# LOGIC-LM++: Multi-Step Refinement for Symbolic Formulations

**Shashank Kirtania, Priyanshu Gupta, Arjun Radhakrishna**
**Microsoft**
{t-skirtania, priyansgupta, arradha} @microsoft.com

## Abstract

In this paper we examine the limitations of Large Language Models (LLMs) for complex reasoning tasks. Although recent works have started to employ formal languages as an intermediate representation for reasoning tasks, they often face challenges in accurately generating and refining these formal specifications to ensure correctness. To address these issues, this paper proposes Logic-LM++, an improvement on Logic-LM (Pan et al., 2023). It uses the ability of LLMs to do pairwise comparisons, allowing the evaluation of the refinements suggested by the LLM. The paper demonstrates that Logic-LM++ outperforms Logic-LM and other contemporary techniques across natural language reasoning tasks on three datasets, FO-LIO, ProofWriter and AR-LSAT, with an average improvement of 18.5% on standard prompting, 12.3% on chain of thought prompting and 5% on Logic-LM.

## 1 Introduction

Large language models (LLMs) have shown proven capability of reasoning (Brown et al., 2020; Chowdhery et al., 2022) but still struggle at complex reasoning problems as seen in real world assessments (Zhong et al., 2021). For complex multi hop reasoning tasks current state of the art approaches (Pan et al., 2023; Ye et al., 2023) leverage formal languages as intermediate representation of these reasoning problems and utilize symbolic reasoners to come up with the right response. A typical workflow of such techniques consist of 3 steps: a natural language prompt which consist of the task information, a response formulation for the problem, final response generated with symbolic executor.

While logic-assisted LLM reasoning techniques are promising, we observe following problems in such systems: Firstly, LLMs are poor at generating intermediate formal specifications. A few techniques try to counter this problem with a refinement loop (Madaan et al., 2023a; Welleck et al., 2022; Shinn et al., 2023) to improve upon the syntactical correctness of the symbolic formulation. Secondly, the LLMs are poor at repairing the formal representations with limited information with error information. For example, in Figure 1 the LLM initially generates a syntactically incorrect formulation. After a turn of refinement, while the LLM is able to generate a response that is syntactically correct, it introduces a *semantic* error in the formulation by incorrectly translating the statement "No young person teaches". These kind of incorrect translations from Natural Language (NL) to intermediate formal specifications is a common problem we observe over the failing cases of refinement. Thirdly, we observe that refinements are not always linear-resolving an error with the symbolic formulation can take multiple steps of careful edits and evaluation. The formulations generated in refinement stage in 1 introduced the wrong interpretation of "No young person teaches" to "All young people teaches".

To address these challenges we propose to add following measures in Logic-LM to enhance it's capabilities resulting in improved variant Logic-LM++.

We leverage the ability of LLMs to do pairwise comparison (Zheng et al., 2023a), this gives us an opportunity to evaluate the refinements suggested by the LLM and do a semantic check with respect to the problem statement to ensure if the edits in the symbolic formulation generated while refinement improve the formulation semantically not just syntactically.

We also improve on the refinement mechanism present in Logic-LM to give more context of the problem statement during refinement stage, this eliminates cases where recommended edits are appalling and do not improve the formulation significantly.
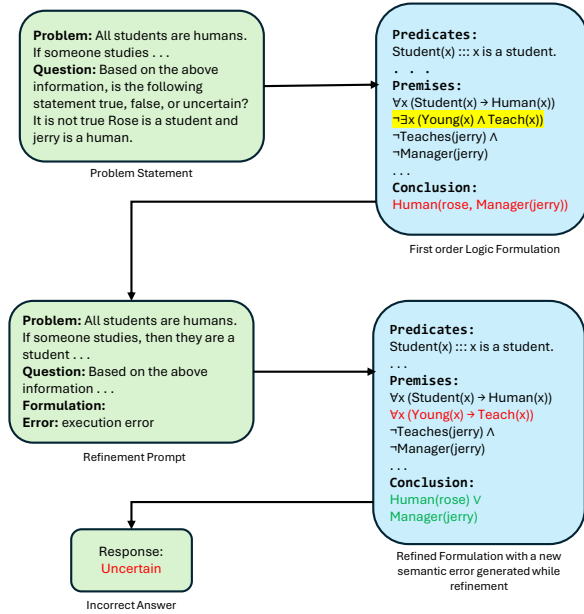
## 2   Related Work



