potential future outcomes (if applicable). Pay attention to the meta knowledge. First, present a detailed comparison, showing every step without skipping any. Then, provide a conclusion, selecting only one answer. Example: ### Input: "Meta knowledge": "To avoid errors when determining the value of variables in algebraic equations, follow these general guidelines: 1. Understand conditions for solutions: Know the conditions required for specific solution types and how they relate to equation parameters. 2. Discriminant analysis: Use the discriminant to determine solution types and ensure it meets conditions for the desired type. 3. Factorization techniques: Use factorization to simplify equations and find relationships between parameters. 4. Systematic approach to finding solutions: Solve related equations systematically and consider all potential factor pairs. 5. Verification: Verify solutions by substituting back to confirm correctness and catch errors.", "Problem": "Consider the quadratic equation $ax^2 + bx + c = 0$ where $a, b,$ and c are integers. For this equation to have integer solutions, the discriminant $b^2 - 4ac$ must be a perfect square. Suppose the quadratic equation $2x^2 + kx + 3 = 0$ has integer solutions. What is the value of k ?",
1323 1324 1325 1326 1327 1328 1329 1330 1331 1332 1333 1334 1335 1336 1337	Compare the provided candidate options, considering both their current attributes and potential future outcomes (if applicable). Pay attention to the meta knowledge. First, present a detailed comparison, showing every step without skipping any. Then, provide a conclusion, selecting only one answer. Example: ### Input: "Meta knowledge": "To avoid errors when determining the value of variables in algebraic equations, follow these general guidelines: 1. Understand conditions for solutions: Know the conditions required for specific solution types and how they relate to equation parameters. 2. Discriminant analysis: Use the discriminant to determine solution types and ensure it meets conditions for the desired type. 3. Factorization techniques: Use factorization to simplify equations and find relationships between parameters. 4. Systematic approach to finding solutions: Solve related equations systematically and consider all potential factor pairs. 5. Verification: Verify solutions by substituting back to confirm correctness and catch errors.", "Problem": "Consider the quadratic equation $ax^2 + bx + c = 0$ where $a, b,$ and c are integers. For this equation to have integer solutions, the discriminant $b^2 - 4ac$ must be a perfect square. Suppose the quadratic equation $2x^2 + kx + 3 = 0$ has integer solutions. What is the value of k ?", "Goal": "Find the value of k .",
1323 1324 1325 1326 1327 1328 1329 1330 1331 1332 1333 1334 1335 1336 1337 1338	Compare the provided candidate options, considering both their current attributes and potential future outcomes (if applicable). Pay attention to the meta knowledge. First, present a detailed comparison, showing every step without skipping any. Then, provide a conclusion, selecting only one answer. Example: ### Input: "Meta knowledge": "To avoid errors when determining the value of variables in algebraic equations, follow these general guidelines: 1. Understand conditions for solutions: Know the conditions required for specific solution types and how they relate to equation parameters. 2. Discriminant analysis: Use the discriminant to determine solution types and ensure it meets conditions for the desired type. 3. Factorization techniques: Use factorization to simplify equations and find relationships between parameters. 4. Systematic approach to finding solutions: Solve related equations systematically and consider all potential factor pairs. 5. Verification: Verify solutions by substituting back to confirm correctness and catch errors.", "Problem": "Consider the quadratic equation $ax^2 + bx + c = 0$ where $a, b,$ and c are integers. For this equation to have integer solutions, the discriminant $b^2 - 4ac$ must be a perfect square. Suppose the quadratic equation $2x^2 + kx + 3 = 0$ has integer solutions. What is the value of k ?", "Goal": "Find the value of k .", "Initial state": "We know that $2x^2 + kx + 3 = 0$ has integer solutions.", "Initial state": "We know that $2x^2 + kx + 3 = 0$ has integer solutions.", "Initial graph": {"Statement": {"sl": " $2x^2 + kx + 3 = 0$ has integer solutions.", "Initial graph": {"sl": "Given
1323 1324 1325 1326 1327 1328 1329 1330 1331 1332 1333 1334 1335 1336 1337 1338 1339	Compare the provided candidate options, considering both their current attributes and potential future outcomes (if applicable). Pay attention to the meta knowledge. First, present a detailed comparison, showing every step without skipping any. Then, provide a conclusion, selecting only one answer. Example: ### Input: "Meta knowledge": "To avoid errors when determining the value of variables in algebraic equations, follow these general guidelines: 1. Understand conditions for solutions: Know the conditions required for specific solution types and how they relate to equation parameters. 2. Discriminant analysis: Use the discriminant to determine solution types and ensure it meets conditions for the desired type. 3. Factorization techniques: Use factorization to simplify equations and find relationships between parameters. 4. Systematic approach to finding solutions: Solve related equation $ax^2 + bx + c = 0$ where <i>a</i> , <i>b</i> , and <i>c</i> are integers. For this equation to have integer solutions, the discriminant $b^2 - 4ac$ must be a perfect square. Suppose the quadratic equation $2x^2 + kx + 3 = 0$ has integer solutions. What is the value of <i>k</i> ?", "Goal": "Find the value of <i>k</i> .", "Initial state": "We know that $2x^2 + kx + 3 = 0$ has integer solutions."}, "Entailment": {"s1": "Given condition"}},
1323 1324 1325 1326 1327 1328 1329 1330 1331 1332 1333 1334 1335 1336 1337 1338 1339 1340	Compare the provided candidate options, considering both their current attributes and potential future outcomes (if applicable). Pay attention to the meta knowledge. First, present a detailed comparison, showing every step without skipping any. Then, provide a conclusion, selecting only one answer. Example: ### Input: "Meta knowledge": "To avoid errors when determining the value of variables in algebraic equations, follow these general guidelines: 1. Understand conditions for solutions: Know the conditions required for specific solution types and how they relate to equation parameters. 2. Discriminant analysis: Use the discriminant to determine solution types and ensure it meets conditions for the desired type. 3. Factorization techniques: Use factorization to simplify equations and find relationships between parameters. 4. Systematic approach to finding solutions: Solve related equations systematically and consider all potential factor pairs. 5. Verification: Verify solutions by substituting back to confirm correctness and catch errors.", "Problem": "Consider the quadratic equation $ax^2 + bx + c = 0$ where $a, b,$ and c are integers. For this equation to have integer solutions, the discriminant $b^2 - 4ac$ must be a perfect square. Suppose the quadratic equation $2x^2 + kx + 3 = 0$ has integer solutions.", "Initial state": "We know that $2x^2 + kx + 3 = 0$ has integer solutions.", "Entailment": {"s1": "Given condition"}, "Initial graph": {"Statement": {"s1": " $2x^2 + kx + 3 = 0$ has integer solutions."}, "Entailment": {"s1": "Given condition"}, "Plan": "The discriminant $b^2 - 4ac$ must be a perfect square. Suppose to favore condition"},
1323 1324 1325 1326 1327 1328 1329 1330 1331 1332 1333 1334 1335 1336 1337 1338 1339 1340 1341	Compare the provided candidate options, considering both their current attributes and potential future outcomes (if applicable). Pay attention to the meta knowledge. First, present a detailed comparison, showing every step without skipping any. Then, provide a conclusion, selecting only one answer. Example: ### Input: "Meta knowledge": "To avoid errors when determining the value of variables in algebraic equations, follow these general guidelines: 1. Understand conditions for solutions: Know the conditions required for specific solution types and how they relate to equation parameters. 2. Discriminant analysis: Use the discriminant to determine solution types and ensure it meets conditions for the desired type. 3. Factorization techniques: Use factorization to simplify equations and find relationships between parameters. 4. Systematic approach to finding solutions: Solve related equations systematically and consider all potential factor pairs. 5. Verification: Verify solutions by substituting back to confirm correctness and catch errors.", "Problem": "Consider the quadratic equation $ax^2 + bx + c = 0$ where a , b , and c are integers. For this equation to have integer solutions, the discriminant $b^2 - 4ac$ must be a perfect square. Suppose the quadratic equation $2x^2 + kx + 3 = 0$ has integer solutions.", "Initial graph": {"Statement": {"s1": " $2x^2 + kx + 3 = 0$ has integer solutions."}, "Entailment": {"s1": "Given condition"}}, "Plan": "The discriminant $b^2 - 4ac$ must be a perfect square. We will set up the discriminant and solve for k under the condition.",
1323 1324 1325 1326 1327 1328 1329 1330 1331 1332 1333 1334 1335 1336 1337 1338 1339 1340 1341 1342	Compare the provided candidate options, considering both their current attributes and potential future outcomes (if applicable). Pay attention to the meta knowledge. First, present a detailed comparison, showing every step without skipping any. Then, provide a conclusion, selecting only one answer. Example: ### Input: "Meta knowledge": "To avoid errors when determining the value of variables in algebraic equations, follow these general guidelines: 1. Understand conditions for solutions: Know the conditions required for specific solution types and how they relate to equation parameters. 2. Discriminant analysis: Use the discriminant to determine solution types and ensure it meets conditions for the desired type. 3. Factorization techniques: Use factorization to simplify equations and find relationships between parameters. 4. Systematic approach to finding solutions: Solve related equations systematically and consider all potential factor pairs. 5. Verification: Verify solutions by substituting back to confirm correctness and catch errors.", "Problem": "Consider the quadratic equation $ax^2 + bx + c = 0$ where a , b , and c are integers. For this equation to have integer solutions, the discriminant $b^2 - 4ac$ must be a perfect square. Suppose the quadratic equation $2x^2 + kx + 3 = 0$ has integer solutions.", "Initial graph": {"Statement": {"s1": " $2x^2 + kx + 3 = 0$ has integer solutions."}, "Entailment": {"s1": "Given condition"}}, "Plan": "The discriminant $b^2 - 4ac$ must be a perfect square. We will set up the discriminant and solve for k under the condition.", "Xatemather": "The discriminant $b^2 - 4ac$ for the quadratic equation.", "State 1": "The discriminant Δ
1323 1324 1325 1326 1327 1328 1329 1330 1331 1332 1333 1334 1335 1336 1337 1338 1339 1340 1341 1342 1343	Compare the provided candidate options, considering both their current attributes and potential future outcomes (if applicable). Pay attention to the meta knowledge. First, present a detailed comparison, showing every step without skipping any. Then, provide a conclusion, selecting only one answer. Example: ### Input: "Meta knowledge": "To avoid errors when determining the value of variables in algebraic equations, follow these general guidelines: 1. Understand conditions for solutions: Know the conditions required for specific solution types and how they relate to equation parameters. 2. Discriminant analysis: Use the discriminant to determine solution types and ensure it meets conditions for the desired type. 3. Factorization techniques: Use factorization to simplify equations and find relationships between parameters. 4. Systematic approach to finding solutions: Solve related equation systematically and consider all potential factor pairs. 5. Verification: Verify solutions by substituting back to confirm correctness and catch errors.", "Problem": "Consider the quadratic equation $ax^2 + bx + c = 0$ where $a, b, and c$ are integers. For this equation to have integer solutions, the discriminant $b^2 - 4ac$ must be a perfect square. Suppose the quadratic equation $2x^2 + kx + 3 = 0$ has integer solutions.", "Initial state": "We know that $2x^2 + kx + 3 = 0$ has integer solutions.", "Entailment": {"s1": "Given condition"}, "That be a perfect square. Suppose the discriminant $b^2 - 4ac$ must be a perfect square. Suppose the k muder the condition"}, "Initial graph": {"Statement": {"s1": " $2x^2 + kx + 3 = 0$ has integer solutions.", "Entailment": {"s1": "Given condition"}, "Plan": "The discriminant $b^2 - 4ac$ for the quadratic equation.", "State 1": "The discriminant Δ for the equation is given by: $\Delta = k^2 - 4 \cdot 2 \cdot 3 = k^2 - 24$.",
