LTRAG: Enhancing autoformalization and self-refinement for logical reasoning with Thought-Guided RAG

Anonymous ACL submission

Abstract

002 Logical reasoning is fundamental to intelligent systems. Large language models (LLMs) have demonstrated promise in natural language (NL) reasoning, especially with techniques like chain-of-thought (CoT) prompting. Neuro-symbolic methods like Logic-LM and LINC further enhance performance on challenging datasets FOLIO and AR-LSAT by integrating formalization with LLMs and symbolic solvers, and possibly refinement with LLMs. However, these methods still struggle with the accurate formalization of com-013 plex NL problems. In this paper, we introduce LTRAG, a framework to enhance autoformalization and self-refinement for logical reasoning with Retrieval-Augmented Gen-017 eration (RAG), by building knowledge bases of thought-guided examples¹. Experimental results on FOLIO and AR-LSAT show that LTRAG consistently outperforms Logic-LM and LINC across different models. On GPT-4 and AR-LSAT, it achieves an accuracy gain of 13% over Logic-LM.

1 Introduction

007

039

Logical reasoning is a fundamental aspect of human intelligence and essential for complex tasks such as problem-solving and decision-making. Recently, logical reasoning over natural language (NL) exploiting large language models (LLMs) has received much attention. Many datasets have been proposed, including synthetic ProofWriter for rule reasoning (Tafjord et al., 2021), humancrafted FOLIO for complex first-order logic (FOL) reasoning (Han et al., 2024), and AR-LSAT extracted from the LSAT exams (Zhong et al., 2021). Various methods have been explored, including prompting (Wei et al., 2022; Kojima et al., 2022), fine-tuning (Zelikman et al., 2022), and neuro-

> ¹https://anonymous.4open.science/r/ dataset-example-LTRAG/

symbolic methods based on search (Kazemi et al., 2023; Hao et al., 2023).

040

041

042

043

045

047

051

054

057

060

061

062

063

064

065

066

067

068

069

071

073

074

075

076

078

A recent endeavor for logical reasoning over NL is neuro-symbolic methods based on autoformalization, by combining translation with LLMs from NLs to formal languages and rigorous reasoning of symbolic solvers. Typical works are Logic-LM (Pan et al., 2023) which introduces self-refinement to use the symbolic solver's error messages to refine the formalization, and LINC (Olausson et al., 2023). On GPT-4, LINC achieves an accuracy of 98% on ProofWriter, and Logic-LM reaches 79% on FOLIO. However, existing methods still face significant challenges, since many difficult NL problems are hard to formalize precisely in an automatic fashion. For example, with GPT-4, on AR-LSAT, Logic-LM only achieves an accuracy of 43%. In particular, these methods rely on a fixed set of examples for autoformalization and self-refinement, thus struggling with diverse and complex inputs, limiting their generalizability. On the other hand, Retrieval-Augmented Generation (RAG) (Lewis et al., 2020) enhances generation by dynamically retrieving relevant information from an external knowledge base (KB).

In this paper, we propose a framework LTRAG to enhance autoformalization and self-refinement for logical reasoning with thought-guided RAG, by building KBs of thought-guided examples for formalization and refinement. Experimental results on FOLIO and AR-LSAT show LTRAG consistently outperforms Logic-LM and LINC across different models. Especially, LTRAG achieves an accuracy gain of 8% over Logic-LM and LINC on GPT-3.5-turbo and FOLIO, and 13% on GPT-4 and AR-LSAT.

2 **Related Work**

Translating a NL into a formal language is challenging due to its fuzziness, ambiguity, and implicit information. Nguyen et al. (2022) proposed a method combining manually translating with deep learning. Yang et al. (2024) introduced LogicLLaMA, combining supervised fine-tuning and reinforcement learning with human feedback.
Chen et al. (2023) developed a framework using LLMs for translation between NL and temporal logic using intermediate languages.

079

080

081

087

100

101

102

104

105

106

107

108

109

110

111

112

113

114

115

116

117

118

119

120

121

122

123

Neuro-symbolic methods for logical reasoning based on autoformalization have recently gained much attention. Pan et al. (2023) proposed Logic-LM, including a Problem Formulator to formalize the problem, a Symbolic Reasoner, and a Self-Refiner module for error correction. Olausson et al. (2023) introduced LINC, which also combines autoformalization and symbolic solving, but uses majority voting to aggregate results from multiple formalizations. Ye et al. (2023b) proposed SATLM, which leverages LLMs to formalize NL inputs into satisfiability (SAT) problems and uses a SAT solver for reasoning. Jiang et al. (2024) introduced LeanReasoner, which fine-tunes with data in the Lean theorem-proving environment, formalizes problems into Lean theorems, and solves them with a tactic generator and proof search. Xu et al. (2024) proposed SymbCoT, which does autoformalization but reasons with CoT prompting based on both NL and FOL inputs.

Lewis et al. (2020) proposed RAG, using document retrieval to improve output precision with external knowledge. Fan et al. (2024) showed RAG reduces hallucinations and improves generation quality. Jiang et al. (2023) introduced FLARE, enabling efficient retrieval during generation.

Example selection is key to in-context learning. Liu et al. (2022) used a sentence encoder to select top-k similar examples for given problems, showing dynamic selection improves LLM performance. Levy et al. (2023) studied compositional generalization in semantic parsing, selecting diverse examples via coverage and diversity based methods. Ye et al. (2023a) proposes CEIL, an example selection method using Determinantal Point Processes and contrastive learning.

3 Framework

The structure of LTRAG is depicted in Figure 1. It
comprises four key components: a Retrieval Module, a Translator LLM, a Solver, and a Fixer LLM.
The Translator LLM (similar to Logic-LM's Problem Formulator) converts NL problems into for-

mal representations, while the Fixer LLM (akin to Logic-LM's Self-Refine) corrects translation errors. The Retrieval Module, built upon FastGPT², dynamically retrieves similar examples from the RAG KBs. These retrieved examples serve as guiding references for both Translator and Fixer LLMs, enhancing the accuracy of formalization and error correction. The detailed method for constructing the KBs can be found in Section 4.2. Once the problem is formalized, the Solver takes over to perform logical reasoning. The Solver is based on Microsoft's Z3 solver ³ (de Moura and Bjørner, 2008), and is capable of handling FOL expressions (in FOLIO) and constraint satisfaction problems (in AR-LSAT). If errors are detected, they are reported to the Fixer LLM for another formalization, and the above process is repeated.

