# LEARNING TO REASON VIA MIXTURE-OF-THOUGHT FOR LOGICAL REASONING

**Anonymous authors**Paper under double-blind review

000

001

002003004

006

008 009

010 011

012

013

014

015

016

017

018

019

021

024

025

026

027

028

029

031

033 034

035

037

040

041

042

043 044

046 047

048

049

050

051

052

#### **ABSTRACT**

Human beings naturally utilize multiple reasoning modalities to learn and solve logical problems, i.e., different representational formats such as natural language, code, and symbolic logic. In contrast, most existing LLM-based approaches operate with a single reasoning modality during training, typically natural language. Although some methods explored modality selection or augmentation at inference time, the training process remains modality-blind, limiting synergy among modalities. To fill in this gap, we propose Mixture-of-Thought (MoT), a framework that enables LLMs to reason across three complementary modalities: natural language, code, and a newly introduced symbolic modality, truth-table, which systematically enumerates logical cases and partially mitigates key failure modes in natural language reasoning. MoT adopts a two-phase design: (1) self-evolving MoT training, which jointly learns from filtered, self-generated rationales across modalities; and (2) **MoT inference**, which fully leverages the synergy of three modalities to produce better predictions. Experiments on logical reasoning benchmarks including FOLIO and ProofWriter demonstrate that our MoT framework consistently and significantly outperforms strong LLM baselines with single-modality chain-of-thought approaches, achieving up to +11.7pp average accuracy gain. Further analyses show that our MoT framework benefits both training and inference stages; that it is particularly effective on harder logical reasoning problems; and that different modalities contribute complementary strengths, with truth-table reasoning helping to overcome key bottlenecks in natural language inference.

# 1 Introduction

Large language models (LLMs) have demonstrated remarkable progress in logical reasoning tasks, especially propelled by methods like Chain-of-Thought (CoT) prompting (Wei et al., 2022). However, these CoT approaches predominantly rely on single reasoning modality, *i.e.*, natural language, even when employing ensemble methods (Li et al., 2023; Wang et al., 2025; 2022; Brown et al., 2024; Snell et al., 2025; Liang et al., 2023). Here we refer to a *modality* as a distinct thought paradigm<sup>1</sup> (*e.g.* natural language, symbolic, or code), which differs in representation and inference process. On the other hand, neuro-symbolic methods (Pan et al., 2023; Olausson et al., 2023; Ryu et al., 2025) utilize LLMs as translators and delegate reasoning to external symbolic solvers. Recent work combines CoT with symbolic reasoning via either selecting a single modality per instance (Xiong et al., 2024) or augmenting one modality with the other—while keeping reasoning confined to symbolic (Xu et al., 2024a) or natural language (Liu et al., 2024). **These methods combine modalities only during inference and ignore the synergy of different modalities during training**, thus failing to fully exploit the complementary strengths of different reasoning modalities.

This limitation contrasts sharply with human cognition: Humans naturally employ multiple reasoning modalities, flexibly switching among natural language explanations, code-based procedural thinking, and formal symbolic manipulation, both when learning complex logical skills and when solving novel problems (Newell et al., 1972; Gentner, 1983; Larkin & Simon, 1987; Goldin, 1998). This cognitive versatility, the ability to represent and process information in diverse formats, is crucial for robust reasoning. Current LLMs, largely confined to single-modality training and inference, lack this

<sup>&</sup>lt;sup>1</sup>We use the terms thought paradigm and reasoning modality interchangeably.

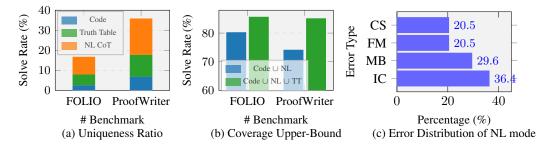


Figure 1: (a) Qwen-2.5-7B-Instruct solves  $\simeq$ 20% of FOLIO and  $\simeq$ 35% of ProofWriter exclusively per paradigm. (b) Code+NL+truth-table yields higher upper-bound coverage than code+NL alone (Xiong et al., 2024). (c) In NL modes, invalid-converse (IC) and missing-branch (MB) errors comprise  $\simeq$ 66% of failures (CS: commonsense injection; FM: factual misquote). Percentages sum to more than 100% because some cases exhibit multiple error types. We provide illustrative examples in Appendix H.1

flexibility. It raises a critical question: Can LLMs achieve more robust and versatile logical reasoning by explicitly learning to operate across multiple complementary reasoning modalities?

Addressing this question requires tackling two challenges: 1) It is still unclear which reasoning modalities should be included; the selected modalities must be complementary to make joint learning worthwhile. 2) Teaching an LLM with multiple modalities is non-trivial, as large aligned reasoning trajectories are scarce. Without those high-quality reasoning trajectories, we cannot teach LLMs to reason via multiple modalities. Our investigation reveals crucial insights for modality selection.

- Natural language bottleneck. Figure 1 (c) shows that nearly two thirds of CoT errors arise from *missing branches* and *invalid converse*, *i.e.*, poor exhaustive enumeration and complex inference (See examples in Appendix H.1). Truth-table reasoning, which systematically lists all possibilities, naturally complements this weakness; therefore, we incorporate a symbolic truth-table paradigm.
- Code–NL complementarity. Inspired by HybridMind (Yue et al., 2024; Xiong et al., 2024), where
  they show preliminary results that a code paradigm could complement NL reasoning, we also
  incorporate code as one reasoning modality into our framework.
- **Paradigm overlap & uniqueness.** Figure 1 (a-b) shows that 35.8% of ProofWriter items and 16.7% of FOLIO items are solved by *exactly one* paradigm, while the union of three reasoning modalities covers up to 85% of all instances—evidence that combining NL, code, and truth-table reasoning is necessary, outperforming the simple combination of code and NL (Xiong et al., 2024).

Building on these insights, we propose Mixture-of-Thought (MoT), a human-inspired framework that enables LLMs to reason via three complementary reasoning modalities: NL, code, and symbolic; an example is shown in table 1 to illustrate each modality. Notably, we introduce a new truth-table-based symbolic reasoning where LLMs ground propositional variables, construct a partial truth table by pruning assignments that violate any premise, and infer the final answer by checking the truth table. Our MoT consists of two parts. One part is training: we propose a self-evolving MoT training algorithm, which improves the model's reasoning ability in each modality through joint iterative optimization (Figure 2 (a)). Another part is inference, where we generate responses under each modality and leverage a voting mechanism to produce the final answer (Figure 2 (b)). This simple strategy allows the model to combine diverse perspectives and make robust predictions.

Empirically, we show that across three base models—Gemma-2-2B-IT, Gemma-2-9B-IT, and Qwen-2.5-7B-Instruct—our MoT consistently surpasses the CoT baseline on ProofWriter (Tafjord et al., 2021) and FOLIO (Han et al., 2024), with an average accuracy gain of up to +11.7pp. Notably, our 9B-parameter MoT matches the results of GPT-4 + Logic-LM on FOLIO. Additional analyses show that 1) MoT training outperforms single-thought training; 2) Mixture-of-Thought sampling yields a higher oracle upper bound than single-thought sampling under the same inference budget 3) The gains grow with problem difficulty: MoT helps most on depth-5 and other harder problems; and 4) A fine-grained error study reveals a key natural-language bottleneck, *i.e.*, missing branches and frequent invalid converse errors, while the truth-table paradigm help resolve some cases of exactly these types.

# 2 THE MIXTURE-OF-THOUGHT FRAMEWORK

In this section, we introduce 1) three complementary reasoning modalities for logical reasoning (Sec. 2.1); 2) our self-evolving training framework that jointly improves these reasoning modalities (Sec.

109

110

111

112

113 114

115

116

117

118

119

120

121

122

123

124

125

126

127

128

129 130 131

132

133

134 135

136

137

138

139

140 141

142

143

144

145

146

147

148

149

150

151 152

153

154

155

156

157

158

159

160

161

Table 1: Illustration of the three complementary reasoning modalities, *i.e.*, natural-language CoT, code-based reasoning, and truth-table reasoning. We provide the corresponding outputs of LLMs in section G.

**Premise:** Peter Parker is either a superhero or a civilian. The Hulk wakes up when he is angry. If he wakes up, he will break a bridge. If a destroyer breaks a bridge, Peter is not a civilian. Peter wears a uniform when he is a superhero. Thor is a god. A god is not a destroyer. Thor will break a bridge when he is happy. If Thor is happy, the Hulk is angry.

Question: If Thor is happy, does Peter Parker wear a uniform? Options: (A) True (B) False (C) Uncertain

```
Natural Language Reasoning
                                                    Code-Based Reasoning (Abstract)
                                                                                                            Truth Table Reasoning
                                                                                                             <truth table>
<nl cot>
                                                     <code>
                                                                                                            Let: T = \text{Thor happy}, H = \text{Hulk angry}, A = \text{wakes} up, B = \text{bridge breaks}, C = \text{Peter is civilian}, S = \text{the happy}
Step 1: given premises "If Thor is happy, the Hulk
                                                    class Hulk:
                                                          def __init__(self, angry): ...
def wakes_up(self): ...
is angry." and "The Hulk wakes up when he is an-
gry.", we can know "If Thor is happy, then hulk
                                                                                                            superhero, U = wears uniform.
                                                          def breaks_bridge(self): ...
                                                                                                            Premises:
                                                    class Thor:
Step 5: given premise "Peter Parker is either a
                                                          def __init__(self, happy): ...
def breaks_bridge(self): ...
                                                                                                             S \lor C, H \to A, A \to B, T \to H, T \to B,
superhero or a civilian." and derivation "If Thor
                                                                                                            B \to \neg C, S \to U
is happy, then Peter Parker is not a civilian", we
                                                     class PeterParker:
                                                          can know "If Thor is happy, then Peter Parker is a
superhero."
                                                                                                            T \Rightarrow H \Rightarrow A \Rightarrow B

B \Rightarrow \neg C \Rightarrow C = \text{False}
Step 6: given premise "Peter Parker wears a uni-
                                                          def wears_uniform(self): ...
                                                    \frac{1}{\text{def apply\_premises(thor, hulk, peter)}} \colon S \vee C \Rightarrow S = \text{True} \Rightarrow U = \text{True}
form when he is a superhero." and derivation "If
Thor is happy, then Peter Parker is a superhero.", we can know "If Thor is happy, then Peter Parker
                                                    {\tt def \ run\_inference(thor, \ hulk, \ peter):} \ \ {\tt Truth \ Table:}
wears a uniform"
<end_of_nl_cot>
                                                    def check_conclusion(...): ...
                                                                                                              T \quad H \quad A \quad B \quad C
                                                                                                                                           S
                                                                                                                                                 U
<answer:
                                                     thor = Thor(happy=True)
Answer: (A)
                                                    hulk = Hulk(angry=False)
                                                                                                             True True True False True True
<end_of_answer>
                                                    peter = PeterParker(...)
                                                                                                            <end of truth table>
                                                     result = check_conclusion(...)
                                                                                                             <answer>Answer: (A)<end of answer>
                                                     <end of code>
                                                     <answer>Answer: (A)<end of answer>
```

2.2); and 3) our mixture-of-thought inference strategy that combines diverse but complementary reasoning paths to derive robust final predictions (Sec. 2.3).

## 2.1 Human-Inspired Complementary Reasoning Modalities

Drawing inspiration from human cognition and error analysis in Figure 1, we argue that no single reasoning modality suffices for all logical challenges. Therefore, we equip a single model with three complementary modalities: natural language CoT, code CoT, and truth table CoT. Specifically, because natural-language CoT often misses branches or makes invalid-converse errors, we design a truth-table approach that explicitly enumerates truth assignments and thus complements these weaknesses. Table 1 illustrates how the three modalities solve a representative problem.

- Natural Language CoT: The model explains its reasoning in plain natural language, decomposing the problem into step-by-step justifications. This format is flexible and interpretable.
- Code CoT: The model first transforms a logical problem to a PYTHON code and then derives the answer based on the PYTHON code. We do not execute the code; instead, we treat it as a way to describe logic in a structured form.
- **Truth Table CoT:** The model first explicitly generates a truth table by defining predicates based on the premises and conclusion, then enumerating possible truth assignments, and finally checking which ones satisfy the conclusion.

These complementary modalities are jointly exploited in our self-evolving training (Sec. 2.2) and majority-vote inference (Sec. 2.3). We now detail the design of the Truth Table CoT approach.

**Truth-Table CoT: Challenges and Design.** Two main challenges arise when enabling LLMs to reason with truth tables: 1) Exponential blow-up: the number of rows grows exponentially with the propositional variables, easily exceeding the context window and compute budget; 2) First-order grounding: practical tasks are given in first-order logic; one must ground variables, select a finite predicate set, and still ensure that the resulting (partial) truth table remains tractable. To address these challenges, we propose a two-step strategy: (i) grounding, which instantiates first-order formulas into a finite set of propositional predicates (Clarke et al., 2001; Wittocx et al., 2010), and (ii) reason to prune, which eliminates rows that violate any premise through reasoning via LLMs, keeping partial truth table (see Table 1 and Appendix G.3). Finally, the LLMs derive the final output with the following rule: True if every surviving assignment satisfies the conclusion, False if none do, and Uncertain otherwise. Moreover, we assign modality-specific tags (e.g., <code> ...

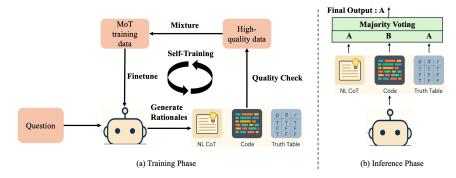


Figure 2: Illustration of our MoT Framework. (a) **Training phase** with three key steps: 1) **Rationale Generation** where given an initial seed dataset, LLM generates rationales across reasoning modalities (NL, Code, and Truth Table); 2) **Quality Checking and Merging** where generated rationales are checked for correctness and format consistency, then merged into high-quality MoT training data; and 3) **Finetuning** where the model is trained using the MoT data. These steps iteratively repeats, forming a self-evolving training cycle. (b) **Inference phase:** the trained model generates outputs for each reasoning modality and applies majority voting to yield the final prediction (e.g., A).

<end\_of\_code>) to explicitly indicate the format during training and inference. The prompts are
detailed in Appendix F.