Figure 1: Refinement of logical formulations in Logic-LM



Figure 2: Improvement in refinement by Logic-LM++

### 2.1   Reasoning with LLMs

Large Language model (LLM) - based reasoning techniques commonly entail the deconstruction of complex questions into a sequence of intermediate steps, often referred to as chains, before reaching the ultimate solution. This technique is a reflection of methods such as Chain of Thought (CoT) prompting and its variations, as shown in various studies (Wei et al., 2022; Kojima et al., 2022). These methodologies require the meticulous segmentation of a problem into a chain of smaller, manageable chunks. Each of these chunks represents a step in the reasoning process, guiding the model towards a comprehensive solution. The concept of the reflection loop, as explored in previous research (Shinn et al., 2023; Madaan et al., 2023b), offers a means of refining the reasoning by identifying and eliminating any flaws that may be introduced by the LLM during a reasoning step. This process enhances the inherent capability of the LLM to self-correct, contributing to more accurate and reliable outcomes. Recent works have further explore the process of self-evaluation at these intermediate steps (Welleck et al., 2022; Paul et al., 2024). This process involves the LLM assessing its reasoning at each step, allowing it to identify any inaccuracies. By rectifying these issues before proceeding to the next step, the LLM can ensure
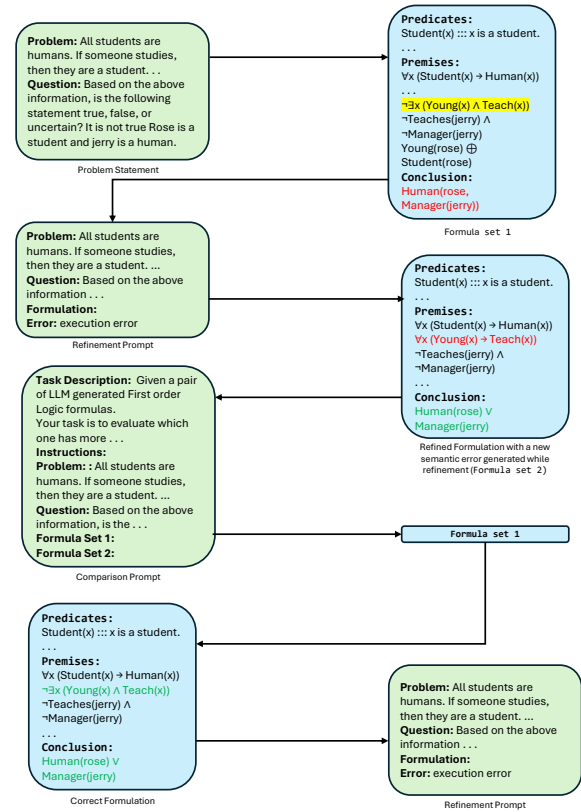
a more accurate and reliable chain of reasoning. Aligned with our objective of capturing the natural language intent of the user from symbolic formulations, recent works (Endres et al., 2024) have also explored the translation of natural language into formal language post conditions. This research investigates how effectively we can convert the often-ambiguous language of human conversation into the precise, unambiguous language of formal logic. This translation process is crucial for the accurate interpretation and execution of user intent, particularly in complex or technical tasks.

### 2.2   Tool-augmented Large Language Models

Language models face inherent limitations, unable to access real-time data, execute actions, or conduct precise mathematical reasoning.To address this, recent research endeavors have sought to augment language models by integrating external resources such as retrievers (Nakano et al., 2021; Shi et al., 2023; Lazaridou et al., 2022), calculators (Cobbe et al., 2021), code interpreters (Wang et al., 2022), planners (Liu et al., 2023), symbolic solvers (Ye et al., 2023; Pan et al., 2023), and other pretrained models (Shen et al., 2023). Notably, in the realm of mathematical reasoning, numerous investigations

| Dataset | GPT-3.5 Turbo | | | | GPT-4 | | | |
|---|---|---|---|---|---|---|---|---|
| | Standard | CoT | Logic-LM | Logic-LM++ | Standard | CoT | Logic-LM | Logic-LM++ |
| FOLIO | 45.09 | 57.35 | **62.80** | *62.25* | 69.11 | 70.58 | 78.92 | **84.80** |
| ProofWriter | 35.50 | 49.17 | 58.33 | **58.83** | 52.67 | 68.11 | 79.66 | 79.66 |
| AR-LSAT | 20.34 | 17.31 | 26.41 | **28.13** | 33.33 | 35.06 | 43.04 | **46.32** |