129

130

131

132

133

134

135

136

137

138

139

140

141

142

143

144

145

146

147

148

149

151

152

153

154

155

156

157

158

159

161

162

163

164

165

166

167

168

169

170

171

172

173

174

175

176

Here is an example (full version in Appendix A.1): One of the input premises is "There are four seasons in a year: Spring, Summer, Fall, and Winter". The full problem is used to retrieve the translation KB, and both the problem and the retrieved examples are provided to the Translator LLM, resulting in an initial formalization:

 $\forall x \ (Season(x) \rightarrow (x = Spring \lor))$

 $x = Summer \lor x = Fall \lor x = Winter)$). The solver will return an error, indicating that using "=" is not allowed. The error is used to retrieve the fixer KB, and both the error and the retrieved examples are provided to the Fixer LLM, resulting in the final correct formalization:

 $\forall x \ (Season(x) \rightarrow (IsSpring(x) \lor IsSummer(x) \lor IsFall(x) \lor IsWinter(x))).$

4 **Experiments**

4.1 Experimental Setup

We evaluate LTRAG on FOLIO and AR-LSAT, comparing it against baselines such as Standard prompting, CoT prompting (using 2 examples on FOLIO and 1 on AR-LSAT), LINC, and Logic-LM, across models including GPT-40 (OpenAI, 2024), DeepSeek v2.5 (DeepSeek-AI, 2024), Llama3.3-70b (Grattafiori et al., 2024), GPT-3.5-turbo, and Gemma2-27b (Rivière et al., 2024).

On FOLIO, as our test set, we use the 182 verified samples retained by LINC after filtering out the problematic ones from the original 204 samples. On AR-LSAT, as Logic-LM, we use the 231 samples from the dev set as our test set.

²https://github.com/labring/FastGPT

³https://github.com/Z3Prover/z3



Figure 1: The framework of LTRAG.

Note that AR-LSAT samples are single-choice questions, and the solver may return multiple answers. When no or multiple answers are obtained after repeated repairs, Logic-LM returns an answer using the CoT method, and we do the same.

4.2 Knowledge Base Construction

177

178

179

182

183

184

185

190

191

192

194

195

196

197

199

205

207

210

212

Our RAG KBs are semi-automatically constructed.
Roughly, the translation KBs are automatically constructed from the training sets with programs or LLMs, while the fixer KBs are constructed by first manually fixing a small set of error cases and then using it to guide LLMs to generate the rest.
Detailed examples are in Appendix A.2. In the following, we give details of our KB construction.

For FOLIO, the translation KB is built from the training set, with the 122 samples where DeepSeek fails to obtain the correct answer using the Translator LLM with an empty KB and the solver. Each such sample is provided with the annotated formulas and translation steps obtained with a program as follows: first, extract all predicates and constants from the problem; second, identify relevant predicates and constants for each sentence; and finally, translate the sentence into FOL. The fixer KB is built from the training set, using samples where DeepSeek fails to obtain an executable formalization. A subset is manually corrected for high-quality references, while the rest are automatically repaired as follows: for each such sample, using the annotated formulas and the manual repair subset as context, let DeepSeek return the thought process for repairing, and retain the sample if the solver returns the correct answer. The fixer KB ends up with 57 samples.

> For AR-LSAT, the translation KB is built as follows: first, we pick 1,500 samples from the train

ing set; second, for each of these samples, we use DeepSeek to formalize it with the RAG KB of Logic-LM's formalizations of the dev set; finally, we retain 538 executable samples. As to the second step, we focus on defining the formalization's syntax and precautions in the prompts; due to the length of AR-LSAT samples, we choose not to let LLMs generate the thought process because long outputs confuse LLMs and easily generate unexecutable formalizations. The fixer KB is built from the training set, using samples where DeepSeek fails to obtain an executable formalization or a unique answer. For each error type, one or more samples are manually selected and corrected with an analysis of the error. The rest samples are semi-automatically repaired as follows: for each such sample, using the manual repair subset, let DeepSeek return the thought process for repairing the formulation, which is manually checked for correctness and decided for being retained or not. The fixer KB ends up with 41 samples.

213

214

215

216

217

218

219

221

222

223

224

225

226

227

228

229

230

231

232

233

234

236

237

238

240

241

242

243

244

245

246

247

248

4.3 Results

The experimental results are summarized in Table 1. To handle longer contexts, we use GPT-40 (128K), a variant of GPT-4 (8K) used in prior work, with minimal performance differences.

On FOLIO, LTRAG consistently outperforms other methods across different models. For GPT-40, it achieves 80.77%, surpassing Logic-LM's 78.92% and LINC's 72.50%. For GPT-3.5-turbo, LTRAG attains 70.88%, significantly outperforming Logic-LM (61.27%) and LINC (62.60%).

On AR-LSAT, LTRAG also consistently improves over the baseline methods. For GPT-40, it outperforms both Logic-LM and CoT by about 13%. LTRAG also enhances DeepSeek v2.5's per-

| Model | FOLIO (Accuracy %) | | | AR-LSAT (Accuracy %) | | | / %) | | |
|---------------|--------------------|----------|-------|----------------------|----------|-------|--------------|-------|----------|
| | LTRAG | Standard | CoT | LINC | Logic-LM | LTRAG | Standard | CoT | Logic-LM |
| GPT-40 | 80.77 | 73.63 | 78.02 | 72.50 | 78.92 | 56.71 | 40.26 | 43.72 | 43.04 |
| DeepSeek v2.5 | 78.57 | 74.73 | 76.37 | - | - | 68.40 | 51.52 | 64.50 | - |
| Llama3.3 | 78.57 | 72.53 | 71.43 | - | - | 59.31 | 40.26 | 39.83 | - |
| GPT-3.5-turbo | 70.88 | 56.59 | 59.34 | 62.60 | 61.27 | 26.84 | 24.24 | 19.48 | 26.41 |
| Gemma2 | 79.67 | 59.89 | 62.09 | - | - | 35.06 | 25.97 | 24.67 | - |

Table 1: Performance comparison on FOLIO and AR-LSAT. The data for Logic-LM and LINC comes from their papers, and '-' denotes that they did not experiment on the model. LINC did not experiment on AR-LSAT.

formance, achieving 68.40% compared to 51.52% under Standard prompting. However, LTRAG attains limited improvements on GPT-3.5-turbo, with the performance gain being less than 3% compared to Standard Prompting and Logic-LM.

249

251

257

260

261

263

264

266

267

269

273

274

275

276

In Table 2, we analyze the executable rate and execution accuracy of LTRAG on AR-LSAT, in comparison to Logic-LM. In terms of executable rate, on both GPT-40 and GPT-3.5-turbo, LTRAG outperforms Logic-LM by about 30%, indicating its superior ability to generate executable programs. It is easy to notice that our execution accuracy is lower compared to Logic-LM. A possible reason is that we get much more executable programs, making the error rate in execution increase.

| Model | Exe_rate | Exe_accuracy |
|-------------------|----------|--------------|
| GPT-40 | 69.26 | 50.00 |
| GPT-4(Logic-LM) | 39.8 | 58.8 |
| GPT-3.5-turbo | 54.11 | 19.20 |
| GPT-3.5(Logic-LM) | 21.8 | 60.3 |
| DeepSeek v2.5 | 71.00 | 44.51 |
| Llama3.3 | 66.67 | 52.60 |
| Gemma2 | 45.02 | 36.54 |

Table 2: Executable rate (Exe_rate) and Execution Accuracy (Exe_accuracy) on AR-LSAT.

