

---

# HYBRIDMIND: Meta Selection of Natural Language and Symbolic Language for Enhanced LLM Reasoning

---

Anonymous Author(s)

Affiliation

Address

email

## Abstract

LLMs approach logical and mathematical reasoning through natural or symbolic languages. While natural language offers human-accessible flexibility but suffers from ambiguity, symbolic reasoning provides precise, machine-executable inferences at the cost of strict domain constraints. We introduce HYBRIDMIND, an adaptive strategy that selects the optimal reasoning approach for each reasoning problem. Through extensive experiments, we evaluate both prompting-based approaches with state-of-the-art LLMs and fine-tuned open-source models. We find that fine-tuning LLaMA-3.1-8B-Instruct as a meta-selector outperforms GPT-4o’s natural language reasoning by 4.4% on FOLIO and 1.3% on MATH. More notably, using GPT-3.5-turbo as a prompted meta-selector yields a 10% improvement on FOLIO’s challenging subset compared to GPT-4o. We will release our code and data to support future research.

## 1 Introduction

Mathematical reasoning with LLMs is typically approached through two paradigms. The first paradigm focuses on designing various prompting strategies to elicit detailed and natural language (NL) reasoning processes. This line of research continues from Chain-of-Thought prompting (Wei et al., 2023; Zhou et al., 2023; Zheng et al., 2024a). The second one leverages LLMs to generate solutions in the form of symbolic language, which can then be executed with external tools to derive the final answer (Olausson et al., 2023; Ye et al., 2023; Zhou et al., 2024a; Gu et al., 2024; Gou et al., 2024; Gao et al., 2023; Chen et al., 2023).

We propose HYBRIDMIND, which leverages the best of both worlds. By dynamically meta-selecting whether to reason in natural language, symbolic language, or a combination of the two, HYBRIDMIND tailors the solution style to the characteristics of each problem. Conceptual or explanatory tasks benefit from the clarity of step-by-step NL reasoning, while problems requiring exact heavy computation are more effectively handled in symbolic form (Gao et al., 2023; Zhou et al., 2024b). In other words, rather than forcing a single method to fit all problem types, HYBRIDMIND adaptively selects the right choice for a problem.

The contributions of our paper are threefold. 1) We propose HYBRIDMIND, a method that dynamically chooses between natural language reasoning and symbolic (Python code or first-order logic formulas) reasoning based on a specific reasoning problem. This “meta-selection” strategy ensures that the language model uses the most effective approach, natural language or symbolic reasoning, depending on each problem’s characteristics. 2) STaR finetuning over LLaMA-3.1-8B as a meta-selector yields a 4.4% accuracy gain on the FOLIO dataset, a logical reasoning dataset measuring complex logical reasoning capabilities of LLMs (Han et al., 2024). Using GPT-3.5-Turbo as a meta-selector yields a 10% improvement on the more challenging subset of FOLIO over GPT-4o. 3) We also identify key statistics beyond final performance and perform case study and prompt ablation study. We

also provide justifications for the performance difference of FOLIO with MATH to interpret the performance improvement of our method, discussing how the difference between Python code and first-order logic formulas could affect the meta-selection results.

## 2 Method

We consider four fundamental approaches involving reasoning with natural language and symbolic language. Natural Language (NL), Symbolic Language (SL), SymbolNL, and NLSymbol. NL represents that the LLM tackles the problem by generating a step-by-step breakdown of the reasoning process in natural language, guiding toward the solution. SL represents that the LLM is instructed to generate a symbolic language solution or symbolic form of the problem and it will be executed to arrive at the final solution. SymbolNL is a two-stage method. In the first stage, the LLM is prompted to write a solution in symbolic language or converts the problem into symbolic language. In the second stage, the LLM step-by-step analyzes the problem in natural language based on the symbolic form to obtain the final answer. NLSymbol is also a two-stage method. The LLM first generates a natural language solution or hints for solving the problem. Then the output will be used to generate symbolic language for solving the problem. Details can be found in Appendix.

## 3 Experiments

Our experiments are designed to test whether large language models can reason more effectively by dynamically choosing between natural language reasoning and symbolic language reasoning (or mixing both) rather than relying on a single approach. We explore whether the model can analyze each problem first and then pick the most suitable method — natural language, symbolic language, or a combination to arrive at a solution. By doing so, we aim to see if different types of problems in math and logical reasoning benefits from different styles of reasoning, and whether a “meta-selection” strategy can yield higher accuracy. Implementation details can be found in Appendix B.4.