Table 1: Accuracy of standard promoting, chain-of-thought (CoT) promoting, Logic-LM and Logic-LM++.

have illustrated the efficacy of incorporating calculators (Cobbe et al., 2021; Imani et al., 2023) or Python interpreters (Gao et al., 2023; Chen et al., 2022) into language models, significantly enhancing performance by offloading numerical computations. Recent studies (Gao et al., 2023; Chen et al., 2022) have showcased improved effectiveness in arithmetic reasoning tasks by generating Python programs that delineate the reasoning process through sequenced chained commands.

## 3 Methodology

### 3.1 Background

Logic-LM (Pan et al., 2023) is a framework to decompose a reasoning problem into three stages:

1. *Problem Formulation*, where given a task description and a problem statement LLM write symbolic formulations that represents the NL problem. In Figure 1 the NL prompt with task description is the problem formulator in Logic-LM.

2. *Symbolic Reasoning*, where we use a symbolic solver like Prover92 (Robinson, 1965)and Z3 theorem prover (Moura and Bjørner, 2008) to solve the formulations generated earlier.

3. *Result interpretation*, where the produced output is mapped to the right answer using regex parsing.

Logic-LM uses a refinement loop to fix errors in symbolic formulation at formulation and reasoning stages. However, Logic-LM still struggles to improve on logical representations, showing almost no improvement after multiple iterations. Authors attribute this to semantic limitations of the formulation. To this end, Logic-LM++ aims to mitigate this limitation by improving the Logic-LM refinement loop.

### 3.2 Self-Refinement Agent

Logic-LM defines the notion of a *Self-Refinement Agent* to implement the refinement loop in the symbolic formulations in cases where the formulations did not yield a successful execution within the system. This agent is characterized by a *refinement prompt* 1. In the original work, the refinement prompt constituted various few shot examples to act as exemplar for the model. While similar techniques have proven useful (Madaan et al., 2023a; Shinn et al., 2023), we anecdotally observe that instead of helping the model it adds extra irrelevant information that distracts the model from fixing the issues relevant to the current formulation, consistent with similar studies in other domains (Pan et al., 2023). To alleviate this, instead of adding few-shots, we add the problem statement and the question to the refinement prompt alongside instructions to self-reflect on the model's failure to generate the right response. As we show later in Section 4, this structure helps better *contextualize* (Shinn et al., 2023) the formulation to the self-reflection agent and help the system generate better refinements.

### 3.3 Backtracking Agent

LLMs has shown remarkable results in automated evaluation benchmarks (Zheng et al., 2023b) and has shown high alignment with the human judgement (Wei et al., 2024). We use this capability of LLMs to assess if the repaired formulation by self-refinement improves the alignment of the human intent with LLM generated formulations. This allows us to get rid of the updates that are not helpful in future iterations and only use those updates where the changes help the model to come to the right formulation. In Figure 1 we can see without the backtracking agent the LLM accepts the semantically incorrect symbolic formulations as the statement "No young person teaches" is translated to "all young people teach" since the code is syntactically correct there is no proof-check on the refinement.

However, In Figure 2 we demonstrate in the same example with the backtracking agent Logic-LM++ is able to generate right formulation by us-
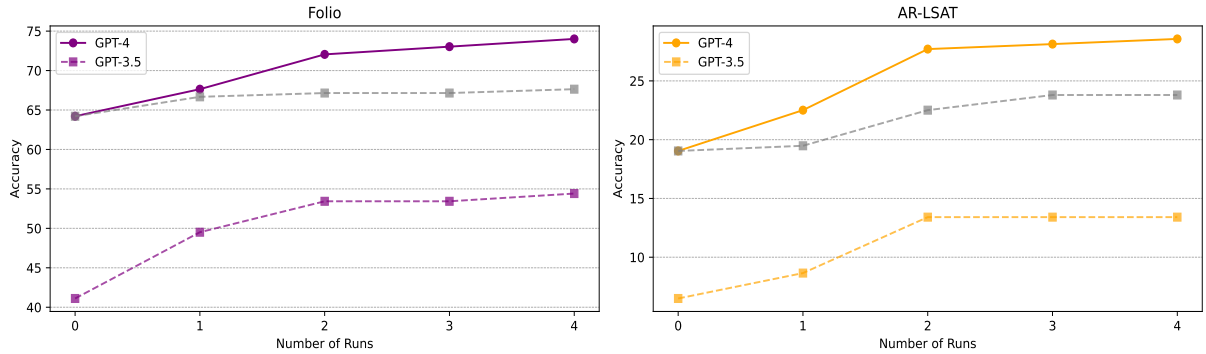
Figure 3: Accuracy in subsequent rounds of refinement. The grey line here represents the accuracy scores on self-refinement without backtracking with GPT-4.

ing the right formulation to represent "No young person teaches" and use the right formulation to describe the question "Rose is a student and Jerry is a human". This showcase how backtracking agent works as funnel to reduce the semantic error that is propagated at the refinement stage. In Figure 2 we show on comparison of the two sets of the formulation, returns a more semantically correct formulation this allows Logic-LM++ to only accept the edits if it improves or preserve the logical structure of the formulation.

## 4 Experiments and Analysis

### 4.1 Dataset

**FOLIO** (Han et al., 2022) is a challenging expert-written dataset for logical reasoning. The problems are aligned with real-world knowledge and use natural wordings, and the questions require complex first-order logic reasoning to solve. We use the FOLIO test set for evaluation with 204 examples.

**AR-LSAT** (Zhong et al., 2022) is a dataset that collects all analytical logic reasoning questions from the Law School Admission Test from 1991 to 2016. We use the test set which has 231 multiple-choice questions. AR-LSAT is particularly challenging, with state-of-the-art model's performance slightly better than random guessing (Liang et al., 2022; Ribeiro et al., 2023).

**ProofWriter** (Tafjord et al., 2021) is another popular dataset on deductive logical reasoning. Since the problems are in more natual language like setting it makes semantic evaluation very relevant in the problem set. Logic-LM use the open-world assumption (OWA) subset in which each example is a (problem, goal) pair and the label is one of PROVED, DISPROVED, UNKNOWN. Logic-LM evaluate the pipeline with the hardest section of

|  |  | GPT-3.5 turbo | | GPT-4 | |
|---|---|---|---|---|---|
| **BT** |  | - | + | - | + |
| **FOLIO** | $E_r$ | 84.3 | 84.3 | 85.8 | **86.7** |
|  | $E_a$ | 64.3 | **66.2** | 79.9 | 85.8 |
| **ProofWriter** | $E_r$ | 95.6 | 95.6 | 99.0 | 99.0 |
|  | $E_a$ | 74.1 | **77.2** | <u>79.6</u> | <u>79.6</u> |
| **AR-LSAT** | $E_r$ | 21.8 | **22.9** | **32.6** | 32.0 |
|  | $E_a$ | 60.3 | **64.1** | 60.0 | **66.2** |

Table 2: Execution rate ($E_r$) and Execution Accuracy ($E_a$) agent with Backtracking (BT).

ProofWriter which contain total of 600 randomly sampled five step multi-hop reasoning questions.

### 4.2 Principal Findings

We report the final results of Logic-LM++ in Table 1. We try to answer 2 major research questions.

**RQ1: Can LLMs conduct pairwise comparisons of symbolic formulations based on their relevance to a natural language task description?** LLMs have demonstrated promising capabilities in pairwise comparisons for NLG evaluations (Kim et al., 2024), even in low-resource languages where their natural language generation abilities remain underdeveloped (Zheng et al., 2023a). As depicted in Table 2, the execution accuracy of the framework employing a backtracking agent is enhanced by approximately 6% with GPT-4 and around 3% with GPT3.5-turbo. Despite the average gain in execution rate being less than 1%, these statistics underscore the empirical improvements in code quality in terms of semantic correctness. Figure 1 provides a working example from the FOLIO dataset. Although the code is syntactically correct after refinement, it misinterprets a logical statement.