4.4 Ablation Experiments

In ablation studies, we investigate the effect of the Fixer LLM with different numbers of examples on different models. On FOLIO, the Fixer LLM improves accuracy by 2–5% for large models, while the improvement is 3–6% for small models. On AR-LSAT, the Fixer LLM improves accuracy by 15–20% for large models, while the improvement is 2–8% for small models. Detailed results are included in Appendix A.3.

4.5 Discussion

We first analyze on FOLIO, why LTRAG achieves better improvements on small models than on

large models. We think the thought-guided examples for the Translator LLM notably benefit small models by mitigating their inherent limitations. Large models produce fewer errors, and small models have limited repair capabilities, leading to limited improvements by the Fixer LLM.

We then analyze on AR-LSAT, why LTRAG achieves better improvements on large models than on small models. AR-LSAT samples are primarily constraint satisfaction problems, having unique answers, making it difficult to provide a systematic translation approach. The samples also involve complex long-text constraints, thus limiting the number of reference examples given to the Fixer LLM. As a result, while both large and small models face difficulties, large models can make effective corrections with limited assistance, whereas small models cannot.

5 Conclusion

In this paper, we propose the LTRAG framework to enhance autoformalization and self-refinement for logical reasoning with thought-guided RAG. The translation KBs are automatically constructed, and the fixer KBs are semi-automatically constructed where a small set of error cases are manually fixed and used to guide LLMs to generate more repairing examples. Empirical results on the challenging datasets FOLIO and AR-LSAT demonstrate that our approach significantly improves refinement capabilities of large models and formalization capabilities of small models. An outstanding advantage of our work is to improve formalization with less computational resources than approaches based on fine-tuning. Future work will focus on more automated construction of thoughtguided RAG KBs, particularly for challenging datasets. A possible approach is to utilize models pre-trained with iterative reasoning strategies to generate the RAG examples.

302

303

304

305

306

307

308

309

310

311

312

313

314

315

277

278

6 Limitations

316

317

318

319

321

323

325

327

329

331

333

334

335

337

341

342

346

347

348

350

354

366

Below, we outline some of the key challenges and constraints associated with our framework:

Difficulty in Constructing Thought Processes for Certain Tasks While structured reasoning steps can be effectively constructed for datasets like FOLIO with short context, other tasks such as AR-LSAT present challenges. AR-LSAT problems often involve complex constraints and relationships that are harder to break down into a thought process. This makes it difficult to provide the same level of guidance for small models, limiting their performance improvements.

Limited Impact of the Fixer LLM on Tasks with Few Syntax Errors The Fixer LLM, which corrects errors flagged by the Solver, shows limited improvement on tasks where syntax errors are rare, such as FOLIO. This is particularly true for large models like GPT-40, DeepSeek, and Llama, which already produce fewer syntax errors due to their advanced reasoning capabilities. As a result, the Fixer LLM's contributions are marginal in such cases, and the primary benefits of LTRAG come from the structured formalization process. Conversely, the Fixer LLM proves more effective on complex tasks like AR-LSAT, where the error types are more varied. Large models, with their superior refinement capabilities, can leverage the Fixer LLM to achieve significant improvements. However, small models, that struggle with both autoformalization and self-refinement, gain less benefit from the Fixer LLM in these scenarios.

Limitations in the Fixer LLM on Semantic Errors The Fixer LLM is primarily designed to address surface-level syntax errors, such as incorrect predicate usage or invalid logical operators. It is not capable of resolving deeper semantic errors, where the logical formalization may be syntactically correct but semantically flawed. On FO-LIO, where the solver provides unique answers, it's hard to detect semantic errors, while AR-LSAT provides extra feedback when the solver returns multiple answers. This limitation highlights the need for more advanced mechanisms that can handle both syntactic and semantic errors.

Dependency on Knowledge Base Quality The performance of LTRAG heavily relies on the quality and comprehensiveness of the KBs. In cases where the KB lacks sufficient examples or contains inaccuracies, the system's abilities may be compromised.

References

- Yongchao Chen, Rujul Gandhi, Yang Zhang, and Chuchu Fan. 2023. NL2TL: transforming natural languages to temporal logics using large language models. In *Proceedings of the 2023 Conference* on Empirical Methods in Natural Language Processing, EMNLP 2023, Singapore, December 6-10, 2023, pages 15880–15903. Association for Computational Linguistics.
- Leonardo Mendonça de Moura and Nikolaj S. Bjørner. 2008. Z3: an efficient SMT solver. In Tools and Algorithms for the Construction and Analysis of Systems, 14th International Conference, TACAS 2008, Held as Part of the Joint European Conferences on Theory and Practice of Software, ETAPS 2008, Budapest, Hungary, March 29-April 6, 2008. Proceedings, volume 4963 of Lecture Notes in Computer Science, pages 337–340. Springer.
- DeepSeek-AI. 2024. Deepseek-v2: A strong, economical, and efficient mixture-of-experts language model. *Preprint*, arXiv:2405.04434.
- Wenqi Fan, Yujuan Ding, Liangbo Ning, Shijie Wang, Hengyun Li, Dawei Yin, Tat-Seng Chua, and Qing Li. 2024. A survey on rag meeting llms: Towards retrieval-augmented large language models. In Proceedings of the 30th ACM SIGKDD Conference on Knowledge Discovery and Data Mining, KDD '24, pages 6491–6501, New York, NY, USA. Association for Computing Machinery.
- Aaron Grattafiori, Abhimanyu Dubey, Abhinav Jauhri, Abhinav Pandey, Abhishek Kadian, Ahmad Al-Dahle, and et al. 2024. The llama 3 herd of models. *Preprint*, arXiv:2407.21783.
- 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 R. Fabbri, Wojciech Kryscinski, Semih Yavuz, Ye Liu, Xi Victoria Lin, Shafiq Joty, Yingbo Zhou, Caiming Xiong, Rex Ying, Arman Cohan, and Dragomir Radev. 2024. FOLIO: natural language reasoning with first-order logic. In Proceedings of the 2024 Conference on Empirical Methods in Natural Language Processing, EMNLP 2024, Miami, FL, USA, November 12-16, 2024, pages 22017-22031. Association for Computational Linguistics.
- Shibo Hao, Yi Gu, Haodi Ma, Joshua Jiahua Hong, Zhen Wang, Daisy Zhe Wang, and Zhiting Hu. 2023. Reasoning with language model is planning with world model. In *Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing, EMNLP 2023, Singapore, December 6-10, 2023*, pages 8154–8173. Association for Computational Linguistics.