1323 1324 1325 1326 1327 1328 1329 1330 1331 1332 1333 1334 1335 1336 1337 1338 1339 1340 1341 1342 1343 1344	Compare the provided candidate options, considering both their current attributes and potential future outcomes (if applicable). Pay attention to the meta knowledge. First, present a detailed comparison, showing every step without skipping any. Then, provide a conclusion, selecting only one answer. Example: ### Input: "Meta knowledge": "To avoid errors when determining the value of variables in algebraic equations, follow these general guidelines: 1. Understand conditions for solutions: Know the conditions required for specific solution types and how they relate to equation parameters. 2. Discriminant analysis: Use the discriminant to determine solution ts solutions and find relationships between parameters. 4. Systematic approach to finding solutions: Solve related equations systematically and consider all potential factor pairs. 5. Verification: Verify solutions by substituting back to confirm correctness and catch errors.", "Problem": "Consider the quadratic equation $ax^2 + bx + c = 0$ where $a, b, and c$ are integers. For this equation to have integer solutions, the discriminant $b^2 - 4ac$ must be a perfect square. Suppose the quadratic equation $2x^2 + kx + 3 = 0$ has integer solutions.", "Initial state": "We know that $2x^2 + kx + 3 = 0$ has integer solutions.", "Initial graph": {"Statement": {"s1": " $2x^2 + kx + 3 = 0$ has integer solutions."}, "Entailment": {"s1": "Given condition"}}, "Plan": "The discriminant $b^2 - 4ac$ must be a perfect square. Suppose the k under the condition.", "Action 1": "Calculate the discriminant $b^2 - 4ac$ for the quadratic equation.", "State 1": "The discriminant Δ for the equation is given by: $\Delta = k^2 - 4 \cdot 2 \cdot 3 = k^2 - 24$.", "Entailment": {"s1": " $\Delta = k^2 - 24$.", "Entailment": {"s1": ["s1": " $\Delta = k^2 - 24$.", "Entailment": {"s1": ["s1": ["s1": ["]]},
1323 1324 1325 1326 1327 1328 1329 1330 1331 1332 1333 1334 1335 1336 1337 1338 1339 1340 1341 1342 1343 1344 1345	Compare the provided candidate options, considering both their current attributes and potential future outcomes (if applicable). Pay attention to the meta knowledge. First, present a detailed comparison, showing every step without skipping any. Then, provide a conclusion, selecting only one answer. Example: ### Input: "Meta knowledge": "To avoid errors when determining the value of variables in algebraic equations, follow these general guidelines: 1. Understand conditions for solutions: Know the conditions required for specific solution types and how they relate to equation parameters. 2. Discriminant analysis: Use the discriminant to determine solution types and ensure it meets conditions for the desired type. 3. Factorization techniques: Use factorization to simplify equations and find relationships between parameters. 4. Systematic approach to finding solutions: Solve related equations systematically and consider all potential factor pairs. 5. Verification: Verify solutions by substituting back to confirm correctness and catch errors.", "Problem": "Consider the quadratic equation $ax^2 + bx + c = 0$ where a , b , and c are integers. For this equation to have integer solutions, the discriminant $b^2 - 4ac$ must be a perfect square. Suppose the quadratic equation $2x^2 + kx + 3 = 0$ has integer solutions.", "Initial graph": {"Statement": {"s1": " $2x^2 + kx + 3 = 0$ has integer solutions."}, "Entailment": {"s1": "Given condition"}}, "Plan": "The discriminant $b^2 - 4ac$ for the quadratic equation $4^2 - 4ac$ is up the discriminant and solve for k under the condition.", "Action 1": "Calculate the discriminant $b^2 - 4ac$ for the quadratic equation.", "State 1": "The discriminant Δ for the equation is given by: $\Delta = k^2 - 4 \cdot 2 \cdot 3 = k^2 - 24$.", "Graph 1": {"Statement": {"s2": " $\Delta = k^2 - 24$.", "Entailment": {"s1": "The discriminant Δ for the equation is given by: $\Delta = k^2 - 4 \cdot 2 \cdot 3 = k^2 - 24$.", "Graph 1": {"Statement": {"s2": " $\Delta = k^2 - 24$.", "Entailment": {"s2": ["s1"]}}, "Action 2"
1323 1324 1325 1326 1327 1328 1329 1330 1331 1332 1333 1334 1335 1336 1337 1338 1339 1340 1341 1342 1343 1344 1345 1346	Compare the provided candidate options, considering both their current attributes and potential future outcomes (if applicable). Pay attention to the meta knowledge. First, present a detailed comparison, showing every step without skipping any. Then, provide a conclusion, selecting only one answer. Example: ### Input: "Meta knowledge": "To avoid errors when determining the value of variables in algebraic equations, follow these general guidelines: 1. Understand conditions for solutions: Know the conditions required for specific solution types and how they relate to equation parameters. 2. Discriminant analysis: Use the discriminant to determine solution types and ensure it meets conditions for the desired type. 3. Factorization techniques: Use factorization to simplify equations and find relationships between parameters. 4. Systematic approach to finding solutions: Solve related equations systematically and consider all potential factor pairs. 5. Verification: Verify solutions by substituting back to confirm correctness and catch errors.", "Problem": "Consider the quadratic equation $ax^2 + bx + c = 0$ where $a, b,$ and c are integers. For this equation $2x^2 + kx + 3 = 0$ has integer solutions. What is the value of k ?", "Goal": "Find the value of k .", "Initial state": "We know that $2x^2 + kx + 3 = 0$ has integer solutions.", "Entailment": {"s1": "Given condition"}}, "Plan": "The discriminant $b^2 - 4ac$ for the quadratic equation.", "Entailment": {"s1": "Given condition.", "Action 1": "Calculate the discriminant $b^2 - 4ac$ for the quadratic equation.", "State 1": "The discriminant Δ for the equation is given by: $\Delta = k^2 - 4 \cdot 2 \cdot 3 = k^2 - 24$.", "Graph 1": ("Statement": {"s2": " $\Delta = k^2 - 24$.", "Entailment": {"s1": "The discriminant Δ for the equation is given by: $\Delta = k^2 - 4 \cdot 2 \cdot 3 = k^2 - 24$.", "Graph 1": ("Statement": {"s2": " $\Delta = k^2 - 24$.", "Entailment": {"s2": ["s1"]}}, "Action 2": "Set the discriminant $k^2 - 24$ to be a perfect square.", "State 2": "For $k^2 - 24$ to be
1323 1324 1325 1326 1327 1328 1329 1330 1331 1332 1333 1334 1335 1336 1337 1338 1339 1340 1341 1342 1343 1344 1345 1346 1347	Compare the provided candidate options, considering both their current attributes and potential future outcomes (if applicable). Pay attention to the meta knowledge. First, present a detailed comparison, showing every step without skipping any. Then, provide a conclusion, selecting only one answer. Example: ### Input: "Meta knowledge": "To avoid errors when determining the value of variables in algebraic equations, follow these general guidelines: 1. Understand conditions for solutions: Know the conditions required for specific solution types and how they relate to equation parameters. 2. Discriminant analysis: Use the discriminant to determine solution types and ensure it meets conditions for the desired type. 3. Factorization techniques: Use factorization to simplify equations and find relationships between parameters. 4. Systematic approach to finding solutions: Solve related equations systematically and consider all potential factor pairs. 5. Verification: Verify solutions by substituting back to confirm correctness and catch errors.", "Problem": "Consider the quadratic equation $ax^2 + bx + c = 0$ where a , b , and c are integers. For this equation $2x^2 + kx + 3 = 0$ has integer solutions. What is the value of k ?", "Goal": "Find the value of k .", "Initial state": "We know that $2x^2 + kx + 3 = 0$ has integer solutions.", "Entailment": {"s1": "Given condition"}}, "Plan": "The discriminant $b^2 - 4ac$ must be a perfect square. Suppose the quadratic diverte the condition.", "Action 1": "Calculate the discriminant $b^2 - 4ac$ for the quadratic equation.", "State 1": "The discriminant Δ for the equation is given by: $\Delta = k^2 - 24 \cdot$, "Entailment": {"s1": "State ent": {"s2": " $\Delta = k^2 - 24$.", "Graph 1": {"Statement": {"s2": " $\Delta = k^2 - 24$.", "Graph 1": {"Statement": {"s2": " $\Delta = k^2 - 24$.", "Graph 1": {"Statement": {"s2": " $\Delta = k^2 - 24$.", "Graph 1": {"Statement": {"s2": " $\Delta = k^2 - 24$.", "Entailment": {"s1": "The discriminant $k^2 - 24$ to be a perfect square.", "State 2": "For $k^2 - $
1323 1324 1325 1326 1327 1328 1329 1330 1331 1332 1333 1334 1335 1336 1337 1338 1339 1340 1341 1342 1343 1344 1345 1344 1345 1346 1347 1348	Compare the provided candidate options, considering both their current attributes and potential future outcomes (if applicable). Pay attention to the meta knowledge. First, present a detailed comparison, showing every step without skipping any. Then, provide a conclusion, selecting only one answer. Example: ### Input: "Meta knowledge": "To avoid errors when determining the value of variables in algebraic equations, follow these general guidelines: 1. Understand conditions for solutions: Know the conditions required for specific solution types and how they relate to equation parameters. 2. Discriminant analysis: Use the discriminant to determine solution types and ensure it meets conditions for the desired type. 3. Factorization techniques: Use factorization to simplify equations and find relationships between parameters. 4. Systematic approach to finding solutions: Solve related equations systematically and consider all potential factor pairs. 5. Verification: Verify solutions by substituting back to confirm correctness and catch errors.", "Problem": "Consider the quadratic equation $ax^2 + bx + c = 0$ where $a, b, and c$ are integers. For this equation to have integer solutions, the discriminant $b^2 - 4ac$ must be a perfect square. Suppose the quadratic equation $2x^2 + kx + 3 = 0$ has integer solutions. What is the value of k ?", "Goal": "Find the value of k .", "Initial graph": {"Statement": {"s1": " $2x^2 + kx + 3 = 0$ has integer solutions.", "Initial graph": {"Statement": {"s1": " $2x^2 + kx + 3 = 0$ has integer solutions."}, "Entailment": {"s1": "Given condition"}, "Action 1": "Calculate the discriminant $b^2 - 4ac$ for the quadratic equation.", "State 1": "The discriminant Δ for the equation is given by: $\Delta = k^2 - 42 \cdot 3 = k^2 - 24$.", "Graph 1": {"Statement": {"s2": "} = $\Delta = k^2 - 24$.", "Entailment": {"s1": "S1": "Given 1": "S2te medicine: $k^2 - n^2 = 24$ ($k - n$)($k + n$) = 24 .", "Graph 2": {"Statement": {"s3": "} = 24 ($k - n$)($k + n$) = 24 .", "Graph 2": {"Statement": {"s3"
1323 1324 1325 1326 1327 1328 1329 1330 1331 1332 1333 1334 1335 1336 1337 1338 1339 1340 1341 1342 1343 1344 1345 1344 1345 1346 1347 1348 1349	Compare the provided candidate options, considering both their current attributes and potential future outcomes (if applicable). Pay attention to the meta knowledge. First, present a detailed comparison, showing every step without skipping any. Then, provide a conclusion, selecting only one answer. Example: ### Input: "Meta knowledge": "To avoid errors when determining the value of variables in algebraic equations, follow these general guidelines: 1. Understand conditions for solutions: Know the conditions required for specific solution types and how they relate to equation parameters. 2. Discriminant analysis: Use the discriminant to determine solution types and ensure it meets conditions for the desired type. 3. Factorization techniques: Use factorization to simplify equations and find relationships between parameters. 4. Systematic approach to finding solutions: Solve related equations systematically and consider all potential factor pairs. 5. Verification: Verify solutions by substituting back to confirm correctness and catch errors.", "Problem": "Consider the quadratic equation $ax^2 + bx + c = 0$ where a , b , and c are integers. For this equation to have integer solutions, the discriminant $b^2 - 4ac$ must be a perfect square. Suppose the quadratic equation $2x^2 + kx + 3 = 0$ has integer solutions.", "Initial graph": {"Statement": {"s1": " $2x^2 + kx + 3 = 0$ has integer solutions.", "Initial graph": {"Statement": {"s1": " $2x^2 + kx + 3 = 0$ has integer solutions.", "Initial graph": {"Statement": {"s1": " $2x^2 + kx + 3 = 0$ has integer solutions.", "The discriminant $b^2 - 4ac$ must be a perfect square. We will set up the discriminant and solve for k under the condition.", "Action 1": "Calculate the discriminant $b^2 - 4ac$ for the quadratic equation.", "State 1": "The discriminant Δ for the equation is given by: $\Delta = k^2 - 24$.", "Entailment": {"s2": " $\Delta = k^2 - 24$.", "Entailment": {"s3": "The discriminant $\Delta^2 - ac = 24$ is $b^2 - 24 = n^2$ for some integer n . This gives us the equ