However, by implementing pairwise comparisons, the LLM can select the semantically correct formulation. This leads to the correct answer in the subsequent refinement iteration.

**RQ2: Does refinement by LLM always positively affect the formulations?**

In Figure 3, we evaluate the refinement process with and without backtracking. Logic-LM's accuracy plateaus with more runs because refined solutions may not represent the intended code. The author's also discuss this as a known limitation of the refinement process in the refinement loop they proposed. Backtracking, which reverts to the initial code if no semantic improvement is found, allows Logic-LM++ to perform consistently better by continually reassessing and correcting refinements for more reliable results.

Figure 4 shows that the backtracking agent significantly improves results in the second round within the FOLIO dataset, with a similar impact in later rounds. This indicates that backtracking is most effective early on since the generated refinement can also degrade the performance of the formulations, enabling Logic-LM++ to achieve substantial better and iterative improvements over time.

### 4.3 Error Analysis

Even though Logic-LM++ shows impressive improvements over standard refinement techniques, it still lacks behind in the cases where the first set of formulation generated is completely different from the ground truth formulation. On analyzing the failure cases in Logic-LM we note that the current pipeline relies a lot on fixing the bugs with current formulation without losing on semantic understanding, however in cases where the generating semantically correct formulations is hard the technique is contingent to initial formulations generated.

### 5 Discussion and Future Work

Figure 3 reveals a significant observation regarding the iteration increase of Logic-LM, which appears to reach convergence substantially earlier than Logic-LM++. Logic-LM associates attributes this to the hard limit of semantically correctness that can be achieved with Logic-LM. The findings stress the importance of semantic accuracy, as the Logic-LM++ exhibits consistently improved outcomes over multiple iterations, contrary to findings by Logic-LM. This outcome is primarily attributed
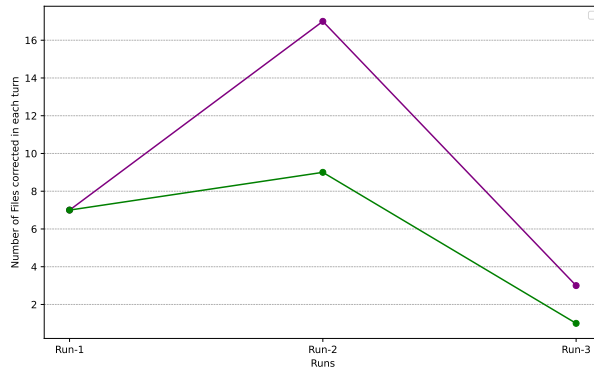


Figure 4: Number of symbolic formulations corrected after each turn of self-refinement with backtracking agent (purple) and without backtracking agent (green) in FOLIO with GPT-4.

to the model's capability to revert to the initial formulation if the refined version does not offer a semantically superior representation. Eventhough, Logic-LM++ show promising results it only focus on symbolic formulations, this effort can be well generalised to other tool augmented techniques that rely on intermediate code representation with semantic improvements during refinement.

### 6 Conclusion

We propose Logic-LM++ which beats state of the art results on natural language reasoning tasks on three datasets. Logic-LM++ takes leverage of LLMs' reasoning capabilities to show significant improvements in efficient use of logic solvers for reasoning, we demonstrate that LLMs show promising results at conducting comparison between symbolic formulations even in cases where generating symbolic formulations is a hard task for LLM.

### Limitation

At present, Logic-LM++ faces constraints in its capacity to effectively capture the semantic intricacies in reasoning tasks. This limitation notably complicates the evaluation process, especially when dealing with smaller LLMs like (Rozière et al., 2023). The understanding required for accurate reasoning poses a significant challenge, particularly in contexts where the model's semantic comprehension may be insufficient. Due to this the assessment of performance becomes notably more complex. This limitation underscores the need for continued advancements in semantic understanding within LLMs to enhance their efficacy across reasoning tasks.