367 368

369

370

371

373

374

375

376

377

378

379

381

382

383

387

388

389

390

391

392

393

394

395

396

398

399

400

401

402

403

404

405

406

407

408

409

410

411

412

413

414

415

416

417

418

419

420

421

422

537

538

539

Dongwei Jiang, Marcio Fonseca, and Shay B. Cohen. 2024. Leanreasoner: Boosting complex logical reasoning with lean. In Proceedings of the 2024 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies (Volume 1: Long Papers), NAACL 2024, Mexico City, Mexico, June 16-21, 2024, pages 7497–7510. Association for Computational Linguistics.

424

425

426

427 428

429

430

431

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

- Zhengbao Jiang, Frank F. Xu, Luyu Gao, Zhiqing Sun, Qian Liu, Jane Dwivedi-Yu, Yiming Yang, Jamie Callan, and Graham Neubig. 2023. Active retrieval augmented generation. In Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing, EMNLP 2023, Singapore, December 6-10, 2023, pages 7969–7992. Association for Computational Linguistics.
- Mehran Kazemi, Najoung Kim, Deepti Bhatia, Xin Xu, and Deepak Ramachandran. 2023. LAMBADA: backward chaining for automated reasoning in natural language. In Proceedings of the 61st Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers), ACL 2023, Toronto, Canada, July 9-14, 2023, pages 6547–6568. Association for Computational Linguistics.
- Takeshi Kojima, Shixiang Shane Gu, Machel Reid, Yutaka Matsuo, and Yusuke Iwasawa. 2022. Large language models are zero-shot reasoners. In Advances in Neural Information Processing Systems 35: Annual Conference on Neural Information Processing Systems 2022, NeurIPS 2022, New Orleans, LA, USA, November 28 - December 9, 2022.
- Itay Levy, Ben Bogin, and Jonathan Berant. 2023. Diverse demonstrations improve in-context compositional generalization. In *Proceedings of the 61st Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 1401– 1422, Toronto, Canada. Association for Computational Linguistics.
- Patrick S. H. Lewis, Ethan Perez, Aleksandra Piktus, Fabio Petroni, Vladimir Karpukhin, Naman Goyal, Heinrich Küttler, Mike Lewis, Wen-tau Yih, Tim Rocktäschel, Sebastian Riedel, and Douwe Kiela. 2020. Retrieval-augmented generation for knowledge-intensive NLP tasks. In Advances in Neural Information Processing Systems 33: Annual Conference on Neural Information Processing Systems 2020, NeurIPS 2020, December 6-12, 2020, virtual.
- Jiachang Liu, Dinghan Shen, Yizhe Zhang, Bill Dolan, Lawrence Carin, and Weizhu Chen. 2022. What makes good in-context examples for GPT-3? In Proceedings of Deep Learning Inside Out (DeeLIO 2022): The 3rd Workshop on Knowledge Extraction and Integration for Deep Learning Architectures, pages 100–114, Dublin, Ireland and Online. Association for Computational Linguistics.
 - Ha-Thanh Nguyen, Wachara Fungwacharakorn, Fumihito Nishino, and Ken Satoh. 2022. A multi-step

approach in translating natural language into logical formula. In Legal Knowledge and Information Systems - JURIX 2022: The Thirty-fifth Annual Conference, Saarbrücken, Germany, 14-16 December 2022, volume 362 of Frontiers in Artificial Intelligence and Applications, pages 103–112. IOS Press.

- Theo Olausson, Alex Gu, Benjamin Lipkin, Cedegao E. Zhang, Armando Solar-Lezama, Joshua B. Tenenbaum, and Roger Levy. 2023. LINC: A neurosymbolic approach for logical reasoning by combining language models with first-order logic provers. In *Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing, EMNLP* 2023, Singapore, December 6-10, 2023, pages 5153– 5176. Association for Computational Linguistics.
- OpenAI. 2024. Gpt-4o system card. *Preprint*, arXiv:2410.21276.
- Liangming Pan, Alon Albalak, Xinyi Wang, and William Wang. 2023. Logic-Im: Empowering large language models with symbolic solvers for faithful logical reasoning. In *Findings of the Association for Computational Linguistics: EMNLP 2023, Singapore, December 6-10, 2023*, pages 3806–3824. Association for Computational Linguistics.
- Morgane Rivière, Shreya Pathak, Pier Giuseppe Sessa, Cassidy Hardin, Surya Bhupatiraju, Léonard Hussenot, Thomas Mesnard, Bobak Shahriari, Alexandre Ramé, Johan Ferret, and et al. 2024. Gemma 2: Improving open language models at a practical size. *CoRR*, abs/2408.00118.
- Oyvind Tafjord, Bhavana Dalvi, and Peter Clark. 2021. Proofwriter: Generating implications, proofs, and abductive statements over natural language. In *Findings of the Association for Computational Linguistics: ACL/IJCNLP 2021, Online Event, August 1-6, 2021*, volume ACL/IJCNLP 2021 of *Findings of ACL*, pages 3621–3634. Association for Computational Linguistics.
- Jason Wei, Xuezhi Wang, Dale Schuurmans, Maarten Bosma, Brian Ichter, Fei Xia, Ed H. Chi, Quoc V. Le, and Denny Zhou. 2022. Chain-of-thought prompting elicits reasoning in large language models. In Advances in Neural Information Processing Systems 35: Annual Conference on Neural Information Processing Systems 2022, NeurIPS 2022, New Orleans, LA, USA, November 28 - December 9, 2022.
- Jundong Xu, Hao Fei, Liangming Pan, Qian Liu, Mong-Li Lee, and Wynne Hsu. 2024. Faithful logical reasoning via symbolic chain-of-thought. In Proceedings of the 62nd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers), ACL 2024, Bangkok, Thailand, August 11-16, 2024, pages 13326–13365. Association for Computational Linguistics.
- Yuan Yang, Siheng Xiong, Ali Payani, Ehsan Shareghi, and Faramarz Fekri. 2024. Harnessing the power of large language models for natural language to

first-order logic translation. In Proceedings of the 62nd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers), pages 6942-6959, Bangkok, Thailand. Association for Computational Linguistics.