	Option 2": "Eactor the right side of the equation: $(k - n)(k + n) = 24$ "
351	"Option 3": "Solve for k by factoring $(k - n)(k + n) = 24$."
352],
353	"Future": {
354	"Future 1": {"State 3": "We need to find integer pairs $(k - n)$ and $(k + n)$ such that their product is 24.
355	Possible pairs are: $(1, 24), (2, 12), (3, 8), (4, 6), (-1, -24), (-2, -12), (-3, -8), (-4, -6).$ Solving for k
356	In each case, we find: $k - n = 2$, $k + n = 12$, $2k = 14$, $k = 7$. Graph 5 : { Statement : { s4 : $k = 7$ },
357	"Future 2": {"State 3": "The factors of 24 are 1 2 3 4 6 8 12 and 24 We can check each of these to see
358	if they can be expressed as the product of two integers with a sum of k and a difference of $n_1 \cdot 24 = 24$.
350	$2 \cdot 12 = 24$, and $3 \cdot 8 = 24$ are the only possibilities. Now, let's check these: $k - n = 1$, $k + n = 24$,
260	k = 12, n = 23, k - n = 2, k + n = 12, k = 7, n = 5, and k - n = 3, k + n = 8, k = 5, n = 5. So, the
361	possible values of k are 5, 7, and 12. However, we need to check if k is an integer.", "Graph 3": {"Statement": {" $s4$ ": " $k = 5, 7, 12$ "}, "Entailment": {" $s4$ ": [" $s3$ "]}}.
362	"Future 3": {"State 3": "We need to find pairs of integers $(k - n)$ and $(k + n)$ that multiply to 24. The
363	possible pairs are $(1, 24), (2, 12), (3, 8), (4, 6), (-1, -24), (-2, -12), (-3, -8), (-4, -6)$. For each pair,
364	solve for k and n.", "Graph 3": {"Statement": {"s4": "Possible pairs for $(k - n)$ and $(k + n)$ are $(1, 24)$,
365	$(2,12), (3,8), (4,6), (-1,-24), (-2,-12), (-3,-8), (-4,-6)."\}$, "Entailment": {"s4": ["s3"]}}
366	}
367	### Output:
898	$\pi\pi\pi$ Output. "Comparison": "All options involve solving for k based on the equation. Option 1 mentions finding pairs
00	Option 2 and 3 mention factorization. All candidates are valid. Given the actions, the future states are
70	different. Future 1 shows $k = 7$. Future 2 shows $k = 5, 7, 12$. Future 3 does not mention k. Future 1 is correct
/U	In Future 2, by solving $k - n = 1$, $k + n = 24$, we obtain $k = 12.5$, $n = 11.5$ rather than $k = 12$, $n = 23$.
71	Thus, Future 2 is incorrect. Future 3 does not mention k .",
72	"Conclusion": "Option 1"
73	
'4	
75	Prompt - Contrastive Process Supervision for State Generation
76	Compare the provided candidate options, considering both their current attributes and potential future
77	outcomes (if applicable). Pay attention to the meta knowledge. First, present a detailed comparison, showing
78	every step without skipping any. Then, provide a conclusion, selecting only one answer.
79	
80	Example:
81	
82	Meta knowledge : To avoid mistakes when solving polynomial multiplication problems, follow these general guidelines: 1. Understand Polynomial Terms and Their Degrees: Pacognize the individual terms of a
83	polynomial and their respective degrees 2 Distributive Property in Polynomial Multiplication: Apply the
84	distributive property correctly by multiplying each term of the first polynomial by each term of the second
25	polynomial. 3. Combine Like Terms: After distributing, combine the like terms, which are terms with the
10	same degree. Be systematic in organizing terms to ensure all like terms are combined correctly 4. Pay
0	sume degree, be systemate in organizing terms to ensure un inte terms are comorned content, in ru,
57	Attention to Signs: Be careful with positive and negative signs during multiplication and when combining like
	Attention to Signs: Be careful with positive and negative signs during multiplication and when combining like terms. Ensure that the signs of the terms are handled correctly during the distribution process.",
38	Attention to Signs: Be careful with positive and negative signs during multiplication and when combining like terms. Ensure that the signs of the terms are handled correctly during the distribution process.", "Problem": "Given the polynomials $p(x) = 2x + 3$ and $q(x) = x^2 - x + 4$, find the polynomial resulting
8	Attention to Signs: Be careful with positive and negative signs during multiplication and when combining like terms. Ensure that the signs of the terms are handled correctly during the distribution process.", "Problem": "Given the polynomials $p(x) = 2x + 3$ and $q(x) = x^2 - x + 4$, find the polynomial resulting from multiplying $p(x)$ and $q(x)$ and express it in standard polynomial form. What is the coefficient of x^2 in
8 9 0	Attention to Signs: Be careful with positive and negative signs during multiplication and when combining like terms. Ensure that the signs of the terms are handled correctly during the distribution process.", "Problem": "Given the polynomials $p(x) = 2x + 3$ and $q(x) = x^2 - x + 4$, find the polynomial resulting from multiplying $p(x)$ and $q(x)$ and express it in standard polynomial form. What is the coefficient of x^2 in the resulting polynomial?",
38 39 90 91	Attention to Signs: Be careful with positive and negative signs during multiplication and when combining like terms. Ensure that the signs of the terms are handled correctly during the distribution process.", "Problem": "Given the polynomials $p(x) = 2x + 3$ and $q(x) = x^2 - x + 4$, find the polynomial resulting from multiplying $p(x)$ and $q(x)$ and express it in standard polynomial form. What is the coefficient of x^2 in the resulting polynomial?", "Goal": "Find the coefficient of x^2 in the resulting polynomial from multiplying $p(x)$ and $q(x)$.",
38 39 90 91 92	Attention to Signs: Be careful with positive and negative signs during multiplication and when combining like terms. Ensure that the signs of the terms are handled correctly during the distribution process.", "Problem": "Given the polynomials $p(x) = 2x + 3$ and $q(x) = x^2 - x + 4$, find the polynomial resulting from multiplying $p(x)$ and $q(x)$ and express it in standard polynomial form. What is the coefficient of x^2 in the resulting polynomial?", "Goal": "Find the coefficient of x^2 in the resulting polynomial from multiplying $p(x)$ and $q(x)$.", "Initial state": "We have the polynomials $p(x) = 2x + 3$ and $q(x) = x^2 - x + 4$.",
38 39 90 91 92 93	Attention to Signs: Be careful with positive and negative signs during multiplication and when combining like terms. Ensure that the signs of the terms are handled correctly during the distribution process.", "Problem": "Given the polynomials $p(x) = 2x + 3$ and $q(x) = x^2 - x + 4$, find the polynomial resulting from multiplying $p(x)$ and $q(x)$ and express it in standard polynomial form. What is the coefficient of x^2 in the resulting polynomial?", "Goal": "Find the coefficient of x^2 in the resulting polynomial from multiplying $p(x)$ and $q(x)$.", "Initial state": "We have the polynomials $p(x) = 2x + 3$ and $q(x) = x^2 - x + 4$.", "Initial graph": {"Statement": {"s1": " $p(x) = 2x + 3$ ", "s2": " $q(x) = x^2 - x + 4$ "}, "Entailment": {"s1": "
68 69 00 01 02 03 04	Attention to Signs: Be careful with positive and negative signs during multiplication and when combining like terms. Ensure that the signs of the terms are handled correctly during the distribution process.", "Problem": "Given the polynomials $p(x) = 2x + 3$ and $q(x) = x^2 - x + 4$, find the polynomial resulting from multiplying $p(x)$ and $q(x)$ and express it in standard polynomial form. What is the coefficient of x^2 in the resulting polynomial?", "Goal": "Find the coefficient of x^2 in the resulting polynomial from multiplying $p(x)$ and $q(x)$.", "Initial state": "We have the polynomials $p(x) = 2x + 3$ and $q(x) = x^2 - x + 4$.", "Initial graph": {"Statement": {"s1": " $p(x) = 2x + 3$ ", "s2": " $q(x) = x^2 - x + 4$ "}, "Entailment": {"s1": "Given condition", "s2": "Given condition"}
88 99 90 91 92 93 94	Attention to Signs: Be careful with positive and negative signs during multiplication and when combining like terms. Ensure that the signs of the terms are handled correctly during the distribution process.", "Problem": "Given the polynomials $p(x) = 2x + 3$ and $q(x) = x^2 - x + 4$, find the polynomial resulting from multiplying $p(x)$ and $q(x)$ and express it in standard polynomial form. What is the coefficient of x^2 in the resulting polynomial?", "Goal": "Find the coefficient of x^2 in the resulting polynomial from multiplying $p(x)$ and $q(x)$.", "Initial state": "We have the polynomials $p(x) = 2x + 3$ and $q(x) = x^2 - x + 4$.", "Initial graph": {"Statement": {"s1": " $p(x) = 2x + 3$ ", "s2": " $q(x) = x^2 - x + 4$ "}, "Entailment": {"s1": "Given condition", "s2": "Given condition"}},
38 39 90 91 92 93 94 95	Attention to Signs: Be careful with positive and negative signs during multiplication and when combining like terms. Ensure that the signs of the terms are handled correctly during the distribution process.", "Problem": "Given the polynomials $p(x) = 2x + 3$ and $q(x) = x^2 - x + 4$, find the polynomial resulting from multiplying $p(x)$ and $q(x)$ and express it in standard polynomial form. What is the coefficient of x^2 in the resulting polynomial?", "Goal": "Find the coefficient of x^2 in the resulting polynomial from multiplying $p(x)$ and $q(x)$.", "Initial state": "We have the polynomials $p(x) = 2x + 3$ and $q(x) = x^2 - x + 4$.", "Initial graph": {"Statement": {"s1": " $p(x) = 2x + 3$ ", "s2": " $q(x) = x^2 - x + 4$ "}, "Entailment": {"s1": "Given condition", "s2": "Given condition"}}, "Plan": "First, multiply the polynomials $p(x)$ and $q(x)$. Then, identify the coefficient of x^2 in the resulting polynomial.",
88 89 90 91 92 93 94 95 96	Attention to Signs: Be careful with positive and negative signs during multiplication and when combining like terms. Ensure that the signs of the terms are handled correctly during the distribution process.", "Problem": "Given the polynomials $p(x) = 2x + 3$ and $q(x) = x^2 - x + 4$, find the polynomial resulting from multiplying $p(x)$ and $q(x)$ and express it in standard polynomial form. What is the coefficient of x^2 in the resulting polynomial?", "Goal": "Find the coefficient of x^2 in the resulting polynomial from multiplying $p(x)$ and $q(x)$.", "Initial state": "We have the polynomials $p(x) = 2x + 3$ and $q(x) = x^2 - x + 4$.", "Initial graph": {"Statement": {"s1": " $p(x) = 2x + 3$ ", "s2": " $q(x) = x^2 - x + 4$ "}, "Entailment": {"s1": "Given condition", "s2": "Given condition"}}, "Plan": "First, multiply the polynomials $p(x)$ and $q(x)$. Then, identify the coefficient of x^2 in the resulting polynomial.", "Action 1": "Multiply the polynomials $p(x)$ and $q(x)$.",
88 89 90 91 92 93 94 95 96 97	Attention to Signs: Be careful with positive and negative signs during multiplication and when combining like terms. Ensure that the signs of the terms are handled correctly during the distribution process.", "Problem": "Given the polynomials $p(x) = 2x + 3$ and $q(x) = x^2 - x + 4$, find the polynomial resulting from multiplying $p(x)$ and $q(x)$ and express it in standard polynomial form. What is the coefficient of x^2 in the resulting polynomial from multiplying $p(x)$ and $q(x)$.", "Goal": "Find the coefficient of x^2 in the resulting polynomial from multiplying $p(x)$ and $q(x)$.", "Initial state": "We have the polynomials $p(x) = 2x + 3$ and $q(x) = x^2 - x + 4$.", "Initial graph": {"Statement": {"s1": " $p(x) = 2x + 3$ ", "s2": " $q(x) = x^2 - x + 4$ "}, "Entailment": {"s1": "Given condition"}}, "Plan": "First, multiply the polynomials $p(x)$ and $q(x)$. Then, identify the coefficient of x^2 in the resulting polynomial.", "Action 1": "Multiply the polynomials $p(x)$ and $q(x)$.", "State 1": "To multiply $p(x) = 2x + 3$ and $q(x) = x^2 - x + 4$, distribute each term of $p(x)$ to each term of
88 89 90 91 92 93 94 95 96 97 98	Attention to Signs: Be careful with positive and negative signs during multiplication and when combining like terms. Ensure that the signs of the terms are handled correctly during the distribution process.", "Problem": "Given the polynomials $p(x) = 2x + 3$ and $q(x) = x^2 - x + 4$, find the polynomial resulting from multiplying $p(x)$ and $q(x)$ and express it in standard polynomial form. What is the coefficient of x^2 in the resulting polynomial from multiplying $p(x)$ and $q(x)$.", "Goal": "Find the coefficient of x^2 in the resulting polynomial from multiplying $p(x)$ and $q(x)$.", "Initial state": "We have the polynomials $p(x) = 2x + 3$ and $q(x) = x^2 - x + 4$.", "Initial graph": {"Statement": {"s1": " $p(x) = 2x + 3$ ", "s2": " $q(x) = x^2 - x + 4$ "}, "Entailment": {"s1": "Given condition"}}, "Plan": "First, multiply the polynomials $p(x)$ and $q(x)$. Then, identify the coefficient of x^2 in the resulting polynomial.", "Action 1": "Multiply the polynomials $p(x)$ and $q(x)$.", "State 1": "To multiply $p(x) = 2x + 3$ and $q(x) = x^2 - x + 4$, distribute each term of $p(x)$ to each term of $q(x)$: $(2x + 3)(x^2 - x + 4) = 2x(x^2) + 2x(-x) + 2x(4) + 3(x^2) + 3(-x) + 3(4) = 2x^3 - 2x^2 + 8x + 4$
88 89 90 91 92 93 94 95 96 97 98 99	Attention to Signs: Be careful with positive and negative signs during multiplication and when combining like terms. Ensure that the signs of the terms are handled correctly during multiplication and when combining like terms. Ensure that the signs of the terms are handled correctly during the distribution process.", "Problem": "Given the polynomials $p(x) = 2x + 3$ and $q(x) = x^2 - x + 4$, find the polynomial resulting from multiplying $p(x)$ and $q(x)$ and express it in standard polynomial form. What is the coefficient of x^2 in the resulting polynomial from multiplying $p(x)$ and $q(x)$.", "Goal": "Find the coefficient of x^2 in the resulting polynomial from multiplying $p(x)$ and $q(x)$.", "Initial state": "We have the polynomials $p(x) = 2x + 3$ and $q(x) = x^2 - x + 4$.", "Initial graph": {"Statement": {"s1": " $p(x) = 2x + 3$ ", "s2": " $q(x) = x^2 - x + 4$ "}, "Entailment": {"s1": "Given condition", "s2": "Given condition"}}, "Plan": "First, multiply the polynomials $p(x)$ and $q(x)$. Then, identify the coefficient of x^2 in the resulting polynomial.", "Action 1": "Multiply the polynomials $p(x)$ and $q(x)$.", "State 1": "To multiply $p(x) = 2x + 3$ and $q(x) = x^2 - x + 4$, distribute each term of $p(x)$ to each term of $q(x)$: $(2x + 3)(x^2 - x + 4) = 2x(x^2) + 2x(-x) + 2x(4) + 3(x^2) + 3(-x) + 3(4) = 2x^3 - 2x^2 + 8x + 3x^2 - 3x + 12$.",
888 899 990 991 992 993 993 995 995 996 997 998 999 000	Attention to Signs: Be careful with positive and negative signs during multiplication and when combining like terms. Ensure that the signs of the terms are handled correctly during the distribution process.", "Problem": "Given the polynomials $p(x) = 2x + 3$ and $q(x) = x^2 - x + 4$, find the polynomial resulting from multiplying $p(x)$ and $q(x)$ and express it in standard polynomial form. What is the coefficient of x^2 in the resulting polynomial from multiplying $p(x)$ and $q(x)$.", "Goal": "Find the coefficient of x^2 in the resulting polynomial from multiplying $p(x)$ and $q(x)$.", "Initial state": "We have the polynomials $p(x) = 2x + 3$ and $q(x) = x^2 - x + 4$.", "Initial graph": {"Statement": {"s1": " $p(x) = 2x + 3$ ", "s2": " $q(x) = x^2 - x + 4$ "}, "Entailment": {"s1": "Given condition", "s2": "Given condition"}}, "Plan": "First, multiply the polynomials $p(x)$ and $q(x)$. Then, identify the coefficient of x^2 in the resulting polynomial.", "Action 1": "Multiply the polynomials $p(x)$ and $q(x)$.", "State 1": "To multiply $p(x) = 2x + 3$ and $q(x) = x^2 - x + 4$, distribute each term of $p(x)$ to each term of $q(x)$: $(2x + 3)(x^2 - x + 4) = 2x(x^2) + 2x(-x) + 2x(4) + 3(x^2) + 3(-x) + 3(4) = 2x^3 - 2x^2 + 8x + 3x^2 - 3x + 12$.", "Graph 1": {"Statement": {"s3": " $p(x)q(x) = 2x^3 - 2x^2 + 8x + 3x^2 - 3x + 12$ "}. "Entailment": {"s3": ["s1": "
88 89 90 91 92 93 94 95 96 97 98 99 99 90 00	Attention to Signs: Be careful with positive and negative signs during multiplication and when combining like terms. Ensure that the signs of the terms are handled correctly during the distribution process.", "Problem": "Given the polynomials $p(x) = 2x + 3$ and $q(x) = x^2 - x + 4$, find the polynomial resulting from multiplying $p(x)$ and $q(x)$ and express it in standard polynomial form. What is the coefficient of x^2 in the resulting polynomial from multiplying $p(x)$ and $q(x)$.", "Goal": "Find the coefficient of x^2 in the resulting polynomial from multiplying $p(x)$ and $q(x)$.", "Initial state": "We have the polynomials $p(x) = 2x + 3$ and $q(x) = x^2 - x + 4$.", "Initial graph": {"Statement": {"s1": " $p(x) = 2x + 3$ and $q(x) = x^2 - x + 4$ "}, "Entailment": {"s1": "Given condition", "s2": "Given condition"}}, "Plan": "First, multiply the polynomials $p(x)$ and $q(x)$. Then, identify the coefficient of x^2 in the resulting polynomial.", "Action 1": "Multiply the polynomials $p(x)$ and $q(x)$.", "State 1": "To multiply $p(x) = 2x + 3$ and $q(x) = x^2 - x + 4$, distribute each term of $p(x)$ to each term of $q(x)$: $(2x + 3)(x^2 - x + 4) = 2x(x^2) + 2x(-x) + 2x(4) + 3(x^2) + 3(-x) + 3(4) = 2x^3 - 2x^2 + 8x + 3x^2 - 3x + 12.",$ "Graph 1": {"Statement": {"s3": " $p(x)q(x) = 2x^3 - 2x^2 + 8x + 3x^2 - 3x + 12"$ }, "Entailment": {"s3": ["s1": "s1": "s1": "s1": "s2": "]}
888 889 890 891 892 893 894 895 896 897 898 899 600 601 602	Attention to Signs: Be careful with positive and negative signs during multiplication and when combining like terms. Ensure that the signs of the terms are handled correctly during the distribution process.", "Problem": "Given the polynomials $p(x) = 2x + 3$ and $q(x) = x^2 - x + 4$, find the polynomial resulting from multiplying $p(x)$ and $q(x)$ and express it in standard polynomial form. What is the coefficient of x^2 in the resulting polynomial from multiplying $p(x)$ and $q(x)$ and express it in standard polynomial form. What is the coefficient of x^2 in the resulting polynomial from multiplying $p(x)$ and $q(x)$.", "Total state": "We have the polynomials $p(x) = 2x + 3$ and $q(x) = x^2 - x + 4$.", "Initial graph": {"Statement": {"s1": " $p(x) = 2x + 3$ ", "s2": " $q(x) = x^2 - x + 4$ "}, "Entailment": {"s1": "Given condition", "s2": "Given condition"}}, "Plan": "First, multiply the polynomials $p(x)$ and $q(x)$. Then, identify the coefficient of x^2 in the resulting polynomial.", "Action 1": "Multiply the polynomials $p(x)$ and $q(x)$.", "State 1": "To multiply $p(x) = 2x + 3$ and $q(x) = x^2 - x + 4$, distribute each term of $p(x)$ to each term of $q(x)$: $(2x + 3)(x^2 - x + 4) = 2x(x^2) + 2x(-x) + 2x(4) + 3(x^2) + 3(-x) + 3(4) = 2x^3 - 2x^2 + 8x + 3x^2 - 3x + 12.",$ "Graph 1": {"Statement": {"s3": " $p(x)q(x) = 2x^3 - 2x^2 + 8x + 3x^2 - 3x + 12$ "}, "Entailment": {"s3": ["s1": ", "s2"]}, "Action 2": "Combine like terms in the resulting polynomial.",
 388 389 390 391 392 393 394 395 396 397 398 399 400 401 402 403 	Attention to Signs Be careful with positive and negative signs during multiplication and when combining like terms. Ensure that the signs of the terms are handled correctly during the distribution process.", "Problem": "Given the polynomials $p(x) = 2x + 3$ and $q(x) = x^2 - x + 4$, find the polynomial resulting from multiplying $p(x)$ and $q(x)$ and express it in standard polynomial form. What is the coefficient of x^2 in the resulting polynomial from multiplying $p(x)$ and $q(x)$.", "Goal": "Find the coefficient of x^2 in the resulting polynomial from multiplying $p(x)$ and $q(x)$.", "Initial state": "We have the polynomials $p(x) = 2x + 3$ and $q(x) = x^2 - x + 4$.", "Initial graph": {"Statement": {"s1": " $p(x) = 2x + 3$ ", "s2": " $q(x) = x^2 - x + 4$ "}, "Entailment": {"s1": "Given condition", "s2": "Given condition"}}, "Plan": "First, multiply the polynomials $p(x)$ and $q(x)$. Then, identify the coefficient of x^2 in the resulting polynomial.", "Action 1": "Multiply the polynomials $p(x)$ and $q(x)$.", "State 1": "To multiply $p(x) = 2x + 3$ and $q(x) = x^2 - x + 4$, distribute each term of $p(x)$ to each term of $q(x)$: $(2x + 3)(x^2 - x + 4) = 2x(x^2) + 2x(-x) + 2x(4) + 3(x^2) + 3(-x) + 3(4) = 2x^3 - 2x^2 + 8x + 3x^2 - 3x + 12.",$ "Graph 1": {"Statement": {"s3": " $p(x)q(x) = 2x^3 - 2x^2 + 8x + 3x^2 - 3x + 12$ "}, "Entailment": {"s3": ["s1": ", "s2"]}}, "Action 2:: "Combine like terms in the resulting polynomial.", "State 2:: {
 388 389 390 391 392 393 394 395 396 397 398 399 400 401 402 403 	Attention to Signs: Be careful with positive and negative signs during multiplication and when combining like terms. Ensure that the signs of the terms are handled correctly during the distribution process.", "Problem": "Given the polynomials $p(x) = 2x + 3$ and $q(x) = x^2 - x + 4$, find the polynomial resulting from multiplying $p(x)$ and $q(x)$ and express it in standard polynomial form. What is the coefficient of x^2 in the resulting polynomial form multiplying $p(x)$ and $q(x)$.", "Goal": "Find the coefficient of x^2 in the resulting polynomial from multiplying $p(x)$ and $q(x)$.", "Initial state": "We have the polynomials $p(x) = 2x + 3$ and $q(x) = x^2 - x + 4$.", "Initial graph": {"Statement": {"s1": " $p(x) = 2x + 3$ ", "s2": " $q(x) = x^2 - x + 4$ "}, "Entailment": {"s1": "Given condition", "s2": "Given condition"}}, "Plan": "First, multiply the polynomials $p(x)$ and $q(x)$. Then, identify the coefficient of x^2 in the resulting polynomial.", "Action 1": "Multiply the polynomials $p(x)$ and $q(x)$.", "State 1": "To multiply $p(x) = 2x + 3$ and $q(x) = x^2 - x + 4$, distribute each term of $p(x)$ to each term of $q(x)$: $(2x + 3)(x^2 - x + 4) = 2x(x^2) + 2x(-x) + 2x(4) + 3(x^2) + 3(-x) + 3(4) = 2x^3 - 2x^2 + 8x + 3x^2 - 3x + 12.",$ "Graph 1": {"Statement": {"s3": " $p(x)q(x) = 2x^3 - 2x^2 + 8x + 3x^2 - 3x + 12"$ }, "Entailment": {"s3": ["s1": ", "S2"]}, "Action 2": "Combine like terms in the resulting polynomial.", "State 2": { "Option 1": " $2x^3 + (2x^2 - 2x^2) + 8x + 3x^2 + (-3x) + 12 = 2x^3 + 3x^2 + 5x + 12"$,

1404 "Option 2": $2x^{3} + (-2x^{2} + 3x^{2}) + (8x - 3x) + 12 = 2x^{3} + x^{2} + 5x + 12$ ", 1405 "Option 3": " $2x^3 - 2x^2 + 8x + 3x^2 + 12 = 2x^3 + x^2 + 8x + 12$ " 1406 }, "Future": { 1407 "Future 1": {"Graph 2": {"Statement": {"s4": " $p(x)q(x) = 2x^3 + 3x^2 + 5x + 12"$ }, "Entailment": {"s4": 1408 ["s3"]}}, 1409 "Future 2": {"Graph 2": {"Statement": {"s4": " $p(x)q(x) = 2x^3 + x^2 + 5x + 12$ "}, "Entailment": {"s4": 1410 ["s3"]}} 1411 "Future 3": {"Graph 2": {"Statement": {"s4": " $p(x)q(x) = 2x^3 + x^2 + 8x + 12$ "}, "Entailment": {"s4": 1412 ["s3"]}} 1413 } 1414 ### Output: 1415 "Comparison": 1416 "All options involve combining like terms to find the coefficient of x^2 . Option 1 combines terms resulting in 1417 $2x^3 + 3x^2 + 5x + 12$, Option 2 simplifies to $2x^3 + x^2 + 5x + 12$, and Option 3 simplifies to $2x^3 + x^2 + 5x + 12$, and Option 3 simplifies to $2x^3 + x^2 + 5x + 12$. 1418 8x + 12. Thus, they yield different coefficients for x^2 and x. Option 1 mentions $2x^3 + (2x^2 - 2x^2) + 8x + 3x^2$ 1419 $3x^{2} + (-3x) + 12$. However, State 1 gives $2x^{3} - 2x^{2} + 8x + 3x^{2} - 3x + 12$. Option 1 mistakes $-2x^{2}$ as 1420 $2x^2 - 2x^2$. On the other hand, Option 2 uses the correct equation $2x^3 + (-2x^2 + 3x^2) + (8x - 3x) + 12$, 1421 which leads to the correct results. The problem of Option 3 is that it omits the term of -3x, leading to a 1422 wrong coefficient for x.", 1423 "Conclusion": "Option 2" 1424

D IMPLEMENTATION DETAILS

1425 1426

1427 1428

SWAP is fine-tuned from a base language model using LoRA. To enable scalability and general-1429 ization in our framework, we fine-tune a single generator and a single discriminator, and repurpose 1430 them to serve as the policy model, world model, and controller. For each dataset, the generator 1431 is fine-tuned on all positive trajectories in the training set that lead to the correct final answer. As 1432 illustrated in Figure 2, the generator contains two LoRAs. The original LoRA is fine-tuned on the 1433 positive trajectories as usual, while the SemEquiv-LoRA is fine-tuned on semantic equivalence data, 1434 which are bootstrapped using GPT-40, for plan, actions and states. Specially, the number of tra-1435 jectories for the generator are as follows: GSM8k (28.3k), MATH (49.3k), ReClor (14.5k), FOLIO 1436 (7.3k), HumanEval (3.1k), and MBPP (1.3k). For each positive trajectory, we random sampled some 1437 steps and generated two alternatives for each step. The number of semantically equivalent pairs we obtained are as follows: GSM8k (8.1k), MATH (24.2k), ReClor (7.1k), FOLIO (3.8k), HumanEval 1438 (1.6k) and MBPP (0.7k). 1439

1440 The discriminator is fine-tuned on contrastive process annotations for every dataset (Figure 3). 1441 Specifically, given a positive trajectory, we randomly search two alternatives for each step and ob-1442 tain their ranking. The number of trajectories for the discriminator are as follows: GSM8k (48.0k), 1443 MATH (98.2k), ReClor (28.7k), FOLIO (14.1k), HumanEval (6.0k), and MBPP (2.5k). We bootstrap the meta-knowledge text for each training question using GPT-40. For inference, we use a 1444 DPR model (Karpukhin et al., 2020) to obtain embeddings for both training questions and the test 1445 query, then calculate the Cosine similarity to select the top 5 matches. Once the relevant knowledge 1446 is extracted from these top 5 training questions, we consider two approaches: 1) using the original 1447 text directly or 2) compressing it into a shorter version. In our experiments, we found that Approach 1448 1 resulted in higher accuracy, whereas Approach 2 offered slightly lower accuracy but faster infer-1449 ence speed. The length of the future trajectory τ^{j} is also determined experimentally. We found that 1450 the optimal strategy for plan discrimination is to include all future steps. For action discrimination, 1451 including only the next state led to an increase in discrimination accuracy, whereas including ad-1452 ditional future steps caused the accuracy to decrease. We attribute this decline to the new errors 1453 introduced by the subsequent steps after analyzing the error cases. Additionally, we observed a class 1454 imbalance issue in the contrastive process annotations. Specifically, when generating discrimination 1455 data (for both actions and states), GPT-40 tends to select the first candidate if all options are similar. To address this, we propose two strategies: 1) Pre-processing: During data generation, we randomly 1456 alter the index of the ground truth. For samples where GPT-40 cannot provide the correct answer, 1457 we supply the ground truth to assist the model. 2) Post-processing: After generating the training

1458 Table 3: Fine-grained performance on MATH across different subsets: Algebra (ALG), Counting 1459 and Probability (CP), Geometry (GEO), Intermediate Algebra (IA), Number Theory (NT), Precal-1460 culus (PRE), and Prealgebra (PALG). The number of test questions for each subset is shown in parentheses. Bold values indicate the best performance per subset and overall. 1461

63					М	ath			
64	Model	ALG	СР	GEO	IA	NT	PRE	PALG	Total
65		(1187)	(474)	(479)	(903)	(540)	(546)	(871)	(5000)
6	LLaMA3-8B (0-shot CoT)	38.6	21.5	20.0	12.2	21.3	15.2	47.5	27.6
	LLaMA3-8B (4-shot CoT)	35.0	19.4	17.7	8.1	17.2	11.9	41.2	23.6
	LLaMA3-8B (0-shot CoT + SC)	38.4	19.9	17.3	10.4	18.3	12.6	46.4	26.0
	LLaMA3-8B (SFT-CoT)	37.5	19.7	18.0	10.1	18.5	11.7	45.0	25.4
	SWAP (w/o discriminator)	51.0	40.2	27.4	18.0	30.9	18.2	60.7	37.3
	SWAP	55.4	43.4	31.8	22.3	37.8	23.1	68.3	42.3

Table 4: Fine-grained performance on MATH across different difficulty levels (Level 1-5). The number of test questions for each level is shown. Bold values indicate the best performance.

	Math					
Model	L1 (437)	L2 (894)	L3 (1131)	L4 (1214)	L5 (1324)	Total (5000)
LLaMA3-8B (0-shot CoT)	65.2	45.4	30.7	18.9	8.4	27.6
LLaMA3-8B (4-shot CoT)	54.7	38.3	27.1	16.5	7.2	23.6
LLaMA3-8B (0-shot CoT + SC)	65.2	45.0	29.3	16.4	6.2	26.0
LLaMA3-8B (SFT-CoT)	64.4	44.5	28.7	15.9	5.7	25.4
SWAP (w/o discriminator)	76.5	59.9	43.3	29.4	11.7	37.3
SWAP	78.8	65.4	50.6	34.2	14.9	42.3

1482 1483

1462 1463

1471 1472

1473

1484 data, we manually change the index of the options and adjust the output accordingly. To further 1485 enhance model robustness, we also apply data augmentation by increasing the training data through 1486 varying the index and description of the options. 1487

To ensure effective training, we also employ specialized strategies such as curriculum learning and 1488 self-improving training. For curriculum learning, we first divide the training questions into groups 1489 based on their difficulty levels. For some datasets, such as MATH, the difficulty level of the problems 1490 is already provided; for other datasets, we determine the difficulty level based on the length of the 1491 solution. In the first round, we use Level 1 problems; in the second round, use both Level 1 and 1492 Level 2 problems, and so on. In each round, we train the model until convergence, using early 1493 stopping to prevent overfitting. We also employ self-improving training to iteratively refine the 1494 model's accuracy. After training, the system is run on the training samples, and the errors it produces 1495 are collected. These errors are then used to further fine-tune the discriminator, while the generator remains fixed. This process is repeated until convergence. 1496

1497

1498 E **FINE-GRAINED RESULTS** 1499

1500 To gain a comprehensive understanding of the model's strengths and weaknesses, we provide fine-1501 grained results on MATH (Table 3 and 4). We choose MATH for this analysis since it categorizes the 1502 test set by both problem types and difficulty levels, facilitating a more detailed evaluation of model 1503 performance across different dimensions. From Table 3, we observe that SWAP consistently outperforms other models across all subsets and overall, surpassing various reasoning methods applied to 1504 LLaMA3-8B-Instruct. This demonstrates that SWAP significantly enhances the overall mathematical reasoning capability compared to the baseline. The inclusion of the discriminator mechanism 1506 enables more accurate reasoning and selection, improving performance across different subsets. 1507

Interestingly, SWAP achieves better results on basic algebra compared to more complex topics like Intermediate Algebra or Precalculus, indicating variability in difficulty across different problem 1509 types. Meanwhile, the differences in performance between different LLaMA3-8B-Instruct reason-1510 ing methods are minor, and direct reasoning appears more effective than few-shot learning for these 1511 mathematical problems. Overall, SWAP demonstrates superior mathematical reasoning, particularly 1512 achieving significant improvements in challenging subsets like Intermediate Algebra and Precalcu-1513 lus, highlighting the effectiveness of our approach. 1514

Similarly, in Table 4, SWAP achieved the best performance across all difficulty levels, particularly 1515 excelling in the most challenging Level 5, where it reached an accuracy of 14.9%, compared to the 1516 best baseline performance of 8.4%. The discriminator mechanism contributes to improved accuracy 1517 on high-difficulty problems, demonstrating its effectiveness in enhancing reasoning capabilities. As 1518 difficulty increases, all models show a significant decline in performance, particularly at Levels 1519 4 and 5, indicating the increased complexity of reasoning required for these problems. Overall, 1520 SWAP consistently outperforms the baseline, especially on higher-difficulty problems, highlighting 1521 its advantage in handling complex reasoning tasks.

1522 1523

1524

F EFFICIENCY STUDY

1525 In this section, we analyze the efficiency of dif-1526 ferent planning methods. The time complexity 1527 of SWAP is O(bNT), where b is the breadth 1528 limit, N is the generation number limit, and 1529 T is the step limit. In contrast, the time complexity of RAP (using Monte Carlo Tree Search 1531 (MCTS)) (Hao et al., 2023) is $O(N_{sim}NT)$, 1532 where $N_{\rm sim}$ is the total simulation number 1533 limit. Typically, a large number of simulations $N_{\rm sim} \gg b$ are required to reliably esti-1534 mate Q(s, a) in MCTS. For ToT (Yao et al., 1535 2023), the time complexity depends on the im-1536 plementation strategy: 1) Breadth-First Search 1537 (BFS): without pruning: $O(N^T)$; with prun-1538 ing: O(bNT). 2) Depth-First Search (DFS): 1539

Table 5: Efficiency-performance trade-off across
different reasoning & planning methods.

	GSM8K	
	Avg token usage	Acc
Llama3-	-8b-Instruct	
Zero-shot CoT Few-shot CoT	175.6 148.0	70.0 72.4
RAP	2214.7 5241.4	75.2 76.0
SWAP (w/o discriminator) SWAP	306.9 3612.0	78.1 82.7

The complexity depends on the state evaluator. The traversal continues until the state evaluator 1540 deems the final state satisfactory, making the complexity tied to the evaluation criteria. In conclu-1541 sion, SWAP is more efficient than RAP and ToT (BFS without pruning version). It is similar to ToT 1542 (BFS with pruning version). 1543

We further evaluated the average number of tokens generated using different methods on the GSM8K 1544 dataset with the Llama-3-8B-Instruct model. The results are summarized in Table 5. We observed 1545 that while the theoretical time complexity of SWAP is comparable to ToT (BFS with pruning), it 1546 generates more tokens in practice due to the incorporation of a world model and the construction of 1547 an entailment graph. On the other hand, SWAP is significantly more efficient than RAP (MCTS), 1548 which involves extensive simulations to reliably estimate the Q value. 1549

1550

EXAMPLES OF SWAP G

GSM8K Example

1551 1552

In this section, we provide example solutions generated by SWAP for all benchmarks used in our 1553 paper, including GSM8K, MATH, FOLIO, ReClor, HumanEval, and MBPP. Our framework, SWAP, 1554 consists of three main components: the policy model M_{π} , the world model $M_{\rm wm}$, and the controller 1555 $M_{\rm c}$ (Figure 1). Starting with the goal G and the initial state $s_0(\mathcal{G}_0)$, the policy model M_{π} generates 1556 an optimized plan PL. Using G, PL, and the current state $s_t(\mathcal{G}_t), M_{\pi}$ proposes the next action a_t 1557 through deliberate planning. After the action is chosen, the world model $M_{\rm wm}$ predicts the next state 1558 s_{t+1} and updates the entailment graph \mathcal{G}_{t+1} . Finally, based on G and the updated state $s_{t+1}(\mathcal{G}_{t+1})$, 1559 the controller M_c decides whether to continue the reasoning process or output the final answer. 1560

1561

1563

1562

Problem: 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?

1564 1565 Solution:

"Goal": "Determine how much profit Josh made after flipping the house."

value by 150%." 1568 "Initial graph": {"Statement": {"s1": "Josh bought a house for \$80,000 and put in \$50,000 in repairs, increasing the value by 150%."}, "Entailment": {"s1": "Given condition"}} 1569 "Plan": "First, calculate the total investment made by adding the purchase price and repair costs. Then 1570 calculate the increased value of the house by multiplying the original value by 1.5. Finally, subtract the total investment from the new value to determine the profit." 1572 "Action 1": "Calculate the total investment by adding \$80,000 (purchase price) and \$50,000 (repair costs)." 1573 "State 1": "Total investment is \$80,000 + \$50,000 = \$130,000." "Graph 1": {"Statement": {"s2": "Total investment is \$130,000."}, "Entailment": {"s2": ["s1"]}} 1574 "Action 2": "Calculate the increased value of the house by multiplying \$80,000 by 1.5." 1575 "State 2": "The increase in value is \$80,000 * 1.5 = \$120,000." 1576 "Graph 2": {"Statement": {"s3": "The increase in value is \$120,000."}, "Entailment": {"s3": ["s1"]}} 1577 "Action 3": "Calculate the new value of the house by adding \$80,000 and \$120,000." 1578 "State 3": "The new value of the house is \$80,000 + \$120,000 = \$200,000." "Graph 3": {"Statement": {"s4": "The new value of the house is \$200,000."}, "Entailment": {"s4": ["s3"]}} 1579 "Action 4": "Subtract the total investment (\$130,000) from the new value of the house (\$200,000) to 1580 determine the profit." 1581 "State 4": "The profit is \$200,000 - \$130,000 = \$70,000." 1582 "Graph 4": {"Statement": {"s5": "The profit is \$70,000."}, "Entailment": {"s5": ["s2", "s4"]}} 1583 "Final answer": "70000" 1585 MATH Example 1586 **Problem:** Find the distance between the vertex of the graph of the equation $f(x) = x^2 - 8x + 15$ and the 1587 point (0, 2). 1588 1589 Solution: 1590 "Goal": "Find the distance between the vertex of the parabola and the point (0, 2)." 1591 "Initial state": "We need to find the vertex of the parabola defined by the function $f(x) = x^2 - 8x + 15$." 1592 "Initial graph": {"Statement": {"s1": "The function is $f(x) = x^2 - 8x + 15$ "}, "Entailment": {"s1": "Given 1593 function"}} "Plan": "To find the vertex of the parabola, we can use the formula for the vertex of a quadratic function in the 1594 form $f(x) = ax^2 + bx + c$. Then, we calculate the distance between this vertex and the point (0, 2)." 1595 "Action 1": "Find the vertex of the quadratic function $f(x) = x^2 - 8x + 15$." 1596 "State 1": "The vertex of a quadratic function $f(x) = ax^2 + bx + c$ is given by the formula $x = -\frac{b}{2a}$. For 1597 $f(x) = x^2 - 8x + 15$: $\ln a = 1, b = -8$, and c = 15. $\ln Substituting these values: <math>\ln x = -\frac{-8}{2.1} = \frac{8}{2} = 4 \ln \lambda$ 1598 nTo find the y-coordinate of the vertex, substitute x = 4 back into f(x):\n 1599 $f(4) = 4^2 - 8 \times 4 + 15 = 16 - 32 + 15 = -1$ \nThus, the vertex of the parabola is at (4, -1)." 1600 "Graph 1": {"Statement": {"s2": "The vertex is (4, -1)"}, "Entailment": {"s2": ["s1"]}} 1601 "Action 2": "Calculate the distance between the vertex (4, -1) and the point (0, 2). Use the distance formula." 1602 "State 2": "The distance d between two points (x_1, y_1) and (x_2, y_2) is given by:\n 1603 $d = \sqrt{(x_2 - x_1)^2 + (y_2 - y_1)^2} \ln Substituting (x_1, y_1) = (4, -1) \text{ and } (x_2, y_2) = (0, 2): \ln (x_1, y_1) = (4, -1)$ $d = \sqrt{(0-4)^2 + (2-(-1))^2} \\ \ln d = \sqrt{(-4)^2 + (2+1)^2} \\ \ln d = \sqrt{16+9} \\ \ln d = \sqrt{25} = 5 \\ \ln The \ distance = 10^{-10} \\ \ln d = \sqrt{16+9} \\ \ln$ 1604 1605 is 5.' "Graph 2": {"Statement": {"s3": "The distance between the vertex and the point (0, 2) is 5"}, "Entailment": {" 1606 s3": ["s2"]} "Final answer": "5" 1609 1610 FOLIO Example 1611 **Problem:** Determine whether the conclusion is true, false, or uncertain based on the premises provided. 1612 1613

"Initial state": "We know that Josh bought a house for \$80,000 and put in \$50,000 in repairs, increasing the

Premises:

- 1614 When the Monkeypox virus occurs in a being, it may get Monkeypox.
- Monkeypox virus can occur in certain animals. 1615
- Humans are mammals. 1616
- Mammals are animals.
- 1617 Symptoms of Monkeypox include fever, headache, muscle pains, and tiredness.
- 1618 People feel tired when they get the flu.
- 1619

1566

1567

Conclusion: There is an animal.

1620	
1621	Solution:
1622	"Goal": "Determine whether the conclusion ("There is an animal") logically follows from the given premises."
1623	"Initial state": "Premises: 1. When the Monkeypox virus occurs in a being, it may get Monkeypox. 2.
1624	Monkeypox virus can occur in certain animals. 3. Humans are mammals. 4. Mammals are animals. 5.
1625	Symptoms of Monkeypox include fever, headache, muscle pains, and tiredness. 6. People feel tired when they
1626	get the flu." "Initial graph", ("Statement", ("Dramica 1", "Dramica 2", "Dramica 4", "Dramica 5", "Dramica 6")
1627	"Finial graph: { Statement: { Premise 1, Premise 2, Premise 3, Premise 3, Premise 3, Premise 0},
1628	", "Premise 4": "Given condition", "Premise 5": "Given condition", "Premise 6": "Given condition"}
1629	"Action 1": "Identify any direct connection between the premises and the conclusion."
1630	"State 1": "The conclusion is stating "There is an animal." From Premises 3 and 4, we have that humans are
1631	mammals and mammals are animals. Therefore, humans are animals."
1632	"Graph 1": {"Statement": {"s1": "Humans are animals."}, "Entailment": {"s1": ["Premise 3", "Premise 4"]}}
1622	"State 2": "Premise 2 states that the Monkeynov virus can occur in certain animals, implying that there are
1624	animals susceptible to Monkeypox. This further supports the idea that animals exist."
1625	"Graph 2": {"Statement": {"s2": "Certain animals can get Monkeypox."}, "Entailment": {"s2": ["Premise 2"]}}
1626	"Action 3": "Draw a conclusion based on the gathered information."
1627	"State 3": "From Premises 3, 4, and 2, we can infer that there are animals (humans being a specific type of
1600	animal and other animals susceptible to Monkeypox). Therefore, the conclusion "There is an animal" is true."
1038	"Final answer": "True"
1639	
1640	
1641	
1642	ReClor Example
1643	Problem: In a business whose owners and employees all belong to one family, the employees can be paid
1644	exceptionally low wages. Hence, general operating expenses are much lower than they would be for other
1645 1646	business ventures, making profits higher. So a family business is a family's surest road to financial prosperity.
1647	0. ignores the fact that in a family business, paying family members low wages may itself reduce the family's
1648	prosperity
1649	1. presumes, without providing justification, that family members are willing to work for low wages in a
1650	family business because they believe that doing so promotes the family's prosperity
1651	2. ignores the fact that businesses that achieve high levels of customer satisfaction are often profitable even if
1652	3. presumes, without providing justification, that only businesses with low general operating expenses can
1653	succeed
1654	
1655	The reasoning in the argument is flawed because the argument
1656	Solution
1657	Solution: "Goal": "Identify the flaw in the reasoning presented in the argument "
1658	"Initial state": "We know that the problem presents an argument suggesting that low wages in a family
1659	business lead to financial prosperity."
1660	"Initial graph": {"Statement": {"s1": "Family businesses can pay low wages to employees.", "s2": "Lower
1661	wages result in lower operating expenses.", "s3": "Lower operating expenses lead to higher profits, implying
1662	Innancial prosperity."}, "Entailment": {"s1": "Given condition", "s2": "Given condition", "s3": "Given condition
1663	"Plan": "We need to evaluate the reasoning used in the argument and determine which option best describes
1664	the flaw."
1665	"Action 1": "Analyze the statement about ignoring the impact of low wages on family prosperity (option 0)."
1666	"State 1": "The argument assumes that paying low wages to family members will lead to prosperity without
1667	considering the potential negative effect on the family's overall financial well–being. This flaw makes option 0
1669	a plausible conclusion." "Graph 1": 5"Statement": 5"s/": "Paying low wagas may nagativaly affact family prosperity ") "Enterlineart":
1000	Staph 1. , Statement . , S4 . 1 aying low wages may negatively affect failing prosperity. }, Entailment : {"s4". ["s4". [
1670	"Action 2": "Analyze the statement about family members' willingness to work for low wages (option 1)."
1070	"State 2": "The argument assumes that family members are willing to accept low wages to promote prosperity
10/1	without providing any justification. This makes option 1 another valid flaw."
16/2	"Graph 2": {"Statement": {"s5": "Family members are assumed to accept low wages without justification."},"

Entailment": {"s5": ["s1"]}} "Action 3": "Analyze the statement about high customer satisfaction leading to profitability (option 2)."

	State 3": "The argument does not address whether high levels of customer satisfaction could also lead to profitability, even if wages are high. This makes option 2 a valid flaw as well." Graph 3": {"Statement": {"s6": "High customer satisfaction could lead to profitability even with high wages "}, "Entailment": {"s6": ["s2"]}} Action 4": "Analyze the statement about low operating expenses being necessary for success (option 3)."
	State 4": "The argument presumes, without justification, that only businesses with low general operating spenses can be successful. This makes option 3 another potential flaw."
j	ustification."}, "Entailment": {"s7": ["s3"]}} Final answer": "0"
ł	IumanEval Example
f	Problem: rom typing import List
c	<pre>lef has_close_elements(numbers: List[float], threshold: float) -> bool: "' Check if in given list of numbers, are any two numbers closer to each other than</pre>
	given threshold.
	$>>>$ has_close_elements([1.0, 2.0, 3.0], 0.5) False
	>>> has_close_elements($[1.0, 2.8, 3.0, 4.0, 5.0, 2.0], 0.3$)
	True
ç	Solution:
	Goal": "Write a function `has_close_elements`."
"	Initial state": "We have an input list `numbers` of floats and a float `threshold`. We need to check if any two
	numbers in the list have a difference smaller than `threshold`."
	"Finitial graph . { Statement : { S1 : The input list numbers of hoats and the noat threshold are given."},
	Plan": "We will sort the list first so that we only need to check adjacent elements. Then, we iterate through
t	he sorted list and check if the absolute difference between any two consecutive numbers is smaller than the `
t T	hreshold . If we find such a pair, we return True. If the iteration completes without finding a pair, we return
1	Action 1": "Sort the input list of numbers."
	State 1": "The list is now sorted, which allows us to only check consecutive elements for small differences."
	Graph 1": {"Statement": {"s2": "The list `numbers` is sorted."}, "Entailment": {"s2": ["s1"]}}
" "	Action 2": "Iterate through the sorted list and calculate the difference between consecutive elements."
t	state 2: we will stop and return true it we find two consecutive numbers whose difference is smaller than he `threshold`. If no such pair is found, we return False after the iteration completes."
	Graph 2": {"Statement": {"s3": "We check consecutive elements in the sorted list for differences smaller than
t	he threshold."}, "Entailment": {"s3": ["s1", "s2"]}}
"	Final answer":
f	pymon rom typing import List
1	tom typing import Dist
c	lef has_close_elements(numbers: List[float], threshold: float) -> bool:
	# Step 1: Sort the numbers
	numbers.sort()
	# Step 2. Iterate through the sorted list and check differences between consecutive elements
	for i in range(len(numbers) $- 1$):
	if abs(numbers[i] – numbers[i + 1]) < threshold:
	return True
	# Stan 3: Paturn Falsa if no close elements are found
	# Step 5. Return False II no close elements are found return False
~	
Ì	MBPP Example

Problem: Write a python function to remove first and last occurrence of a given character from the string. Your code should pass these tests:

assert remove_Occ("hello", "I") == "heo"
assert remove_Occ("abcda", "a") == "bcd"
assert remove_Occ("PHP", "P") == "H"
Solution:
"Initial state": "We have a string `s` and a character `ch` We need to find the first and last occurrences of `ch
in 's' and remove them."
"Initial graph": {"Statement": {"s1": "Input string `s` and character `ch` are provided."}, "Entailment": {"s1":
"Given condition"}}
"Plan": "We'll first find the index of the first occurrence of `ch` in the string and then the index of the last
"Action 1": "Find the indices of the first and last occurrences of ch` in `s` "
"State 1": "We now have the indices of the first and last occurrences of `ch` in `s`."
"Graph 1": {"Statement": {"s2": "The indices of the first and last occurrences are found."}, "Entailment": {"s2
": ["s1"]}}
"Action 2": "Remove the characters at these two indices from the string."
"State 2": "The first and last occurrences of "ch" have been removed from the string."
(1) (1) (1) (1) (1) (1) (1) (1) (1) (1)
"Final answer":
```python
def remove_Occ(s: str, ch: str) $\rightarrow$ str:
# Step 1: Find the first occurrence of `ch` first index = a find(ab)
$\lim_{t \to 0} st_{t} = s_{t} \lim_{t \to 0} u(cn)$
# Step 2: Find the last occurrence of `ch`
last_index = s.rfind(ch)
# Step 3: Remove the first occurrence of `ch`
If Inst_Index $!=-1$ : s = s[-first index] + s[first index + 1]
# If the first occurrence was before the last, adjust last_index
if last_index > first_index:
$last_index -= 1$
# Step 4: Remove the last occurrence of `ch`
if last index $!= -1$ and last index $!=$ first index:
$s = s[:last_index] + s[last_index + 1:]$
return s