# References

Tom B. Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared Kaplan, Prafulla Dhariwal, Arvind Neelakantan, Pranav Shyam, Girish Sastry, Amanda Askell, Sandhini Agarwal, Ariel Herbert-Voss, Gretchen Krueger, Tom Henighan, Rewon Child, Aditya Ramesh, Daniel M. Ziegler, Jeffrey Wu, Clemens Winter, Christopher Hesse, Mark Chen, Eric Sigler, Mateusz Litwin, Scott Gray, Benjamin Chess, Jack Clark, Christopher Berner, Sam McCandlish, Alec Radford, Ilya Sutskever, and Dario Amodei. 2020. Language models are few-shot learners. In *Proceedings of the Annual Conference on Neural Information Processing Systems (NeurIPS)*.

Wenhu Chen, Xueguang Ma, Xinyi Wang, and William W. Cohen. 2022. Program of thoughts prompting: Disentangling computation from reasoning for numerical reasoning tasks. *CoRR*, abs/2211.12588.

Aakanksha Chowdhery, Sharan Narang, Jacob Devlin, Maarten Bosma, Gaurav Mishra, Adam Roberts, Paul Barham, Hyung Won Chung, Charles Sutton, Sebastian Gehrmann, et al. 2022. PaLM: Scaling Language Modeling with Pathways. *Preprint*, arXiv:2204.02311. *arXiv preprint arXiv:2204.02311*. https://arxiv.org/abs/2204.02311.

Karl Cobbe, Vineet Kosaraju, Mohammad Bavarian, Jacob Hilton, Reiichiro Nakano, Christopher Hesse, and John Schulman. 2021. Training verifiers to solve math word problems. *CoRR*, abs/2110.14168.

Madeline Endres, Sarah Fakhoury, Saikat Chakraborty, and Shuvendu Lahiri. 2024. Can large language models transform natural language intent into formal method postconditions? In *The ACM International Conference on the Foundations of Software Engineering (FSE)*. ACM. Https://2024.esec-fse.org/details/fse-2024-research-papers/51/Can-Large-Language-Models-Transform-Natural-Language-Intent-into-Formal-Method-Postco.

Luyu Gao, Aman Madaan, Shuyan Zhou, Uri Alon, Pengfei Liu, Yiming Yang, Jamie Callan, and Graham Neubig. 2023. PAL: program-aided language models. In *Proceedings of the International Conference on Machine Learning (ICML)*, volume 202, pages 10764–10799.

Simeng Han, Hailey Schoelkopf, Yilun Zhao, Zhenting Qi, Martin Riddell, Luke Benson, Lucy Sun, Ekaterina Zubova, Yujie Qiao, Matthew Burtell, David Peng, Jonathan Fan, Yixin Liu, Brian Wong, Malcolm Sailor, Ansong Ni, Linyong Nan, Jungo Kasai, Tao Yu, Rui Zhang, Shafiq R. Joty, Alexander R. Fabbri, Wojciech Kryscinski, Xi Victoria Lin, Caiming Xiong, and Dragomir Radev. 2022. FOLIO: natural language reasoning with first-order logic. *CoRR*, abs/2209.00840.

Shima Imani, Liang Du, and Harsh Shrivastava. 2023. Mathprompter: Mathematical reasoning using large language models. In *Proceedings of the 61st Annual Meeting of the Association for Computational Linguistics (ACL), Industry Track*, pages 37–42.

Seungone Kim, Juyoung Suk, Shayne Longpre, Bill Yuchen Lin, Jamin Shin, Sean Welleck, Graham Neubig, Moontae Lee, Kyungjae Lee, and Minjoon Seo. 2024. Prometheus 2: An open source language model specialized in evaluating other language models. *Preprint*, arXiv:2405.01535.

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*, volume 35, pages 22199–22213. Curran Associates, Inc.

Angeliki Lazaridou, Elena Gribovskaya, Wojciech Stokowiec, and Nikolai Grigorev. 2022. Internet-augmented language models through few-shot prompting for open-domain question answering. *CoRR*, abs/2203.05115.

Percy Liang, Rishi Bommasani, Tony Lee, Dimitris Tsipras, Dilara Soylu, Michihiro Yasunaga, Yian Zhang, Deepak Narayanan, Yuhuai Wu, Ananya Kumar, Benjamin Newman, Binhang Yuan, Bobby Yan, Ce Zhang, Christian Cosgrove, Christopher D. Manning, Christopher Ré, Diana Acosta-Navas, Drew A. Hudson, Eric Zelikman, Esin Durmus, Faisal Ladhak, Frieda Rong, Hongyu Ren, Huaxiu Yao, Jue Wang, Keshav Santhanam, Laurel J. Orr, Lucia Zheng, Mert Yüksekgönül, Mirac Suzgun, Nathan Kim, Neel Guha, Niladri S. Chatterji, Omar Khattab, Peter Henderson, Qian Huang, Ryan Chi, Sang Michael Xie, Shibani Santurkar, Surya Ganguli, Tatsunori Hashimoto, Thomas Icard, Tianyi Zhang, Vishrav Chaudhary, William Wang, Xuechen Li, Yifan Mai, Yuhui Zhang, and Yuta Koreeda. 2022. Holistic evaluation of language models. *CoRR*, abs/2211.09110.

Bo Liu, Yuqian Jiang, Xiaohan Zhang, Qiang Liu, Shiqi Zhang, Joydeep Biswas, and Peter Stone. 2023. LLM+P: empowering large language models with optimal planning proficiency. *CoRR*, abs/2304.11477.

Aman Madaan, Niket Tandon, Prakhar Gupta, Skyler Hallinan, Luyu Gao, Sarah Wiegreffe, Uri Alon, Nouha Dziri, Shrimai Prabhumoye, Yiming Yang, Shashank Gupta, Bodhisattwa Prasad Majumder, Katherine Hermann, Sean Welleck, Amir Yazdanbakhsh, and Peter Clark. 2023a. Self-refine: Iterative refinement with self-feedback. In *Advances in Neural Information Processing Systems*, volume 36, pages 46534–46594. Curran Associates, Inc.

Aman Madaan, Niket Tandon, Prakhar Gupta, Skyler Hallinan, Luyu Gao, Sarah Wiegreffe, Uri Alon, Nouha Dziri, Shrimai Prabhumoye, Yiming Yang, Shashank Gupta, Bodhisattwa Prasad Majumder, Katherine Hermann, Sean Welleck, Amir Yazdanbakhsh, and Peter Clark. 2023b. Self-refine: Iterative refinement with self-feedback. In *Advances in Neural Information Processing Systems*, volume 36, pages 46534–46594. Curran Associates, Inc.

Leonardo Mendonça de Moura and Nikolaj S. Bjørner. 2008. Z3: An efficient smt solver. In *Proceedings of the 14th International Conference of Tools and Algorithms for the Construction and Analysis of Systems (TACAS)*, volume 4963 of *Lecture Notes in Computer Science*, pages 337–340.

Reiichiro Nakano, Jacob Hilton, Suchir Balaji, Jeff Wu, Long Ouyang, Christina Kim, Christopher Hesse, Shantanu Jain, Vineet Kosaraju, William Saunders, Xu Jiang, Karl Cobbe, Tyna Eloundou, Gretchen Krueger, Kevin Button, Matthew Knight, Benjamin Chess, and John Schulman. 2021. Webgpt: Browser-assisted question-answering with human feedback. *CoRR*, abs/2112.09332.

Liangming Pan, Alon Albalak, Xinyi Wang, and William Wang. 2023. Logic-LM: Empowering large language models with symbolic solvers for faithful logical reasoning. In *Findings of the Association for Computational Linguistics: EMNLP 2023*, pages 3806–3824, Singapore. Association for Computational Linguistics.

Debjit Paul, Mete Ismayilzada, Maxime Peyrard, Beatriz Borges, Antoine Bosselut, Robert West, and Boi Faltings. 2024. Refiner: Reasoning feedback on intermediate representations. *Preprint*, arXiv:2304.01904.

Danilo Neves Ribeiro, Shen Wang, Xiaofei Ma, Henghui Zhu, Rui Dong, Deguang Kong, Juliette Burger, Anjelica Ramos, Zhiheng Huang, William Yang Wang, George Karypis, Bing Xiang, and Dan Roth. 2023. STREET: A multi-task structured reasoning and explanation benchmark. In *Proceedings of the Eleventh International Conference on Learning Representations (ICLR)*.