#### 2.2 Self-Evolving Mixture-of-Thought Training

Explicitly learning to reason across multiple complementary modalities, such as natural language, code, and symbolic truth tables, is non-trivial. A key challenge lies in the lack of annotated reasoning trajectories for each modality, especially for our newly introduced truth-table approach. Collecting labeled CoT traces for all of these modalities is also costly and often infeasible. To address this, we propose a self-evolving MoT training approach, which enables the model to operate across multiple complementary reasoning modalities by iteratively learning from its own generated reasoning traces.

Given the policy M, our goal is to maximize the following objective across the problems x and modalities  $\mathcal{T} \in \{\text{NL}, \text{Code}, \text{TruthTable}\}$ :

$$\mathbb{E}_{(x_i, y_i) \sim \mathcal{D}, t \sim \mathcal{T}, (z_i^t, \hat{y}_i^t) \sim M(\cdot | x_i, t, \mathcal{E}_t)} [R(z_i^t, \hat{y}_i^t, y_i; t)], \tag{1}$$

where  $\mathcal{D} = \{(x_i, y_i)\}_{i=1}^D$  is the dataset with problem  $x_i$  and corresponding ground-truth  $y_i$ ;  $z_i^t$  and  $\hat{y}_i^t$  be model-generated reasoning trace/answer with modality t for i-th problem. To elicit the reasoning modality t, we design a small few-shot example set  $\mathcal{E}_t$  for each t, and prepend the exemplar from the set to each problem  $x_i$ . Conditioned on  $(x_i, t, \mathcal{E}_t)$ ,  $z_i$  is sampled from policy M, followed by the prediction of the final answer  $\hat{y}_i$ . R is the reward function and the design is detailed in the following.

**Reward Function** R. In preliminary experiments, we observe mismatch between tags and reasoning traces. This error leads to performance degradation, as different modalities negatively interfere with each other. Notably, this error is especially prevalent in the code modality. We define the reward as:

$$R(z, \hat{y}, y; t) = \begin{cases} 1, & y = \hat{y} \land \text{isValid}(z, t), \\ 0, & \text{otherwise,} \end{cases}$$
 (2)

where the isValid function checks the format consistency by two standards: a) each trace should correctly include its modality's structural tag (e.g., <end\_of\_nl\_cot> for nl) and b) for code traces, ensuring the presence of both a valid function definition (def) and a class definition (class); Following (Zelikman et al., 2022), we also filter out the traces with incorrect answer. We do not perform step-level verification of intermediate reasoning, as this would require additional external tools, e.g., LLMs judge (Zheng et al., 2023) and greatly slow down training. Instead, we find that simple checks (e.g., modality tags, basic code structures) already yield significant performance gains.

**Training.** We conduct multiple rounds of self-evolving training until performance saturates.  $M_n$  is used to denote the policy in the n-th round with trainable parameters  $\theta_n$ . Leveraging the

217

218

219

220

221

222

223

224

225

226

227

228

229

230

231

232

233

234

235

236

237

238

239

240

241

242

243

244

245246

247 248

249

250

251

252253254

255256

257

258

259

260

261

262

263

264

265 266

267

268

269

policy-gradient (Sutton et al., 1999) trick, we can easily obtain the gradient of eq. (1) as  $\nabla J = \mathbb{E}_{\substack{(x_i,y_i) \sim \mathcal{D}, t \sim \mathcal{T} \\ (z_i^t, \hat{y}_i^t) \sim M_{n-1}(\cdot | x_i, t, \mathcal{E}_t)}} \left[ R(z_i^t, \hat{y}_i^t, y_i; t) \, \nabla_{\theta_{n-1}} \log M_{n-1}(z_i^t, \hat{y}_i^t \mid x_i, t, \mathcal{E}_t) \right]. \tag{3}$ 

In our current setting, the reward is binary (Eq. 2). This reduces the framework to a degenerate case of RL that can be equivalently seen as rejection sampling plus supervised finetuning, closely related to selfevolving or STaR training (Zelikman et al., 2022). Algorithm 1 and Figure 2 illustrate our multi-round training procedure. At round n, we prompt the model  $M_{n-1}$ to generate a reasoning trace  $z_i^t$  and a predicted answer  $\hat{y}_i^t$  for each  $x_i$  across all reasoning modalities  $t \in \mathcal{T}$  (Line 4-9). It is worth noting that we use few-shot prompting only in the first round (Line 7); once the model has bootstrapped its own reasoning ability, all subsequent rounds run in zero-shot mode without additional exemplars (Line 9). We retain a sample only if it passes the quality filter (Line 11-13) and merge all surviving traces into  $\mathcal{D}^{\mathrm{gen}}_{\mathrm{all},n}$ (Line 16). The updated model  $M_n$ , which is finetuned from  $M_{n-1}$  on the filtered dataset  $\mathcal{D}_{\text{all},n}^{\text{gen}}$  (Line 17). Unlike (Zelikman et al., 2022), which restarts from the base model each round, our training proceeds on-policy—learning from its own validated outputs. We demonstrate the effectiveness of this change in Appendix E.3.

```
Input: an LLM M; dataset \mathcal{D} = \{(x_i, y_i)\}_{i=1}^D; reasoning modality \mathcal{T} = \{\text{NL}, \text{Code}, \text{TruthTable}\}, Sampling times S, few-shot examples \mathcal{E} = \{\mathcal{E}_{\text{NL}}, \mathcal{E}_{\text{Code}}, \mathcal{E}_{\text{TruthTable}}\}
```

**Output:** Mixture-of-Thought enhanced model  $M_N$ 

**Algorithm 1** Self-Evolving MoT Training

```
1: M_0 \leftarrow M
  2: for n = 1 to N do
                      Initialize \mathcal{D}_{\mathrm{all,n}}^{\mathrm{gen}} \leftarrow \emptyset; \mathcal{D}_{\mathrm{NL,n}}^{\mathrm{gen}} \leftarrow \emptyset; \mathcal{D}_{\mathrm{Code,n}}^{\mathrm{gen}} \leftarrow
  3:
            \emptyset; \mathcal{D}^{\text{gen}}_{\text{TruthTable},n} \leftarrow \emptyset
                       for all t \in \mathcal{T} do
   4:
   5:
                                  for i = 1 to D do
                                             if n = 1 then
   6:
                                                       z_i^t, \hat{y}_i^t \leftarrow M_{n-1}(x_i; t; \mathcal{E}_t; S)
   7:
   8:
                                                        z_i^t, \hat{y}_i^t \leftarrow M_{n-1}(x_i; t; S)
   9:
 10:
                                             end if
                                            \begin{array}{l} \textbf{if } R(z_i^t, \hat{y}_i^t, y_i; t) \text{=} 1 \textbf{ then} \\ \mathcal{D}_{t,n}^{\text{gen}} \leftarrow \mathcal{D}_{t,n}^{\text{gen}} \cup \{(x_i, z_i^t, y_i)\} \\ \textbf{end if} \end{array}
11:
12:
13:
                                  end for
14:
15:
                       \begin{aligned} \mathcal{D}_{\text{all},n}^{\text{gen}} \leftarrow \text{Mix}(\mathcal{D}_{\text{NL},n}^{\text{gen}}, \mathcal{D}_{\text{Code},n}^{\text{gen}}, \mathcal{D}_{\text{TruthTable},n}^{\text{gen}}) \\ M_n \leftarrow \text{Train}(M_{n-1}, \mathcal{D}_{\text{all},n}^{\text{gen}}) \end{aligned} 
16:
17:
18: end for
19: return M_N
```

# 2.3 MIXTURE-OF-THOUGHT INFERENCE

To leverage the complementary strengths of three modalities, for each problem, we have three outputs corresponding to three modalities elicited by tagging, *i.e.*, <nl\_cot>, <code>, and <truth\_table>, then we apply majority voting over outputs to decide the final answer. In case of ties, we randomly pick up the answer from one reasoning modality. We further explore the test-time scaling of MoT, and analyze its effectiveness in section 3.4.

# 3 EMPIRICAL EVALUATIONS

# 3.1 EXPERIMENTAL SETUPS

**Models.** To validate the effectiveness of our MoT, we select three widely-used LLMs across different sizes and model families: Qwen-2.5-7B-Instruct (Yang et al., 2024) and Gemma-2-2B-It/Gemma-2-9B-It (Team et al., 2024) as base models.

**Baselines.** Our approach is a kind of chain-of-thought approach. To this end, we select baselines from two folds: 1) neuro-symbolic approach and 2) chain-of-though approach. In the first category, we select Logic-LM (Pan et al., 2023) as a comparison. For the CoT approach, we select CoT (Wei et al., 2022) as a comparison. Since these approaches heavily rely on strong instruction-following capabilities, we directly cite their performance results from the original papers based on GPT-4.

**Dataset.** We select two logical reasoning benchmarks: ProofWriter (Tafjord et al., 2021) and FOLIO (Han et al., 2024) for evaluation. For ProofWriter, we select the hardest subset, which consists of 600 questions with reasoning depths of 5, the same as Pan et al. (2023). FOLIO is recognized for its high-quality export-made realistic test cases with more diverse reasoning depths ranging from 1-8. It consists of 203 questions. We utilize accuracy and pass@k as metrics.

Table 2: Accuracy (%) on the FOLIO and ProofWriter benchmarks. Our MoT training consistently improves the performance of each base model. Applying MoT inference further enhances performance across both benchmarks and all models. @3 denotes Self-Consistency approach (Wang et al., 2022) with three votes. We provide full results of extra baselines (e.g., LoT) in Appendix E.1 & E.2.

Model	Method Type	Reasoning Modality	FOLIO	ProofWriter	Avg				
(A) Prior SOTA Approach									
GPT-4	Logic-LM	-	78.9	79.7	79.3				
GP1-4	CoT (Vanilla)	-	70.6	68.1	69.4				
	(B) Base Model:	Gemma-2-2B-It							
Gemma-2-2B-It (3-Shot)	Single-Thought	Best (nl)	42.4	39.8	41.1				
Gemma-2-2B-It @ 3 (3-Shot)	Single-Thought	Best (nl)	45.3	38.8	42.1				
MoT (0-Shot)	Single-Thought	Best	61.1	62.7	61.9				
MoT (0-Shot)	Mixture-of-Thought	All	62.6	65.0	63.8				
	(C) Base Model:	Gemma-2-9B-It							
Gemma-2-9B-It (3-shot)	Single-Thought	Best (nl)	69.5	61.2	65.4				
Gemma-2-9B-It @ 3 (3-shot)	Single-Thought	Best (nl)	72.9	62.7	67.8				
MoT (0-shot)	Single-Thought	Best	76.9	69.5	73.2				
MoT (0-shot)	Mixture-of-Thought	All	<b>78.9</b>	70.7	74.8				
	(D) Base Model: Q	ven2.5-7B-Instruct							
Qwen2.5-7B-Instruct (3-shot)	Single-Thought	Best (nl)	71.9	60.5	66.2				
Qwen2.5-7B-Instruct @ 3 (3-shot)	Single-Thought	Best (nl)	73.4	65.8	69.6				
MoT (0-shot)	Single-Thought	Best	75.9	69.2	72.6				
MoT (0-shot)	Mixture-of-Thought	All	78.3	71.8	75.1				

**Training/Inference Details.** For each dataset, we collect 1000 training samples from the training set. We perform 2 or 3 rounds of self-evolving training. In each round, the model is fine-tuned for two epochs using a learning rate of 2e-5 and a batch size of 128. During the trajectory collection phase, the temperature, max\_tokens, and sample count are set to 1.0, 2048, and 10, respectively. Notably, we use a high temperature (1.0) to encourage diverse trajectories. We sample each problem 10 times during trajectory collection to maximize coverage. Of all the generated traces, only the first single trajectory that satisfies our quality criteria is retained for the final training set. For evaluation, the temperature and max\_tokens are configured to 0.7 and 2048, respectively. We do not perform hyperparameter tuning, so further optimization may yield even better performance. We run all experiments on 4 H100 GPUs. We employ vLLM engine (Kwon et al., 2023) to improve inference efficiency.

## 3.2 MAIN RESULTS

Table 2 displays the results on FOLIO and ProofWriter benchmarks. First, our Mixture-of-Thought (MoT) training with Single-Thought inference outperforms the corresponding base models by an average of **11.7pp** (from 41.1% to 61.9% for Gemma-2-2b-It, from 65.4% to 73.2% for Gemma-2-9b-It and from 66.2% to 72.6% for Qwen-2.5-7b-Instruct), demonstrating the effectiveness of our training strategy. When we further apply MoT inference, the MoT-trained model yields consistent additional gains of up to **2.0pp**. Notably, our 9B model achieves 78.9% accuracy on FOLIO, matching the performance of Logic-LM, which uses an external solver and focuses on close-sourced SoTA LLMs. We provide a detailed performance of both base models and the corresponding MoT models, as well as stronger baselines, in Appendix E.1 & E.2.

#### 3.3 MIXTURE-OF-THOUGHT TRAINING VS. SINGLE-THOUGHT TRAINING

In this section, we try to answer the key question: *Does MoT training truly offer benefits over Single-Thought training?* We have two baselines: 1) models trained on single-thought data and 2) models trained on partially MoT data, e.g., Code + NL. We evaluate both in-mode accuracy and cross-mode generalization. To enhance model's format following ability, we use 3-shot prompting to make model output the specific reasoning modality. Table 3 illustrates the results on FOLIO.

**SoT vs. MoT.** First, MoT training achieves the highest average accuracy across all three modalities, beating single-thought trained model, which indicates that our MoT training can jointly improve reasoning ability across all modalities. Second, MoT training can further push the performance boundary for each reasoning modality. For example, by using two of the three modalities, *i.e.*, Code and NL\_CoT, the trained models outperform all single-thought baselines. This clearly indicates

Table 3: Accuracy (%) of different training strategies across reasoning modalities (Same Round). Shaded cells denote in-domain evaluation, i.e., testing on the same modalities during training. Avg. refers to the average performance using three modalities while Ensemble means the majority vote results on three modalities. Values underlined indicate that the model did not follow the instruction (*e.g.*, when asked to use Code, it still used NL).

Training Approach	Param	Data	Code	NL_CoT	Truth Table	Avg.	Ensemble
w/o Training							
-	9B	N/A	56.7	69.5	63.6	63.3	66.0
Single-Thought Training							
Single-Thought (Code)	9B	-	61.6	59.1	64.0	61.6	70.4
Single-Thought (NL_CoT)	9B	-	52.7	73.9	69.0	65.2	73.4
Single-Thought (Truth Table)	9B	-	53.2	69.0	71.9	64.7	71.9
Single-Thought (Three Models Combined)	3x9B	$\sim 3 \times$	61.6	73.9	71.9	69.1	77.3
Mixture-of-Thought Training							
Mixture-of-Thought (NL_CoT + Truth Table)	9B	$\sim 2 \times$	65.5	72.9	69.5	69.3	72.9
Mixture-of-Thought (Truth Table + Code)	9B	$\sim 2 \times$	70.0	71.4	62.1	67.8	72.4
Mixture-of-Thought (Code + NL_CoT)	9B	$\sim 2 \times$	70.9	70.0	74.4	71.8	74.9
Mixture-of-Thought (Default, All)	9B	$\sim 3 \times$	73.9	76.9	70.0	73.6	78.9

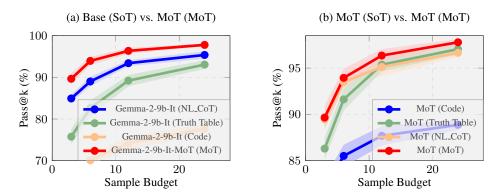


Figure 3: Pass@k vs. Sample Budget on FOLIO. (a) MoT-trained model with MoT sampling outperforms the base model (Gemma-2-9b-It) with SoT sampling. (b) Within the MoT-trained model, MoT sampling yields higher Pass@k than SoT sampling (NL\_CoT, Truth Table, Code).

synergy between these three complementary modalities during training. Third, deploying one model for each modality is resource-expensive. In contrast, MoT training enables a single model to seamlessly switch among reasoning modalities based on prompts.

**Partial MoT vs. MoT.** Our default Mixture-of-Thought setting yields the best average performance and achieves the best accuracy by using two combined reasoning paradigms, which indicates that all the modalities are useful. This superiority is further reflected in the ensemble accuracy, where MoT achieves 78.9%. We provide more evidence in Sec. 4.2 and Appendix E.7.

**Additional Ablations for MoT Training.** We further give more analysis to show 1) robust and optimal design of the MoT framework (Appendix E.3); 2) MoT training is better than single-thought training with distillation data (Appendix E.4) and 3) MoT data outperform an equivalent amount of diverse single-thought CoT data (Appendix E.5). These results underscore the practical and broader value of our MoT framework.

# 3.4 TEST-TIME SCALING ACROSS REASONING MODALITIES

We investigate how different single-thought and MoT inference scale with an increased test-time budget. To do this, we first generate 128 responses from each model with each modality. Then we evaluate two sampling strategies: 1) Single-Thought Sampling: We randomly select k responses from the 128 generated responses. and 2) MoT Sampling: Assuming there are  $N_T$  reasoning modalities, we sample  $\frac{k}{N_T}$  responses from each modality (so that the total number of responses is k). We choose k ranging from 3 to 24 and have 10 runs for each setting.

**MoT framework vs. Single-thought Baseline.** We compare our Gemma-2-9b-It-MoT with Gemma-2-9b-It. Figure 3 (a) shows our MoT model with MoT sampling consistently outperforms Gemma-2-9b-It with single-thought sampling. When the sample budget is less than 20, the performance gap is significant. It suggests that our MoT approach significantly increases the response diversity, leading to a more efficient use of inference compute. We observe a consistent phenomenon in terms of averaged accuracy (Appendix E.6, Figure 6).

Comparison of different modalities. We further plot the scaling curves of our MoT model (based on Gemma-2-9B-It) under three reasoning modalities in Figure 3 (b). Here are insights: 1) While NL significantly outperforms the truth-table paradigm at low k, their theoretical upper bounds converge as k increases; 2) The code paradigm exhibits the lowest upper bound among the three; 3) Across all values of k, our MoT framework consistently achieves the highest pass@k and attains the largest upper bound, indicating the largest potential of MoT trained models in test-time scaling.

## 4 FURTHER ANALYSIS

## 4.1 MIXTURE-OF-THOUGHT INFERENCE BENEFITS MORE FOR DIFFICULT PROBLEMS

We further identify the types of problems that benefit most from the proposed MoT inference approach. Specifically, we focus on problem difficulty, which can be effectively measured by the depth of reasoning. We conduct analysis on FOLIO and ProverQA (Qi et al., 2025). Figure 4 shows the performance of our MoT model with different reasoning modalities across reasoning depths.

We can see that MoT inference benefits more in solving more difficult logical problems. Our MoT model with MoT inference achieves an accuracy of 73.0% on challenging logical tasks with reasoning depths ranging from 5 to 8, outperforming each modality by a significant margin, with an average improvement of 9% points. However, such performance gains turn into slight degradation when

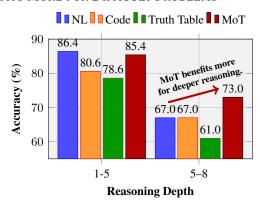


Figure 4: Performance comparison of different thought paradigms across reasoning depths. On FO-LIO, MoT inference exhibits better performance on difficult problems.

dealing with easy problems. Similar phenomenon can be observed in ProverQA (See Figure 7).

#### 4.2 Complementary, Uniqueness and Error Analysis

In this section, we quantify the complementary and uniqueness of our reasoning modalities and the training dynamics of our self-evolving MoT training. We focus on three metrics: 1) *Unique coverage*, *i.e.*, examples solved by exactly one modality; 2) *Complementarity coverage*, *i.e.*, examples solved by at least two modalities; and 3) *Oracle upper bound*, *i.e.*, examples solved by at least one modality.

Figure 1(a),(b) shows each modality's solve rate and oracle upper bound on ProofWriter and FOLIO. We further give a detailed unique and complementarity coverage and oracle upper bound in Table 8 in the Appendix. First, although our approach slightly reduces unique coverage compared to the baseline, both methods still achieve strong performance in this metric. Second, in terms of complementarity, our method increases the number of examples solved by multiple modalities—particularly on ProofWriter—demonstrating enhanced synergy. Third, by incorporating the truth-table paradigm alongside Code and NL, our model attains a higher oracle upper bound than prior work using only Code+NL, underscoring the benefit and necessity of the truth-table paradigm.

**Bottleneck of NL reasoning modality.** We perform a human evaluation of model outputs generated by natural language reasoning on the FOLIO dataset. We identify two major error patterns in the incorrectly solved cases: 1) failure to consider multiple cases when handling disjunction operations, such as "either/or"; 2) failure to utilize the transposition inference rule. For example, given  $A \rightarrow B$ , the model might sometimes incorrectly produce  $\neg A \rightarrow \neg B$ . Motivated by these observations and error types identified in prior work (Han et al., 2024; Olausson et al., 2023), we define four error categories: (i) invalid converse; (ii) missing branch; (iii) factual misquote; and (iv) incorporation of

commonsense knowledge and design an automatic pipeline to assess model rationales. Figure 1(c) presents the results, showing that invalid converse and missing-branch errors together account for nearly 66% of all errors. These findings underscore the value of the Truth Table thought paradigm.

**Scenarios that Truth Table excels in.** We manually analyze all 13 examples (Table 8) that were solved only using the truth table paradigm and find that 1) 5 out of 13 problems require transposition; 2) 5 out of 13 problems contain disjunction or similar operations (e.g., 'Rock can fly, or is a bird, or cannot breathe') and 3) 2 out of 13 problems contain both. This indicates that Truth Table may indeed complement the NL paradigm to some extent. We give two examples in Appendix H.2.

#### 5 RELATED WORK

432

433

434

435

436

437

438

439

440 441

442 443

444

445

446

447

448

449

450

451

452

453

454

455

456

457

458

459

460

461

462

463

464

465

466

467

468

469

470

471

472

473

474

475

476 477 478

479 480

481

482

483

484

485

LLMs for Symbolic Reasoning. Prior work has explored adapting LLMs to symbolic reasoning. One common approach treats LLMs as nl-to-fol translators, and then use an external symbolic prover to derive the final answer (Pan et al., 2023; Olausson et al., 2023; Callewaert et al., 2025; Ryu et al., 2025). While effective, this pipeline largely bypasses the model's internal reasoning capabilities, which our work seeks to fully leverage. To alleviate this problem, another line of work seeks to directly leverage LLMs' reasoning ability via CoT prompting (Wei et al., 2022). However, natural language remains inherently flexible and sometimes insufficient for structured reasoning. To bridge the gap between flexibility and formal rigor, recent work has explored combining natural and symbolic reasoning (Xu et al., 2024a; Liu et al., 2024; Xiong et al., 2024). These approaches often either rely on a primary reasoning modality (e.g., symbolic or NL), augmented with auxiliary signals from other representations (Xu et al., 2024a; Liu et al., 2024) or select one from multiple reasoning modalities (Xiong et al., 2024) at inference time. In contrast, our work 1) explicitly defines three kinds of reasoning paradigms covering natural language, symbolic and code-based reasoning. 2) goes beyond modality selection by jointly learning and inferring with all modalities, via a self-evolving MoT training and inference framework.

**Encouraging Diverse Thinking in Chain-of-Thoughts.** Previous work diversifies the CoT to further improve reasoning performance. A common strategy is to sample multiple outputs with higher temperatures (Wang et al., 2022; Brown et al., 2024), but this cannot guarantee true diversity (Wang et al., 2025). To address this, some work uses varied prompts—by task type (Wang et al., 2025), difficulty (Li et al., 2022), or strategy (Li et al., 2023; He et al., 2024)—and agent-based prompting via multi-agent debate (Liang et al., 2023; Hegazy, 2024) or self-reflection (Zhang et al., 2024) to elicit diverse CoTs. These methods diversify within one modality (NL or code). In contrast, we systematically introduce modality-level diversity—truth table, natural language, and code reasoning—which better aligns with the structural requirements of symbolic tasks and complements existing approaches. Recent work has also explored training smaller models on diverse CoTs generated by large LLMs (Ho et al., 2022; Puerto et al., 2024), though these approaches are limited to single-modality supervision and rely on external teacher models. In contrast, our method introduces modality-level diversity and requires no external supervision. We demonstrate that inter-modality diversity yields greater benefits for self-training than intra-modality diversity in Appendix E.5. Concurrent to our work, Chain-of-Reasoning (CoR) (Yu et al., 2025) also studies synergy of multiple reasoning modalities. However, CoR focuses on sequential synergy and targets mathematical reasoning. In contrast, our approach (i) focuses on logical reasoning and introduces truth table—based reasoning, (ii) develops a self-evolving Mixture-of-Thought training algorithm to bootstrap multi-modality capabilities, and (iii) exploits parallel synergy across modalities during both training and inference.

#### 6 Conclusion

We presented **Mixture-of-Thought** (MoT), a unified framework for improving logical reasoning by enabling LLMs to reason through natural language, code-based, and symbolic (truth table) paradigms within a single system. Unlike previous work, our approach combines a self-evolving training process that fosters cross-paradigm synergy with an inference-time voting mechanism that aggregates complementary reasoning strategies. Extensive experiments on two challenging logical reasoning benchmarks, FOLIO and ProofWriter, demonstrate that MoT substantially outperforms strong baselines, particularly on complex, high-depth problems.

# **ETHICS STATEMENT**

Our Mixture-of-Thought (MoT) framework is designed to improve logical reasoning by integrating multiple complementary modalities (natural language, code, truth tables). We do not foresee any risks related to ethics issues.

# REPRODUCIBILITY STATEMENT

We include detailed information about our approach in Section 2 and include experimental settings in Section 3 and Appendix D. We have give details in Section 3 and Appendix D to reproduce our results and we will open-source the code as soon as we collect our scripts for an easy way to reproduce our results

#### REFERENCES

- Bradley Brown, Jordan Juravsky, Ryan Ehrlich, Ronald Clark, Quoc V Le, Christopher Ré, and Azalia Mirhoseini. Large language monkeys: Scaling inference compute with repeated sampling. *arXiv* preprint arXiv:2407.21787, 2024.
- Benjamin Callewaert, Simon Vandevelde, and Joost Vennekens. Verus-lm: a versatile framework for combining llms with symbolic reasoning. *arXiv preprint arXiv:2501.14540*, 2025.
- Edmund Clarke, Armin Biere, Richard Raimi, and Yunshan Zhu. Bounded model checking using satisfiability solving. *Formal methods in system design*, 19:7–34, 2001.
- Dedre Gentner. Structure-mapping: A theoretical framework for analogy. *Cognitive science*, 7(2): 155–170, 1983.
- Gerald A. Goldin. Representational systems, learning, and problem solving in mathematics. *The Journal of Mathematical Behavior*, 17(2):137–165, 1998. ISSN 0732-3123. doi: https://doi.org/10.1016/S0364-0213(99)80056-1. URL https://www.sciencedirect.com/science/article/pii/S0364021399800561. Representations and the Psychology of Mathematics Education: Part II.
- Xinyu Guan, Li Lyna Zhang, Yifei Liu, Ning Shang, Youran Sun, Yi Zhu, Fan Yang, and Mao Yang. rstar-math: Small llms can master math reasoning with self-evolved deep thinking. *arXiv* preprint *arXiv*:2501.04519, 2025.
- Simeng Han, Hailey Schoelkopf, Yilun Zhao, Zhenting Qi, Martin Riddell, Wenfei Zhou, James Coady, David Peng, Yujie Qiao, Luke Benson, Lucy Sun, Alexander Wardle-Solano, Hannah Szabó, Ekaterina Zubova, Matthew Burtell, Jonathan Fan, Yixin Liu, Brian Wong, Malcolm Sailor, Ansong Ni, Linyong Nan, Jungo Kasai, Tao Yu, Rui Zhang, Alexander Fabbri, Wojciech Maciej Kryscinski, Semih Yavuz, Ye Liu, Xi Victoria Lin, Shafiq Joty, Yingbo Zhou, Caiming Xiong, Rex Ying, Arman Cohan, and Dragomir Radev. FOLIO: Natural language reasoning with first-order logic. In Yaser Al-Onaizan, Mohit Bansal, and Yun-Nung Chen (eds.), *Proceedings of the 2024 Conference on Empirical Methods in Natural Language Processing*, pp. 22017–22031, Miami, Florida, USA, November 2024. Association for Computational Linguistics. doi: 10.18653/v1/2024. emnlp-main.1229. URL https://aclanthology.org/2024.emnlp-main.1229/.
- Zhiwei He, Tian Liang, Wenxiang Jiao, Zhuosheng Zhang, Yujiu Yang, Rui Wang, Zhaopeng Tu, Shuming Shi, and Xing Wang. Exploring human-like translation strategy with large language models. *Transactions of the Association for Computational Linguistics*, 12:229–246, 2024.
- Mahmood Hegazy. Diversity of thought elicits stronger reasoning capabilities in multi-agent debate frameworks. *ArXiv*, abs/2410.12853, 2024. URL https://api.semanticscholar.org/CorpusID:273403906.
- Namgyu Ho, Laura Schmid, and Se-Young Yun. Large language models are reasoning teachers. *arXiv* preprint arXiv:2212.10071, 2022.

Arian Hosseini, Xingdi Yuan, Nikolay Malkin, Aaron Courville, Alessandro Sordoni, and Rishabh Agarwal. V-star: Training verifiers for self-taught reasoners. *arXiv preprint arXiv:2402.06457*, 2024.

- Woosuk Kwon, Zhuohan Li, Siyuan Zhuang, Ying Sheng, Lianmin Zheng, Cody Hao Yu, Joseph E. Gonzalez, Hao Zhang, and Ion Stoica. Efficient memory management for large language model serving with pagedattention. In *Proceedings of the ACM SIGOPS 29th Symposium on Operating Systems Principles*, 2023.
- Jill H Larkin and Herbert A Simon. Why a diagram is (sometimes) worth ten thousand words. *Cognitive science*, 11(1):65–100, 1987.
- Chengpeng Li, Mingfeng Xue, Zhenru Zhang, Jiaxi Yang, Beichen Zhang, Xiang Wang, Bowen Yu, Binyuan Hui, Junyang Lin, and Dayiheng Liu. Start: Self-taught reasoner with tools. *arXiv* preprint arXiv:2503.04625, 2025.
- Xin-Ye Li, Jiang-Tian Xue, Zheng Xie, and Ming Li. Think outside the code: Brainstorming boosts large language models in code generation. *arXiv* preprint arXiv:2305.10679, 2023.
- Yujia Li, David Choi, Junyoung Chung, Nate Kushman, Julian Schrittwieser, Rémi Leblond, Tom Eccles, James Keeling, Felix Gimeno, Agustin Dal Lago, et al. Competition-level code generation with alphacode. *Science*, 378(6624):1092–1097, 2022.
- 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 multi-agent debate. *arXiv preprint arXiv:2305.19118*, 2023.
- Haohan Lin, Zhiqing Sun, Sean Welleck, and Yiming Yang. Lean-star: Learning to interleave thinking and proving. *arXiv preprint arXiv:2407.10040*, 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. *arXiv* preprint arXiv:2409.17539, 2024.
- Allen Newell, Herbert Alexander Simon, et al. *Human problem solving*, volume 104. Prentice-hall Englewood Cliffs, NJ, 1972.
- Theo X. Olausson, Alex Gu, Ben Lipkin, Cedegao E. Zhang, Armando Solar-Lezama, Joshua B. Tenenbaum, and Roger P. Levy. LINC: A neurosymbolic approach for logical reasoning by combining language models with first-order logic provers. In *The 2023 Conference on Empirical Methods in Natural Language Processing*, 2023. URL https://openreview.net/forum?id=h00GHjWDEp.
- Liangming Pan, Alon Albalak, Xinyi Wang, and William Yang Wang. Logic-lm: Empowering large language models with symbolic solvers for faithful logical reasoning. *arXiv preprint arXiv:2305.12295*, 2023.
- XIANGYU PENG, Congying Xia, Xinyi Yang, Caiming Xiong, Chien-Sheng Wu, and Chen Xing. Regenesis: LLMs can grow into reasoning generalists via self-improvement. In *The Thirteenth International Conference on Learning Representations*, 2025. URL https://openreview.net/forum?id=YUYJsHOf3c.
- Gabriel Poesia, Kanishk Gandhi, Eric Zelikman, and Noah Goodman. Certified deductive reasoning with language models. *Transactions on Machine Learning Research*, 2024. ISSN 2835-8856. URL https://openreview.net/forum?id=yXnwrs2T16.
- Haritz Puerto, Tilek Chubakov, Xiaodan Zhu, Harish Tayyar Madabushi, and Iryna Gurevych. Fine-tuning with divergent chains of thought boosts reasoning through self-correction in language models. *ArXiv*, abs/2407.03181, 2024. URL https://api.semanticscholar.org/CorpusID:270924195.

- Chengwen Qi, Ren Ma, Bowen Li, He Du, Binyuan Hui, Jinwang Wu, Yuanjun Laili, and Conghui He. Large language models meet symbolic provers for logical reasoning evaluation. In *The Thirteenth International Conference on Learning Representations*, 2025. URL https://openreview.net/forum?id=C25SqeXWjE.
  - 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 *The Thirteenth International Conference on Learning Representations*, 2025. URL https://openreview.net/forum?id=09FiNmvNMw.
  - Charlie Victor Snell, Jaehoon Lee, Kelvin Xu, and Aviral Kumar. Scaling LLM test-time compute optimally can be more effective than scaling parameters for reasoning. In *The Thirteenth International Conference on Learning Representations*, 2025. URL https://openreview.net/forum?id=4FWAwZtd2n.
  - Richard S Sutton, David McAllester, Satinder Singh, and Yishay Mansour. Policy gradient methods for reinforcement learning with function approximation. In S. Solla, T. Leen, and K. Müller (eds.), *Advances in Neural Information Processing Systems*, volume 12. MIT Press, 1999. URL https://proceedings.neurips.cc/paper\_files/paper/1999/file/464d828b85b0bed98e80ade0a5c43b0f-Paper.pdf.
  - Oyvind Tafjord, Bhavana Dalvi, and Peter Clark. ProofWriter: Generating implications, proofs, and abductive statements over natural language. In Chengqing Zong, Fei Xia, Wenjie Li, and Roberto Navigli (eds.), *Findings of the Association for Computational Linguistics: ACL-IJCNLP 2021*, pp. 3621–3634, Online, August 2021. Association for Computational Linguistics. doi: 10.18653/v1/2021.findings-acl.317. URL https://aclanthology.org/2021.findings-acl.317/.
  - Gemma Team, Morgane Riviere, Shreya Pathak, Pier Giuseppe Sessa, Cassidy Hardin, Surya Bhupatiraju, Léonard Hussenot, Thomas Mesnard, Bobak Shahriari, Alexandre Ramé, et al. Gemma 2: Improving open language models at a practical size. *arXiv preprint arXiv:2408.00118*, 2024.
  - Lewis Tunstall, Edward Beeching, Nathan Lambert, Nazneen Rajani, Shengyi Huang, Kashif Rasul, Alvaro Bartolome, Alexander M. Rush, and Thomas Wolf. The Alignment Handbook. URL https://github.com/huggingface/alignment-handbook.
  - Danqing Wang, Jianxin Ma, Fei Fang, and Lei Li. Typedthinker: Diversify large language model reasoning with typed thinking. In *The Thirteenth International Conference on Learning Representations*, 2025. URL https://openreview.net/forum?id=VIUisLx8lQ.
  - Tianlu Wang, Ilia Kulikov, Olga Golovneva, Ping Yu, Weizhe Yuan, Jane Dwivedi-Yu, Richard Yuanzhe Pang, Maryam Fazel-Zarandi, Jason Weston, and Xian Li. Self-taught evaluators. *arXiv* preprint arXiv:2408.02666, 2024.
  - Xuezhi Wang, Jason Wei, Dale Schuurmans, Quoc Le, Ed Chi, Sharan Narang, Aakanksha Chowdhery, and Denny Zhou. Self-consistency improves chain of thought reasoning in language models. *arXiv* preprint arXiv:2203.11171, 2022.
  - 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.
  - Johan Wittocx, Maarten Mariën, and Marc Denecker. Grounding fo and fo (id) with bounds. *Journal of Artificial Intelligence Research*, 38:223–269, 2010.
  - Xuyuan Xiong, Simeng Han, Ziyue Zhou, and Arman Cohan. Hybridmind: Meta selection of natural language and symbolic language for enhanced llm reasoning. 2024. URL https://api.semanticscholar.org/CorpusID:273501516.
  - Jundong Xu, Hao Fei, Liangming Pan, Qian Liu, Mong-Li Lee, and Wynne Hsu. Faithful logical reasoning via symbolic chain-of-thought. 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)*, pp. 13326–13365, Bangkok, Thailand, August 2024a. Association for Computational Linguistics. doi: 10.18653/v1/2024.acl-long.720. URL https://aclanthology.org/2024.acl-long.720/.

- Jundong Xu, Hao Fei, Liangming Pan, Qian Liu, Mong-Li Lee, and Wynne Hsu. Faithful logical reasoning via symbolic chain-of-thought. *arXiv* preprint arXiv:2405.18357, 2024b.
- An Yang, Baosong Yang, Beichen Zhang, Binyuan Hui, Bo Zheng, Bowen Yu, Chengyuan Li, Dayiheng Liu, Fei Huang, Haoran Wei, et al. Qwen2. 5 technical report. *arXiv preprint arXiv:2412.15115*, 2024.
- Yiyao Yu, Yuxiang Zhang, Dongdong Zhang, Xiao Liang, Hengyuan Zhang, Xingxing Zhang, Ziyi Yang, Mahmoud Khademi, Hany Awadalla, Junjie Wang, et al. Chain-of-reasoning: Towards unified mathematical reasoning in large language models via a multi-paradigm perspective. *arXiv* preprint arXiv:2501.11110, 2025.
- Murong Yue, Jie Zhao, Min Zhang, Liang Du, and Ziyu Yao. Large language model cascades with mixture of thought representations for cost-efficient reasoning. In *The Twelfth International Conference on Learning Representations*, 2024. URL https://openreview.net/forum?id=6okaSfANzh.
- Eric Zelikman, Yuhuai Wu, Jesse Mu, and Noah Goodman. Star: Bootstrapping reasoning with reasoning. *Advances in Neural Information Processing Systems*, 35:15476–15488, 2022.
- Eric Zelikman, Georges Harik, Yijia Shao, Varuna Jayasiri, Nick Haber, and Noah D Goodman. Quiet-star: Language models can teach themselves to think before speaking. *arXiv preprint arXiv:2403.09629*, 2024a.
- Eric Zelikman, Eliana Lorch, Lester Mackey, and Adam Tauman Kalai. Self-taught optimizer (stop): Recursively self-improving code generation. In *First Conference on Language Modeling*, 2024b.
- Wenqi Zhang, Yongliang Shen, Linjuan Wu, Qiuying Peng, Jun Wang, Yueting Zhuang, and Weiming Lu. Self-contrast: Better reflection through inconsistent solving perspectives. 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)*, pp. 3602–3622, Bangkok, Thailand, August 2024. Association for Computational Linguistics. doi: 10.18653/v1/2024. acl-long.197. URL https://aclanthology.org/2024.acl-long.197/.
- Lianmin Zheng, Wei-Lin Chiang, Ying Sheng, Siyuan Zhuang, Zhanghao Wu, Yonghao Zhuang, Zi Lin, Zhuohan Li, Dacheng Li, Eric Xing, et al. Judging llm-as-a-judge with mt-bench and chatbot arena. *Advances in neural information processing systems*, 36:46595–46623, 2023.
- Tong Zheng, Hongming Zhang, Wenhao Yu, Xiaoyang Wang, Xinyu Yang, Runpeng Dai, Rui Liu, Huiwen Bao, Chengsong Huang, Heng Huang, et al. Parallel-r1: Towards parallel thinking via reinforcement learning. *arXiv* preprint arXiv:2509.07980, 2025.

# A THE USE OF LLMS

We only use LLMs to polish the paper writing.

#### B LIMITATIONS AND FUTURE WORK

While our Mixture-of-Thought (MoT) framework demonstrates strong performance on logical reasoning tasks, we have not evaluated its effectiveness on other types of reasoning tasks, such as mathematical or commonsense reasoning. Additionally, our test-time scaling experiments suggest promising directions—such as dynamic mixture-of-thought sampling under budget constraints—but our current work still has not fully explored the benefits of complementary reasoning modalities. Further exploring these aspects could be important to further push the performance boundary of open-source models on reasoning.

We plan to further explore them in the two aspects:

- Extended to boarder tasks: currently our work cannot directly applied to reasoning tasks out of logical reasoning. This is because the reasoning modality we define in our work is specific for logical reasoning, *e.g.*, Truth Table. Therefore, we plan to define more general but complementary reasoning modality that can be applied to more general broader of reasoning tasks and further show how our MoT framework can further improve performance of reasoning tasks beyond logical reasoning.
- Adaptive parallel thinking via mixture-of-thoughts. An interesting question is *How can we fully leverage the benefits of complementary reasoning modalities during inference?* Recent work has shown that reinforcement learning can instill parallel thinking in LLMs (Zheng et al., 2025), enabling models to dynamically activate and coordinate multiple reasoning paths. This opens up a natural opportunity to integrate such adaptive parallel thinking with our MoT framework: treating each reasoning modality as an atomic unit, and allowing the model to dynamically trigger, combine, or prune modalities during inference. We believe that this unified view, blending MoT's modality-level diversity with dynamic parallelism, can lead to more flexible and efficient reasoning systems.

# C Broader Related Work: Self-evolving training.

Self-evolving training techniques have been widely adopted to improve reasoning ability in LLMs, especially when there is lack of reasoning trajectories. Notably, Zelikman et al. (2022) propose a bootstrapping framework that iteratively generates and verifies reasoning trajectories based on the derived final answer, then fine-tunes the model on these self-labeled examples to improve reasoning performance with minimal human supervision. Following this idea, several works adapt self-evolving training to a wider range of tasks (Hosseini et al., 2024; Xiong et al., 2024; Zelikman et al., 2024b; Wang et al., 2024; Lin et al., 2024; Zelikman et al., 2024a; Poesia et al., 2024; Guan et al., 2025; Li et al., 2025). Additionally, researchers also explore improving the high-quality of rationales during STaR algorithm (Poesia et al., 2024; Guan et al., 2025; PENG et al., 2025), incorporating techniques such as formal verification, monte carlo tree search, and abstract-to-concrete prompting. While previous work primarily focuses on generating higher-quality reasoning paths within a single modality, our work explores a complementary direction: how to jointly evolve and coordinate reasoning across multiple thought paradigms.

# D DETAILED EXPERIMENTAL SETTINGS

#### D.1 DATASETS

In this work, we adopt three logical reasoning datasets: 1) FOLIO (Han et al., 2024), 2) ProofWriter (Tafjord et al., 2021), and 3) ProverQA (Qi et al., 2025), to evaluate the effectiveness of our MoT framework.

**FOLIO** (Han et al., 2024). FOLIO provides both the training and validation subsets, consisting of 1003 and 203 samples, respectively. There are two subsets with different difficulties: 1) HybLogic:

contains 100 complex logical problems (5-8 reasoning steps) and 2) WikiLogic: contains 103 simper logical problems (1-5 reasoning steps). In this work, we sample 1000 training samples from the FOLIO training set as seed dataset for our self-evovling MoT training and evaluate both baselines and our trained model on the FOLIO validation set.

**ProofWriter** (**Tafjord et al., 2021**). ProofWriter is a synthetic dataset designed for evaluating the logical reasoning abilities of language models. It consists of multiple subsets, each containing logical reasoning problems of varying reasoning depths—from depth 0 (direct inference) up to depth 5 (requiring multi-step logical deductions). Following Pan et al. (2023), we select the most challenging subset (reasoning depth 5) to construct our training and test data. Specifically, we sample 1,000 instances from the training set provided by Pan et al. (2023) as our training data and adopt their original test set directly for fair evaluation.

**ProverQA** (**Qi et al., 2025**). ProverQA is a recently proposed logical reasoning benchmark, notable for its large scale, high quality, and diversity. It consists of three subsets, each corresponding to a different reasoning difficulty level (*i.e.*, reasoning depth). We select these subsets to evaluate the performance of our MoT framework across varying levels of reasoning complexity.

# D.2 TRAINING DETAILS

We conduct all experiments on 4 H100 GPUs with Alignment Handbook (Tunstall et al.). For each dataset, we sample 1,000 training examples and perform 2–3 rounds of self-evolving training. In each round, the model is fine-tuned for 2 epochs with a learning rate of 2e-5 and a batch size of 128. We do not perform hyperparameter tuning. Further tuning may lead to better performance. All experiments are run with a fixed seed, *i.e.*, 42, for reproducibility.

# D.3 INFERENCE DETAILS

We employ vLLM (Kwon et al., 2023) for efficient inference. During trajectory collection, we generate 10 reasoning traces per example using temperature 1.0, max\_tokens 2048, and sampling count 10. To maximize coverage while ensuring quality, we retain only the first generated trace that passes our quality check. For evaluation, we set the temperature to 0.7 and max\_tokens to 2048. All experiments are run with a fixed seed, *i.e.*, 42, for reproducibility.

# E ADDITIONAL EXPERIMENTAL RESULTS

# E.1 EVALUATING LLM PERFORMANCE ACROSS REASONING MODALITIES ON FOLIO AND PROOFWRITER

Table 4: Performance of three models across reasoning modalities on FOLIO and ProofWriter.

Model	FOLIO			ProofWriter			
	NL	Code	Truth Table	NL	Code	Truth Table	
Gemma-2-2B-It	42.4	38.4	36.5	39.8	40.8	37.5	
+ MoT training	61.1 (18.7\(\dagger))	61.1 (22.7\(\dagger))	58.6 (22.1\(\dagger)\)	62.7 (22.9\(\dagger))	61.7 (20.9\(\dagger))	60.2 (22.7\(\dagger))	
Gemma-2-9B-It	69.5	56.7	63.6	61.2	39.5	55.8	
+ MoT training	76.9 (7.4\(\dagger))	73.9 (17.2\(\dagger)\)	70.0 (6.4\(\dagger)\)	68.5 (7.3\(\dagger))	69.5 (30.0\(\dagger)\)	66.7 (10.9\(\dagger))	
Qwen-2.5-7B-Instruct	71.9	62.1	69.0	60.5	42.3	53.0	
+ MoT training	75.9 (4.0†)	68.5 (6.4†)	71.9 (2.9\(\dagger))	69.2 (8.7†)	66.7 (24.4†)	64.3 (11.3†)	

Table 4 displays detailed results of baselines across reasoning modalities on FOLIO and ProofWriter. We can observe that LLMs owns uneven ability across these reasoning modalities. This also highlights the necessary of our self-evolving MoT training, which can equip LLMs with three complementary reasoning modalities. After self-evolving MoT training, all modalities show joint improvements. This effect is especially significant in smaller models, *i.e.*, Gemma-2-2B-It achieves up to a more than 20% increase in accuracy on average.

# E.2 COMPARISON WITH MORE BASELINES ON FOLIO AND PROOFWRITER

Table 5: Comparison with more baselines on FOLIO and ProofWriter

Method	Base Model	FOLIO (Acc %)	ProofWriter (Acc %)
HybridMind (Xiong et al., 2024)	GPT-3.5	76.6	_
LINC (Olausson et al., 2023)	GPT-4	72.5	=
Symbolic CoT (Xu et al., 2024b)	GPT-4	83.3	82.5
Logic-of-Thoughts @ 5	GPT-3.5	81.5	65.9
Logic-of-Thoughts @ 5	GPT-4	88.2	72.0
MoT	Gemma-2-2b-It	62.6	65.0
MoT	Gemma-2-9b-It	78.9	70.7
MoT	Qwen2.5-7B-Instruct	78.3	71.8

Table 5 presents a comparison between our approach and prior state-of-the-art systems. It demonstrates that our open-source MoT models nearly match the performance of leading closed-source prompting methods (*e.g.*, GPT-3.5 and GPT-4). This indicates that enabling LLMs to learn complementary reasoning modalities is a promising direction.

Table 6: Ablation studies on (1) policy strategy; and (2) mixing strategy.

Ablation	Setting	FOLIO Accuracy (%)				
1101001011	500000	NL	Code	Truth Table	MoT	
1. Policy Strategy	Off-policy MoT	55.2	54.7	53.7	56.7	
	On-policy MoT (default)	<b>61.1</b>	<b>61.1</b>	<b>58.6</b>	<b>62.6</b>	
2. Mixing Strategy	Random single-modality per question	49.8	50.3	48.3	53.7	
	Direct mixing (default)	<b>61.1</b>	<b>61.1</b>	<b>58.6</b>	<b>62.6</b>	

#### E.3 ABLATION STUDIES

We perform ablation studies on three core components: 1) policy strategy, *i.e.*, on-policy vs. off-policy (Zelikman et al., 2022) and 2) mixing approach, *i.e.*, direct mixture vs. mixture by unique conclusion (randomly select single-modality per question).

Table 6 reports FOLIO accuracies under each setting. We make two key observations:

- On-policy training yields consistent gains. Switching from off-policy to on-policy increases single-modality CoT accuracy by approximately 5–6 pp (e.g., NL CoT from 55.2% to 61.1%) and raises MoT's final accuracy from 56.7% to 62.6%. This demonstrates the importance of updating the model with its most recent outputs.
- **Direct mixing outperforms random single-modality sampling.** Presenting all three modalities together boosts accuracy by 8–10 pp compared to randomly picking one modality per question (MoT: 62.6% vs. 53.7%). This indicates that joint exposure to multiple modalities provides stronger complementary signals than isolated examples.

# E.4 IMPACT OF QUALITY OF INITIAL TRAINING DATA: DISTILLATION + SINGLE-MODAL TRAINING VS. RAW DATA + MOT TRAINING

Intuitively, the first-round data are crucial and have a strong impact on the efficacy of self-evolving training. Therefore, we are interested in the following question: Can self-evolving single-thought training enhanced by first-round distillation outperform our self-evolving mixture-of-thought training without any distillation? To answer this question, we compare the following settings: 1) Self-evolving single-thought (nl) training but with distillation data from o4-mini for first round training, which can provide a better initialization; 2) our MoT training without any distillation data; and 3) Self-evolving single-thought (nl) training without any distillation data. Figure 5 displays the results of Gemma-2-2b-It on FOLIO benchmark.

The key observations are: adding distillation data from stronger LLMs is beneficial for improving performance and convergence rate (blue line vs. orange line), but still lags behind our self-evolving

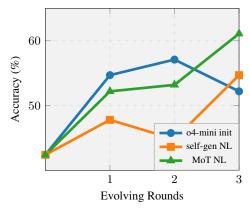


Figure 5: Accuracy (%) over three self-evolving rounds on the FOLIO benchmark for: distilled NL-CoT (first-round only), raw NL-CoT (no distillation), and MoT (no distillation). The performance is evaluated with NL-based reasoning.

MoT training (blue line vs. green line). This suggests the advantages of our self-evolving MoT training: 1) It requires no reliance on stronger—often more expensive—LLMs; 2) It provides a higher upper bound accuracy.

# E.5 FINETUNING WITH DIVERSE SINGLE-MODALITY COT VS. FINETUNING WITH MOT

Ho et al. (2022); Puerto et al. (2024) have explored that finetuning LLMs with diverse CoT can further improve the performance. A natural question then is: given a fixed budget of training examples, which strategy yields better results? (1) self-training with 3N natural-language CoT samples, or (2) self-training with a total of 3N samples composed of N examples from each of three modalities (NL, Code, Truth-Table)?

We consider two settings to answer this question: 1) Self-evolving training with 3N natural-language CoT samples for 2 epochs per round over 3 rounds. We sample 10 reasoning traces per question with temperature of 1.0 and keep the 3 reasoning traces that satisfy our filtering criteria; 2) Self-evolving training with a total of 3N samples comprising N examples from each of the three modalities (NL, Code, Truth-Table) for 2 epochs per round over 3 rounds. We evaluate those trained model with natural language modality on FOLIO dataset.

Table 7: Accuracy (%) of Gemma-2-2b-It under three self-evolving regimes, with budgets of N or 3N training samples. The accuracy is evaluated with NL-based reasoning on FOLIO benchmark. We can see self-evolving training with MoT achieves the best accuracy, demonstrating the benefit of modality-level diversity.

#	Setting	Training Samples	Accuracy (%)
1	NL_CoT	N	54.7
2	NL_CoT	3N	57.1
3	MoT data	3N	61.1

Table 7 shows the results. We can have the following observations: 1) finetuning with diverse NL CoT can indeed improve the performance (#1 vs. #2), which is consistent with findings from Ho et al. (2022); Puerto et al. (2024). 2) Finetuning with MoT data is more efficient than finetuning with same amount of diverse NL CoT data (#2 vs. #3). This indicates that the diversity of single-modality CoT data obtained by sampling with high temperature is not sufficient. By contrast, our MoT data, which leverages the complementarity of truth table, code and nl, can produce more diversity, and therefore improve the training efficiency.

#### E.6 ADDITIONAL RESULTS ON TEST-TIME SCALING ACROSS REASONING MODALITIES

**MoT With Different Thought Paradigms** Table 3 (b) illustrates the scaling behavior of our MoT model across different thought paradigms under varying sample budgets. We observe that code-based reasoning consistently lags behind all other paradigms, indicating its relatively poor performance and limited scalability.

Another interesting phenomenon is that natural language-based reasoning achieves relatively strong performance when the sample budget is small (e.g., k < 20), outperforming the truth table-based paradigm in this regime. However, as the sample budget increases (e.g., k > 20), truth table reasoning begins to match even outperform NL-based reasoning—highlighting its greater potential when more inference resources are available.

Notably, our MoT (ALL) approach offers a favorable trade-off between these two paradigms: it achieves strong performance under low-budget conditions, while delivering better performance when the sample budget is large.

Accuracy vs. Sample Budget Figure 6 presents accuracy-vs-sample-budget curves across different reasoning paradigms. We find that our MoT (ALL) model—trained and inferred under the mixture-of-thought setting—consistently achieves the highest accuracy, outperforming all other approaches regardless of budget size. Additionally, our MoT model can benefit better from increased sample budget compared with all other approaches. Among individual paradigms, NL-CoT performs best under majority voting, while truth table reasoning is more stable but shows limited improvement with increased budget. Code-based reasoning remains the least effective. These results reinforce the value of our MoT framework.

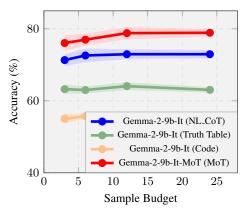


Figure 6: Accuracy vs. Sample Budget for different modes

#### E.7 DETAILED COMPLEMENTARY, UNIQUENESS ANALYSIS

Across both ProofWriter and FOLIO benchmarks, our Mixture-of-Thought (MoT) model shifts away from single-paradigm reliance and toward multi-paradigm collaboration. First, the number of examples solved exclusively by the NL paradigm drops by over 50% (ProofWriter: from 109 to 55; FOLIO: from 18 to 8), and "Only TT correct" cases likewise decrease, indicating that MoT reduces brittle, single-mode reasoning. Second, pairwise overlaps (NL  $\cap$  Code, NL  $\cap$  TT, Code  $\cap$  TT) all increase substantially—NL ∩ Code on ProofWriter rises by 76% (172  $\rightarrow$  304), and similar gains appear on FO-LIO—showing that MoT effectively combines different reasoning formats on the same instance. Finally, the overall coverage (Code  $\cup$  NL  $\cup$  TT) improves from 511 to 544 (+6.5%) on ProofWriter and from 174 to 181 (+4%) on FOLIO, demonstrating that MoT recovers difficult

Table 8: Prediction Category Distribution on Two Benchmarks (Qwen-2.5-7B-Instruct vs Qwen-2.5-7B-Instruct-MoT).

Category	ProofW	riter	FOLIO		
category	Baseline	Our	Baseline	Our	
Single-paradigm only	,				
Only NL correct	109	55	18	8	
Only Code correct	40	32	5	6	
Only TT correct	66	33	11	13	
Pairwise overlap only	v				
$NL \cap Code$ only	172	304	109	125	
$NL \cap TT$ only	210	289	117	125	
Code $\cap$ TT only	170	297	110	112	
$Code \cup NL$	445	511	163	168	
$Code \cup NL \cup TT$	511	544	174	181	

cases missed by the baseline. The consistent trends across two datasets confirm that encouraging multi-paradigm synergy yields more robust and comprehensive logical reasoning performance.

# E.8 Performance vs. Difficulty on ProverQA (QI et al., 2025)

Figure 7 displays the results.

973

974

975

976 977

978

979 980

981 982

983

984

985

986 987

988 989

990

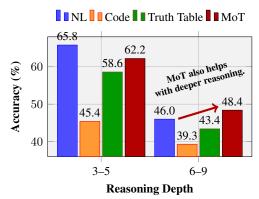


Figure 7: Performance comparison of different thought paradigms across reasoning depths. On ProverQA, MoT inference exhibits better performance on difficult problems.

# F FULL PROMPTS FOR MOT

The full prompts we utilized in this work are illustrated as follows:

```
991
992
993
            You are a rigorous and logically precise AI assistant. Your
994
                task is to answer a logical reasoning problem strictly
995
                following one of three modes, as explicitly specified in
996
                the input. Only one mode will be present in the input.
997
                Follow that mode exclusively.
998
            - Code Mode (<code> ... <end_of_code> <answer> ... <
999
                end_of_answer>)
1000
              - If the input contains <code>, translate the problem into
1001
                  Python code.
1002
              - Execute the logic and derive the answer.
1003
            - Natural Language Chain-of-Thought Mode (<nl_cot>
1004
                end_of_nl_cot> <answer> ... <end_of_answer>)
              - If the input contains <nl_cot>, solve the problem step by
                  step in natural language.
1007
            - Truth Table Mode (<truth_table> ... <end_of_truth_table> <
1008
                answer> ... <end_of_answer>)
1009
              - If the input contains <truth_table>, construct a truth
1010
                  table and derive the answer from it.
1011
1012
             ### Rules
            - Only use the mode specified in the input. Do not switch modes
1013
1014
             - Generate output strictly in the specified mode and format,
1015
                with no additional text.
1016
             - Enclose all reasoning strictly within the corresponding mode
1017
                tags.
              The final answer must be strictly enclosed in <answer> ... <
1018
                end_of_answer>.
1019
            - Do not provide any reasoning or explanations outside of the
1020
                designated mode tags.
            The following is the problem you need to solve.
1023
             premises>
1024
             {premises}
1025
```

```
1026
1027
             </premises>
             <conclusion>
1029
              {conclusion}
1030
              </conclusion>
1031
1032
             <question>
             Is the following statement true, false, or uncertain? {
1033
                 conclusion}
1034
              </question>
1035
1036
             <options>
1037
              (A) True
              (B) False
1038
              (C) Uncertain
1039
             </options>
1040
1041
             <{taq}>
1042
1043
```

# Full prompt used for Error Detection

1044 1045 1046

1047

1048 1049

1050

1051

1052

1053

1054

1055

1056

1057

1058

1059

1060

1061

1062

1063

1064

1065

1066

1067

1068

1069

1070

1071

1072

1073

1074

1075

1076

1077

1078

```
"You must determine whether a rationale faithfully justifies
   the truth value of a conclusion given a set of premises.\n\
      "Faithful means all and only the steps actually used in
         deriving the conclusion:\n"
      "- are grounded in the given premises or prior derived
         steps, \n"
      "- apply valid inference rules (no illicit converse or
         contraposition), \n"
      "- cover every disjunction branch or quantifier case, \n"
      "- use no unstated assumptions, external knowledge, or
         background commonsense, \n"
      "- and correctly assess whether the conclusion is
         supported or contradicted by the premises.\n\n"
      "You must also diagnose where and how the rationale fails
          when it is unfaithful, allowing trivial unused
         remarks to be overridden.\n\n"
      "Error Types:\n"
      "- Missing Branch: Failing to exhaustively consider all
         branches of a disjunction, conditionals, or
         quantified cases.\n"
      "- Invalid Converse: Illicitly reversing the direction of
          a conditional (e.g., mistaking 'A \rightarrowB' for 'B \rightarrowA').\
         n"
      "- Commonsense Injection: Using external background
         knowledge or commonsense not entailed or implied by
         the premises.\n"
      "- Factual Misquote: Misrepresenting, distorting, or
         misquoting the explicit content of the premises.\n\n"
      "Input (JSON):\n"
      "{\n"
      ' "premises": "<string>",\n'
      ' "conclusion": "<string>",\n'
      ' "rationale": "<string>",\n'
      ' "label": "<string>",\n'
      ' "predict": "<string>"\n'
      "}\n\n"
```

```
1080
1081
                   "Output (JSON):\n"
1082
                   "{\n"
1083
                     "faithful": true | false, \n'
1084
                   ' "error_type": "<missing branch | invalid converse |</pre>
1085
                       commonsense injection | factual misquote>",\n'
                   ' "error_location": "<e.g., Step 3, Clause 2>",\n'
1086
                   ' "override": true | false, \n'
1087
                   ' "analysis": "<brief summary explaining why the
1088
                       reasoning is faithful or unfaithful, citing specific
1089
                       logical failures>"\n'
1090
                   "}\n\n"
1091
                   "Notes:\n"
1092
                   "- If multiple error types apply, list them all separated
1093
                        by commas.\n"
1094
                   "- Always identify the first point in the rationale where
1095
                        the faithfulness failure occurs.\n"
1096
                   "- Be concise, precise, and consistent in your labeling.\
1097
                      n\n"
                   "Input:\n"
1098
1099
```

# G More Examples Generated By Our MoT

In this section, we show more examples generated by our MoT with different thought paradigms.

#### G.1 NATURAL LANGUAGE BASED REASONING

1100 1101

11021103

1104 1105

```
1107
1108
1109
1110
            The following is the problem you need to solve.
1111
             premises>
1112
            All rabbits that can be spotted near the campus are cute.
1113
            Some turtles can be spotted near the campus.
1114
            The only animals that can be spotted near the campus are
1115
                rabbits and squirrels.
            If something is skittish, then it is not calm.
1116
            All the squirrels that can be spotted near the campus are
1117
                skittish.
1118
            Rockie can be spotted near the campus, and it is calm.
1119
            </premises>
1120
            <conclusion>
1121
            Rockie is a turtle or cute.
1122
            </conclusion>
1123
1124
            <question>
1125
            Is the following statement true, false, or uncertain? Rockie is
                 a turtle or cute.
1126
             </question>
1127
1128
            <options>
1129
             (A) True
1130
             (B) False
             (C) Uncertain
1131
             </options>
1132
1133
```

```
1134
1135
            <nl_cot>
            Step 1: We know from the premise "The only animals that can be
                spotted near the campus are rabbits and squirrels." that
1137
                Rockie, being spotted near the campus, must be either a
1138
                rabbit or a squirrel.
1139
            Step 2: We also know from the premise "All rabbits that can be
1140
                spotted near the campus are cute" that if Rockie is a
                rabbit, then Rockie is cute.
1141
            Step 3: The premise "All the squirrels that can be spotted near
1142
                 the campus are skittish" and "If something is skittish,
1143
                then it is not calm" tell us that if Rockie is a squirrel,
1144
                then Rockie is not calm.
1145
            Step 4: The premise "Rockie can be spotted near the campus, and
                 it is calm" contradicts the conclusion from Step 3. This
1146
                means Rockie cannot be a squirrel.
1147
            Step 5: Therefore, Rockie must be a rabbit.
1148
            Step 6: Since Rockie is a rabbit and all rabbits spotted near
1149
                the campus are cute, we can conclude that Rockie is cute.
1150
1151
            <end_of_nl_cot>
            <answer>
1152
            The final answer is (A).
1153
            <end_of_answer>
1154
1155
```

#### G.2 CODE BASED REASONING

```
1158
1159
1160
1161
            The following is the problem you need to solve.
1162
1163
            <premises>
1164
            Everyone at the mixer is a Grand Slam champion or an Oscar-
                nominated actor.
1165
            Every Grand Slam champion at the mixer is a professional tennis
                 player.
1167
            All Oscar-nominated actors at the mixer are celebrities.
1168
            All professional tennis players at the mixer are athletes.
1169
            If a person at the mixer is a celebrity, then they are well
                paid.
1170
            If a person at the mixer is an athlete, then they are famous.
1171
            All well-paid people at the mixer live in tax havens.
1172
            Djokovic is at the mixer: if Djokovic is a famous athlete, then
1173
                 Djokovic is well-paid.
1174
            </premises>
1175
            <conclusion>
1176
            Djokovic does not live in a tax haven.
1177
            </conclusion>
1178
1179
             <question>
             Is the following statement true, false, or uncertain? Djokovic
1180
                does not live in a tax haven.
1181
             </question>
1182
1183
             <options>
1184
             (A) True
             (B) False
1185
             (C) Uncertain
1186
             </options>
1187
```

```
1188
1189
1190
             <code>
            from typing import Optional
1191
1192
            class Person:
1193
                def __init
                            (self,
1194
                          is_grand_slam_champion: Optional[bool] = None,
                          is_oscar_nominated_actor: Optional[bool] = None,
1195
                          is_professional_tennis_player: Optional[bool] =
1196
1197
                          is_celebrity: Optional[bool] = None,
1198
                          is_athlete: Optional[bool] = None,
1199
                          is_well_paid: Optional[bool] = None,
                          lives_in_tax_haven: Optional[bool] = None,
1200
                          is_famous: Optional[bool] = None,
1201
                          name: str = "Djokovic"):
1202
                   self.is_grand_slam_champion = is_grand_slam_champion
1203
                   self.is_oscar_nominated_actor = is_oscar_nominated_actor
1204
                   self.is_professional_tennis_player =
1205
                       is_professional_tennis_player
                   self.is_celebrity = is_celebrity
1206
                   self.is_athlete = is_athlete
1207
                   self.is_well_paid = is_well_paid
1208
                   self.lives_in_tax_haven = lives_in_tax_haven
1209
                   self.is_famous = is_famous
1210
            def apply_premises(person: Person) -> bool:
1211
                changed = False
1212
1213
                # Premise 1
1214
                if person.is_grand_slam_champion is None and person.
                    is_oscar_nominated_actor is None:
1215
                   person.is_grand_slam_champion = True
1216
                   changed = True
1217
1218
                # Premise 2
1219
                if person.is_grand_slam_champion is True and person.
                   is_professional_tennis_player is None:
1220
                   person.is_professional_tennis_player = True
1221
                   changed = True
1222
1223
                # Premise 3
1224
                if person.is_oscar_nominated_actor is True and person.
                   is_celebrity is None:
1225
                   person.is_celebrity = True
1226
                   changed = True
1227
1228
                # Premise 4
1229
                if person.is_professional_tennis_player is True and person.
                    is_athlete is None:
1230
                   person.is_athlete = True
1231
                   changed = True
1232
1233
                # Premise 5
1234
                if person.is_celebrity is True and person.is_well_paid is
1235
                   None:
                   person.is_well_paid = True
1236
                   changed = True
1237
1238
                # Premise 6
1239
                if person.is_athlete is True and person.is_famous is None:
                   person.is_famous = True
1240
1241
```

```
1242
1243
                   changed = True
1244
                # Premise 7
1245
                if person.is_well_paid is True and person.lives_in_tax_haven
1246
                     is None:
1247
                   person.lives_in_tax_haven = True
1248
                   changed = True
1249
                # Premise 8
1250
                if person.name == "Djokovic" and person.is_famous is True
1251
                    and person.is_well_paid is None:
1252
                   person.is_well_paid = True
1253
                   changed = True
1254
                return changed
1255
1256
            def run_inference(person: Person):
1257
                while apply_premises(person):
1258
                   pass
1259
            def check_conclusion(person: Person) -> str:
1260
                run_inference(person)
1261
                if person.name == "Djokovic" and person.lives_in_tax_haven
1262
                    is True:
                   return "False"
1263
                else:
1264
                   return "True"
1265
1266
            def func():
1267
               person = Person(
1268
                   is_grand_slam_champion=None,
                   is_oscar_nominated_actor=None,
1269
                   is_professional_tennis_player=None,
1270
                   is_celebrity=None,
1271
                   is_athlete=None,
1272
                   is_well_paid=None,
1273
                   lives_in_tax_haven=None,
                   is_famous=None,
1274
                   name="Djokovic"
1275
1276
                return check_conclusion(person)
1277
            if __name__ == '__main__':
1278
               result = func()
1279
                print ("Conclusion: Djokovic does not live in a tax haven?",
1280
                    result)
1281
            <end_of_code>
1282
             <answer>
1283
            The final answer is (B).
             <end_of_answer>
1284
1285
1286
```

# G.3 TRUTH TABLE BASED REASONING

1287 1288

1293 1294

```
1296
1297
             People in this club who perform in school talent shows often
1298
                 attend and are very engaged with school events.
             People in this club either perform in school talent shows often
1299
                  or are inactive and disinterested community members.
1300
             People in this club who chaperone high school dances are not
1301
                  students who attend the school.
1302
             All people in this club who are inactive and disinterested
                 members of their community chaperone high school dances.
1303
             All young children and teenagers in this club who wish to
1304
                  further their academic careers and educational
1305
                 opportunities are students who attend the school.
1306
             Bonnie is in this club and she either both attends and is very
1307
                 engaged with school events and is a student who attends the
                   school or is not someone who both attends and is very
1308
                  engaged with school events and is not a student who attends
1309
                   the school.
1310
             </premises>
1311
1312
             <conclusion>
1313
             Bonnie performs in school talent shows often.
             </conclusion>
1314
1315
              <auestion>
1316
             Is the following statement true, false, or uncertain? Bonnie
1317
                 performs in school talent shows often.
              </question>
1318
1319
             <options>
1320
              (A) True
1321
              (B) False
1322
              (C) Uncertain
              </options>
1323
1324
             <truth_table>
1325
             Let:
1326
                 T: Bonnie performs in school talent shows often.
1327
                 E: Bonnie attends and is very engaged with school events.
                 S: Bonnie is a student who attends the school.
                 I: Bonnie is an inactive and disinterested community member.
1329
                 C: Bonnie chaperones high school dances.
1330
1331
             Premises:
1332
             1. E \rightarrow S
             2. T VI
1333
             3. C \rightarrow \neg S
1334
             4. I \rightarrowC
1335
             5. S \rightarrowE
1336
             6. (E \wedgeS) \vee (\negE \wedge\negS)
1337
             Logical chain (assuming T = False):
1338
                 From (2): T = False \rightarrow I = True
1339
                 From (4): I = True \rightarrow C = True
1340
                 From (3): C = True \rightarrow \neg S = True \rightarrow S = False
1341
                 From (6): S = False \rightarrow True \rightarrow E = False
1342
                 From (1): E = False \rightarrow E \rightarrow S holds vacuously
                 -> No contradiction if T = False
1343
1344
             Logical chain (assuming T = True):
1345
                 From (2): T = True \rightarrow I = False
1346
                 From (4): I = False \rightarrow C = False
1347
                From (3): C = False \rightarrow \neg S = False \rightarrow S = True
                 From (5): S = True \rightarrow E = True
1348
1349
```

```
1350
1351
                From (6): E = True \land S = True \rightarrow (E \land S) = True
                -> No contradiction if T = True
1353
             Truth Table (remaining consistent rows):
1354
             | T | E | S | I | C |
1355
                                      -- | ----- | ----- |
                  ---- | ----- | ---
1356
             | True | True | True | False | False |
             | False | False | True | True |
1357
1358
             <end_of_truth_table>
1359
             <answer>
1360
             The final answer is (C).
1361
             <end_of_answer>
1362
```

#### 1363 1364 1365 1366 The following is the problem you need to solve. 1368 <premises> 1369 All fir trees are evergreens. 1370 Some objects of worship are fir trees. </premises> 1371 1372 <conclusion> 1373 Some evergreens are not objects of worship. 1374 </conclusion> 1375 <question> 1376 Is the following statement true, false, or uncertain? Some 1377 evergreens are not objects of worship. 1378 </question> 1379 1380 <options> (A) True 1381 (B) False (C) Uncertain 1383 </options> 1384 1385 <truth\_table> Let: 1386 - F: An object is a fir tree. 1387 - E: An object is an evergreen. 1388 - W: An object is an object of worship. 1389 1390 Premises: 1. All fir trees are evergreens. 1391 $F \rightarrow E$ 1392 2. Some objects of worship are fir trees. 1393 $\exists x (W(x) \land F(x))$ 1394 1395 Conclusion to Evaluate: Some evergreens are not objects of worship. 1396 $\exists x (E(x) \land \neg W(x))$ 1397 1398 Logical Chain: 1399 1. From (1), $F \rightarrow E$ . 1400 2. From (2), there exists some object a such that $W(a) \wedge F(a)$ . 1401 3. Since F(a) and $F \rightarrow E$ , it follows that E(a). 4. Therefore, we have E(a) and W(a). 1402 5. We need $\exists x (E(x) \land \neg W(x))$ . 1403

```
1404
1405
             6. The premises do not provide any information about whether
                there exists an evergreen that is not an object of worship.
             7. Therefore, it is uncertain whether some evergreens are not
1407
                objects of worship.
1408
1409
             Truth Table (Key Possibilities):
1410
             | F | E | W | E \\¬W |
1411
             | True | True | True | False |
1412
             | True | True | False| True |
1413
             | False | True | False |
1414
             | False| True | False| True |
1415
             Rows omitted: any with E= False (irrelevant to E \land \neg W) or with F
1416
                = True and E= False (violates F \rightarrow E).
1417
1418
             <end_of_truth_table>
1419
             <answer>
1420
            The final answer is (C).
1421
             <end_of_answer>
1422
1423
```

# H ERROR ANALYSIS

1424

1425 1426

1427

# H.1 EXAMPLES OF ERROR TYPE

```
1428
1429
1430
1431
1432
            <premises>
1433
            All employees who schedule a meeting with their customers will
1434
                go to the company building today.
            Everyone who has lunch in the company building schedules
1435
                meetings with their customers.
1436
            Employees have lunch either in the company building or at home.
1437
            If an employee has lunch at home, they are working remotely
1438
                from home.
1439
            All employees who are in other countries work remotely from
                home.
1440
            No managers work remotely from home.
1441
            James will appear in the company today if and only if he is a
1442
                manager.
1443
            </premises>
1444
            <conclusion>
1445
            James has lunch in the company.
1446
            </conclusion>
1447
1448
            <question>
1449
            Is the following statement true, false, or uncertain? James has
                 lunch in the company.
1450
            </question>
1451
1452
            <options>
1453
             (A) True
1454
             (B) False
1455
            (C) Uncertain
            </options>
1456
1457
```

```
1458
1459
            <nl_cot>
1460
            Step 1: From the premise \"All employees who schedule a meeting
                 with their customers will go to the company building today
1461
                 ,\" we know that if an employee schedules a meeting with
1462
                their customers, they will be in the company building today.
1463
1464
            Step 2: From the premise \"Everyone who has lunch in the
1465
                company building schedules meetings with their customers, \"
1466
                 we can infer that if an employee has lunch in the company
1467
                building, they must schedule a meeting with their customers.
1468
1469
            Step 3: From the premise \"Employees have lunch either in the
1470
                company building or at home, \" we know that an employee can
1471
                 only choose one of these two options for lunch.
1472
1473
            Step 4: From the premise \"If an employee has lunch at home,
1474
                they are working remotely from home, \" we can deduce that
1475
                having lunch at home means they are not in the company
                building and are working remotely.
1476
1477
            Step 5: From the premise \"All employees who are in other
1478
                countries work remotely from home, \" we know that employees
1479
                 in other countries do not have lunch in the company
                building.
1480
1481
            Step 6: From the premise \"No managers work remotely from home
1482
                ,\" we know that managers do not have lunch at home and
1483
                thus do not work remotely.
1484
            Step 7: From the premise \"James will appear in the company
1485
                today if and only if he is a manager, \" we can infer that
1486
                if James is not a manager, he will not appear in the
1487
                company today.
1488
1489
            Step 8: Combining the information from steps 1-7, we can
                conclude that if James is a manager, he will be in the
1490
                company today. Since James will only be in the company
1491
                today if and only if he is a manager, and managers do not
1492
                work remotely, James cannot be having lunch at home.
1493
            \RightarrowIncorrect. It only considers the case where James is a
1494
                manger, missing another branch: James is not a manger.
1495
            Step 9: Therefore, if James is not having lunch at home, he
1496
                must be having lunch in the company building, and by step 2,
1497
                 he must be scheduling meetings with his customers, which
1498
                means he will be in the company building today.
1499
            <end_of_nl_cot>
1500
            <answer>
1501
            The final answer is (A).
1502
            <end_of_answer>
1503
1504
            <Label>
            (C) Uncertain
1505
            </Label>
1506
1507
1508
```

1512 1513 1514 1515 1516 premises> People in this club who perform in school talent shows often 1517 attend and are very engaged with school events. 1518 People in this club either perform in school talent shows often 1519 or are inactive and disinterested community members. 1520 People in this club who chaperone high school dances are not 1521 students who attend the school. All people in this club who are inactive and disinterested 1522 members of their community chaperone high school dances. 1523 All young children and teenagers in this club who wish to 1524 further their academic careers and educational 1525 opportunities are students who attend the school. Bonnie is in this club and she either both attends and is very 1526 engaged with school events and is a student who attends the school or is not someone who both attends and is very 1528 engaged with school events and is not a student who attends 1529 the school. 1530 </premises> 1531 <conclusion> 1532 Bonnie performs in school talent shows often. 1533 </conclusion> 1534 1535 <question> Is the following statement true, false, or uncertain? Bonnie 1536 performs in school talent shows often. 1537 </question> 1538 1539 <options> 1540 (A) True (B) False 1541 (C) Uncertain 1542 </options> 1543 1544 <nl cot> 1545 Step 1: Analyze the given premises and the statement about Bonnie. 1546 - Premise 1: People in this club who perform in school 1547 talent shows often attend and are very engaged with 1548 school events. 1549 - Premise 2: People in this club either perform in school 1550 talent shows often or are inactive and disinterested community members. 1551 - Premise 3: People in this club who chaperone high school 1552 dances are not students who attend the school. 1553 - Premise 4: All people in this club who are inactive and 1554 disinterested members of their community chaperone high 1555 school dances. - Premise 5: All young children and teenagers in this club 1556 who wish to further their academic careers and 1557 educational opportunities are students who attend the 1558 school. 1559 - Bonnie's statement: Bonnie is in this club and she either 1560 both attends and is very engaged with school events and is a student who attends the school or is not someone 1561 who both attends and is very engaged with school events 1562 and is not a student who attends the school. 1563 Step 2: Determine Bonnie's status in the club. 1564

```
1566
1567
               - According to Bonnie's statement, she is either a very
1568
                   engaged and attending member or not such a member. This
                   statement is a disjunction (or) of two conditions, which
1569
                    means Bonnie can only fit into one of the two
1570
                   categories described.
1571
            Step 3: Consider the implications for Bonnie based on her
1572
                status in the club.
                - If Bonnie is a very engaged and attending member, she must
1573
                    perform in school talent shows often (Premise 1).
1574
               ⇒Incorrect (Invalid Converse). It incorrectly utilizes
1575
                   "Bonnie is a very engaged and attending member ->
1576
                   perform".
1577
                - If Bonnie is not a very engaged and attending member, she
                   must be an inactive and disinterested community member (
1578
                   Premise 2). Inactive and disinterested members must
1579
                   chaperone high school dances (Premise 4), which means
1580
                   they cannot be students who attend the school (Premise
1581
1582
               \RightarrowIncorrect (Missing Branches). It fails to enumerate all
1583
                   possible cases stated in Premise.
            Step 4: Evaluate the given statement.
1584
                 The statement claims that Bonnie performs in school talent
1585
                    shows often.
1586
               - Given the conditions and implications, if Bonnie is not a
1587
                   very engaged and attending member, she must be an
                   inactive and disinterested member, which means she would
1588
                    be a chaperone and not a student, contradicting Premise
1589
                    5. Therefore, Bonnie must be a very engaged and
1590
                   attending member.
1591
            Step 5: Conclude based on the analysis.
1592
                 Since Bonnie must be a very engaged and attending member,
                   she performs in school talent shows often.
1593
            <end_of_nl_cot>
1594
            <answer>
1595
            The final answer is (A).
1596
            <end_of_answer>
1597
1598
            <Label>
1599
            (C) Uncertain
            </Label>
1601
1602
```

#### Example 3 of Missing Branches

1608

1609

1610

1611

1612

1613

1614

1615

1616

1617

1620 1621 <conclusion> A design by Max is evocative and dreamy. </conclusion> 1623 1624 <question> 1625 Is the following statement true, false, or uncertain? A design 1626 by Max is evocative and dreamy. </question> 1627 1628 <options> 1629 (A) True 1630 (B) False 1631 (C) Uncertain 1632 </options> 1633 <nl\_cot> 1634 Step 1: Analyze the premise \"All of Zaha Hadid's design styles 1635 that Max adores have interesting geometries. \" This means 1636 if a design by Max adores is a Zaha Hadid design, then it 1637 has interesting geometries. 1638 Step 2: Analyze the premise \"No brutalist buildings that Max adores have interesting geometries. \" This means if a 1640 design by Max adores is a brutalist building, it does not 1641 have interesting geometries. 1642 Step 3: Analyze the premise \"Every style that Max adores is 1643 either Zaha Hadid's design style or Kelly Wearstler's 1644 design style.\" This means all designs by Max adores must 1645 be either Zaha Hadid's or Kelly Wearstler's. 1646 Step 4: Analyze the premise \"All of Kelly Wearstler's design 1647 styles that Max adores are evocative.\" This means if a 1648 design by Max adores is a Kelly Wearstler design, then it 1649 is evocative. 1650 1651 Step 5: Analyze the premise \"All of Kelly Wearstler's design styles that Max adores are dreamy. \" This means if a design 1652 by Max adores is a Kelly Wearstler design, then it is 1653 dreamy. 1654 1655 Step 6: Analyze the premise \"If a design by Max that he adores has interesting geometries, then the design is a brutalist 1656 building and evocative. \" This means if a design by Max 1657 adores has interesting geometries, then it is both a 1658 brutalist building and evocative. 1659 1660 Step 7: From Step 1 and Step 6, if a design by Max adores is a 1661 Zaha Hadid design and has interesting geometries, it must also be a brutalist building and evocative. 1662 1663 Step 8: From Step 2, if a design by Max adores is a brutalist 1664 building, it does not have interesting geometries. So, a 1665 design by Max adores cannot be both a brutalist building 1666 and have interesting geometries. 1667 Step 9: Since a design by Max adores cannot have both 1668 interesting geometries (which would make it a brutalist and evocative) and not have interesting geometries (which 1670 would make it not a brutalist), the design by Max adores 1671 cannot have interesting geometries. 1672

```
1674
1675
            \RightarrowIncorrect (Missing Branches). It fails to enumerate all
1676
                possible cases stated in Premise 3.
1677
            Step 10: Since the design by Max adores does not have
1678
                interesting geometries, according to Step 4 and Step 5, it
1679
                cannot be a Kelly Wearstler design. Therefore, it must be a
1680
                 Zaha Hadid design.
1681
            Step 11: Since the design by Max adores is a Zaha Hadid design,
1682
                 it must have interesting geometries (from Step 1), but
1683
                this contradicts the conclusion that it cannot have
1684
                interesting geometries (from Step 8).
1685
            Step 12: Given the premises, it is impossible for a design by
1686
                Max adores to be both evocative and dreamy, as it would
1687
                require having interesting geometries and not having them
1688
                at the same time.
1689
            <end_of_nl_cot>
1690
            <answer>
            The final answer is (B).
1691
            <end_of_answer>
1692
1694
1695
            <Label>
            (A) True
1696
             </Label>
1697
1698
```

## H.2 EXAMPLES UNIQUELY SOLVED BY TRUTH TABLE

```
1701
1702
1703
1704
1705
             <premises>
             "Stranger Things" is a popular Netflix show.
1706
             If a Netflix show is popular, Karen will binge-watch it.
1707
             If and only if Karen binge-watches a Netflix show, she will
1708
                 download it.
1709
             Karen does not download "Black Mirror."
             "Black Mirror" is a Netflix show.
1710
             If Karen binge-watches a Netflix show, she will share it with
1711
                 Lisa.
1712
             </premises>
1713
1714
             <conclusion>
             "Black Mirror" is popular.
1715
             </conclusion>
1716
1717
             <question>
1718
             Is the following statement true, false, or uncertain? "Black
1719
                 Mirror" is popular.
             </question>
1720
1721
             <options>
1722
             (A) True
1723
             (B) False
1724
             (C) Uncertain
1725
             </options>
1726
             <nl_cot>
1727
```

```
1728
1729
             Step 1: Given that \"Stranger Things\" is a popular Netflix
1730
                 show and if a Netflix show is popular, Karen will binge-
                 watch it, we know that Karen would binge-watch \"Stranger
1731
                 Things\" if it were popular. ⇒Correct.
1732
1733
             Step 2: However, the popularity of \TStranger Things\T does not
1734
                  provide any direct information about the popularity of \"
                 Black Mirror.\" ⇒Correct.
1735
1736
             Step 3: We know that \"Black Mirror\" is a Netflix show, but we
1737
                  do not have any information that directly links its
1738
                 popularity to Karen's actions. ⇒Correct.
1739
             Step 4: The statement that if Karen binge-watches a Netflix
1740
                 show, she will download it, and the fact that Karen does
1741
                 not download \"Black Mirror, \" imply that Karen does not
1742
                 binge-watch \"Black Mirror.\" ⇒Correct.
1743
1744
             Step 5: Since Karen does not binge-watch \"Black Mirror, \" and
1745
                 if a show is popular, Karen will binge-watch it, we cannot
                 conclude that \"Black Mirror\" is popular based on Karen's
1746
                 actions. \RightarrowIncorrect. Fail to utilize
1747
                 Transposition.
1748
1749
             Step 6: The conclusion \"Black Mirror\" is popular cannot be
                 definitively supported or refuted with the given premises.
1750
                 \RightarrowIncorrect.
1751
1752
             <end_of_nl_cot>
1753
             <answer>
1754
             The final answer is (C). \RightarrowIncorrect.
             <end_of_answer>"
1755
1756
             <truth_table>
1757
             Let:
1758
                P: \"Stranger Things\" is a popular Netflix show.
1759
                Q: \"Black Mirror\" is a popular Netflix show.
                B: Karen will binge-watch a Netflix show.
                D: Karen will download a Netflix show.
1761
                S: Karen will share a Netflix show with Lisa.
1762
             Premises:
1763
                P = True (Given) \Rightarrow Correct.
1764
                If P, then B (Given) \RightarrowCorrect.
                If B, then D (Given) \RightarrowCorrect.
1765
                If and only if B, then D (Given) \RightarrowCorrect.
1766
                \"Black Mirror\" is a Netflix show (Given) ⇒Correct.
1767
                If B, then S (Given) \RightarrowCorrect.
1768
             Conclusion to Evaluate:
1769
                Q (Is \"Black Mirror\" popular?) ⇒Correct.
             Logical Chain:
1770
                From (5), \"Black Mirror\" is a Netflix show. ⇒Correct.
1771
                From (2) and (3), if \Black Mirror\ is popular (Q), then
1772
                    Karen will binge-watch it (B). \RightarrowCorrect.
1773
                From (4), if Karen binge-watches \"Black Mirror, \" then she
1774
                    will download it (D). \RightarrowCorrect.
                However, from (6), we know that Karen does not download \"
1775
                    Black Mirror,\" so D = False. \RightarrowCorrect.
1776
                Since B implies D (from (3) and (4)), and D is False, B must
1777
                     also be False. \RightarrowCorrect.
1778
                Since B is False (Karen does not binge-watch \"Black Mirror
1779
                    \"), Q (popularity of \"Black Mirror\") cannot be True.
                    \RightarrowCorrect.
1780
1781
```

```
1782
1783
            Truth Table (Key Possibilities):
             | Q (Q: \"Black Mirror\" is popular) | B (B: Karen will binge-
                watch \"Black Mirror\") | D (D: Karen will download \"Black
1785
                 Mirror\") |
1786
             | False | False | False |
1787
             | True | True | True |
1788
            Given that D is False, Q must be False. \RightarrowCorrect.
            <end_of_truth_table>
1789
             <answer>
1790
            The final answer is (B). \RightarrowCorrect.
1791
            <end_of_answer>
1792
1793
             <code>
1794
            from typing import Optional
1795
1796
            class Show:
1797
                def __init__(self, name: str, is_popular: Optional[bool] =
1798
                    None, is_downloaded: Optional[bool] = None, is_binged:
1799
                    Optional[bool] = None):
                   self.name = name
1800
                   self.is_popular = is_popular
1801
                   self.is_downloaded = is_downloaded
1802
                   self.is_binged = is_binged
1803
            def apply_premises(show: Show) -> bool:
1804
                changed = False
1805
1806
                # Premise 1: "Stranger Things" is a popular Netflix show.
1807
                if show.name == "Stranger Things": ⇒Correct.
1808
                   if show.is_popular is None:
                      show.is_popular = True
1809
                      changed = True
1810
1811
                # Premise 2: If a Netflix show is popular, Karen will binge-
1812
                    watch it.
1813
                if show.is_popular is True and show.name != "Black Mirror":
1814
                    ⇒Correct.
                   if show.is_binged is None:
1815
                      show.is_binged = True
1816
                      changed = True
1817
1818
                # Premise 3: If and only if Karen binge-watches a Netflix
                    show, she will download it.
1819
                if show.is_binged is True: ⇒Correct.
1820
                   if show.is_downloaded is None:
1821
                      show.is_downloaded = True
1822
                      changed = True
1823
                # Premise 4: "Black Mirror" is a Netflix show. ⇒Incorrect.
1824
                   1. Do not encode Transposition. 2. It directly set
1825
                    show.is_popular = False., so even if it set it correctly, as
1826
                    it do not include Transposition, code approach still
1827
                    cannot work correctly.
1828
                if show.name == "Black Mirror":
                   if show.is_popular is None:
1829
                      show.is_popular = False
1830
                      changed = True
1831
1832
                # Premise 5: If Karen binge-watches a Netflix show, she will
1833
                     share it with Lisa.
                if show.is_binged is True: ⇒Correct.
1834
1835
```

```
1836
1837
                   if show.is_downloaded is True:
1838
                       if show.name == "Black Mirror":
                          if show.is_popular is False:
1839
                             show.is_popular = True
1840
                             changed = True
1841
1842
                return changed
1843
             def run_inference(show: Show):
1844
                while apply_premises(show):
1845
                   pass
1846
1847
             def check_conclusion(show: Show) -> str:
                run_inference(show)
1848
                if show.name == "Black Mirror":
1849
                   if show.is_popular is True:
1850
                      return "True"
1851
                   elif show.is_popular is False:
1852
                      return "False"
1853
                   else:
                      return "Uncertain"
1854
                else:
1855
                   return "Uncertain"
1856
1857
             def func():
                show = Show(name="Black Mirror", is_popular=None,
1858
                    is_downloaded=None, is_binged=None)
1859
                return check_conclusion(show)
1860
1861
            if __name__ == '__main__':
1862
                result = func()
                print("Conclusion: 'Black Mirror' is popular?", result)
1863
1864
             <end_of_code>
1865
             <answer>
1866
            The final answer is (C). \RightarrowIncorrect.
1867
             <end_of_answer>
1868
1869
```

#### Examples uniquely solved by Truth Table

```
1872
1873
            <premises>
1874
            All of Peter's pets that are mammals are also animals.
1875
            All monkeys are mammals.
1876
            Peter's pets are all either monkeys or birds.
            Peter's birds can fly.
1877
            All animals that are Peter's pets can breathe.
1878
            If Peter's pet can fly, then it has wings.
1879
            Rock is Peter's pet.
1880
            Rock can fly, or Rock is a bird, or Rock cannot breathe.
1881
            </premises>
1882
            <conclusion>
1883
            Rock is a monkey.
1884
            </conclusion>
1885
1886
            <question>
            Is the following statement true, false, or uncertain? Rock is a
1887
                 monkey.
1888
            </question>
1889
```

```
1890
1891
1892
            <options>
             (A) True
1893
             (B) False
1894
             (C) Uncertain
1895
             </options>
1896
            <nl cot>
1897
            Step 1: From the premise \"All of Peter's pets that are mammals
1898
                 are also animals, \" we know that any mammal pet of Peter's
1899
                 is an animal.
1900
1901
            Step 2: From the premise \"All monkeys are mammals,\" we know
                that if something is a monkey, it is also a mammal.
1902
1903
            Step 3: From the premise \"Peter's pets are all either monkeys
1904
                or birds, \" we know that Rock, being a pet of Peter's, must
1905
                 be either a monkey or a bird.
1906
            Step 4: From the premise \"Peter's birds can fly,\" we know
1907
                that if something is a bird, it can fly.
1908
1909
            Step 5: From the premise \"All animals that are Peter's pets
1910
                can breathe, \" we know that any pet of Peter's, whether a
1911
                bird or a mammal, can breathe.
1912
            Step 6: From the premise \"If Peter's pet can fly, then it has
1913
                wings,\" we know that if a pet of Peter's can fly, it must
1914
                have wings.
1915
1916
            Step 7: From the premise \mbox{"Rock can fly, or Rock is a bird, or}
                Rock cannot breathe, \" we have a disjunction. We need to
1917
                evaluate each part of this disjunction:\n- If Rock can fly,
1918
                 then it must have wings (Step 6), and it must be a bird (
1919
                Step 4), because only birds can fly among Peter's pets.\n-
1920
                If Rock is a bird, it can fly (Step 4), and it can breathe
1921
                (Step 5).\n- If Rock cannot breathe, this contradicts the
                fact that all of Peter's pets can breathe (Step 5).
1922
                \RightarrowInCorrect. If Rock can fly, then it must have wings
1923
                (Step 6), and it must be a bird (Step 4): this step
1924
                incorrectly utilizes rule fly -> birds to derive rock must
1925
                be a bird.
1926
            Step 8: Given the disjunction in Step 7, the only consistent
1927
                option is that Rock is a bird, because if Rock can fly, it
1928
                must be a bird, and if Rock cannot breathe, it contradicts
1929
                the premises. Therefore, Rock is a bird.\n\nStep 9: Since
1930
                Rock is a bird, and all of Peter's pets are either monkeys
1931
                or birds (Step 3), Rock cannot be a monkey.
1932
            <end_of_nl_cot>
1933
            <answer>
1934
            The final answer is (B).
1935
            <end_of_answer>
1936
1937
1938
1939
            <truth_table>
1940
            Let:
1941
               M: Rock is a monkey
               B: Rock is a bird
1942
1943
```

```
1944
1945
               F: Rock can fly
                A: Rock is an animal
                W: Rock has wings
1947
                R: Rock can breathe
1948
1949
            Premises:
1950
                1. All of Peter's pets that are mammals are also animals. (
                   Not directly about Rock)
1951
                2. All monkeys are mammals. (Not directly about Rock)
1952
                3. Peter's pets are all either monkeys or birds. (Not
1953
                   directly about Rock)
1954
                4. Peter's birds can fly. (B \rightarrowF)
1955
                5. All animals that are Peter's pets can breathe. (A \rightarrowR)
                6. If Peter's pet can fly, then it has wings. (F \rightarrowW)
1956
                7. Rock is Peter's pet.
1957
                8. Rock can fly, or Rock is a bird, or Rock cannot breathe.
1958
                    (F VB V¬R)
1959
1960
            Conclusion to Evaluate:
1961
                Rock is a monkey (M)
1962
            \Rightarrow Correct. The model correctly produce two cases that lead to
1963
                different conclusions, i.e., True and False. Therefore, it
1964
                 finally produces Uncertain as the final answer.
1965
            Truth Table (Key Possibilities):
1966
1967
             | M (Monkey) | B (Bird) | F (Can Fly) | A (Animal) | W (Wings)
1968
                | R (Can Breathe) |
1969
             | True | False | True | True | True | True |
1970
             | False | True | True | True | True |
             <end_of_truth_table>
1971
             <answer>
1972
            The final answer is (C).
1973
            </end_of_answer>
1974
1975
            <code>
            from typing import Optional
1977
            class Pet:
1978
                def __init__(self,
1979
                          is_monkey: Optional[bool] = None,
1980
                          is_bird: Optional[bool] = None,
                          can_fly: Optional[bool] = None,
1981
                          has_wings: Optional[bool] = None,
1982
                          is_animal: Optional[bool] = None):
1983
                   self.is_monkey = is_monkey
1984
                   self.is_bird = is_bird
1985
                   self.can_fly = can_fly
                   self.has_wings = has_wings
1986
                   self.is_animal = is_animal
1987
1988
            def apply_premises(pet: Pet) -> bool:
1989
                changed = False
1990
                # Premise 1: All of Peter's pets that are mammals are also
1991
                   animals.
1992
                # Since all of Peter's pets are either monkeys or birds, and
1993
                    monkeys are mammals,
1994
                # this premise does not directly affect the state of the pet
1995
1996
1997
```

```
1998
1999
                # Premise 2: All monkeys are mammals.
                if pet.is_monkey is True and pet.is_animal is not True:
                   pet.is_animal = True
2001
                   changed = True
2002
                # Premise 3: Peter's pets are all either monkeys or birds.
2004
                if pet.is_monkey is None and pet.is_bird is None:
                   pet.is_monkey = False
                   pet.is_bird = False
2006
                   changed = True
2007
2008
                # Premise 4: Peter's birds can fly.
2009
                if pet.is_bird is True and pet.can_fly is not True:
                   pet.can_fly = True
2010
                   changed = True
2011
2012
                # Premise 5: All animals that are Peter's pets can breathe.
2013
               if pet.is_animal is True and pet.can_fly is not True:
2014
                   pet.can_fly = True
2015
                   changed = True
2016
                # Premise 6: If Peter's pet can fly, then it has wings.
2017
               if pet.can_fly is True and pet.has_wings is not True:
2018
                   pet.has_wings = True
2019
                   changed = True
2020
                # Premise 7: Rock is Peter's pet.
2021
                # This does not change the state of the pet but indicates
2022
                   that the pet is Rock.
2023
2024
                # Premise 8: Rock can fly, or Rock is a bird, or Rock cannot
                     breathe.
2025
                if pet.name == "Rock":
2026
                   if pet.can_fly is None:
2027
                      pet.can_fly = True
2028
                      changed = True
2029
                   if pet.is_bird is None:
                      pet.is_bird = True
2030
                      changed = True
2031
                   if pet.can_fly is True and pet.has_wings is not True:
2032
                      pet.has_wings = True
2033
                      changed = True
2034
               ⇒ Incorrect. The model struggles to correctly interpret
2035
                   premises that involve multiple possible branches of
                    reasoning, such as disjunctions (V). It often commits
2036
                   prematurely to one or more branches instead of
2037
                   preserving logical uncertainty.
2038
2039
               return changed
2040
            def run_inference(pet: Pet):
2041
                while apply_premises(pet):
2042
                  pass
2043
2044
            def check_conclusion(pet: Pet) -> str:
2045
                run_inference(pet)
                if pet.name == "Rock" and pet.is_monkey is True:
2046
                   return "True"
2047
                else:
2048
                   return "False"
2049
            def func():
2050
2051
```

```
2052
2053
                pet = Pet(
2054
                   is_monkey=None,
                   is_bird=None,
2055
                   can_fly=None,
2056
                   has_wings=None,
2057
                   is_animal=None
2058
                pet.name = "Rock"
2059
                return check_conclusion(pet)
2060
2061
             if __name__ == '__main__':
2062
                result = func()
2063
                print("Conclusion: Rock is a monkey?", result)
2064
             <end_of_code>
             <answer>
2065
             The final answer is (B).
2066
             </end_of_answer>
2067
2068
2069
```