### 3.1 Experiments on FOLIO

Method	WikiLogic	HybLogic	Avg.
<i>Baselines</i>			
Random	75.68	60.00	67.70
NL	80.18	60.87	70.35
SL	72.07	69.57	70.80
NL (8 MV)	<u>81.98</u>	66.09	73.89
<i>Finetune</i>			
2-class (base)	81.08	63.48	72.12
2-class SFT	81.08	60.87	70.35
STaR	<b>82.88</b>	66.69	74.78
<i>Prompting (0-shot)</i>			
GPT-3.5-Turbo	80.18	59.13	69.74
GPT-4o-mini	75.68	60.00	67.70
GPT-4o	80.18	66.96	73.45
o3-mini	80.18	62.61	71.24
<i>Prompting (2-shot)</i>			
GPT-4o-mini	81.08	66.09	73.45
GPT-4o	81.08	67.83	74.34
o3-mini	80.18	70.43	75.22
HYBRIDMIND (GPT-3.5-Turbo)	79.28	<b>73.91</b>	<b>76.55</b>
<i>Upperbound</i>			
Best-of-two (NL/SL)	92.79	83.48	88.05

Table 1: Model performance on FOLIO *test* set. **Bold** and underlined numbers indicate the best and second-best performance in each category, respectively. MV: Majority Vote. Best-of-2 (NL/SL): At least one of NL/SL is correct.

For FOLIO, NL and SL denote solving the questions based on Chain-of-Thought (Wei et al., 2023) and LINC (Olausson et al., 2023) using GPT-4o. We only sample once for both methods for efficiency.

**Baseline and upperbound.** Although state-of-the-art LLMs achieve better performance at reasoning in natural language compared to symbolic language on both of the two subsets we tested (Table 1), the accuracy of at least one of the methods being correct (best-of-2) is around 17% higher than NL or SL reasoning. This score is higher than performing reasoning solely in natural language or symbolic language, demonstrating the complementary strengths of these different reasoning approaches and the potential benefits of selecting between one of them effectively.

Method	Algebra	Counting & Probability	Geometry	Number & Theory	Intermediate & Algebra	Precalculus	Prealgebra	Avg.
<i>Baselines</i>								
NL	<b>92.67</b>	78.69	<b>61.38</b>	81.48	58.91	60.07	<b>89.78</b>	76.98
SL	72.37	78.48	46.56	71.48	44.30	28.94	80.60	62.00
SymbolNL	79.87	74.05	56.99	81.48	58.69	57.51	81.63	71.43
NLSymbol	64.95	82.07	50.31	81.11	51.72	39.38	85.07	65.24
MV (4 methods)	88.29	<b>84.60</b>	57.62	85.56	59.36	53.85	<b>89.78</b>	75.98
<i>Finetune</i>								
2-class (base)	72.37	78.48	46.56	71.78	44.30	28.94	80.60	62.00
2-class	<b>92.67</b>	78.69	<b>61.38</b>	<b>81.48</b>	58.91	<b>60.07</b>	<b>89.78</b>	<b>76.98</b>
4-class (base)	72.37	78.48	46.56	71.78	44.30	28.94	80.60	62.00
STaR (fine-tuning)	90.65	79.32	<b>61.17</b>	77.04	58.03	58.24	86.91	75.20
HYBRIDMIND (4-class)	<u>92.42</u>	79.54	<b>61.38</b>	<b>82.78</b>	<u>60.35</u>	<b>60.99</b>	<b>89.78</b>	<b>77.50</b>
<i>Prompting (0-shot)</i>								
GPT-3.5-Turbo	84.58	78.27	58.04	75.93	55.26	47.44	85.30	71.28
GPT-4o	91.58	<b>81.65</b>	60.96	79.26	<b>60.69</b>	57.69	88.98	76.64
o3-mini	87.87	79.75	58.87	76.30	56.81	50.18	87.26	73.24
<i>Upperbound</i>								
Best-of-2 (NL/SL)	95.79	87.55	66.81	89.63	68.77	64.10	93.57	82.84
Best-of-4	97.47	92.41	72.86	94.07	76.63	72.71	94.60	87.30

Table 2: Model performance on MATH *test* set. **Bold** and underlined numbers indicate the best and second-best performance in each category, respectively. MV: Majority Vote. Best-of-2 (NL/SL): At least one of NL/SL is correct. Best-of-4: At least one of NL/SL/SymbolNL/NLSymbol is correct.