John Alan Robinson. 1965. A machine-oriented logic based on the resolution principle. *Journal of the ACM (JACM)*, 12(1):23–41.

Baptiste Rozière, Jonas Gehring, Fabian Gloeckle, Sten Sootla, Itai Gat, Xiaoqing Ellen Tan, Yossi Adi, Jingyu Liu, Tal Remez, Jérémy Rapin, Artyom Kozhevnikov, Ivan Evtimov, Joanna Bitton, Manish Bhatt, Cristian Canton Ferrer, Aaron Grattafiori, Wenhan Xiong, Alexandre Défossez, Jade Copet, Faisal Azhar, Hugo Touvron, Louis Martin, Nicolas Usunier, Thomas Scialom, and Gabriel Synnaeve. 2023. Code Llama: Open Foundation Models for Code. *arXiv preprint arXiv:2308.12950*.

Yongliang Shen, Kaitao Song, Xu Tan, Dongsheng Li, Weiming Lu, and Yueting Zhuang. 2023. Hugginggpt: Solving AI tasks with chatgpt and its friends in huggingface. *CoRR*, abs/2303.17580.

Weijia Shi, Sewon Min, Michihiro Yasunaga, Minjoon Seo, Rich James, Mike Lewis, Luke Zettlemoyer, and Wen-tau Yih. 2023. REPLUG: retrieval-augmented black-box language models. *CoRR*, abs/2301.12652.

Noah Shinn, Federico Cassano, Ashwin Gopinath, Karthik Narasimhan, and Shunyu Yao. 2023. Reflexion: language agents with verbal reinforcement learning. In *Advances in Neural Information Processing Systems*, volume 36, pages 8634–8652. Curran Associates, Inc.

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*, pages 3621–3634, Online. Association for Computational Linguistics.

Xingyao Wang, Sha Li, and Heng Ji. 2022. Code4struct: Code generation for few-shot structured prediction from natural language. *CoRR*, abs/2210.12810.

Fangyun Wei, Xi Chen, and Lu Luo. 2024. Rethinking generative large language model evaluation for semantic comprehension.

Jason Wei, Xuezhi Wang, Dale Schuurmans, Maarten Bosma, brian ichter, Fei Xia, Ed 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*, volume 35, pages 24824–24837. Curran Associates, Inc.

Sean Welleck, Ximing Lu, Peter West, Faeze Brahman, Tianxiao Shen, Daniel Khashabi, and Yejin Choi. 2022. Generating sequences by learning to self-correct. *Preprint*, arXiv:2211.00053.

Xi Ye, Qiaochu Chen, Isil Dillig, and Greg Durrett. 2023. Satlm: Satisfiability-aided language models using declarative prompting. In *Proceedings of NeurIPS*.

Lianmin Zheng, Wei-Lin Chiang, Ying Sheng, Siyuan Zhuang, Zhanghao Wu, Yonghao Zhuang, Zi Lin, Zhuohan Li, Dacheng Li, Eric Xing, Hao Zhang, Joseph E Gonzalez, and Ion Stoica. 2023a. Judging llm-as-a-judge with mt-bench and chatbot arena. In *Advances in Neural Information Processing Systems*, volume 36, pages 46595–46623. Curran Associates, Inc.

Lianmin Zheng, Wei-Lin Chiang, Ying Sheng, Siyuan Zhuang, Zhanghao Wu, Yonghao Zhuang, Zi Lin, Zhuohan Li, Dacheng Li, Eric Xing, Hao Zhang, Joseph E. Gonzalez, and Ion Stoica. 2023b. Judging LLM-as-a-judge with MT-bench and chatbot arena. In *Thirty-seventh Conference on Neural Information Processing Systems Datasets and Benchmarks Track*.

Wanjun Zhong, Siyuan Wang, Duyu Tang, Zenan Xu, Daya Guo, Yining Chen, Jiahai Wang, Jian Yin, Ming Zhou, and Nan Duan. 2022. Analytical reasoning of text. In *Findings of the 2022 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies (NAACL-HLT)*, pages 2306–2319.

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.