540

541

543

545 546

547

551

552

553

554

555

556 557

559

560

561

566

- Jiacheng Ye, Zhiyong Wu, Jiangtao Feng, Tao Yu, and Lingpeng Kong. 2023a. Compositional exemplars for in-context learning. In International Conference on Machine Learning, ICML 2023, 23-29 July 2023, Honolulu, Hawaii, USA, volume 202 of Proceedings of Machine Learning Research, pages 39818-39833. PMLR.
- Xi Ye, Qiaochu Chen, Isil Dillig, and Greg Durrett. 2023b. Satlm: Satisfiability-aided language models using declarative prompting. Advances in Neural Information Processing Systems, 36:45548–45580.
- Eric Zelikman, Yuhuai Wu, Jesse Mu, and Noah D. Goodman. 2022. Star: Bootstrapping reasoning with reasoning. In Advances in Neural Information Processing Systems 35: Annual Conference on Neural Information Processing Systems 2022, NeurIPS 2022, New Orleans, LA, USA, November 28 - December 9, 2022.
- Wanjun Zhong, Siyuan Wang, Duyu Tang, Zenan Xu, Daya Guo, Jiahai Wang, Jian Yin, Ming Zhou, and Nan Duan. 2021. AR-LSAT: investigating analytical reasoning of text. CoRR, abs/2104.06598.

Appendix Α

An example of solving FOLIO problems A.1

567

568

598

599

600

615

616

| Premises | 569 570 |
|--|------------|
| 1.[BG] There are four seasons in a | 571 |
| → year: Spring, Summer, Fall, and | 572 |
| \hookrightarrow Winter. | 573 |
| 2.All students who want to have a long | 574 |
| \hookrightarrow vacation love summer the most. | 575 |
| 3.Emma's favorite season is summer. | 576 |
| 4.Mia's favorite season is not the same | 577 |
| ↔ as Emma's. | 578 |
| 5.James wants to have a long vacation. | 579 |
| | 580 |
| Conclusion: | 581 |
| James's favorite season is summer. | 583 |

Listing 1: Problem

| Premises: | 584 585 |
|---|------------|
| 1. $\forall x(Season(x) \rightarrow (x = Sprinq \lor$ | 586 |
| $x = Summer \lor x = Fall \lor x = Winter))$ | 587 |
| 2. $\forall x (WantsLongVacation(x) \rightarrow$ | 588 |
| FavoriteSeason(x, Summer)) | 589 |
| 3. FavoriteSeason(Emma, Summer) | 590 |
| 4. $\neg(FavoriteSeason(Mia, Emma))$ | 591 |
| FavoriteSeason(Emma, Mia)) | 592 |
| 5. WantsLongVacation(James) | 593 |
| | 594 |
| Conclusion: | 595 |
| FavoriteSeason(James, Summer) | 596 |

Listing 2: Translated Formulas

The solver will return an error, indicating that using "=" is not allowed on premise 1. Instead, predicates should be used to represent relationships.

| | 601 |
|---|-----|
| Premises: | 602 |
| 1. $\forall x(Season(x) \rightarrow (IsSpring(x)) \lor$ | 603 |
| $IsSummer(x) \lor IsFall(x) \lor IsWinter(x)))$ | 604 |
| 2. $\forall x (WantsLongVacation(x) \rightarrow$ | 605 |
| FavoriteSeason(x, Summer)) | 606 |
| 3. FavoriteSeason(Emma, Summer) | 607 |
| 4. $\neg(FavoriteSeason(Mia, Emma))$ | 608 |
| FavoriteSeason(Emma, Mia)) | 609 |
| 5. WantsLongVacation(James) | 610 |
| | 611 |
| Conclusion: | 612 |
| FavoriteSeason(James, Summer) | 613 |

Listing 3: Fixed Formulas

Examples from the RAG KBs A.2

A.2.1 FOLIO

| Input: | 617 618 |
|---|------------|
| Premises: | 619 |
| 1.All pets are animals. | 620 |
| 2.Pets can be either a dog or a cat. | 621 |
| 3.If a person has a pet, they care for | 622 |
| ↔ that pet. | 623 |
| 4.Dogs and Cats can be naughty. | 624 |
| 5.Pets who are naughty are not liked as | 625 |
| \hookrightarrow much. | 626 |

```
Conclusion:
              Leo is an animal.
              Output:
              We can get these predicates:
633
              Pet(x),Animal(x),Dog(x),Cat(x),HasPet(x),
              Cares(x,y),Naughty(x),Liked(x,y)
636
              We can get these constants:
637
              charlie,leo
              Let's translate one by one.
638
              Premises:
              1.Text:All pets are animals.
641
              Predicates:
              Animal(x), Pet(x)
642
643
              Fol: \forall x(Pet(x) \rightarrow Animal(x))
644
              2.Text:Pets can be either a dog or a
                  \hookrightarrow cat.
              Predicates:
646
647
              Cat(x), Pet(x), Dog(x)
              Fol: \forall x(Pet(x) \rightarrow Dog(x) \oplus Cat(x))
649
              3.Text: If a person has a pet, they care
                   \hookrightarrow for that pet.
650
651
              Predicates:
              Cares(x, y), HasPet(x)
              Fol: \forall x \forall y (HasPet(x) \rightarrow Cares(x, y))
653
              4. Text: Dogs and Cats can be naughty.
655
              Predicates:
              Cat(x), Naughty(x), Dog(x)
657
              Fol: \exists x ((Cat(x) \land Naughty(x))) \lor
658
              (Dog(x) \land Naughty(x)))
              5.Text:Pets who are naughty are not
                  \hookrightarrow liked as much.
              Predicates:
              Naughty(x), Pet(x), Liked(x,y)
              Fol: \forall x \forall y (Pet(x) \land Naughty(x) \rightarrow \neg Liked(x, y))
              6.Text:Charlie has a naughty pet dog

→ named Leo.

              Predicates:
667
              Pet(x), Naughty(x), HasPet(x), Dog(x)
              Constants:
              charlie,leo
670
              Fol: HasPet(charlie) \land Pet(leo)
              \land Dog(leo) \land Naughty(leo)
671
              Conclusion:
              Text:Leo is an animal.
673
674
              Predicates:
              Animal(x)
676
              Constants:
677
              leo
              Fol: Animal(leo)
679
              Listing 4: An Example from the Translation KB for
              FOLIO
              Input:
              Premises:
              1. All commercial aircraft are produced
                  \hookrightarrow by either Boeing or Airbus.
              origin: \forall x (Commercial Aircraft(x))
              \rightarrow (ProducedByBoeing(x) \lor ProducedByAirbus(x)))
              true answer: \forall x (CommercialAircraft(x))
              \rightarrow ProducedBy(x, boeing) \oplus ProducedBy(x, airbus))
              2. All American Airlines jets are

→ commercial aircraft.

              origin: \forall x (AmericanAirlinesJet(x))
```