### 3.2 Experiments on MATH

The results of experiments on MATH are in Table 2. NL baseline achieves the highest performance, while SL performs the worst among the four baseline methods<sup>1</sup>. NLSymbol shows apparently better results, indicating that outlining a reasoning path prior to generating code can improve reasoning with code, particularly in models with limited code capabilities.

In this section, we demonstrate that a smaller-scale model can be fine-tuned as a meta-selector to choose the most effective strategy for a given mathematical reasoning problem, achieving stronger performance than state-of-the-art baselines. Detailed analyses and meta-selector statistics on FOLIO and MATH (provided in the Appendix) further show that our approach produces a more balanced and principled selection between SL and NL.

## 4 Conclusion and Future Work

In this paper, we proposed HYBRIDMIND, meta selection of natural language, symbolic language or a mix of both for enhancing LLM math and logical reasoning. Extensive experiments on MATH and FOLIO show that dynamically selecting between reasoning with natural language and reasoning with symbolic language improves the reasoning performance of state-of-the-art LLMs, especially on the harder subset of FOLIO. We performed extensive analysis to understand the performance improvement achieved by HYBRIDMIND.

<sup>1</sup>While SL achieves better performance than CoT in Gao et al. (2023), our adopted dataset and base model are both different from the ones used in their study.

## References

- Janice Ahn, Rishu Verma, Renze Lou, Di Liu, Rui Zhang, and Wenpeng Yin. Large language models for mathematical reasoning: Progresses and challenges. *arXiv preprint arXiv:2402.00157*, 2024.
- Wenhu Chen, Xueguang Ma, Xinyi Wang, and William W. Cohen. Program of thoughts prompting: Disentangling computation from reasoning for numerical reasoning tasks. *Transactions on Machine Learning Research*, 2023. ISSN 2835-8856. URL <https://openreview.net/forum?id=YfZ4ZPt8zd>.
- Luyu Gao, Aman Madaan, Shuyan Zhou, Uri Alon, Pengfei Liu, Yiming Yang, Jamie Callan, and Graham Neubig. Pal: Program-aided language models, 2023. URL <https://arxiv.org/abs/2211.10435>.
- Peizhong Gao, Ao Xie, Shaoguang Mao, Wenshan Wu, Yan Xia, Haipeng Mi, and Furu Wei. Meta reasoning for large language models, 2024. URL <https://arxiv.org/abs/2406.11698>.
- Zhibin Gou, Zhihong Shao, Yeyun Gong, yelong shen, Yujiu Yang, Minlie Huang, Nan Duan, and Weizhu Chen. ToRA: A tool-integrated reasoning agent for mathematical problem solving. In *The Twelfth International Conference on Learning Representations*, 2024. URL <https://openreview.net/forum?id=Ep0TtjVoap>.
- Alex Gu, Baptiste Rozière, Hugh Leather, Armando Solar-Lezama, Gabriel Synnaeve, and Sida I. Wang. Cruxeval: A benchmark for code reasoning, understanding and execution, 2024. URL <https://arxiv.org/abs/2401.03065>.
- Shashank Gupta, Vaishnavi Shrivastava, Ameet Deshpande, Ashwin Kalyan, Peter Clark, Ashish Sabharwal, and Tushar Khot. Bias runs deep: Implicit reasoning biases in persona-assigned llms. *arXiv preprint arXiv:2311.04892*, 2023.
- 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/>.
- Dan Hendrycks, Collin Burns, Saurav Kadavath, Akul Arora, Steven Basart, Eric Tang, Dawn Song, and Jacob Steinhardt. Measuring mathematical problem solving with the math dataset, 2021. URL <https://arxiv.org/abs/2103.03874>.
- Shima Imani, Liang Du, and Harsh Shrivastava. Mathprompter: Mathematical reasoning using large language models, 2023. URL <https://arxiv.org/abs/2303.05398>.
- Takeshi Kojima, Shixiang Shane Gu, Machel Reid, Yutaka Matsuo, and Yusuke Iwasawa. Large language models are zero-shot reasoners, 2023. URL <https://arxiv.org/abs/2205.11916>.
- Tianyu Liu, