ENHANCING COMPLEX SYMBOLIC LOGICAL REASONING OF LARGE LANGUAGE MODELS VIA SPARSE MULTI-AGENT DEBATE

Anonymous authors

000

001

002

004

006

008 009 010

011 012

013

014

015

016

017

018

019

021

023

024

025

026

027

028

029

031

032

034

037

040

041

042 043

044

046

047

048

051

052

Paper under double-blind review

ABSTRACT

Large language models (LLMs) struggle with complex logical reasoning. Existing logical question-answering methods can be briefly summarized into two categories: (1) translating natural language (NL) to symbolic language (SL) then reasoning via external solvers, and (2) adopting LLMs to reason directly in NL based on prompting or fine-tuning. Previous methods mainly focus on using a single agent in logical reasoning based on one of the above approaches. However, we point out that both methods have unignorable limitations. For example, the former (SL-based) methods suffer from imperfect translation, and the latter (NL-based) methods are prone to hallucinations. By recognizing the complementarity between SL and NL reasoning, as well as among different SL systems and NL reasoning techniques, we propose a multi-agent debate approach to absorb the strengths of various methods. In the translation stage, multiple agents translate the NL into different SL and refine translations through debate. In the reasoning stage, multiple agents based on symbolic language (obtained by the corresponding solver) and natural language debate multiple rounds, with the final answer determined by majority vote. In addition, to address the inefficiency of multi-agent debates, we introduce an adaptive sparse communication mechanism that prunes unnecessary interactions based on agent confidence and information gains. Extensive experiments on three datasets show that our method enhances logical QA performance while reducing computational cost.

1 Introduction

Large language models (LLMs) have demonstrated exceptional capabilities across a wide range of tasks. However, they still face significant challenges when performing complex logical reasoning, limiting their applicability in real-world scenarios (Cheng et al., 2025). Previous methods for logical question-answering (QA) can be broadly divided into two categories: (1) they translate natural language (NL) problems into symbolic language (SL), such as logic programming (LP), first-order logic (FOL), or Boolean satisfiability (SAT) format, based on the LLM prompting strategy, and then perform reasoning using these symbolic representations based on a external logical solver, i.e., reasoning in SL (Ye et al., 2023; Ryu et al., 2025), and (2) they use prompting (Liu et al., 2025; Xu et al., 2025) or fine-tuning strategies (Morishita et al., 2024; Wan et al., 2024) to enable LLMs to answer such questions from NL directly, i.e., reasoning in NL.

Previous works have primarily explored single-agent methods based on one of these two methods, whose performance is fundamentally limited by the capabilities of a single model (Olausson et al., 2023; Xu et al., 2024). In addition, both types of methods have unignorable drawbacks. Specifically, for the former, though they can enable a rigorous reasoning process through rule-based symbolic operations, they will get wrong results (or even be unable to run the solver) when the translation process is imperfect (Pan et al., 2023; Feng et al., 2024). For the latter, they can conduct flexible reasoning and tolerate for inaccurate expressions in the texture data by leveraging the powerful semantic understanding and generation capabilities of LLMs, but there will be uncontrollable hallucinations in reasoning and an inability to reason strictly based on the logical rules (Yao et al., 2023).

We note that reasoning in symbolic language and reasoning in natural language are not merely alternatives but complementary. In addition, when reasoning in SL, multiple SL systems (LP, FOL, SAT, etc.) have their own advantages and disadvantages, and when reasoning in NL, there also exist

various ways, such as Chain-of-thought (Wei et al., 2022) and Plan-and-Solve Wang et al. (2023), to obtain the final answer. Thus, we propose a multi-agent debate approach for logical reasoning in terms of both SL and NL, aiming to absorb the strengths of various methods. Specifically, in the translation stage, we employ multiple agents, with each agent responsible for translating NL to a specific symbol system, and then refining the translation through debate to enhance the accuracy. In the reasoning stage, we assign multiple agents using symbolic language (obtained by the corresponding solver) and natural language to debate multiple rounds. Through the exchange and argumentation in each round, agents can draw on the strengths of different reasoning methods and perspectives, and determine the final answer through a majority vote among agents for robust reasoning.

Moreover, deploying a multi-agent debate framework suffers computational overhead and token consumption (Du et al., 2023), particularly when debates involve repetitive exchanges or redundant information sharing. To address this inefficiency, we propose an adaptive sparse communication mechanism that prunes unnecessary communication by assessing the agent's confidence and information gains, allowing each agent to selectively retain only the most valuable outputs from others.

The main contribution of this paper can be summarized as follows:

- We analyze the complementarity between SL and NL reasoning paradigms, as well as the complementarity within various SL systems and NL reasoning approaches.
- We are the first to propose a multi-agent approach with an adaptive sparse communication mechanism for logical reasoning, which not only enables the absorption of advantages from multiple reasoning methods through debate but also optimizes computational efficiency and cost.
- Extensive experiments on three datasets demonstrate our method can improve the performance of logical QA while reducing the computational cost.

2 RELATED WORK

Logical Question Answering. The field of logical question answering seeks to enhance the reasoning capabilities of large language models and is generally pursued through three main approaches: solver-based, fine-tuning, and prompt-based methods (Cheng et al., 2025). Solver-based methods operate by converting natural language queries into formal symbolic expressions before utilizing specialized solvers for inference (Lyu et al., 2023; Ye et al., 2023; Olausson et al., 2023; Ryu et al., 2025). Fine-tuning techniques employ a dual strategy of creating synthetic datasets with explicit reasoning processes and augmenting training corpora with structured logical knowledge to embed reasoning abilities directly within model parameters (Wan et al., 2024; Morishita et al., 2024; Feng et al., 2024). Prompt-based methods explore a variety of strategies, with some generating explicit reasoning chains to guide inference (Wei et al., 2022; Yao et al., 2023; Besta et al., 2024; Zhang et al., 2023; 2024), while others direct models to produce symbolic forms for stepwise verification (Xu et al., 2024; 2025; Liu et al., 2025; Li et al., 2024a; Wang et al., 2024). While existing research has predominantly focused on single-agent systems, our work introduces a multi-agent debate framework to synergize the complementary advantages of both SL and NL reasoning.

Multi-Agent Debate in LLMs. Within this domain, multi-agent debate (MAD) (Du et al., 2023) is a strategy where agents engage in iterative rounds of discussion to improve their final responses through a process of collective refinement. Research on agent roles has explored distinct reasoning modes and functional assignments, such as a proposer, a critic, a planner, and an executor, to increase diversity and reliability (Liang et al., 2024; Park et al., 2023; Li et al., 2023). The inclusion of an independent judge has been shown to enhance the factual accuracy and stability of results across tasks (Du et al., 2023; Estornell & Liu, 2024; Khan et al., 2024; Chan et al.). Additionally, collaboration among heterogeneous models aims for a more robust consensus through opinion aggregation, with methods like Reconcile adding confidence-weighted voting to integrate varying viewpoints (Wang et al.; Chen et al., 2024). To address the inherent cost of these frameworks, some methods, such as SparseMAD, reduce communication by pruning the topology to a static sparse graph where agents read from fixed neighbors (Li et al., 2024b), while CortexDebate constructs a sparse debate graph with equal participation and learns edge weights using the McKinsey Trust Formula (Sun et al., 2025). Our work builds on these efforts by proposing a multi-agent debate framework that combines both symbolic and natural language reasoning, and we introduce a novel adaptive sparse communication mechanism to significantly enhance efficiency.

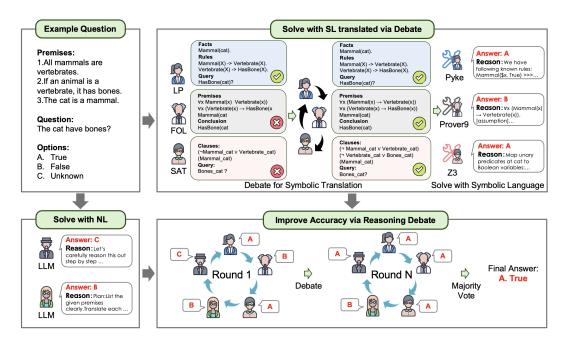


Figure 1: Overview of our sparse multi-agent debate framework for logical reasoning.

3 LOGICAL QUESTION ANSWERING PROBLEM SETUP

The task of logical question answering requires determining if a conclusion can be validly inferred from a provided set of facts and rules. For this type of problem, the model's objective is to classify the statement as true, false, or unknown. This challenge is illustrated by the example below, taken from the ProofWriter dataset (Tafjord et al., 2021):

Premises:

The bear chases the squirrel. The bear is not cold. The bear visits the cat. The bear visits the lion. The cat needs the squirrel. The lion needs the cat. The squirrel needs the lion. If something visits the lion then it visits the squirrel. If something chases the cat then the cat visits the lion.

Rules:

- If something visits the squirrel and it needs the lion then the lion does not chase the bear.
- If something is round and it visits the lion then the lion is not cold.
- If something visits the squirrel then it chases the cat.
- If the cat does not chase the bear then the cat visits the bear.
- If something visits the squirrel then it is not nice.
- If the bear is big then the bear visits the squirrel.

Question: Based on the above information, is the following statement true, false, or unknown? The squirrel does not need the lion.

Options: A) True B) False C) Unknown

Answer: B

Even with recent advancements, models continue to face considerable difficulties with logical reasoning, as evidenced by their limited performance; for instance, prior work achieved only approximately 80% accuracy on ProofWriter (Xu et al., 2025).

4 PROPOSED METHOD

4.1 OVERVIEW

To address the limitations of existing single-agent logical reasoning methods based on SL or NL, we propose a multi-agent debate framework. Specifically, as shown in Figure 1, we first translate NL logical questions into multiple SL, such as logic programming (LP), first-order logic (FOL), and Boolean satisfiability (SAT). Agents debate to refine their translations, ensuring translation accuracy for subsequent SL-based solving with solvers such as Pyke, Prover9, and Z3. Meanwhile, we adopt

LLMs to directly solve the NL logical question based on the Chain-of-though and Plan-and-Solve techniques. Finally, agents based on results from SL and NL perform debates in multiple rounds to absorb the strengths of various methods, and a majority vote among agents is used to determine the final answer. Additionally, a sparse communication mechanism is proposed to optimize the efficiency and cost of multi-agent interactions.

4.2 Debate for Symbolic Translation of Logical QA

To perform logical reasoning in a structured and unambiguous format, we begin by converting the raw natural language question into a formal symbolic expression. As illustrated in the top of Figure 1, this process first translates a logical reasoning question into three distinct symbolic languages (LP, FOL, and SAT) in parallel, then leverages a multi-agent debate to refine the final translations to improve the translation accuracy. In the following, we briefly introduce LP, FOL and SAT with their mutually different advantages and shortcomings, which motivates us to use them simultaneously.

Logic Programming (LP). Logic programming is tailored for rule-based deduction, providing a systematic framework for forward or backward inference chains. For example, a rule could be represented as has_parent(x, y) \land has_parent(y, z) \rightarrow has_grandparent(x, z). While LP excels in its *brief and efficient deduction*, its *expressiveness is constrained to rule-based problems*.

First-Order Logic (FOL). First-Order Logic provides a highly expressive framework of representing complex relations and universal quantifiers. A typical expression might be $\forall x \forall y (\text{Loves}(x,y) \rightarrow \neg \text{Hates}(x,y))$. FOL's power lies in its ability to model *intricate logical structures*, but *limited to the computational complexity for large-scale problems*.

Boolean Satisfiability (SAT). SAT formalizes a problem as a set of Boolean variables and constraints, solvable by highly optimized solvers. An example is $A = Write(Cat), B = Write(Deer), C = Black(Cat), (A \lor B) \land (\neg A \lor C)$. This approach is extremely efficient for constraint-based problems, though its limited expressiveness makes it unsuitable for complex, non-Boolean logical relationships.

4.3 Debate for Reasoning in Symbolic Language and Natural Language

Reasoning via Corresponding Logical Solvers. Given the translated symbolic languages such as LP, FOL, or SAT, solver-based reasoning methods use external logical solvers to perform logical reasoning. Despite the strong symbolic reasoning capabilities of these solvers, their effectiveness is highly sensitive to the accuracy of translation from natural to symbolic language, as even minor errors can distort solver outputs (Li et al., 2024a; Liu et al., 2025), and information loss during translation often prevents execution, rendering the problem unsolvable (Feng et al., 2024).

Example: LP Reasoning Extracted from Pyke Solver

```
We have following known rules from the context:
```

rule1: Sees(\$x, cat, True) Green(\$x, False) » Sees(\$x, cow, True)

rule2: Kind(rabbit, True) Sees(rabbit, squirrel, True) » Needs(squirrel, rabbit, True)

... ...

Now begin reasoning to obtain all implied facts:

Use rule1: Sees(\$x, cat, True) Green(\$x, False) » Sees(\$x, cow, True)

Use rule2: Kind(rabbit, True) Sees(rabbit, squirrel, True) » Needs(squirrel, rabbit, True)

... ...

All newly implied Facts: Cold('cat', True), Cold('cow', True), Eats('squirrel', 'cow', True), Rough('cat', True), Round('cat', False), Round('cow', False), Round('squirrel', False), Sees('cat', 'rabbit', True), Sees('cow', 'rabbit', True), Sees('squirrel', 'rabbit', True)

Reasoning Pipelines in Natural Language. Prompt-based reasoning methods guide LLMs to explicitly construct logical chains during question answering, thereby producing step-by-step natural language reasoning (Wei et al., 2022; Yao et al., 2023; Zhang et al., 2023; 2024). By reasoning directly in natural language, these methods avoid rigid failures caused by symbolic translation errors, thus exhibiting high robustness. However, their reasoning ability is limited by the intrinsic capacity of LLMs, making them prone to errors on complex tasks, while repeated multi-step calls to the model incur substantial computational costs (Yang et al., 2023).

Multi-agents' Debate to Improve the Accuracy of Reasoning. Solver-based methods, which exhibit strong reasoning ability but low robustness, and prompt-based methods, which are highly robust but weaker in reasoning, are inherently complementary. This motivates our proposal of a multi-agent debate strategy for mutual benefit between these two paradigms, ultimately enhancing reasoning accuracy. Specifically, our approach begins by generating a set of initial natural language reasoning narratives. For the solver-based method, its symbolic reasoning process, encompassing the rules, steps, and conclusions, is visualized as a comprehensive natural language description, exemplified by a Logic Programming (LP) reasoning text from a Pyke solver. Concurrently, the prompt-based method directly outputs a narrative documenting its thought process. Subsequently, the process enters an iterative refinement loop driven by LLM. In each round, the LLM is prompted to rewrite each reasoning narrative, using all other narratives as the provided context to inform and guide its revision. This procedure is repeated for N rounds (a predefined hyper-parameter), to facilitate deep interaction and mutual calibration. The final answer is then determined by a majority vote on the conclusions from all refined narratives.

4.4 IMPROVING EFFICIENCY VIA SPARSE DEBATE FRAMEWORK

To reduce the computational cost, we further introduce a sparse communication strategy, in which communication between agents is dynamically pruned based on a preference score. This metric assesses the potential benefit of an interaction between two LLMs in each turn by jointly considering the relative confidence of the agents and the information gains from the opponents.

4.4.1 Multi-Turn Dynamic Interaction Preference Between LLMs

We establish a sparse communication topology to improve the efficiency in multi-turn interactions by a dynamic pruning mechanism, which allows source agent i to communicate its output to the receiving agent j at round d. Specifically, we propose a preference score quantifying the potential utility of the information in the communication, which is defined as:

$$\operatorname{Pre}_{i \to j}^d = \frac{C_i^d}{C_i^d} + \lambda (1 - \cos(A_j^d, A_i^d || A_j^d)).$$

This score comprises two key components. The first is C_i^d/C_j^d , representing the ratio of confidence scores between the source agent i and the receiving agent j at round d. The second is $1 - \cos(A_j^d, A_i^d)$, measuring the difference of two outputs, regarded as information gain.

To guarantee efficiency, we propose a dynamic strategy to determine which agent should be communicated with. Specifically, in round d, we use this average preference score $\overline{\text{Pre}}_{i \to j}^{d-1}$ as the adaptive threshold, we define a binary communication gate $O_{i \to j}^{d}$. Communication from i to j is permitted only if the current preference score is greater than or equal to the historical average, indicating that the current interaction is at least as beneficial as the average past interaction between this pair. The indicator of whether agent i benefits agent j at round d is formally defined as:

$$O_{i \to j}^d = \begin{cases} 1, & \operatorname{Pre}_{i \to j}^d \ge \alpha \cdot \overline{\operatorname{Pre}_{i \to j}^{d-1}} \\ 0, & \operatorname{Pre}_{i \to j}^d < \alpha \cdot \overline{\operatorname{Pre}_{i \to j}^{d-1}} \end{cases}.$$

4.4.2 SELECTIVE MEMORY UPDATING VIA SPARSE COMMUNICATION

The sparse communication mechanism directly informs how each agent updates its internal state or memory across debate rounds. Each agent maintains a personalized memory that aggregates valuable insights from others. At the beginning of the first round (d=1), all agents start with an empty memory $M_s^1 \leftarrow \varnothing$ and communication is fully connected $(O_{i \to j}^d = 1 \text{ for all pairs})$. From the second round, the sparse communication gate $O_{i \to j}^d$ is activated. At the end of each round d, every agent s updates its memory for the next round d_s^{d+1} , by selectively incorporating the outputs d_s^d from only those agents d for which the communication channel was open (i.e., $d_{i \to j}^d = 1$). After the memory updated, agent d generates its output for the next round d_s^{d+1} , by querying the symbolic question and d is newly updated, personalized memory. After d rounds of debate, the final outputs from all agents d_s^{d+1} , ..., d_s^{d+1} , are aggregated via a majority vote to determine the final answer. The complete sparse communication Algorithm 1 is detailed in Appendix E.

270 271 272

274

275

Table 1: Accuracy comparison across three logical reasoning benchmarks

2	7	6
2	7	7
2	7	8
2	7	9
2	8	0
2	8	1

283

284

287 288 289

290

291

292 293

295

296

297

298 299 300 301 302

304 305 306

307

308

303

309 310 311

312 313 314 315 316

317 318 319 320 321 322 323

GPT-4 Claude 3.7 Sonnet DeepSeek-V3 Methods ProntoQA ProofWriter LogiDeduct ProntoQA ProofWriter LogiDeduct ProntoQA ProofWriter LogiDeduct 93.40% 79.17% 87.00% 91.80% 76.17% 94.00% 83.20% 80.50% 93.33% LogicLM 84.33% LINC 90.40% 80.67% 82.33% 91.20% 83.83% 87.67% 91.00% 84.00% 1-shot COT 69.67% 81.50% 71.83% 81.20% 67.17% 87.20% 82.33% 85.00% 83.00% 87.00% 65.67% 98.20% 89.67% 75.33% 85.17% Aristotle 95.80% 94.40% 72.33% SymbCOT 96.00% 82.33% 86.33% 97.20% 92.33% 94.00% 98.00% 85.83% 94.00% CR 93.20% 71.67% 80.33% 96.80% 82.83% 86.67% 95.40% 80.33% 83.67% DetermLR 97 80% 77 33% 85 00% 98 00% 84 33% 88 33% 96 80% 82.17% 88 33% 99.80% 89.50% SparseMAD 88.67% 99.80% 92.83% 99.83% 98.00% 92.50% 95.33% CortexDebate 99.60% 90.83% 92.33% 99.80% 96.17% 99.67% 99.80% 93.00% 99.67% 99.40% 100.00% 97.00% 99.67% 92.83% Ours (w/o sparse) 90.17% 94.00% 99.80% 100.00% Ours (w/ sparse) 92.00% 96.83% 100.00% 93.33%

Table 2: Impact of different debate components on performance

	GPT-4			Claude 3.7 Sonnet			DeepSeek-V3		
Method	ProntoQA	ProofWriter	LogiDeduct	ProntoQA	ProofWriter	LogiDeduct	ProntoQA	ProofWriter	LogiDeduct
w/o MA Trans.	99.40%	89.17%	90.00%	100.00%	96.00%	97.33%	99.60%	92.67%	97.33%
w/o MA Rea. via SL	95.60%	79.33%	84.67%	98.00%	83.33%	91.00%	96.00%	86.17%	93.00%
w/o MA Rea. via NL	/ / · · · · · · · · · · · · · · · · · ·	90.67%	94.00%	100.00%	96.67%	100.00%	99.20%	90.00%	98.00%
Ours	100.00%	92.00%	94.33%	100.00%	96.83%	100.00%	100.00%	93.33%	100.00%

EXPERIMENTS

5.1 EXPERIMENTAL SETUP

Datasets. We evaluate our method on three logical reasoning benchmarks: (1) ProntoQA (Saparov & He, 2022), a synthetic dataset for testing deductive reasoning over ontological knowledge with 500 test examples; (2) **ProofWriter** (Tafjord et al., 2021), We use the test set following (Pan et al., 2023), which is a set of randomly sampled 600 examples from the most challenging depth-5 subset; and (3) Logical Deduction (Srivastava et al., 2023), a dataset from BIG-Bench focusing on complex deductive reasoning with ordering constraints, containing 300 test examples. These benchmarks collectively assess different aspects of logical reasoning, from basic syllogistic inference to complex multi-hop deduction, and constraint-based reasoning.

Models. We conduct experiments on three LLMs: GPT-4 (OpenAI, 2023), Claude 3.7 Sonnet (Anthropic, 2025), and DeepSeek-V3 (Wu et al., 2024). All models are accessed via their respective APIs with temperature set to 0 to ensure deterministic outputs and reproducible results.

Baselines. We compare against nine representative methods spanning different approaches: (1) Solver-based methods: LogicLM (Pan et al., 2023) and LINC (Olausson et al., 2023), which translate natural language to symbolic forms for external solver processing; (2) Prompt-based methods: one-shot COT (Wei et al., 2022), Aristotle (Xu et al., 2025), SymbCOT (Xu et al., 2024), CR (Cumulative Reasoning) (Zhang et al., 2023), and DetermLR (Sun et al., 2024); (3) Multi-agent methods: SparseMAD (Li et al., 2024b) with 2 out of 5 agents communicating, and CortexDebate (Sun et al., 2025). All baseline results are obtained using the same model versions and temperature settings.

Implementation Details. Our framework employs five agents in the reasoning debate stage (three symbolic reasoning agents using LP/Pyke, FOL/Prover9, and SAT/Z3 solvers respectively, plus two natural language reasoning agents using COT and Plan-and-Solve prompting). The translation debate stage uses three agents, each specializing in one symbolic language. We set the debate rounds D=3 for translation and D=4 for reasoning stages based on our parameter analysis (Sections 5.4). The hyperparameter λ for balancing confidence and information gain is set to 1.0. When symbolic solvers fail to execute, we employ the "Simulate" strategy (detailed in Appendix B) where agents fall back to LLM-based reasoning while maintaining their symbolic perspective. The complete prompt used is detailed in the Appendix F. We use Sentence-BERT (Reimers & Gurevych, 2019) to encode agent outputs into dense embeddings for computing cosine similarity.

Evaluation Metrics. We report **Accuracy**, the percentage of correctly answered logical questions, as our evaluation metrics.

Table 3: Effect of agent diversity and composition

			GPT-4		C	laude 3.7 Sor	nnet		DeepSeek-V	3
SL reasoning	NL reasoning	ProntoQA	ProofWriter	LogiDeduct	ProntoQA	ProofWriter	LogiDeduct	ProntoQA	ProofWriter	LogiDeduct
FOL	COT	97.00%	85.50%	81.67%	99.20%	93.83%	97.67%	98.00%	91.00%	92.00%
SAT+FOL	COT	97.20%	86.17%	93.00%	99.60%	94.00%	99.67%	98.40%	92.50%	99.67%
SAT+FOL+LP	COT	100.00%	91.67%	94.00%	100.00%	96.17%	100.00%	99.60%	92.83%	100.00%
SAT+FOL+LP	COT+P&S	100.00%	92.00%	94.33%	100.00%	96.83%	100.00%	100.00%	93.33%	100.00%

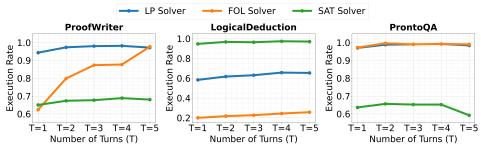


Figure 2: Relation between debate rounds and solver execution rate (GPT-4). Execution rate peaks at 2-3 rounds then declines.

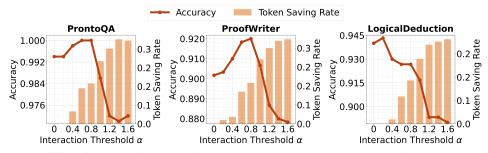


Figure 3: Effect of communication gating threshold on accuracy and token saving rate on GPT-4.

5.2 MAIN RESULTS

Table 1 presents our experimental results across three logical reasoning benchmarks. Our sparse multi-agent debate framework achieves state-of-the-art performance, with average accuracies of 95.44% (GPT-4), 98.94% (Claude 3.7), and 97.78% (DeepSeek-V3).

Overall Performance. Our method with sparse debate consistently outperforms all baselines across benchmarks and models. Compared to single-agent methods, we achieve substantial improvements over LogicLM, LINC, Aristotle, SymbCOT, CR, and DetermLR. Against multi-agent baselines, we surpass both SparseMAD and CortexDebate while maintaining computational efficiency (detailed token cost comparisons are provided in Appendix A).

Sparse vs. Full Communication. Notably, our sparse variant consistently outperforms the fully-connected version: 95.44% vs. 94.53% (GPT-4), 98.94% vs. 98.89% (Claude 3.7), and 97.78% vs. 97.54% (DeepSeek-V3), indicating that selective communication filtering not only reduces computational costs but also mitigates noise from redundant agent interactions, leading to more effective debates. The improvements across diverse base models demonstrate the robustness of our approach.

5.3 ABLATION STUDY

To understand the contribution of each component in our framework, we conduct comprehensive ablation studies on both the debate stages and the agent composition.

Impact of Debate Stages. Table 2 ablates three debate components: (1) translation debate during NL-to-SL conversion, (2) symbolic reasoning agents, and (3) natural language reasoning agents. Removing symbolic reasoning causes the largest performance drop (8.91% on GPT-4), followed by translation debate (2.59%), confirming that formal logical reasoning is most critical while accurate symbolic translation and natural language reasoning provide complementary benefits—validating our multi-stage debate design.

Table 4: Case Study of Translation Debate: Agents collaboratively refine their NL-to-SL translations through debate, \checkmark indicates correct translation, \nearrow indicates incorrect translation.

Round	Agent Translation	Key Points during Debate
	LP X: Predicates: Fruit (\$x), MoreExpensive (\$x,\$y), LessExpensive (\$x,\$y), ThirdMostExpensive (\$x) Facts: LessExpensive (kiwi,plum,True), ThirdMostExpensive (pear,True)	Initial translation with multiple predicates for comparison. Uses separate predicates for each ranking position.
1	FOL X: Rank (fruit, pos) where pos ∈ {one,two,three,four,five} Cheaper (x, y). Premises: Rank (watermelon, one), Rank (pears, three), Rank (kiwis, four)	Uses positional ranking with one=most expensive. Encodes ordering relationships between fruits.
	SAT X: fruits=EnumSort([]) price_rank=IntSort([1,2,3,4,5]) pos=Function([fruits]->[price_rank]) Constraints: pos(Kiwis)==2,pos(Pears)==3,pos(Watermelons)==5	Models prices as integer positions 1-5. Uses constraints like pos (Kiwis) <pos (plums).<="" td=""></pos>
	LP X: Modified to Rank (\$x, \$n, bool) where \$n: 5=most expensive, 1=least expensive Added rules for deriving complete ordering	"I need to be more precise about the ranking system." Realizes ambiguity in ranking direction needs clarification. Adds missing derivation rules for complete ordering.
2	FOL X : Same predicates, but added ordering axioms: $\forall X \forall Y \; (\text{Rank}\; (X, \text{one}) \; \land \; \text{Rank}\; (Y, \text{two}) \; \rightarrow \; \text{Cheaper}\; (Y, X)\;)$ Plus completeness: each fruit must have some rank	"I need to add the ordering relationships between ranks." Adds 10 ordering axioms to fully specify rank relations.
	SAT X: Same structure but notes critical error Realizes rank 1 should be cheapest, not most expensive	"I made a critical error in my ranking system." Identifies that ranking direction was inverted.
	LP√: Final version with both LessExpensive and MoreExpensive Complete rules for rank-based comparisons Rank (\$x, \$n, True) && Rank (\$y, \$m, True) && \$n<\$m >> LessExpensive()	Maintains own symbolic system while incorporating insights. Final translation is syntactically correct and complete.
3	FOL √: Complete with all 10 ordering axioms Asymmetry constraint: ∀X∀Y (Cheaper (X, Y) → ¬Cheaper (Y, X)) Each fruit and rank uniqueness constraints maintained	Final version includes all necessary constraints. Ensures logical consistency of ordering relations.
	SAT \(\sigma : \text{ Corrected ranking: 1=cheapest, 5=most expensive} \) Distinct ([f:fruits], pos(f)) for unique ranks All constraints properly oriented: pos(Kiwis) \(\text{pos}(Plums) \)	Successfully corrected the ranking direction. Final translation aligns with problem semantics. Maintains Z3 syntax requirements.

Impact of Agent Diversity. Table 3 examines how different combinations of reasoning agents affect performance. We progressively add agents: starting from a single FOL agent with COT reasoning, we incrementally incorporate SAT, LP, and Plan&Solve agents. The results reveal steady improvements with each addition ($88.06\% \rightarrow 92.12\% \rightarrow 95.22\% \rightarrow 95.44\%$ on GPT-4), demonstrating that both symbolic reasoning diversity (FOL, SAT, LP) and natural language reasoning diversity (COT, Plan&Solve) are essential for robust logical reasoning.

5.4 HYPERPARAMETER ANALYSIS

Translation Debate. Figure 2 shows executable rates of translated symbolic expressions peak at 2-3 debate rounds before degrading—a pattern consistent across all models (see Appendices C and D for other models). This degradation beyond round 3 indicates excessive debate introduces noise through over-correction of initially accurate translations. The finding validates our choice of D=3 rounds, optimally balancing translation quality improvement against over-refinement risks.

Accuracy-Communication Sparsity Trade-off. We investigate the impact of the communication threshold α on both accuracy and computational efficiency, measured as token saving rate: $(\text{Tokens}_{\text{w/o sparse}} - \text{Tokens}_{\text{w/o sparse}})/\text{Tokens}_{\text{w/o sparse}}$. Higher α values enforce stricter communication filtering, resulting in sparser interaction graphs. Figure 3 illustrates this trade-off for GPT-4 (see Appendices C and D for other models). A notable pattern emerges: as α increases, accuracy often improves while simultaneously reducing token costs by 10-30%. This suggests that moderate sparsity filters out redundant inter-agent communications that can harm reasoning quality. For instance, on GPT-4, accuracy peaks at $\alpha \approx 0.7$ with 95.44% accuracy and 20% token reduction.

Reasoning Debate. Figure 4 shows accuracy saturates after 2-3 debate rounds across three benchmarks, then plateauing or slightly degrading. This pattern suggests agents quickly reach consensus on logical problems, with further rounds introducing noise through overthinking or redundant arguments. The consistent 3-round optimum across datasets validates our choice of D=4, balancing reasoning quality with computational efficiency.

5.5 CASE STUDY: MULTI-AGENT TRANSLATION DEBATE DYNAMICS

We present a case study from LogicalDeduction where agents collaboratively translate a fruit pricing problem into symbolic languages. Table 4 shows how through three debate rounds, agents identify and correct critical errors: SAT discovers its inverted ranking (treating 1 as most expensive rather than cheapest), LP clarifies ranking ambiguity, and FOL adds missing ordering axioms. This peer

Table 5: Case Study of agent reasoning across debate rounds. In this case, A is the correct answer, C is the incorrect answer.

Round	Agent	Answer	Key Reasoning						
	LP A From Needs (Lion, Bear) and Rule 2: Needs (\$x, Bear) → Needs (\$x, Lion), we deduce Needs (Lion, Lion). Then Rule 3: Needs (\$x, Lion) → Visits (\$x, Lion) gives us Visits (Lion)								
	FOL	С	Both proving and denying Visits (Lion, Lion) yield no contradiction. Neither can be conclusively established as true or false. Problem is underdefined.						
1	SAT C Z3 solver finds no explicit fact that implies Visits (Lion, Lion). Rules do not derive it. Status: Unknown/indete								
	P&S A Lion needs bear (Fact 6). By Rule 2: lion needs lion. By Rule 3: lion visits lion. Clear logical chain establishes truth.								
	COT C No direct rule or statement establishes whether lion visits itself. No chain of reasoning leads to this conclusi								
	LP	A FOL's claim of "underdefined" is incorrect. Rules 2→3 are deterministic, not "underdefined".							
	FOL	C	Still claiming Unknown despite clear logical chain shown by others. Insisting problem structure yields no definitive resolution.						
	SAT	A	After reviewing Plan&Solve's argument, I agree. The deduction chain is valid. Changing from Unknown to True.						
2	P&S	A	Standing by deduction. Facts and rules suffice to establish truth. The logical steps are valid and complete.						
	СОТ	OT A I acknowledging my oversight before. Plan&Solve's Rule 2+3 chain is valid.Changing answer to True.							
	LP	A	Chain of reasoning is airtight. Therefore, lion visits itself.						
	FOL	A	After reassessing, I concur. The Rule 2→Rule 3 chain is valid and well-supported. Changing from Unknown to True.						
	SAT A Maintaining revised conclusion.								
3	P&S A Reasoning follows structured approach: Fact 6 + Rule 2 + Rule 3 = definitive True.								
5	СОТ	COT A Agreement confirmed. Chain of logic using Fact 6, Rule 2, and Rule 3 is sound.							
4			All agents reach consensus: Answer A - "The lion visits the lion" is conclusively true						

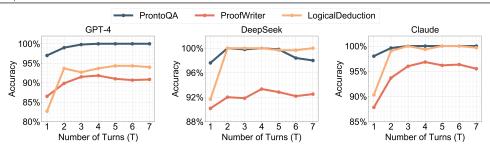


Figure 4: Relation between turns and final accuracy.

review process yields syntactically correct and semantically accurate translations across all three symbolic languages, demonstrating how multi-agent debate enhances translation quality—crucial for symbolic reasoning accuracy. Full question and dialogue details are in Appendix G.3.

5.6 CASE STUDY: MULTI-AGENT REASONING DEBATE DYNAMICS

To illustrate how our multi-agent debate framework achieves consensus through collaborative reasoning, we present a case study from the ProofWriter dataset shown in Table 5. The problem requires determining whether "The lion visits the lion" is true (A), false (B), or unknown (C) based on given logical rules and facts, with ground truth answer being (A). The debate showcases effective peer correction: agents with incorrect initial answers recognize their logical oversights through examining others' reasoning chains and converge to the correct solution, validating multi-agent debate's error-correction capability. Full question and dialogues for this case can be found in Appendix G.3.

6 Conclusion

This paper mitigates the important limitations of large language models (LLMs) in complex logical reasoning. We first analyze the complementarity between symbolic language (SL) and natural language (NL) reasoning paradigms, as well as the complementarity within various SL systems and NL reasoning approaches. Different from previous works, which have primarily been based on a single-agent approach, using one of SL-based or NL-based methods, we are the first to propose a multi-agent approach, which enables the absorption of advantages from multiple reasoning methods through debate. Additionally, we propose a sparse communication mechanism to optimize the efficiency and cost of these multi-agent interactions. Extensive experiments on three datasets show that our method enhances logical QA performance while reducing computational cost.

REFERENCES

- Anthropic. Claude 3.7 sonnet system card, 2025. Accessed 2025-09-25.
- Maciej Besta, Nils Blach, Ales Kubicek, Robert Gerstenberger, Michal Podstawski, Lukas Gianinazzi, Joanna Gajda, Tomasz Lehmann, Hubert Niewiadomski, Piotr Nyczyk, et al. Graph of thoughts: Solving elaborate problems with large language models. In *Proceedings of the AAAI conference on artificial intelligence*, volume 38, pp. 17682–17690, 2024.
- Chi-Min Chan, Weize Chen, Yusheng Su, Jianxuan Yu, Wei Xue, Shanghang Zhang, Jie Fu, and Zhiyuan Liu. Chateval: Towards better llm-based evaluators through multi-agent debate. In *The Twelfth International Conference on Learning Representations*.
- Justin Chen, Swarnadeep Saha, and Mohit Bansal. Reconcile: Round-table conference improves reasoning via consensus among diverse llms. In *Proceedings of the 62nd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pp. 7066–7085, 2024.
- Fengxiang Cheng, Haoxuan Li, Fenrong Liu, Robert van Rooij, Kun Zhang, and Zhouchen Lin. Empowering llms with logical reasoning: A comprehensive survey. *International Joint Conference on Artificial Intelligence, Survey Track*, 2025.
- Yilun Du, Shuang Li, Antonio Torralba, Joshua B Tenenbaum, and Igor Mordatch. Improving factuality and reasoning in language models through multiagent debate. In *Forty-first International Conference on Machine Learning*, 2023.
- Andrew Estornell and Yang Liu. Multi-Ilm debate: Framework, principals, and interventions. *Advances in Neural Information Processing Systems*, 37:28938–28964, 2024.
- Jiazhan Feng, Ruochen Xu, Junheng Hao, Hiteshi Sharma, Yelong Shen, Dongyan Zhao, and Weizhu Chen. Language models can be deductive solvers. In *Findings of the Association for Computational Linguistics: NAACL 2024*, pp. 4026–4042, 2024.
- Akbir Khan, John Hughes, Dan Valentine, Laura Ruis, Kshitij Sachan, Ansh Radhakrishnan, Edward Grefenstette, Samuel R Bowman, Tim Rocktäschel, and Ethan Perez. Debating with more persuasive llms leads to more truthful answers. *Proceedings of Machine Learning Research*, 235: 23662–23733, 2024.
- Guohao Li, Hasan Hammoud, Hani Itani, Dmitrii Khizbullin, and Bernard Ghanem. Camel: Communicative agents for mind exploration of large language model society. *Advances in Neural Information Processing Systems*, 36:51991–52008, 2023.
- Qingchuan Li, Jiatong Li, Tongxuan Liu, Yuting Zeng, Mingyue Cheng, Weizhe Huang, and Qi Liu. Leveraging llms for hypothetical deduction in logical inference: A neuro-symbolic approach. *arXiv preprint arXiv:2410.21779*, 2024a.
- Yunxuan Li, Yibing Du, Jiageng Zhang, Le Hou, Peter Grabowski, Yeqing Li, and Eugene Ie. Improving multi-agent debate with sparse communication topology. In *Findings of the Association for Computational Linguistics: EMNLP 2024*, pp. 7281–7294, 2024b.
- Tian Liang, Zhiwei He, Wenxiang Jiao, Xing Wang, Yan Wang, Rui Wang, Yujiu Yang, Shuming Shi, and Zhaopeng Tu. Encouraging divergent thinking in large language models through multiagent debate. In *EMNLP*, 2024.
- Tongxuan Liu, Wenjiang Xu, Weizhe Huang, Yuting Zeng, Jiaxing Wang, Xingyu Wang, Hailong Yang, and Jing Li. Logic-of-thought: Injecting logic into contexts for full reasoning in large language models. In *Proceedings of the 2025 Conference of the Nations of the Americas Chapter of the Association for Computational Linguistics: Human Language Technologies*, 2025.
- Qing Lyu, Shreya Havaldar, Adam Stein, Li Zhang, Delip Rao, Eric Wong, Marianna Apidianaki, and Chris Callison-Burch. Faithful chain-of-thought reasoning. In *The 13th International Joint Conference on Natural Language Processing and the 3rd Conference of the Asia-Pacific Chapter of the Association for Computational Linguistics (IJCNLP-AACL 2023)*, 2023.

- Terufumi Morishita, Gaku Morio, Atsuki Yamaguchi, and Yasuhiro Sogawa. Enhancing reasoning capabilities of Ilms via principled synthetic logic corpus. *Advances in Neural Information Processing Systems*, 37:73572–73604, 2024.
- Theo Olausson, Alex Gu, Ben Lipkin, Cedegao Zhang, Armando Solar-Lezama, Joshua Tenenbaum, and Roger Levy. Linc: A neurosymbolic approach for logical reasoning by combining language models with first-order logic provers. In *Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing*, pp. 5153–5176, 2023.
- OpenAI. Gpt-4 technical report. 2023.

- Liangming Pan, Alon Albalak, Xinyi Wang, and William Wang. Logic-lm: Empowering large language models with symbolic solvers for faithful logical reasoning. In *Findings of the Association for Computational Linguistics: EMNLP 2023*, pp. 3806–3824, 2023.
- Joon Sung Park, Joseph O'Brien, Carrie Jun Cai, Meredith Ringel Morris, Percy Liang, and Michael S Bernstein. Generative agents: Interactive simulacra of human behavior. In *Proceedings of the 36th annual acm symposium on user interface software and technology*, pp. 1–22, 2023.
- Nils Reimers and Iryna Gurevych. Sentence-bert: Sentence embeddings using siamese bertnetworks. In *Proceedings of the 2019 Conference on Empirical Methods in Natural Language Processing and the 9th International Joint Conference on Natural Language Processing (EMNLP-IJCNLP)*, pp. 3982–3992, 2019.
- Hyun Ryu, Gyeongman Kim, Hyemin S. Lee, and Eunho Yang. Divide and translate: Compositional first-order logic translation and verification for complex logical reasoning. In *ICLR*, 2025.
- Abulhair Saparov and He He. Language models are greedy reasoners: A systematic formal analysis of chain-of-thought. *arXiv preprint arXiv:2210.01240*, 2022.
- Aarohi Srivastava, Abhinav Rastogi, Abhishek Rao, Abu Awal Shoeb, Abubakar Abid, Adam Fisch, Adam R Brown, Adam Santoro, Aditya Gupta, Adri Garriga-Alonso, et al. Beyond the imitation game: Quantifying and extrapolating the capabilities of language models. *Transactions on machine learning research*, 2023.
- Hongda Sun, Weikai Xu, Wei Liu, Jian Luan, Bin Wang, Shuo Shang, Ji-Rong Wen, and Rui Yan. DetermLR: Augmenting LLM-based logical reasoning from indeterminacy to determinacy. In Lun-Wei Ku, Andre Martins, and Vivek Srikumar (eds.), *Proceedings of the 62nd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, Bangkok, Thailand, August 2024. Association for Computational Linguistics. doi: 10.18653/v1/2024.acl-long.531.
- Yiliu Sun, Zicheng Zhao, Sheng Wan, and Chen Gong. Cortexdebate: Debating sparsely and equally for multi-agent debate. In *Findings of the Association for Computational Linguistics: ACL 2025*, pp. 9503–9523, 2025.
- Oyvind Tafjord, Bhavana Dalvi, and Peter Clark. Proofwriter: Generating implications, proofs, and abductive statements over natural language. In *Findings of the Association for Computational Linguistics: ACL-IJCNLP 2021*, pp. 3621–3634, 2021.
- Yuxuan Wan, Wenxuan Wang, Yiliu Yang, Youliang Yuan, Jen-tse Huang, Pinjia He, Wenxiang Jiao, and Michael R Lyu. Logicasker: Evaluating and improving the logical reasoning ability of large language models. In *EMNLP*, 2024.
- Junlin Wang, WANG Jue, Ben Athiwaratkun, Ce Zhang, and James Zou. Mixture-of-agents enhances large language model capabilities. In *The Thirteenth International Conference on Learning Representations*.
- Lei Wang, Wanyu Xu, Yihuai Lan, Zhiqiang Hu, Yunshi Lan, Roy Ka-Wei Lee, and Ee-Peng Lim. Plan-and-solve prompting: Improving zero-shot chain-of-thought reasoning by large language models. In *Proceedings of the 61st Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pp. 2609–2634, 2023.

- Zhongsheng Wang, Jiamou Liu, Qiming Bao, Hongfei Rong, and Jingfeng Zhang. Chatlogic: Integrating logic programming with large language models for multi-step reasoning. In 2024 International Joint Conference on Neural Networks (IJCNN), pp. 1–8. IEEE, 2024.
 - 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.
 - F. Wu, Z. Ren, Z. Sha, Z. Fu, Z. Xu, Z. Huang, Z. Zhang, Z. Xie, Z. Hao, Z. Gou, Z. Ma, Z. Yan, Z. Shao, Z. Xu, Z. Wu, Z. Zhang, Z. Li, Z. Gu, Z. Liu, Z. Li, Z. Xie, Z. Song, Z. Gao, Z. Pan, et al. Deepseek-v3 technical report. 2024. URL https://arxiv.org/abs/2412.19437.
 - Jundong Xu, Hao Fei, Liangming Pan, Qian Liu, Mong-Li Lee, and Wynne Hsu. Faithful logical reasoning via symbolic chain-of-thought. In *Proceedings of the 62nd Annual Meeting of the Association for Computational Linguistics*, 2024.
 - Jundong Xu, Hao Fei, Meng Luo, Qian Liu, Liangming Pan, William Yang Wang, Preslav Nakov, Mong-Li Lee, and Wynne Hsu. Aristotle: Mastering logical reasoning with a logic-complete decompose-search-resolve framework. In *Proceedings of the 62nd Annual Meeting of the Association for Computational Linguistics*, 2025.
 - Sen Yang, Xin Li, Leyang Cui, Lidong Bing, and Wai Lam. Neuro-symbolic integration brings causal and reliable reasoning proofs. *arXiv* preprint arXiv:2311.09802, 2023.
- Shunyu Yao, Dian Yu, Jeffrey Zhao, Izhak Shafran, Tom Griffiths, Yuan Cao, and Karthik Narasimhan. Tree of thoughts: Deliberate problem solving with large language models. *Advances in neural information processing systems*, 36:11809–11822, 2023.
- Xi Ye, Qiaochu Chen, Isil Dillig, and Greg Durrett. Satlm: Satisfiability-aided language models using declarative prompting. *Advances in Neural Information Processing Systems*, 36:45548–45580, 2023.
- Yifan Zhang, Jingqin Yang, Yang Yuan, and Andrew Chi-Chih Yao. Cumulative reasoning with large language models. *arXiv preprint arXiv:2308.04371*, 2023.
- Yifan Zhang, Yang Yuan, and Andrew Chi-Chih Yao. On the diagram of thought. *arXiv preprint* arXiv:2409.10038, 2024.

Table 6: Per-dataset cost-effectiveness comparison. Tokens are prefill tokens per question (\downarrow) , accuracy is in % (\uparrow) .

Model	Method	ProntoQA		ProofWriter		LogicalDeduct	
		Tokens ↓	Acc (%) ↑	Tokens ↓	Acc (%) ↑	Tokens ↓	Acc (%) ↑
	SparseMAD	37,784	99.80	41,678	89.50	43,635	88.67
CDT 4	CortexDebate	35,973	99.60	37,554	90.83	45,487	92.33
GPT-4	Ours (w/o sparse)	46,502	99.40	51,358	90.17	54,857	94.00
	Ours (w/ sparse)	35,854	100.00	42,617	92.00	54,171	94.33
	SparseMAD	79,456	99.80	48,190	92.83	44,245	99.83
Claude 3.7	CortexDebate	68,023	99.80	41,897	96.17	45,962	99.67
Claude 3.7	Ours (w/o sparse)	106,015	100.00	63,105	97.00	68,636	99.67
	Ours (w/ sparse)	52,923	100.00	62,204	96.83	47,121	100.00
	SparseMAD	40,200	98.00	18,527	92.50	53,257	95.33
DeepSeek-V3	CortexDebate	35,381	99.80	19,349	93.00	47,388	99.67
	Ours (w/o sparse)	57,366	99.80	25,059	92.83	70,464	100.00
	Ours (w/ sparse)	36,702	100.00	24,115	93.33	46,107	100.00

Table 7: Aggregate performance across three benchmarks. Token Saving and ΔAcc are relative to *Ours (w/o sparse)*.

Model	Method	Avg. Acc (%) ↑	Avg. Tokens	Token Saving (%) ↑	ΔAcc (pp) ↑
	SparseMAD	92.66	41,032	19.40	-1.87
GPT-4	CortexDebate	94.39	39,671	22.07	-0.14
GP1-4	Ours (w/o sparse)	94.52	50,906	0.00	+0.00
	Ours (w/ sparse)	95.44	44,214	13.15	+0.92
	SparseMAD	97.49	57,297	27.70	-1.40
Claude 3.7	CortexDebate	98.61	51,961	34.44	-0.28
Claude 3.7	Ours (w/o sparse)	98.89	79,252	0.00	+0.00
	Ours (w/ sparse)	98.94	54,082	31.76	+0.05
	SparseMAD	95.28	37,328	26.75	-2.27
DeepSeek-V3	CortexDebate	97.49	34,039	33.21	-0.05
	Ours (w/o sparse)	97.54	50,963	0.00	+0.00
	Ours (w/ sparse)	97.78	35,641	30.06	+0.24

USAGE OF AI

In this work, we made limited use of LLMs as an assistive writing tool. Specifically, we used LLMs to replace synonyms, restructure sentences, and brainstorm alternative ways of expressing ideas within paragraphs. All conceptual contributions, research design, experiments, analyses, and final writing decisions were made by the authors. The authors take full responsibility for the accuracy and originality of the content.

A Cost-Effectiveness Analysis

We evaluate the cost-effectiveness of our sparse communication approach by measuring token consumption and accuracy across three LLMs and three benchmarks. Following our evaluation protocol, we report *prefill tokens per question* as a reproducible cost proxy and accuracy as effectiveness; lower tokens are better (\downarrow) , higher accuracy is better (\uparrow) . We do not report wall-clock time due to API jitter; tokens serve as a stable, reproducible proxy for runtime and dollar cost.

In our experiments, *Ours* (*w/o sparse*) approximates a fully-connected debate topology where all agents communicate in each round, while *Ours* (*w/ sparse*) uses our adaptive sparse communication gate to selectively prune interactions based on confidence and information gains.

Our adaptive sparse gate achieves the highest accuracy while keeping token costs comparable to strong baselines. As shown in Table 6 and Table 7, our sparse method consistently outperforms the

fully-connected baseline on all three models, achieving both higher accuracy (+0.92pp on GPT-4, +0.05pp on Claude 3.7, +0.24pp on DeepSeek-V3) and substantial token savings (13–36%). Remarkably, it also surpasses existing multi-agent baselines (SparseMAD and CortexDebate) in accuracy while maintaining competitive token efficiency. This demonstrates that our confidence-based pruning mechanism not only reduces computational overhead but also improves reasoning quality by filtering redundant inter-agent communications.

B HANDLING SYMBOLIC SOLVER FAILURES

During the symbolic reasoning stage, solvers may occasionally fail to execute the translated logical expressions due to syntax errors, incompatible formula structures, or computational timeouts. Since our multi-agent framework relies on symbolic solvers (Pyke, Prover9, and Z3) to provide formal reasoning, handling these execution failures appropriately is crucial for maintaining system robustness.

Table 8 presents the impact of different failure handling strategies on final accuracy across three benchmarks using GPT-4. We evaluate three strategies:

- **Random**: When a solver fails, the agent randomly selects an answer from the available options. This serves as a baseline strategy.
- **Discard**: Failed solver agents are excluded from the debate, and only successfully executed agents participate in subsequent rounds and final voting.
- **Simulate**: When a solver fails, we prompt the corresponding agent to simulate the solver's reasoning process using the LLM's inherent logical capabilities, effectively falling back to natural language reasoning while maintaining the agent's role in the debate.

Table 8: Final accuracy (%) under different handling strategies when a symbolic solver fails (GPT-4).

Strategy	ProntoQA	ProofWriter	LogicalDeduction
Random Discard	99.20% 99.80%	89.83% 91.33%	91.33% 93.67%
Simulate	100.00%	91.33% 92.00 %	93.07% 94.33 %

The results demonstrate that the *Simulate* strategy consistently achieves the best performance across all benchmarks. This approach leverages the LLM's ability to approximate symbolic reasoning when formal execution fails, maintaining full agent participation while providing reasonable fallback reasoning. The *Discard* strategy performs better than random selection but loses valuable perspectives from failed agents. These findings suggest that maintaining agent diversity through simulation is more beneficial than excluding agents, even when their primary symbolic reasoning mechanism fails.

C ADDITIONAL EXPERIMENTAL RESULTS ON DEEPSEEK-V3

This section presents additional experimental results for DeepSeek-V3 that show similar patterns to the GPT-4 results discussed in the main paper.

C.1 COMMUNICATION THRESHOLD ANALYSIS

Figure 5 shows the effect of communication gating threshold on accuracy and token saving rate for DeepSeek-V3. The results demonstrate patterns consistent with GPT-4, achieving token reduction while maintaining high accuracy.

C.2 TRANSLATION QUALITY ANALYSIS

Figure 6 shows the relationship between debate rounds and solver execution rates for DeepSeek-V3. Consistent with our GPT-4 findings, the execution rate increases during the first 1-2 rounds and then shows diminishing returns.

D ADDITIONAL EXPERIMENTAL RESULTS ON CLAUDE 3.7 SONNET

This section provides supplementary experimental results for Claude 3.7 Sonnet.

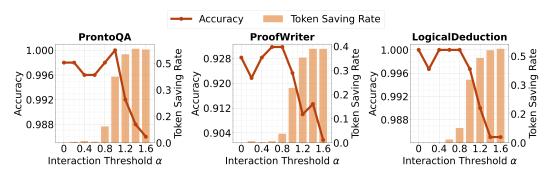


Figure 5: Effect of communication gating threshold on accuracy and token saving rate on DeepSeek-V3.

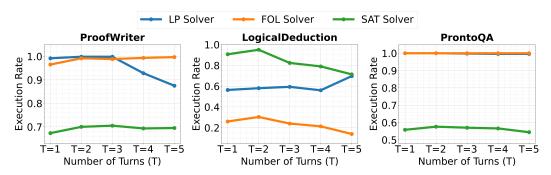


Figure 6: Relation between debate rounds and solver execution rate for DeepSeek-V3.

D.1 COMMUNICATION THRESHOLD ANALYSIS

Figure 7 presents the accuracy-efficiency trade-off for Claude 3.7 Sonnet. Similar to GPT-4 and DeepSeek-V3, Claude 3.7 maintains high accuracy while achieving significant token savings through sparse communication.

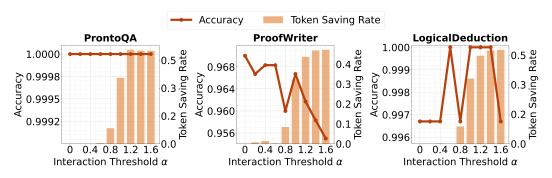


Figure 7: Effect of communication gating threshold on accuracy and token saving rate on Claude 3.7 Sonnet.

D.2 TRANSLATION QUALITY ANALYSIS

Figure 8 illustrates the translation quality dynamics for Claude 3.7 Sonnet. The pattern is consistent with other models: execution rates improve significantly within the first 2-3 debate rounds, validating our multi-agent debate approach for translation refinement.

E MULTI-TURN INTERACTION ALGORITHM FOR SPARSE COMMUNICATION

The sparse communication mechanism directly informs how each agent updates its internal state or memory across debate rounds. Each agent maintains a personalized memory that aggregates

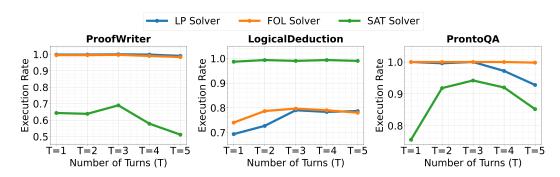


Figure 8: Relation between debate rounds and solver execution rate for Claude 3.7 Sonnet.

Algorithm 1: Multi-Turn Interaction Algorithm for Enhancing LLMs' Logical Reasoning

Input: Communication rounds D, Agent number n, hyper-parameter λ ; 1 Translate raw logical question Q to symbolic expression Sym(Q);

```
2 M_1^{d=1},\ldots,M_n^{d=1}\leftarrow\varnothing; 3 for d\in\{1,\ldots,D\} do
                        O_{i \to j}^d = 1 \text{ for all } i, j \in \{1, \dots, n\};
                        Compute \operatorname{Pre}_{i \to j}^d = \frac{C_i^d}{C_i^d} + \lambda (1 - \cos(A_j^d, A_i^d)) for all i \neq j;
                        \text{Compute } \overline{\text{Pre}_{i \to j}^d} = \frac{1}{d} (\overline{\text{Pre}_{i \to j}^{d-1}} \cdot (d-1) + \frac{C_i^d}{C_i^d} + \lambda (1 - \cos(A_j^d, A_i^d))) \text{ for all } i \neq j;
                       if \operatorname{Pre}_{i 	o j}^d < \alpha \cdot \overline{\operatorname{Pre}_{i 	o j}^{d-1}} then
833
                          O_{i\to j}^d = 0;
                        for s \in \{1, ..., n\} do
836
                                // Memory update of the s{	ext{-}}{	ext{th}} agent at round d
                               \begin{array}{l} M_s^{d+1} \leftarrow M_s^d \cup \{A_i^d \mid i \in \{1,\dots,n\}, O_{i \to s}^d = 1\}; \\ // \text{ Output of the } s\text{-th agent at round } d \text{ using personalized} \end{array}
            10
838
                                A_s^{d+1} \leftarrow \text{LLM}_s(\text{Sym}(Q)||M_s^{d+1});
            11
```

12 Majority vote among the *n* agents $A_1^{D+1}, \ldots, A_n^{D+1}$;

valuable insights from others. At the beginning of the first round (d = 1), all agents start with an empty memory $M_s^1 \leftarrow \emptyset$ and communication is fully connected $(O_{i \to j}^d = 1 \text{ for all pairs})$. From the second round, the sparse communication gate $O_{i \to j}^d$ is activated. At the end of each round d, every agent s updates its memory for the next round M_s^{d+1} , by selectively incorporating the outputs A_i^d from only those agents i for which the communication channel was open (i.e., $O_{i \to i}^d = 1$). After the memory updated, agent s generates its output for the next round A_i^{d+1} , by querying the symbolic question and i's newly updated, personalized memory. After D rounds of debate, the final outputs from all agents $A_1^{D+1}, \ldots, A_n^{D+1}$, are aggregated via a majority vote to determine the final answer.

PROMPT TEMPLATES

TRANSLATION DEBATE

Translation Prompt

Task. You are given a logic problem in natural language including a context and a question as follows: Context: \${context} Question: \${question}

Discussion Rules

- 864 865
- 866 867
- 868 870
- 871 872
- 873 874
- 875 876 877
- 879

- 882 883
- 884 885
- 887 888
- 889 890 891
- 892 893
- 894 895
- 896 897
- 899 900 901

902 903 904

905

910 911 912

913

914

915 916 917 1. Syntax Verification: Carefully review previous discussions to understand others' translations. While maintaining your own symbolic language system, check and correct any syntax errors in your translation (e.g., unclosed parentheses, malformed expressions).

- 2. Completeness Check: Review others' translations to understand their interpretation of the natural language problem. While keeping your own symbolic language system, verify and correct the information completeness of your translation (no missing/extra facts, rules, predicates, or statements from the original problem).
- 3. Language Independence: When referencing others' translations, you must maintain your own symbolic language system. Do not adopt symbols or syntax from other languages.

Discussion history

```
${chat_history}
```

Role-specific description

\${role_description}

Now it's your turn to speak. Please speak as concisely and clearly as possible

Role-specific description — LP translator

Your task is to translate the logic problem in natural language into LP logic formulas:

- 1. define all the predicates in the problem
- 2. parse the problem into logic rules based on the defined predicates
- 3. write all the facts mentioned in the problem
- 4. parse the question into the logic form (Use && to represent AND, and you cannot use NOT or other negations in LP)

Example

```
Context: Each jompus is fruity.
(... more context here ...)
Rompuses are zumpuses. Alex is a tumpus.
Question: True or false: Alex is not shy.
Predicates:
Jompus(x, bool) ::: Does x belong to Jompus?
(... more predicates here ...)
Zumpus(x, bool) ::: Does x belong to Zumpuses?
Tumpuses (Alex, True)
Rules:
Jompus ($x, True) >>> Fruity ($x, True)
(... more rules here ...)
Dumpus($x, True) >>> Rompus($x, True)
Query:
Shy(Alex, False)
```

Role-specific description — FOL translator

Your task is to translate the logic problem in natural language into first-order logic formulas. The grammar of first-order logic is defined as follows:

```
918
          logical conjunction:
                                        \operatorname{expr}_1 \wedge \operatorname{expr}_2
919
          logical disjunction:
                                        \operatorname{expr}_1 \vee \operatorname{expr}_2
920
          logical exclusive disjunction:
                                        \operatorname{expr}_1 \oplus \operatorname{expr}_2
921
                                        \neg expr_1
          logical negation:
922
          expr_1 implies expr_2:
                                        expr_1 \rightarrow expr_2
          expr_1 iff expr_2:
923
                                        \operatorname{expr}_1 \leftrightarrow \operatorname{expr}_2
          logical universal quantification:
                                        \forall x
924
          logical existential quantification:
                                        \exists x
925
926
          Output format: logic form ::: description
927
          Example
928
          Context: All people who regularly drink coffee are
929
          dependent on caffeine.
930
          (... more context here ...)
          If Rina is not a person dependent on caffeine and a
931
          student, then Rina is either a person dependent
932
          on caffeine and a student, or a person dependent
933
          on caffeine nor a student, or neither a person
934
          dependent on caffeine nor a student.
935
936
          Question: Based on the above information, is the
937
          following statement true, false, or uncertain?
          Rina is either a person who jokes about being
938
          addicted to caffeine or is unaware that caffeine
939
          is a drug.
940
941
          Predicates:
          Dependent(x) ::: x is a person dependent on caffeine
942
          (... more predicates here ...)
943
          Student(x) ::: x is a student
944
945
          Premises:
946
          \pi x (Drinks(x) \pi x (Drinks(x)) $\rightarrow$ Dependent(x)) ::: All people who
          regularly drink coffee are dependent on caffeine.
947
          (... more premises here ...)
948
          \pi \ (Jokes(x) \pi \ \rightarrow$ \pi \ \neg$Unaware(x)) ::: No one who
949
               jokes
950
          about being addicted to caffeine is unaware that
951
          caffeine is a drug.
952
          Conclusion:
953
                         Unaware(rina) ::: Rina is either a person who jokes about
          Jokes (rina)
954
               being addicted to caffeine
955
          or is unaware that caffeine is a drug.
956
957
```

Role-specific description — SAT translator

Your task is to parse the logic problem in natural language as a SAT problem using Z3 syntax, defining declarations, constraints, and options.

- 1. Always include all three section headers in order: # Declarations, # Constraints, # Options
- 2. Declarations must follow exact patterns:
 - name = EnumSort([items, ...]) for non-numeric items
 - name = IntSort([numbers, ...]) for numeric items
 - name = Function([types] -> [return_type])
- 3. Constraints support:

958

959 960

961

962

963

964

965

966

969 970

- Direct expressions with ==, !=, <=, >=, <, >, Implies(), And(), Or(), Not()
- ForAll([var:type, ...], expr) and Exists([var:type, ...], expr)
- Count([var:type], condition)

```
972
           • Distinct([var:type], expr)
973
        4. Options must use predefined functions:
           • is_valid(), is_sat(), is_unsat()
975
        5. Add explanation with :::
976
        6. Avoid:
977
           • Add # in any other places apart from three section headers
978
           · Add any other unnecessary comment or dashes
979
        Example
980
        Context: Bob is cold. Bob is quiet. Bob is red. Bob is smart. Charlie
981
            is kind. Charlie is quiet. Charlie is red. Charlie is rough. Dave
982
            is cold. Dave is kind. Dave is smart. Fiona is quiet. If something
983
             is quiet and cold then it is smart. Red, cold things are round.
984
            If something is kind and rough then it is red. All quiet things
            are rough. Cold, smart things are red. If something is rough then
985
            it is cold. All red things are rough. If Dave is smart and Dave is
986
             kind then Dave is quiet.
987
        Question: True or false: Charlie is kind.
988
989
        # Declarations
        objects = EnumSort([Bob, Charlie, Dave, Fiona])
990
        attributes = EnumSort([cold, quiet, red, smart, kind, rough, round])
991
        has_attribute = Function([objects, attributes] -> [bool])
992
993
        # Constraints
994
        has_attribute(Bob, cold) == True ::: Bob is cold.
        has_attribute(Bob, quiet) == True ::: Bob is quiet.
995
        has_attribute(Bob, red) == True ::: Bob is red.
996
        has_attribute(Bob, smart) == True ::: Bob is smart.
997
        has_attribute(Charlie, kind) == True ::: Charlie is kind.
998
        has_attribute(Charlie, quiet) == True ::: Charlie is quiet.
999
        has_attribute(Charlie, red) == True ::: Charlie is red.
        has_attribute(Charlie, rough) == True ::: Charlie is rough.
1000
        has_attribute(Dave, cold) == True ::: Dave is cold.
1001
        has_attribute(Dave, kind) == True ::: Dave is kind.
1002
        has_attribute(Dave, smart) == True ::: Dave is smart.
1003
        has_attribute(Fiona, quiet) == True ::: Fiona is quiet.
1004
        ForAll([x:objects], Implies(And(has_attribute(x, quiet) == True,
            has_attribute(x, cold) == True), has_attribute(x, smart) == True))
1005
             ::: If something is quiet and cold then it is smart.
1006
        ForAll([x:objects], Implies(And(has_attribute(x, red) == True,
1007
            has_attribute(x, cold) == True), has_attribute(x, round) == True))
1008
             ::: Red, cold things are round.
1009
        ForAll([x:objects], Implies(And(has_attribute(x, kind) == True,
            has_attribute(x, rough) == True), has_attribute(x, red) == True))
1010
            ::: If something is kind and rough then it is red.
1011
        ForAll([x:objects], Implies(has_attribute(x, quiet) == True,
1012
            has_attribute(x, rough) == True)) ::: All quiet things are rough.
1013
        ForAll([x:objects], Implies(And(has_attribute(x, cold) == True,
1014
            has_attribute(x, smart) == True), has_attribute(x, red) == True))
            ::: Cold, smart things are red.
1015
        ForAll([x:objects], Implies(has_attribute(x, rough) == True,
1016
            has_attribute(x, cold) == True)) ::: If something is rough then it
1017
             is cold.
1018
        ForAll([x:objects], Implies(has_attribute(x, red) == True,
1019
            has_attribute(x, rough) == True)) ::: All red things are rough.
        Implies (And (has_attribute (Dave, smart) == True, has_attribute (Dave,
1020
            kind) == True), has_attribute(Dave, quiet) == True) ::: If Dave is
1021
             smart and Dave is kind then Dave is quiet.
1022
1023
        # Options
1024
        is_valid(has_attribute(Charlie, kind) == True) ::: Charlie is kind is
            True (A).
1025
```

1026 is_unsat(has_attribute(Charlie, kind) == True) ::: Charlie is kind is 1027 False (B). 1028 1029 1030 F.2 REASONING DEBATE 1031 **Final Debate Prompt** 1032 1033 You are given a logic problem that contains a context, a question, and options: 1034 Context: \${context} Question: \${question} 1035 Options: \${options} 1036 Role description: \${Role-specific description} 1037 Your initial answer is \$ {predict }. Your initial reasoning is: \${reasoning}. 1039 1040 You are now in a collaborative debate with other reasoning agents. Your goal is to reach the correct answer through discussion. 1041 1042 **Important Debate Rules** 1043 1. Review other agents' arguments in the discussion history first. 1044 2. Identify specific points of agreement or disagreement. 3. Challenge weak reasoning with concrete counterexamples. 1046 No need to repeat your whole reasoning if your argument remains unchanged. 1047 5. Acknowledge other arguments when you find them correct, even if they contradict your initial position. 1048 6. If you change your answer, always explain why you changed. 1049 7. Be willing to change your answer if convinced by other arguments. 1050 8. Reference specific agents and their arguments when responding. 1051 9. Be interactive and engaging with other agents! 1052 Discussion history 1053 \${chat_history} 1054 **Turn-specific instruction** 1055 \${turn_specific_instruction} 1056 1057 Role-specific description — LP supporter 1058 1059 You are a supporter of the Logic Programming (LP) approach. **Strengths:**

- Systematic rule-based reasoning with clear steps
- Handle complex relations via predicates and rules
- Transparent reasoning process verifiable step by step
- Strong foundation in formal logic and theorem proving

In the debate, you should:

1061

1062

1063

1064

1065

1067

1068 1069

1070

1071

1072

1074

1077

1078

1079

- Emphasize rigor and reliability of LP reasoning
- Highlight systematic application of logical rules
- · Defend transparency and verifiability
- · Challenge others when lacking formal logical foundation

Role-specific description — FOL supporter

You are a supporter of First-Order Logic (FOL). Strengths:

- Mathematical precision with quantifiers and operators
- · Express complex relationships precisely
- Sound theoretical foundation
- Handle nested structures and implications

In the debate, you should:

- · Emphasize rigor and expressiveness of FOL
- · Highlight formal completeness and soundness
- · Defend against criticisms of complexity
- Challenge others when lacking precision

Role-specific description — SAT supporter

You are a supporter of the SAT/SMT (Z3) approach. Strengths:

- Formal representation using Z3 syntax
- · Complete and sound reasoning
- Ability to prove validity/satisfiability
- Handle quantifiers and complex formulas

In the debate, you should:

1080

1081 1082

1083

1084

1086

1087

1088

1089

1090 1091

1093

1094

1095

1098 1099 1100

1101

1102

1103

1104

1105

1106

11071108

1109

1110

1111 1112

1113

1114

1115

1116

1117

1118

1119

1120

1121

1122

1123 1124

1125

1126

1127

1128

1130

1131

1132

1133

- · Emphasize formal correctness and completeness
- Highlight Z3's power in constraint solving
- Defend clarity of declarative specification
- Challenge others on missed edge cases

Role-specific description — Plan-and-Solve supporter

You are a supporter of the Plan-and-Solve approach. Strengths:

- · First devise a plan, then execute it
- · Structured reasoning framework

In the debate, you should:

- Defend the reasoning process and solution steps
- · Justify the plan structure and accuracy

Role-specific description — COT Solver supporter

You are the Chain-of-Thought (CoT) supporter. **Strengths:**

· Step-by-step reasoning from facts to conclusions

In the debate, you should:

- Defend solutions generated via CoT prompting
- Validate each inference in the reasoning chain
- · Argue for soundness of the conclusion

G CASE STUDIES

G.1 REASONINGS EXTRACTED FROM SYMBOLIC SOLVERS

Case Study of Reasonings Extracted

Context. The cat eats the squirrel. The cat sees the squirrel. The cow eats the squirrel. The cow sees the cat. The rabbit is round. The rabbit sees the cat. The squirrel eats the rabbit. The squirrel is cold. The squirrel needs the rabbit. The squirrel sees the cat. If someone sees the cat and they are not green then they see the cow. If the rabbit is kind and the rabbit sees the squirrel then the squirrel needs the rabbit. Rough people are cold. If someone sees the rabbit then they are not round. If someone sees the squirrel and they are not green then they need the squirrel. If someone eats the cow then they see the rabbit. If someone eats the squirrel then they are rough. If someone is cold then they eat the cow.

Question. Based on the above information, is the following statement true, false, or unknown? *The cat is not round.*

Options. A) True B) False C) Unknown

LP Reasoning Extracted from Pyke Solver:

```
We first define following predicates and corresponding natural language explanations:

Eats($x, $y, bool) ::: Does x eat y?

Sees($x, $y, bool) ::: Does x see y?

Round($x, bool) ::: Is x round?

Cold($x, bool) ::: Is x cold?

Needs($x, $y, bool) ::: Does x need y?

Green($x, bool) ::: Is x green?

Kind($x, bool) ::: Is x round?

Rough($x, bool) ::: Is x round?

We have following known facts from the context:

Eats(cat, squirrel, True)
```

```
1134
1135
              Sees(cat, squirrel, True)
              Eats(cow, squirrel, True)
1136
              Sees(cow, cat, True)
1137
              Round(rabbit, True)
              Sees (rabbit, cat, True)
1138
              Eats(squirrel, rabbit, True)
              Cold(squirrel, True)
1139
              Needs(squirrel, rabbit, True)
1140
              Sees(squirrel, cat, True)
1141
            We have following known rules from the context:
              rule1: Sees($x, cat, True) && Green($x, False) >>> Sees($x, cow, True)
1142
              rule2: Kind(rabbit, True) && Sees(rabbit, squirrel, True) >>> Needs(squirrel, rabbit,
1143
                 True)
              rule3: Rough($x, True) >>> Cold($x, True)
1144
              rule4: Sees($x, rabbit, True) >>> Round($x, False)
1145
              rule5: Sees($x, squirrel, True) && Green($x, False) >>> Needs($x, squirrel, True)
              rule6: Eats($x, cow, True) >>> Sees($x, rabbit, True)
1146
              rule7: Eats($x, squirrel, True) >>> Rough($x, True)
1147
              rule8: Cold($x, True) >>> Eats($x, cow, True)
            Now begin reasoning to obtain all implied facts:
1148
            Use rule1: Sees($x, cat, True) && Green($x, False) >>> Sees($x, cow, True)
1149
            Bind $x to 'cow'
            Unbind $x
1150
            Bind $x to 'rabbit'
1151
            Unbind $x
            Bind $x to 'squirrel'
1152
            Unbind $x
1153
            Finish implied with rule1
            Use rule2: Kind(rabbit, True) && Sees(rabbit, squirrel, True) >>> Needs(squirrel,
1154
                rabbit, True)
1155
            Finish implied with rule2
            Use rule3: Rough($x, True) >>> Cold($x, True)
1156
            Finish implied with rule3
1157
            Use rule4: Sees($x, rabbit, True) >>> Round($x, False)
            Finish implied with rule4
1158
            Use rule5: Sees($x, squirrel, True) && Green($x, False) >>> Needs($x, squirrel, True)
            Bind $x to 'cat
1159
            Unbind $x
1160
            Finish implied with rule5
1161
            Use rule6: Eats($x, cow, True) >>> Sees($x, rabbit, True)
            Finish implied with rule6
1162
            Use rule7: Eats(x, squirrel, True) >>> Rough(x, True)
            Bind $x to 'cat'
1163
            Obtain a new implied fact: Rough ('cat', True)
1164
            Reuse rule3: Rough($x, True) >>> Cold($x, True)
1165
            Bind $x to 'cat'
            Obtain a new implied fact: Cold('cat', True)
1166
            Use rule8: Cold($x, True) >>> Eats($x, cow, True)
1167
            Finish implied with rule8
            Finish implied with rule3
1168
            Reuse rule3: Rough(x, True) >>> Cold(x, True)
1169
            Bind $x to 'cow'
            Obtain a new implied fact: Cold('cow', True)
1170
            Reuse rule8: Cold($x, True) >>> Eats($x, cow, True)
1171
            Finish implied with rule8
            Finish implied with rule3
1172
            Finish implied with rule7
1173
            Reuse rule8: Cold($x, True) >>> Eats($x, cow, True)
            Bind $x to 'squirrel'
1174
            Obtain a new implied fact: Eats('squirrel', 'cow', True)
1175
            Reuse rule6: Eats($x, cow, True) >>> Sees($x, rabbit, True)
            Bind $x to 'squirrel'
1176
            Obtain a new implied fact: Sees('squirrel', 'rabbit', True)
1177
            Reuse rule1: Sees(x, cat, True) && Green(x, False) >>> Sees(x, cow, True)
            Bind $x to 'squirrel'
1178
            Finish implied with rule1
1179
            Reuse rule2: Kind(rabbit, True) && Sees(rabbit, squirrel, True) >>> Needs(squirrel,
                 rabbit, True)
1180
            Finish implied with rule2
            Reuse rule4: Sees($x, rabbit, True) >>> Round($x, False)
1181
            Bind $x to 'squirrel'
1182
            Obtain a new implied fact: Round('squirrel', False)
1183
            Finish implied with rule4
            Reuse rule5: Sees($x, squirrel, True) && Green($x, False) >>> Needs($x, squirrel, True)
1184
            Bind $x to 'squirrel'
            Finish implied with rule5
1185
            Finish implied with rule6
1186
            Reuse rule7: Eats($x, squirrel, True) >>> Rough($x, True)
1187
```

```
1188
1189
             Bind $x to 'squirrel'
             Finish implied with rule7
1190
             Reuse rule6: Eats($x, cow, True) >>> Sees($x, rabbit, True)
1191
             Bind $x to 'cat'
             Obtain a new implied fact: Sees('cat', 'rabbit', True)
1192
             Reuse rule1: Sees(\$x, cat, True) && Green(\$x, False) >>> Sees(\$x, cow, True)
1193
             Bind $x to 'cat'
             Finish implied with rule1
1194
             Reuse rule2: Kind(rabbit, True) && Sees(rabbit, squirrel, True) >>> Needs(squirrel,
1195
                   rabbit, True)
             Finish implied with rule2
1196
             Reuse rule4: Sees($x, rabbit, True) >>> Round($x, False)
1197
             Bind $x to 'cat'
             Obtain a new implied fact: Round('cat', False)
1198
             Finish implied with rule4
             Reuse rule5: Sees($x, squirrel, True) && Green($x, False) >>> Needs($x, squirrel, True)
1199
             Bind $x to 'cat'
1200
             Finish implied with rule5
             Finish implied with rule6
1201
             Reuse rule7: Eats($x, squirrel, True) >>> Rough($x, True)
1202
             Bind $x to 'cat'
1203
             Finish implied with rule7
             Reuse rule6: Eats($x, cow, True) >>> Sees($x, rabbit', True)
1204
             Bind $x to 'cow'
1205
             Obtain a new implied fact: Sees('cow', 'rabbit', True)
             Reuse rule1: Sees($x, cat, True) && Green($x, False) >>> Sees($x, cow, True)
1206
             Bind $x to 'cow'
1207
             Finish implied with rule1
             Reuse rule2: Kind(rabbit, True) && Sees(rabbit, squirrel, True) >>> Needs(squirrel,
1208
                  rabbit, True)
1209
             Finish implied with rule2
             Reuse rule4: Sees($x, rabbit, True) >>> Round($x, False)
1210
             Bind $x to 'cow'
1211
             Obtain a new implied fact: Round('cow', False)
             Finish implied with rule4
1212
             Reuse rule5: Sees(\$x, squirrel, True) && Green(\$x, False) >>> Needs(\$x, squirrel, True)
             Bind $x to 'cow'
1213
             Finish implied with rule5
1214
             Finish implied with rule6
1215
             Reuse rule7: Eats(x, squirrel, True) >>> Rough(x, True)
             Bind $x to 'cow'
1216
             Finish implied with rule7
1217
             Finish implied with rule8
             All newly implied Facts: Cold('cat', True), Cold('cow', True), Eats('squirrel', 'cow', True), Rough('cat', True), Round('cat', False), Round('cow', False), Round('squirrel', False), Sees('cat', 'rabbit', True), Sees('cow', 'rabbit', True), Sees ('squirrel', 'rabbit', True)
1218
1219
1220
             Finish reasoning
1221
             FOL Reasoning Extracted from Prover9 Solver:
1222
             prove original conclusion:
1223
             3 (all x (Rough(x) \rightarrow Cold(x))). [assumption].
1224
             4 (all x (Sees(x, Rabbit) \rightarrow -Round(x))). [assumption].
             6 (all x (Eats(x,Cow) -> Sees(x,Rabbit))). [assumption].
1225
             7 (all x (Eats(x, Squirrel) \rightarrow Rough(x))). [assumption].
1226
             8 (all x (Cold(x) \rightarrow Eats(x,Cow))). [assumption].
             9 -Round(Cat). [goal].
10 -Eats(x,Cow) | Sees(x,Rabbit). [clausify(6)].
1227
1228
             11 Eats(Cat,Squirrel). [assumption].
14 -Eats(x,Squirrel) | Rough(x). [clausify(7)].
1229
             15 -Cold(x) \mid Eats(x, Cow). [clausify(8)].
1230
             16 -Sees(x, Rabbit) | -Round(x). [clausify(4)].
1231
             18 Round (Cat). [deny (9)].
             19 -Cold(x) | Sees(x, Rabbit). [resolve(15, b, 10, a)].
1232
             21 -Rough(x) \mid Cold(x). [clausify(3)].
1233
             27 -Sees(Cat, Rabbit). [resolve(18, a, 16, b)].
             29 Sees(x,Rabbit) | -Rough(x). [resolve(19,a,21,b)].
1234
             31 Rough(Cat). [resolve(14,a,11,a)].
33 -Rough(Cat). [resolve(29,a,27,a)].
1235
             34 $F. [resolve(33,a,31,a)].
1236
1237
```

1242 G.2 Translate Debate 1243 Case Study from logical_deduction 1244 1245 **Context:** A fruit stand sells five fruits: mangoes, kiwis, plums, pears, and watermelons. The 1246 kiwis are less expensive than the plums. The pears are the third-most expensive. The kiwis 1247 are the second-cheapest. The watermelons are the most expensive. 1248 **Question:** Which of the following is true? 1249 A) The mangoes are the third-most expensive. 1250 B) The kiwis are the third-most expensive. 1251 C) The plums are the third-most expensive. 1252 D) The pears are the third-most expensive. 1253 E) The watermelons are the third-most expensive. Answer: D 1255 Round 1 1256 LP translator: 1257 I'll translate this logic problem about fruit prices into LP format: **Predicates:** 1259 Fruit(\$x, bool) ::: \$x is one of the five fruitsMoreExpensive(\$x, \$y, bool) ::: Fruit \$x is more expensive than fruit \$y 1261 LessExpensive(\$x, \$y, bool) ::: Fruit \$x is less expensive than fruit \$y 1262 MostExpensive(\$x, bool) ::: Fruit \$x is the most expensive 1263 SecondMostExpensive(\$x, bool) ::: Fruit \$x is the second-most expensive 1264 ThirdMostExpensive(\$x, bool) ::: Fruit \$x is the third-most expensive FourthMostExpensive(\$x, bool) ::: Fruit \$x is the fourth-most expensive (second-cheapest) 1265 LeastExpensive(\$x, bool) ::: Fruit \$x is the least expensive (cheapest) 1266 **Facts:** 1267 Fruit(mango, True) 1268 Fruit(kiwi, True) 1269 Fruit(plum, True) 1270 Fruit(pear, True) 1271 Fruit(watermelon, True) 1272 LessExpensive(kiwi, plum, True) ::: Kiwis are less expensive than plums ThirdMostExpensive(pear, True) ::: Pears are the third-most expensive 1274 FourthMostExpensive(kiwi, True) ::: Kiwis are the second-cheapest (fourth-most expensive) MostExpensive(watermelon, True) ::: Watermelons are the most expensive 1276 1277 LessExpensive(\$x, \$y, True) >> MoreExpensive(\$y, \$x, True) 1278 MoreExpensive(\$x, \$y, True) >> LessExpensive(\$y, \$x, True) 1279 Query: 1280 ThirdMostExpensive(mango, True) ::: Option A 1281 ThirdMostExpensive(kiwi, True) ::: Option B 1282 ThirdMostExpensive(plum, True) ::: Option C 1283 ThirdMostExpensive(pear, True) ::: Option D 1284 ThirdMostExpensive(watermelon, True) ::: Option E 1285 FOL translator: 1286 I'll translate this logical deduction problem into first-order logic. 1287 **Predicates:** Rank(fruit, pos) ::: fruit has price position pos, where pos $\in \{one, two, three, four, five\}$; one = most expensive, five = cheapest Cheaper(x, y) ::: x is cheaper than y 1290 1291 **Premises:** Rank(watermelon, one) ::: Watermelons are the most expensive Rank(pears, three) ::: Pears are the third-most expensive 1293 Rank(kiwis, four) ::: Kiwis are the second-cheapest (fourth-most expensive) 1294 Cheaper(kiwis, plums) ::: Kiwis are less expensive than plums

```
1296
1297
            \forall F \forall P \forall Q ((Rank(F, P) \land Rank(F, Q)) \rightarrow (P = Q)) ::: Each fruit has exactly one rank
1298
            \forall P \forall F \forall G((Rank(F, P) \land Rank(G, P)) \rightarrow (F = G)) ::: Each rank has exactly one fruit
1299
            Conclusion:
1300
            Rank(mangoes, three) ::: Option A
1301
            Rank(kiwis, three) ::: Option B
            Rank(plums, three) ::: Option C
1302
            Rank(pears, three) ::: Option D
1303
            Rank(watermelon, three) ::: Option E
1304
1305
            SAT translator:
            I'll formulate this problem by modeling the fruit prices as positions in an ordering.
            # Declarations
            fruits = EnumSort([Mangoes, Kiwis, Plums, Pears, Watermelons])
            price_rank = IntSort([1, 2, 3, 4, 5])
1309
            pos = Function([fruits] -> [price_rank])
1310
            # Constraints
1311
            Distinct([f:fruits], pos(f)) ::: Each fruit has a unique price rank (1=cheapest, 5=most
1312
            expensive)
1313
            pos(Kiwis) < pos(Plums) ::: The kiwis are less expensive than the plums
1314
            pos(Pears) == 3 ::: The pears are the third-most expensive
1315
            pos(Kiwis) == 2 ::: The kiwis are the second-cheapest (rank 2)
1316
            pos(Watermelons) == 5 ::: The watermelons are the most expensive (rank 5)
1317
            # Options
1318
            is_valid(pos(Mangoes) == 3) ::: A) The mangoes are the third-most expensive
            is_valid(pos(Kiwis) == 3) ::: B) The kiwis are the third-most expensive
1319
            is valid(pos(Plums) == 3) ::: C) The plums are the third-most expensive
1320
            is_valid(pos(Pears) == 3) ::: D) The pears are the third-most expensive
1321
            is_valid(pos(Watermelons) == 3) ::: E) The watermelons are the third-most expensive
1322
            Round 2
1324
            LP translator:
            Looking at the other translations, I notice some inconsistencies in my approach that need
1326
            correction:
            1. Price ranking clarification: I need to be more precise about the ranking system. In my
1328
               translation, I should clarify that rank 1 is least expensive and rank 5 is most expensive (or
               vice versa) to avoid ambiguity.
1330
            2. Missing derivation rules: I should add rules that allow us to derive the complete ordering
               of fruits based on the given information.
1331
            Let me correct my translation:
1332
            Predicates:
1333
            Fruit($x, bool) ::: $x is one of the five fruits
1334
            Rank($x, $n, bool) ::: Fruit $x has price rank $n (5=most expensive, 1=least expensive)
1335
            LessExpensive($x, $y, bool) ::: Fruit $x is less expensive than fruit $y
1336
            Facts:
1337
            Fruit(mango, True)
1338
            Fruit(kiwi, True)
1339
            Fruit(plum, True)
1340
            Fruit(pear, True)
1341
            Fruit(watermelon, True)
            LessExpensive(kiwi, plum, True) ::: Kiwis are less expensive than plums
            Rank(pear, 3, True) ::: Pears are the third-most expensive
            Rank(kiwi, 2, True) ::: Kiwis are the second-cheapest (rank 2)
1344
            Rank(watermelon, 5, True) ::: Watermelons are the most expensive (rank 5)
1345
            Rules:
            LessExpensive($x, $y, True) && Rank($y, $n, True) >> Rank($x, $m, True) && $m < $n
1347
            ::: If $x is less expensive than $y, then $x's rank is lower than $y's rank
1348
```

```
1350
1351
             Rank($x, $n, True) && Rank($y, $m, True) && $n < $m >> LessExpensive($x, $y, True)
1352
             ::: If $x's rank is lower than $y's, then $x is less expensive than $y
             Ouerv:
1353
             Rank(mango, 3, True) ::: Option A
1354
             Rank(kiwi, 3, True) ::: Option B
1355
             Rank(plum, 3, True) ::: Option C
1356
             Rank(pear, 3, True) ::: Option D
1357
             Rank(watermelon, 3, True) ::: Option E
1358
             FOL translator
1359
             I'll translate this problem into first-order logic, focusing on the ordering of fruits by price.
1360
             Predicates:
             Rank(fruit, pos) ::: fruit has price position pos, where pos \in \{one, two, three, four, five\};
             one = most expensive, five = cheapest
1363
             Cheaper(x, y) ::: x is cheaper (less expensive) than y
1364
             Premises:
1365
             Rank(watermelon, one) ::: Watermelons are the most expensive
             Rank(pears, three) ::: Pears are the third-most expensive
1367
             Rank(kiwis, four) ::: Kiwis are the second-cheapest
             Cheaper(kiwis, plums) ::: Kiwis are cheaper than plums
1369
             \forall F \ \forall P \ \forall Q((Rank(F,P) \land Rank(F,Q)) \rightarrow (P=Q)) :::  One rank per fruit
             \forall P \forall F \forall G((Rank(F, P) \land Rank(G, P)) \rightarrow (F = G)) ::: One fruit per rank
1370
             Rank(mangoes, one) \vee Rank(mangoes, two) \vee Rank(mangoes, three) \vee Rank(mangoes,
1371
             four) ∨ Rank(mangoes, five) ::: Mangoes must have some rank
1372
             Rank(plums, one) ∨ Rank(plums, two) ∨ Rank(plums, three) ∨ Rank(plums, four) ∨
1373
             Rank(plums, five) ::: Plums must have some rank
1374
             I need to add the ordering relationships between ranks:
1375
             \forall X \ \forall Y \ (Rank(X, one) \land Rank(Y, two) \rightarrow Cheaper(Y, X)) ::: "higher rank \rightarrow more
1376
             expensive"
             \forall X \, \forall Y \, (Rank(X, two) \land Rank(Y, three) \rightarrow Cheaper(Y, X))
             \forall X \, \forall Y \, (Rank(X, three) \land Rank(Y, four) \rightarrow Cheaper(Y, X))
             \forall X \, \forall Y \, (Rank(X, four) \land Rank(Y, five) \rightarrow Cheaper(Y, X))
1380
             \forall X \, \forall Y \, (Cheaper(X,Y) \rightarrow \neg Cheaper(Y,X)) ::: "cheaper" is asymmetric
1381
             Conclusion:
             Rank(mangoes, three) ::: Option A
1382
             Rank(kiwis, three) ::: Option B
             Rank(plums, three) ::: Option C
1384
             Rank(pears, three) ::: Option D
1385
             Rank(watermelon, three) ::: Option E
1386
1387
             SAT translator
             I need to correct my SAT translation. Looking at the problem more carefully:
1388
             # Declarations
1389
             fruits = EnumSort([Mangoes, Kiwis, Plums, Pears, Watermelons])
1390
             price_rank = IntSort([1, 2, 3, 4, 5])
1391
             pos = Function([fruits] -> [price_rank])
1392
             # Constraints
1393
             Distinct([f:fruits], pos(f)) ::: Each fruit has a unique price rank
1394
             pos(Kiwis) < pos(Plums) ::: The kiwis are less expensive than the plums
1395
             pos(Pears) == 3 ::: The pears are the third-most expensive
             pos(Kiwis) == 2 ::: The kiwis are the second-cheapest
             pos(Watermelons) == 5 ::: The watermelons are the most expensive
             # Options
1398
             is_valid(pos(Mangoes) == 3) ::: A) The mangoes are the third-most expensive
1399
             is_valid(pos(Kiwis) == 3) ::: B) The kiwis are the third-most expensive
1400
             is valid(pos(Plums) == 3) ::: C) The plums are the third-most expensive
1401
             is_valid(pos(Pears) == 3) ::: D) The pears are the third-most expensive
1402
             is_valid(pos(Watermelons) == 3) ::: E) The watermelons are the third-most expensive
1403
```

```
1404
1405
            I made a critical error in my ranking system. In the problem, rank 1 should be cheapest and
1406
            rank 5 most expensive, not the other way around. Let me fix this.
            Round 3
1407
1408
            LP translator:
1409
            Predicates:
1410
            Fruit(x, bool) ::: x is one of the five fruits
1411
            Rank($x, $n, bool) ::: Fruit $x has price rank $n (5=most expensive, 1=least expensive)
1412
            LessExpensive($x, $y, bool) ::: Fruit $x is less expensive than fruit $y
1413
            MoreExpensive($x, $y, bool) ::: Fruit $x is more expensive than fruit $y
1414
            Facts:
1415
            Fruit(mango, True)
1416
            Fruit(kiwi, True)
1417
            Fruit(plum, True)
            Fruit(pear, True)
1418
            Fruit(watermelon, True)
1419
            LessExpensive(kiwi, plum, True) ::: Kiwis are less expensive than plums
1420
            Rank(pear, 3, True) ::: Pears are the third-most expensive
1421
            Rank(kiwi, 2, True) ::: Kiwis are the second-cheapest (rank 2)
1422
            Rank(watermelon, 5, True) ::: Watermelons are the most expensive (rank 5)
1423
            Rules:
1424
            LessExpensive($x, $y, True) >> MoreExpensive($y, $x, True) ::: If $x is less expensive
1425
            than $y, then $y is more expensive than $x
1426
            MoreExpensive($x, $y, True) >> LessExpensive($y, $x, True) ::: If $x is more expensive
1427
            than $y, then $y is less expensive than $x
1428
            Rank($x, $n, True) && Rank($y, $m, True) && $n < $m >> LessExpensive($x, $y, True)
            ::: If $x's rank is lower than $y's, then $x is less expensive than $y
1429
            LessExpensive($x, $y, True) && Rank($y, $n, True) && $n > 1 >> Rank($x, $m, True)
1430
            && $m < $n ::: If $x is less expensive than $y and $y has rank greater than 1, then $x has a
1431
            lower rank than $y
1432
             Query:
1433
            Rank(mango, 3, True) ::: Option A
1434
            Rank(kiwi, 3, True) ::: Option B
1435
            Rank(plum, 3, True) ::: Option C
1436
            Rank(pear, 3, True) ::: Option D
1437
            Rank(watermelon, 3, True) ::: Option E
1438
             FOL translator:
1439
            Predicates:
1440
            Rank(fruit, pos) ::: fruit has price position pos, where pos \in \{one, two, three, four, five\};
1441
            one = most expensive, five = cheapest
1442
            Cheaper(x, y) ::: x is cheaper (less expensive) than y
1443
            Premises:
1444
            Rank(watermelon, one) ::: Watermelons are the most expensive
1445
            Rank(pears, three) ::: Pears are the third-most expensive
1446
            Rank(kiwis, four) ::: Kiwis are the second-cheapest
            Cheaper(kiwis, plums) ::: Kiwis are cheaper than plums
1447
            \forall F \forall P \forall Q ((Rank(F, P) \land Rank(F, Q)) \rightarrow (P = Q)) ::: One rank per fruit
1448
            \forall P \forall F \forall G ((Rank(F, P) \land Rank(G, P)) \rightarrow (F = G)) ::: One fruit per rank
1449
            Rank(mangoes, one) \vee Rank(mangoes, two) \vee Rank(mangoes, three) \vee Rank(mangoes,
1450
            four) ∨ Rank(mangoes, five) ::: each still-unknown fruit occupies some rank
1451
            Rank(plums, one) \vee Rank(plums, two) \vee Rank(plums, three) \vee Rank(plums, four) \vee
1452
            Rank(plums, five) ::: each still-unknown fruit occupies some rank
1453
            \forall X \ \forall Y \ (Rank(X, one) \land Rank(Y, two) \rightarrow Cheaper(Y, X)) ::: "higher rank \rightarrow more
1454
            expensive" (10 ordered pairs)
1455
            \forall X \ \forall Y \ (Rank(X, one) \land Rank(Y, three) \rightarrow Cheaper(Y, X)) ::: \text{"higher rank} \rightarrow \text{more}
1456
            expensive" (10 ordered pairs)
1457
```

```
1458
1459
              \forall X \ \forall Y \ (Rank(X, one) \land Rank(Y, four) \rightarrow Cheaper(Y, X)) ::: \text{"higher rank} \rightarrow \text{more}
1460
              expensive" (10 ordered pairs)
              \forall X \ \forall Y \ (Rank(X, one) \land Rank(Y, five) \rightarrow Cheaper(Y, X)) ::: \text{"higher rank} \rightarrow \mathsf{more}
1461
              expensive" (10 ordered pairs)
1462
              \forall X \ \forall Y \ (Rank(X, two) \land Rank(Y, three) \rightarrow Cheaper(Y, X)) ::: "higher rank \rightarrow more
1463
              expensive" (10 ordered pairs)
1464
              \forall X \ \forall Y \ (Rank(X, two) \land Rank(Y, four) \rightarrow Cheaper(Y, X)) ::: \text{"higher rank} \rightarrow \text{more}
1465
              expensive" (10 ordered pairs)
1466
              \forall X \ \forall Y \ (Rank(X, two) \land Rank(Y, five) \rightarrow Cheaper(Y, X)) :::  "higher rank \rightarrow more
1467
              expensive" (10 ordered pairs)
1468
              \forall X \ \forall Y \ (Rank(X, three) \land Rank(Y, four) \rightarrow Cheaper(Y, X)) ::: \text{"higher rank} \rightarrow \text{more}
1469
              expensive" (10 ordered pairs)
1470
              \forall X \ \forall Y \ (Rank(X, three) \land Rank(Y, five) \rightarrow Cheaper(Y, X)) ::: \text{"higher rank} \rightarrow \text{more}
1471
              expensive" (10 ordered pairs)
              \forall X \, \forall Y \, (Rank(X, four) \land Rank(Y, five) \rightarrow Cheaper(Y, X)) ::: \text{"higher rank} \rightarrow \text{more}
1472
              expensive" (10 ordered pairs)
              \forall X \ \forall Y \ (Cheaper(X,Y) \rightarrow \neg Cheaper(Y,X)) ::: "cheaper" is asymmetric
1474
              Conclusion:
1475
              Rank(mangoes, three) ::: Option A
1476
              Rank(kiwis, three) ::: Option B
1477
              Rank(plums, three) ::: Option C
1478
              Rank(pears, three) ::: Option D
1479
              Rank(watermelon, three) ::: Option E
1480
              SAT translator:
1481
              # Declarations
1482
              fruits = EnumSort([Mangoes, Kiwis, Plums, Pears, Watermelons])
1483
              price_rank = IntSort([1, 2, 3, 4, 5])
1484
              pos = Function([fruits] -> [price_rank])
1485
              # Constraints
1486
              Distinct([f:fruits], pos(f)) ::: Each fruit has a unique price rank (1=cheapest, 5=most
1487
              expensive)
1488
              pos(Kiwis) < pos(Plums) ::: The kiwis are less expensive than the plums
              pos(Pears) == 3 ::: The pears are the third-most expensive
1489
              pos(Kiwis) == 2 ::: The kiwis are the second-cheapest
1490
              pos(Watermelons) == 5 ::: The watermelons are the most expensive
1491
              # Options
1492
              is_valid(pos(Mangoes) == 3) ::: A) The mangoes are the third-most expensive
1493
              is valid(pos(Kiwis) == 3) ::: B) The kiwis are the third-most expensive
1494
              is valid(pos(Plums) == 3) ::: C) The plums are the third-most expensive
1495
              is valid(pos(Pears) == 3) ::: D) The pears are the third-most expensive
1496
              is_valid(pos(Watermelons) == 3) ::: E) The watermelons are the third-most expensive
1497
```

G.3 REASONING DEBATE

1498

1499

1500 1501 1502

1503

1506

1507

1509

1510

1511

Case Study from ProofWriter Dataset

Context. The bear is round. The bear likes the cow. The bear visits the lion. The cow is big. The cow is rough. The lion needs the bear. The mouse likes the lion. If someone visits the lion then they are round. If someone needs the bear then they need the lion. If someone needs the lion then they visit the lion. If someone visits the lion and the lion is round then they like the lion. If someone is rough then they need the lion. If someone likes the lion and they need the lion then the lion visits the mouse. If someone is rough then they visit the lion. If someone likes the lion and the lion is big then they are kind.

Question. Based on the above information, is the following statement true, false, or unknown? *The lion visits the lion*.

Options. A) True B) False C) Unknown

After reviewing the reasoning and deductions, I conclude that the statement "The lion visits the lion" is **True**. From rule 2 (Needs (\$x, Bear) >> Needs (\$x, Lion)) and the given fact Needs (Lion, Bear), we deduce Needs (Lion, Lion). Applying rule 3 (Needs (\$x, Lion) >> Visits (\$x, Lion)), this implies Visits (Lion, Lion). Therefore, the statement holds true. **Key Points:**

- 1. Needs (Lion, Bear) is a given fact.
- 2. Rule 2 extends this to Needs (Lion, Lion).
- 3. Rule 3 further concludes Visits (Lion, Lion).

Thus, my answer remains **True**, and my confidence in this answer is **1.0**, based on the rigorous application of logical rules. Confidence: 1.0

Based on my analysis using First-Order Logic (FOL), I stand by my initial conclusion of **Unknown**. Here's why:

- 1. Logic Consistency: Both attempting to prove Visits (Lion, Lion) and its negation resulted in no contradiction from the given predicates. This means neither can be conclusively established as true or false.
- 2. Expressiveness of FOL: FOL captured all the nuanced dependencies and conditions in the problem (such as implications involving Needs, Visits, and Likes). However, the information provided still leaves the truth of Visits (Lion, Lion) indetermi-
- 3. Evidence Limitations: Although key facts like Rough (Cow) -> Visits (Cow, Lion) and conditions involving Needs and Likes were derived, none directly establish the behavior of Visits (Lion, Lion), nor does denying Visits (Lion, Lion) provide a definitive contradiction.

In summary, this is a clear example of the limitations of the provided context, not of FOL itself. The problem is underdefined based on the assumptions and given predicates, leading to the result of **Unknown**.

To analyze whether "The lion visits the lion" is true, false, or unknown, I model the logical context using Z3. Here's a formal representation of the problem with the relevant predicates:

- 1. Individuals: {bear, cow, lion, mouse}
 - Likes(x, y)
 - Visits(x, y)
 - Needs (x, y)
 - Kind(x)

1550

1551

1552

1553

1554

1555

1556

1558

1559

1560

1561

1563

1564

1565

Facts encoded: Round (bear), Likes (bear, cow), Visits (bear, lion), Big(cow), Rough(cow), Needs(lion, bear), Likes(mouse, lion)

Rules encoded: (Visits(x, lion) \rightarrow Round(x)), (Needs(x, bear) \rightarrow Needs(x, lion)), $(Needs(x, lion) \rightarrow Visits(x, lion)),$ (Visits(x, lion) \land Round(lion) \rightarrow Likes(x, lion)), (Rough (x) \rightarrow Needs(x, lion)), (Likes $(x, lion) \land Needs(x, lion) \rightarrow$ Visits(lion, mouse)), (Rough(x) \rightarrow Visits(x, lion)), (Likes(x, lion) \land Big(lion) \rightarrow Kind(x))

To verify "The lion visits the lion", I encode: ?Visits(lion, lion).

Solving this using Z3 for satisfiability under the constraints:

- No explicit fact implies Visits (lion, lion).
- Rules do not derive Visits (lion, lion).
- Unknown/indeterminate due to lack of evidence.

1566 1567 Thus, the statement remains UNKNOWN. 1568 Confidence: 1.0 1569 Plan-and-Solve supporter: 1570 The Plan-and-Solve approach confidently concludes that "The lion visits the lion" is true. 1571 This conclusion hinges on two key deductions: 1572 1. The lion needs the bear (Fact 6). 1573 Using **Rule 2** ("If someone needs the bear, then they need the lion"), we deduce that the lion also needs the lion. 1574 2. The lion needs the lion. 1575 Applying **Rule 3** ("If someone needs the lion, then they visit the lion"), it follows directly 1576 that the lion visits the lion. This reasoning is logically consistent and follows the facts and rules provided. My confi-1578 dence in this conclusion remains high unless another agent proves an inconsistency or over-1579 looked factor. 1580 Confidence: 1.0 1581 COT Solver supporter: 1582 I conclude that the answer is C (Unknown) because there is no information in the context 1583 that establishes whether the lion visits itself. None of the given rules or statements necessar-1584 ily imply that the lion visits itself. While others visiting the lion or liking/needing the lion 1585 are explained, no direct condition or chain of reasoning leads to the conclusion about the 1586 lion visiting itself. Therefore, this remains undetermined based on the given information. 1587 Confidence: 0.9 1588 Round 2 LP supporter: 1590 After reviewing FOL supporter's argument, I must directly address the critical error in their 1591 1592 1. FOL state that Needs (Lion, Lion) "could" imply Visits (Lion, Lion) via Rule 3. There is no "could" here - Rule 3 is a universal implication: Needs (\$x, 1594 Lion) \rightarrow Visits (\$x, Lion). This is deterministic, not probabilistic. 2. FOL's claim of "no explicit or derived fact ensures Visits (Lion, Lion) conclusively" ignores that we do derive this fact through the logical chain: Needs (Lion, Bear) \rightarrow Needs (Lion, Lion) \rightarrow Visits (Lion, Lion). 1598 The statement "The lion visits the lion" is conclusively **True**. Confidence: 1.0 1599 *FOL supporter:* I maintain that the answer is **C** (**Unknown**). My reasoning is as follows: 1602 1. Logical Consistency: Attempts to both prove Visits (Lion, Lion) and its nega-1603 tion resulted in no contradictions, which indicates that the problem is underdefined with 1604 the current rules and facts. 2. Dependence: While Needs (Lion, Bear) implies Needs (Lion, Lion) (via 1605 Rule 2), and Needs (Lion, Lion) could imply Visits (Lion, Lion) (via Rule 1606 3), this chain of reasoning isn't fully supported because no explicit or derived fact ensures Visits (Lion, Lion) conclusively. 3. No Contradiction Found: The problem's structure does not yield a definitive resolution 1609 to the truth of the statement. 1610 Thus, the status remains Unknown based on the provided information and logical frame-1611 work. 1612 Confidence: 1.0 1613 SAT supporter: 1614 After reviewing the arguments provided by others, I align with Plan-and-Solve 1615 **Supporter**'s reasoning. The chain of deductions from Rule 2 (Needs (Lion, 1616 Bear) \rightarrow Needs(Lion, Lion)) and then Rule 3 (Needs(Lion, Lion) \rightarrow 1617 Visits (Lion, Lion)) unequivocally establishes that the lion visits itself. This rea-

soning directly satisfies the logical derivation that I would also construct using Z3.

1618

1621 I now change my answer from **Unknown** to **True**, as the required logical steps are valid and 1622 complete. Confidence: 1.0 1623 1624 Plan-and-Solve supporter: 1625 I stand by my original answer that "The lion visits the lion" is True, supported by valid 1626 logical deductions from the provided rules and facts: 1627 1. **Fact 6:** "The lion needs the bear." Applying **Rule 2** ("If someone needs the bear, then they need the lion"), we deduce that 1628 the lion needs the lion. 1629 2. **Rule 3:** "If someone needs the lion, then they visit the lion." 1630 Since the lion needs the lion, it follows that the lion visits the lion. This argument addresses the claim that we lack sufficient information. The facts and rules provided suffice to logically establish the truth of the statement. 1633 Confidence: 1.0 1634 COT Solver supporter: 1635 After reviewing the ongoing arguments, I acknowledge that the argument provided by **Plan**and-Solve supporter is valid: "The lion needs the bear," and by applying Rule 2 and Rule 1637 3, it follows logically that "The lion visits the lion." This was an oversight in my earlier 1638 evaluation. 1639 Thus, I change my answer to A (True) as the reasoning is clear and follows directly from 1640 the rules and the given facts. 1641 Confidence: 1.0 1642 Round 3 1643 LP supporter: 1644 I remain confident that the statement "The lion visits the lion" is True. The reasoning 1645 provided earlier, and supported by others, systematically deduces this conclusion: 1646 1. From Fact 6, Needs (Lion, Bear) is given. 1647 2. Using Rule 2, this implies Needs (Lion, Lion). 1648 3. By Rule 3, Needs (Lion, Lion) directly leads to Visits (Lion, Lion). 1649 This chain of reasoning is airtight and rigorously validated. Therefore, the lion indeed visits 1650 itself. 1651 Confidence: 1.0 1652 FOL supporter: After carefully reassessing the arguments, I now concur with Plan-and-Solve Supporter 1654 and others who assert that "The lion visits the lion" is True. The reasoning that: 1655 1. Fact 6: "The lion needs the bear" leads to "The lion needs the lion" via Rule 2. 1656 2. The lion needing the lion implies "The lion visits the lion" through Rule 3. This chain of reasoning is both valid and well-supported by the premises. I change my initial 1657 conclusion from **Unknown** to **True**. 1658 Confidence: 1.0 1659 SAT supporter: 1661 I maintain my revised conclusion with confidence: "The lion visits the lion" is True. Fol-1662 lowing the reasoning provided: 1663 1. Fact: Needs (Lion, Bear) is given. 2. Rule 2: Needs (Lion, Bear) \rightarrow Needs (Lion, Lion), so Needs (Lion, 1664 Lion) is deduced. 1665 3. Rule 3: Needs (Lion, Lion) \rightarrow Visits (Lion, Lion), so Visits (Lion, Lion) is established. This reasoning aligns perfectly with the constraints I would also encode in Z3. Therefore, 1668 the statement is logically valid and proven to be **True**. 1669 Confidence: 1.0 *Plan-and-Solve supporter:* 1671 Thank you. My position remains that the statement "The lion visits the lion" is True. The 1672 reasoning follows directly from the structured Plan-and-Solve approach: 1673

1674 1675 1. **Fact 6:** The lion needs the bear. 1676 2. **Rule 2:** If someone needs the bear, they need the lion. From Fact 6, this gives **Needs(Lion, Lion)**. 1677 3. **Rule 3:** If someone needs the lion, they visit the lion. 1678 Using Needs(Lion, Lion), this directly implies Visits(Lion, Lion). 1679 Thus, the statement is logically derived as **True** based on a clear application of the given rules. No conflicting interpretations invalidate this deduction. 1681 Confidence: 1.0 1682 COT Solver supporter: 1683 After reviewing the provided arguments, I now agree with the reasoning by **Plan-and-Solve** 1684 **supporter** and others who concluded that the statement is **True**. The chain of logic using 1685 Fact 6, Rule 2, and Rule 3 is sound, showing that "The lion visits the lion" must hold. I initially missed this chain of reasoning. 1687 **Revised Answer: A (True)** 1688 Confidence: 1.0 1689 Round 4 (Final answers) *LP supporter:* <answer>A</answer> FOL supporter: <answer>A</answer> 1692 SAT supporter: <answer>A</answer> 1693 *Plan-and-Solve supporter:* <answer>A</answer> COT Solver supporter: <answer>A</answer> 1695 Case Study from Logical Deduction Dataset 1698 1699 **Context.** The following paragraphs each describe a set of five objects arranged in a fixed 1700 order. The statements are logically consistent within each paragraph. 1701 In an antique car show, there are five vehicles: a minivan, a limousine, a sedan, a tractor, and 1702 a hatchback. The tractor is the second-newest. The limousine is newer than the hatchback. The limousine is older than the sedan. The minivan is newer than the sedan. 1703 1704 **Question.** Which of the following is true? 1705 B) The limousine is the second-oldest. **Options.** A) The minimal is the second-oldest. 1706 C) The sedan is the second-oldest. D) The tractor is the second-oldest. E) The hatch-1707 back is the second-oldest. 1708 Round 1 1709 LP supporter: 1710 Using logic programming (LP), I argue that the correct option is A) the minivan is the 1711 **second-oldest**. My reasoning is as follows. First, the tractor is explicitly stated to be the 1712 second-newest, so exactly one vehicle is newer than the tractor. Next, from the pairwise 1713 constraints we have S > L > H (i.e., the sedan is newer than the limousine, which is 1714 newer than the hatchback), and we also have M > S (the minivan is newer than the sedan). 1715 Intuitively, this pushes the trio (H, L, S) toward the older side in that order while keeping 1716 the minivan ahead of the sedan; with the tractor fixed near the new end (second-newest), the remaining older-side slot that best satisfies all relations is occupied by the **minivan**, making 1717 it the second-oldest. Therefore, I select A. Confidence: 0.8 1718 1719 *FOL supporter:* Let's formalize the problem using First-Order Logic (FOL) for precision: 1721 **Definitions:** 1722 • Let Newer (x, y) mean "x is newer than y." 1723 • Objects: Minivan (M), Limousine (L), Sedan (S), Tractor (T), Hatchback (H) 1724 **Given information:** 1725

1. Newer (S, L) (The sedan is newer than the limousine.)

2. Newer (L, H) (The limousine is newer than the hatchback.)

1726

3. T is the second-newest.

1728 1729

1780

1781

```
1730
            4. Newer (M, S) (The minimum is newer than the sedan.)
            Logical Analysis: We need to determine the ordering from newest to oldest, consistent with
1731
            the above statements.
1732
            Step 1: Identify the newest vehicle. From T being second-newest, exactly one vehicle is
1733
            newer than T. Since M > S and S > L > H, neither S, L, nor H can be newest. Thus,
1734
            the minivan is the newest.
1735
            Step 2: Order the limousine, sedan, and hatchback. From statement 1: Newer (S, L)
1736
            \Rightarrow S > L. From statement 2: Newer (L, H) \Rightarrow L > H. Thus, the order for these three is
1737
            S > L > H.
1738
            Step 3: Position the tractor. From statement 3, the tractor must be the second-newest
1739
            vehicle. Thus, it follows the minimum: M > T.
            Final order (newest \rightarrow oldest): M > T > S > L > H.
1740
            Evaluation of "second-oldest": The second-oldest is the vehicle fourth in the order. From
1741
            the final order, this is Limousine.
1742
            Answer: B) The limousine is the second-oldest.
1743
            Confidence: 1.0
1744
1745
            SAT supporter:
            I will represent the problem using Z3 to rigorously verify the correct ordering.
1746
            The five vehicles (minivan, limousine, sedan, tractor, hatchback) are given, and we define
1747
            their ages using integer variables. Let:
1748
             • minivan = x1
1749
              limousine = x2
1750
             • sedan = x3
1751
             • tractor = x4
1752
             • hatchback = x5
1753
            We impose the following constraints from the context:
1754
            1. x4 = second newest implies there are three vehicles older than tractor and one
1755
            2. x2 > x5: The limousine is newer than the hatchback.
1756
            3. x2 < x3: The limousine is older than the sedan.
1757
            4. x1 > x3: The minimum is newer than the sedan.
1758
            To solve this:
1759
             • Define constraints for ages to ensure strict ordering, i.e., no two vehicles have the same
1760
               Enumerate all possibilities for which vehicle can be "second-oldest" and validate each
1762
               against the constraints.
1763
            Upon solving: The second-oldest is indeed the limousine, satisfying all conditions.
1764
            Agreement/Disagreement with Reasoning: No counterarguments provided yet.
1765
            Conclusion: The correct answer is B) The limousine is the second-oldest.
1766
            Confidence: 1.0
1767
            Plan-and-Solve supporter:
1768
            I have reviewed the arguments provided. The Plan-and-Solve method has clearly analyzed
1769
            the four clues systematically, arranged the vehicles in order of their relative ages, and iden-
1770
            tified the second-oldest vehicle as the limousine. All steps were justified and consistent with
1771
            the given information. No challenges or alternative reasoning have been introduced that con-
1772
            tradict this deduction. Since the solution is logically sound, supported by verification of all
            clues, and accurately identifies the second-oldest vehicle, I maintain that the answer is B)
1773
            The limousine is the second-oldest.
1774
            Confidence: 1.0
1775
1776
            COT Solver supporter:
1777
            I stand by my initial answer of C (the sedan is the second-oldest). Here's my reasoning in
            summary:
1778
            1. The tractor is the second-newest, so it cannot be the second-oldest.
1779
```

2. The **limousine** is older than the sedan and newer than the hatchback, making it older than

at least two vehicles, but not the second-oldest.

- 3. The **minivan** is newer than the sedan, so it cannot be the second-oldest.
- 4. The **hatchback** is the oldest or close to it since the limousine is newer than it, thus eliminating it as the second-oldest.
 - 5. The **sedan** is older than the minivan but not the oldest, making it the second-oldest.

Hence, the sedan fits the criteria for being the second-oldest. My confidence remains high as the logic follows directly from the constraints provided.

Confidence: 0.9

1788 1789 1790

1791

1792

1793

1795

1796

1797

1799

1801

1803

1805

1806

1807

1808

1809 1810

1811

1812

1813

1814

1815

1816

1817

1818 1819

1820

1821

1824

1825

1826

1827

1830

1831

1786

1787

Round 2

LP supporter:

After reviewing the arguments in this collaborative debate, I realize a critical discrepancy in my initial reasoning. The process of logical deduction and the formal order articulated by the **FOL supporter**, **SAT supporter**, and **Plan-and-Solve supporter** has established the following sequence of ages:

Newest \rightarrow **Oldest:** M > T > S > L > H.

This makes the **limousine** the second-oldest vehicle. My earlier labeling of multiple second-oldest vehicles contradicted the constraints and was erroneous. I now agree that the correct option is **B) The limousine is the second-oldest**.

Confidence: 1.0

FOL supporter:

After carefully reviewing the arguments from all agents and reapplying First-Order Logic, I disagree with **COT_Solver_supporter**'s deduction of the sedan as the second-oldest. According to the formal logical constraints derived, the correct ordering from newest to oldest is: M > T > S > L > H. Thus, the second-oldest vehicle is clearly **Limousine**, not the sedan.

This conclusion is consistent with all constraints: the tractor is second-newest, and the ordering ensures all relative relationships hold true.

Final Answer: B) The limousine is the second-oldest.

Confidence: 1.0

SAT supporter:

I disagree with COT_Solver_supporter claiming the **sedan** is the second-oldest. Z3's SMT solving clearly resolves all constraints without ambiguity and confirms the order as M > T > S > L > H. Hence, the **limousine** is fourth in the sequence, or, equivalently, the "second-oldest." The claim that the sedan is the second-oldest is incorrect because it is explicitly positioned third in the correct sequence (newer than limousine and hatchback but older than tractor and minivan).

Thus, the correct answer remains **B)** The limousine is the second-oldest.

Confidence: 1.0

Plan-and-Solve supporter:

After reviewing the discussion and the constraints again, my initial reasoning and solution remain valid. The Plan-and-Solve approach systematically deduced that the order from **oldest** to **newest** is hatchback > limousine > sedan > tractor > minivan. Therefore, the second-oldest vehicle is the **limousine**, consistent with all clues. The argument for the sedan being the second-oldest is incorrect because it conflicts with the clue stating that the limousine is older than the sedan. This directly eliminates the sedan as a contender for second-oldest.

Final Answer: B) The limousine is the second-oldest.

Confidence: 1.0

COT Solver supporter:

I now realize that I made an error in my earlier reasoning. Upon reviewing the arguments and steps from the other agents, it is evident that the proper ordering is M>T>S>L>H (newest to oldest). From this, the **limousine** is the second-oldest vehicle, not the sedan. Thus, I agree with the corrected logical deductions provided by the other agents. The correct answer is **B**) The limousine is the second-oldest.

Confidence: 1.0

1834 1835

Round 3

1884 1885

1836 1837 LP supporter: 1838 Thank you. After reviewing all the reasoning provided by the agents and scrutinizing my initial position, I fully align with the consensus reached by the FOL supporter, SAT 1839 **supporter**, and others. The logical deductions clearly establish that the correct order is 1840 M > T > S > L > H, making the **limousine** the second-oldest vehicle. My initial misla-1841 beling has been corrected. 1842 Final Answer: B) The limousine is the second-oldest. 1843 Confidence: 1.0 1844 FOL supporter: 1845 Given the thorough reasoning and agreement across multiple approaches, including clear 1846 formal deductions, the correct ordering is M > T > S > L > H (newest to oldest). 1847 Therefore, the second-oldest is the **limousine**. 1848 Final Answer: B) The limousine is the second-oldest. 1849 Confidence: 1.0 1850 SAT supporter: 1851 I maintain that the correct answer is **B) The limousine is the second-oldest**. Using Z3 to represent the constraints ensures a complete and rigorous solution, confirming the order: M > T > S > L > H. No contradictions remain in this conclusion, and all logical 1854 relationships are satisfied. 1855 Confidence: 1.0 *Plan-and-Solve supporter:* 1857 The correct answer remains **B)** The limousine is the second-oldest. The Plan-and-Solve approach has systematically deduced the correct ordering and verified it against all constraints. 1859 Other agents have corroborated this conclusion, and alternative claims have been refuted. 1860 The reasoning is consistent, logical, and validated. 1861 Confidence: 1.0 1862 COT Solver supporter: 1863 After reviewing and assessing the arguments presented by other agents, I acknowledge my 1864 earlier error. The systematic deductions using constraints demonstrate that the correct order 1865 is M > T > S > L > H. From this order, the **limousine** is undoubtedly the second-oldest 1866 vehicle. Final Answer: B) The limousine is the second-oldest. 1867 Confidence: 1.0 1868 Round 4 (Final answers) LP supporter: <answer>B</answer> 1870 FOL supporter: <answer>B</answer> 1871 SAT supporter: <answer>B</answer> 1872 Plan-and-Solve supporter: <answer>B</answer> 1873 COT Solver supporter: <answer>B</answer> 1874