6.Charlie has a naughty pet dog named

→ Leo.

```
\rightarrow CommercialAircraft(x))
```

```
true answer: \forall x (AmericanAirlinesAircraft(x))
```

```
\rightarrow CommercialAircraft(x))
3. Airbus made more in revenue than
     \hookrightarrow Boeing last year.
origin: \exists y(Revenue(airbus, y) \land
\exists z (Revenue(boeing, z) \land y > z))
true answer: MoreinRevenue(airbus, boeing)
Conclusion:
There does not exist a United Airlines
     \hookrightarrow plane produced by Boeing.
Error:
\exists y (Revenue(airbus, y) \land
\exists z(Revenue(boeing, z) \land y > z)) error :
     \hookrightarrow ("Illegal comparisons: Matched
\hookrightarrow 'y > z' from position 50 to 55.
     \hookrightarrow You can't use x > something, try
     \leftrightarrow to use predicate(x, something)
     ↔ to replace it.",) Invalid
↔ symbols in `=<>` at position 52
Output:
Premises:
1. Original:
    \forall x (Commercial Aircraft(x))
    \rightarrow (ProducedByBoeing(x))
    \lor ProducedByAirbus(x)))
    Explanation: The use of \lor (logical
         ↔ disjunction) is incorrect
         \hookrightarrow here because it implies that
         \hookrightarrow an aircraft can be produced
         \hookrightarrow by both Boeing and Airbus,
         \hookrightarrow which is not the case. We
         \hookrightarrow should use \oplus (logical
         \hookrightarrow exclusive disjunction) to
         \hookrightarrow indicate that an aircraft is
         \hookrightarrow produced by either Boeing or
         \hookrightarrow Airbus, but not both.
    Corrected:
    \forall x (Commercial Aircraft(x))
    \rightarrow (ProducedByBoeing(x))
    \oplus ProducedByAirbus(x)))
2. Original:
    \forall x (AmericanAirlinesJet(x))
    \rightarrow CommercialAircraft(x))
    Explanation: There is no need to
         \hookrightarrow change.
3. Original:
    \exists y (Revenue(airbus, y))
    \land \exists z (Revenue(boeing, z) \land y > z))
    Explanation: The use of '>' (greater
         \hookrightarrow than) is invalid according to
         \hookrightarrow the grammar rules. We need to
         \hookrightarrow replace it with a predicate
         \hookrightarrow that represents the

→ relationship between the

         \hookrightarrow revenues of Airbus and
         \hookrightarrow Boeing. We can use a
         \rightarrow predicate
         \hookrightarrow MoreInRevenue(airbus, boeing) to
         ↔ indicate that Airbus made
         \hookrightarrow more revenue than Boeing.
    Corrected:
    MoreInRevenue(airbus, boeing)
Conclusion:
```

694

695

696

697

698

699

700

701

702

703

704

705

706

707

708

709

710

711 712

713

714

715

716

717

718

719

720

721

722

723

724

725

726

727

728 729

730

731

732

733

734

735

736

737

738

739

740 741

742

743

744

745

746

747

748

749

750

751

752

754

755

756

757

758

759

760

761

763

Original:

```
681
683
```

627

| 764 765 | $\neg \exists x (United Airlines Plane(x))$ |
|------------|---|
| 766 | Explanation: There is no need to change. |
| 767 | |
| 768 | Final Corrected Formulas: |
| 769 | |
| 770 | Premises: |
| 771 | 1. $\forall x (Commercial Aircraft(x))$ |
| 772 | \rightarrow (ProducedByBoeing(x) \oplus ProducedByAirbus(x))) |
| 773 | \hookrightarrow |
| 774 | 2. $\forall x (AmericanAirlinesJet(x))$ |
| 775 | $\rightarrow CommercialAircraft(x))$ |
| 776 | 3. MoreInRevenue(airbus, boeing) |
| 777 | |
| 778 | Conclusion: |
| 779 | $\neg \exists x (United Airlines Plane(x))$ |
| 789 | $\land ProducedByBoeing(x))$ |

Listing 5: An Example from the Fixer KB for FOLIO

A.2.2 AR-LSAT

| <pre>Input: # Context: Charlie makes a soup by adding exactly</pre> |
|--|
| → added at some time before the → tomatoes or the kale, but not → before both. # Question: If the lentils are added last, then → which one of the following must |
| → be true? # Options: A) At least one of the foods is added → at some time before the zucchini. B) At least two of the foods are added → at some time before the kale. C) The mushrooms are added third. D) The zucchini is added third. E) The tomatoes are added fourth. |
| <pre>Output: # Declarations foods = EnumSort([kale, lentils,</pre> |

| ForAll([f1:foods, f2:foods], Implies(f1 |
|---|
| \hookrightarrow != f2, order(f1) != order(f2)) |
| \hookrightarrow ::: **Implied Information**: All |
| \hookrightarrow food items have different |
| \hookrightarrow positions. |
| <pre>Implies(order(mushrooms) == 3,</pre> |
| → order(lentils) == 6) ::: If the |
| \hookrightarrow mushrooms are added third, then |
| \hookrightarrow the lentils are added last |
| <pre>Implies(order(zucchini) == 1,</pre> |
| ↔ order(lentils) < order(onions)) |
| \hookrightarrow ::: If the zucchini is added |
| \hookrightarrow first, then the lentils are |
| \hookrightarrow added at some time before the |
| \hookrightarrow onions |
| And(order(tomatoes) != 5, order(kale) |
| \rightarrow != 5) ::: Neither the tomatoes |
| \hookrightarrow nor the kale is added fifth |
| Exists([f:foods], |
| \rightarrow And(Or(order(mushrooms) < |
| \hookrightarrow order(tomatoes), |
| \leftrightarrow order(mushrooms) < order(kale)), |
| ↔ Not(And(order(mushrooms) < |
| \rightarrow order(tomatoes), |
| \rightarrow order (mushrooms) < |
| \leftrightarrow order(kale))))) ::: The |
| \rightarrow mushrooms are added at some time |
| \rightarrow before the comatoes of the kale, |
| \rightarrow but not before both order(lentile) == 6 If the lentile |
| () are added last |
| - ale added last |
| <pre># options is valid(Evists([f:foods] order(f) <</pre> |
| \hookrightarrow order(zucchini))) ··· (A) |
| is valid(Count([f:foods] order(f) < |
| \rightarrow order(kale)) >= 2) ··· (B) |
| is valid(order(mushrooms) == 3) ··· (C) |
| is valid(order(zucchini) == 3) ··· (D) |
| is valid(order(tomatoes) == 4) (F) |
| |

Listing 6: An Example from the Translation KB for AR-LSAT

Input: ### Context: Eight camp counselorsFran, George, \hookrightarrow Henry, Joan, Kathy, Lewis, \hookrightarrow Nathan, and Olgamust each be \hookrightarrow assigned to supervise exactly \hookrightarrow one of three activitiesswimming, \hookrightarrow tennis, and volleyball. The \rightarrow assignment of counselors must \rightarrow conform to the following \hookrightarrow conditions: Each activity is \hookrightarrow supervised by at least two, but \hookrightarrow not more than three, of the eight counselors. Henry \hookrightarrow \hookrightarrow supervises swimming. Neither \hookrightarrow Kathy nor Olga supervises \hookrightarrow tennis. Neither Kathy nor Nathan \rightarrow supervises the same activity as Joan. If George supervises \hookrightarrow \hookrightarrow swimming, both Nathan and Olga \hookrightarrow supervise volleyball. ### Question: Which one of the following is a pair of \hookrightarrow counselors who could be two of \hookrightarrow three counselors assigned to \hookrightarrow supervise swimming?

```
900
901
902
903
904
905
906
907
908
909
910
911
912
913
914
915
916
917
918
919
921
922
924
925
926
927
928
929
930
931
933
936
937
938
939
941
942
944
947
951
952
953
955
958
959
960
961
962
963
964
965
967
```

```
### Options:
A) George and Nathan
B) George and Olga
C) Joan and Kathy
D) Joan and Nathan
E) Joan and Olga
### Logic program
```plaintext
Declarations
counselors = EnumSort([Fran, George,
 \hookrightarrow Henry, Joan, Kathy, Lewis,
 \hookrightarrow Nathan, Olga])
activities = EnumSort([swimming,
 \hookrightarrow tennis, volleyball])
assigned = Function([counselors] ->
 \hookrightarrow [activities])
Constraints
ForAll([a:activities],
 ↔ And(Count([c:counselors],
 \hookrightarrow assigned(c) == a) >= 2,
 ↔ Count([c:counselors],
 → assigned(c) == a) <= 3)) :::</p>
 \hookrightarrow Each activity is supervised by
 \hookrightarrow at least two, but not more than
 \hookrightarrow three, of the eight counselors
assigned(Henry) == swimming ::: Henry

→ supervises swimming

assigned(Kathy) != tennis ::: Neither
 ↔ Kathy nor Olga supervises tennis
assigned(Olga) != tennis ::: Neither
 \hookrightarrow Kathy nor Olga supervises tennis
assigned(Kathy) != assigned(Joan) :::
 \hookrightarrow Neither Kathy nor Nathan
 \hookrightarrow supervises the same activity as
 → Joan
assigned(Nathan) != assigned(Joan) :::
 \hookrightarrow Neither Kathy nor Nathan
 \hookrightarrow supervises the same activity as
 → Joan
Implies(assigned(George) == swimming,
 ↔ And(assigned(Nathan) ==
 → volleyball, assigned(Olga) ==
 \hookrightarrow volleyball)) ::: If George
 \hookrightarrow supervises swimming, both Nathan
 \hookrightarrow and Olga supervise volleyball
Options
is_accurate_list([assigned(George) ==
 \hookrightarrow swimming, assigned(Nathan) ==
 \hookrightarrow swimming]) ::: (A)
is_accurate_list([assigned(George) ==

→ swimming, assigned(Olga) ==

 \hookrightarrow swimming]) ::: (B)
is_accurate_list([assigned(Joan) ==
 ↔ swimming, assigned(Kathy) ==
 \hookrightarrow swimming]) ::: (C)
is_accurate_list([assigned(Joan) ==
 \hookrightarrow swimming, assigned(Nathan) ==
 \hookrightarrow swimming]) ::: (D)
is_accurate_list([assigned(Joan) ==
 → swimming, assigned(Olga) ==
 \hookrightarrow swimming]) ::: (E)
. . .
Error
No correct option problem. All options
 \hookrightarrow in this program are judged as
 \hookrightarrow incorrect, but this question is
 \hookrightarrow a single-choice question. It
 \hookrightarrow seems that some information has
```

```
\hookrightarrow been overlooked, and the program
 \hookrightarrow needs to be re-examined.
Output:
Analysis
1. There is no statement that each
 \hookrightarrow counselor can only be assigned
 \hookrightarrow to one activity; this constraint
 \hookrightarrow needs to be added.
2. The question asks which two people
 \hookrightarrow can simultaneously be assigned
 \hookrightarrow to swimming, so `is_sat` should
 \hookrightarrow be used.
The final program is modified as
 \hookrightarrow follows:
```plaintext
# Declarations
counselors = EnumSort([Fran, George,
    \hookrightarrow Henry, Joan, Kathy, Lewis,
    \hookrightarrow Nathan, Olga])
activities = EnumSort([swimming,
    \hookrightarrow tennis, volleyball])
assigned = Function([counselors] ->
    \hookrightarrow [activities])
# Constraints
ForAll([a:activities],
    ↔ And(Count([c:counselors],
    \hookrightarrow assigned(c) == a) >= 2,
    \hookrightarrow Count([c:counselors],
    \hookrightarrow assigned(c) == a) <= 3)) :::
    \hookrightarrow Each activity is supervised by
    \hookrightarrow at least two, but not more than
    \hookrightarrow three, of the eight counselors
ForAll([c:counselors],

→ Count([a:activities],

    \hookrightarrow assigned(c) == a) == 1) ::: A
    \hookrightarrow person can only participate in
    \hookrightarrow one activity.
assigned(Henry) == swimming ::: Henry
    ↔ supervises swimming
assigned(Kathy) != tennis ::: Neither
    \hookrightarrow Kathy nor Olga supervises tennis
assigned(Olga) != tennis ::: Neither
    \hookrightarrow Kathy nor Olga supervises tennis
assigned(Kathy) != assigned(Joan) :::
    \hookrightarrow Neither Kathy nor Nathan
    \hookrightarrow supervises the same activity as
    → Joan
assigned(Nathan) != assigned(Joan) :::
    \hookrightarrow Neither Kathy nor Nathan
    \hookrightarrow supervises the same activity as
    → Joan
Implies(assigned(George) == swimming,
    ↔ And(assigned(Nathan) ==
    → volleyball, assigned(Olga) ==
    ↔ volleyball)) ::: If George
    \hookrightarrow supervises swimming, both Nathan
    ↔ and Olga supervise volleyball
# Options
is_sat(And(assigned(George) ==
    \hookrightarrow swimming, assigned(Nathan) ==
    \hookrightarrow swimming)) ::: (A)
is_sat(And(assigned(George) ==
    \hookrightarrow swimming, assigned(Olga) ==
    \hookrightarrow swimming)) ::: (B)
is_sat(And(assigned(Joan) == swimming,
    ↔ assigned(Kathy) == swimming))
    \hookrightarrow ::: (C)
```

968

969

970

971

972

973

974

975 976

977

978

979

980

981

982

983

984

985

986

987

988

989

990

991

992

993 994

995

996

997

998

999

1000

1001

1002

1003

1005

1007

1008

1009

1010

1011

1012

1013

1014

1015

1016

1017

1018

1019

1020

1021

1022

1023

1024

1025

1026

1027

1028

1029

1030

1031

1032

1033

1034

1036

```
1038
1039
1040
1041
1042
1043
```

1044

1046

1047

1048

1050

1051

1052

1054

1055

1056

1057

1058

1059

1061

1062

1063

1064

1065

1066

1067

1068

1069

Listing 7: An Example from the Fixer KB for AR-LSAT

A.3 Ablation Studies

We tested the performance of each model at various temperatures from 0.1 to 0.3. Due to the limited space, for the translation process, we show the results of each model at the optimal temperature; for the repair process, we show the results where the input is the optimal temperature and the number of examples of the translation process.

A.3.1 Ablation Studies on FOLIO

In this section, we analyze the translation stage results under two evaluation settings: **Exact Match** and **Error as Unknown**. In the Exact Match setting, an output is considered correct only if the predicted formula strictly matches the ground truth label. In contrast, in the Error as Unknown setting, formulas that contain errors are treated as "Unknown", which may allow some outputs to be counted as correct by chance.

Tables 3 and 4 show that varying the number of in-context examples leads to stable performance across models. In some cases, increasing the number of examples slightly improves accuracy, while in others it introduces noise, suggesting that an optimal example count exists for each model.

Table 3: Pre-Fix Accuracy on FOLIO (TranslationStage) – Exact Match

| Model | Ex#=1 | Ex#=2 | Ex#=3 |
|---------------|-------|-------|-------|
| DeepSeek v2.5 | 76.92 | 75.27 | 75.82 |
| GPT-40 | 74.73 | 74.73 | 75.82 |
| Llama3.3 | 74.18 | 73.08 | 74.18 |
| GPT-3.5-turbo | 60.99 | 64.29 | 63.74 |
| Gemma2 | 73.08 | 76.37 | 75.27 |

Table 4: Pre-Fix Accuracy on FOLIO (TranslationStage) – Error as Unknown

| Model | Ex#=1 | Ex#=2 | Ex#=3 |
|---------------|-------|-------|-------|
| DeepSeek v2.5 | 78.57 | 78.02 | 78.57 |
| GPT-40 | 75.82 | 77.47 | 78.02 |
| Llama3.3 | 75.82 | 75.82 | 78.02 |
| GPT-3.5-turbo | 65.38 | 68.68 | 68.13 |
| Gemma2 | 75.27 | 78.57 | 78.02 |

Next, we present the results after applying the Fixer LLM. For each model, the best result from the table 4 is chosen as the baseline for the fix stage. Comparing the best results from the translation stage with post-fix results (Tables 5 and 6), the Fixer LLM significantly boosts the exact match accuracy. This improvement indicates that the module effectively corrects superficial syntax and formatting errors, reducing the incidence of chancecorrect answers ("lucky guesses"). In contrast, the improvement in the "Error as Unknown" setting is relatively minor, suggesting that the Fixer LLM primarily enhances the strict correctness of the outputs. 1070

1071

1072

1073

1074

1075

1076

1078

1079

1080

1081

1082

1083

1084

1087

1089

1090

1091

1093

1094

Table 5: Post-Fix Accuracy on FOLIO (After Applying Fixer LLM) – Exact Match

| Model | Ex#=1 | Ex#=2 | Ex#=3 |
|---------------|-------|-------|-------|
| DeepSeek v2.5 | 78.57 | 78.02 | 78.02 |
| GPT-40 | 80.77 | 80.77 | 80.77 |
| Llama3.3 | 78.02 | 78.02 | 78.02 |
| GPT-3.5-turbo | 69.78 | 70.33 | 69.78 |
| Gemma2 | 78.57 | 79.12 | 78.02 |

Table 6: Post-Fix Accuracy on FOLIO (After Applying Fixer LLM) – Error as Unknown

| Model | Ex#=1 | Ex#=2 | Ex#=3 |
|---------------|-------|-------|-------|
| DeepSeek v2.5 | 78.57 | 78.57 | 78.57 |
| GPT-40 | 80.77 | 80.77 | 80.77 |
| Llama3.3 | 78.57 | 78.57 | 78.57 |
| GPT-3.5-turbo | 70.88 | 70.88 | 70.88 |
| Gemma2 | 79.67 | 79.67 | 79.67 |

A.3.2 Ablation Studies on AR-LSAT

To present the results more intuitively, we use a slightly different statistical method than in the main text, focusing solely on the proportion of problems correctly solved by the Solver. (Alternatively, this value can be considered as Exe_rate * Exe_accuracy.)

Table 7 shows the best results of each model at different example counts during the Translation Stage. Table 8 shows the results of each model after applying the Fixer LLM.

Table 7: Model Performance on Different ExampleNumbers (Translation Stage)

| Model | Ex#=3 | Ex#=5 | Ex#=7 |
|---------------|-------|-------|-------|
| DeepSeek v2.5 | 12.12 | 16.02 | 13.85 |
| GPT-40 | 19.48 | 15.58 | 16.02 |
| Llama-3.3 | 16.02 | 16.88 | 14.72 |
| GPT-3.5-turbo | 6.06 | 8.23 | 6.06 |
| Gemma-2 | 5.19 | 7.79 | 3.03 |

Table 8: Model Performance (After applying of theFixer LLM)

| Model | Performance |
|---------------|-------------|
| DeepSeek v2.5 | 31.60 |
| GPT-40 | 34.63 |
| Llama-3.3 | 35.06 |
| GPT-3.5-turbo | 10.39 |
| Gemma-2 | 16.45 |

In Table 7, we observe that before the appli-1095 cation of the Fixer LLM, large models such as 1096 DeepSeek, GPT-40, and Llama3.3 achieved their 1097 highest accuracy at around 17%, with GPT-40 per-1098 forming best at 19.5%. In contrast, small models 1099 like GPT-3.5-turbo and Gemma-2 reached a maxi-1100 mum accuracy of only 8.23%. We then applied the 1101 1102 Fixer LLM based on the best performance results for each model during the Translation Stage. As 1103 shown in Table 8, the accuracy improvement for 1104 DeepSeek, Llama3.3, and GPT-40 exceeded 15%, 1105 nearly doubling their original performance. How-1106 1107 ever, for GPT-3.5-turbo and Gemma-2, the improvements were below 10%, with GPT-3.5 show-1108 ing only a modest 2% increase. 1109

> This further supports our conclusion. For large models, the Fixer LLM has a more significant effect on reasoning tasks involving long texts, where it is difficult to focus on details for inference.

1110

1111

1112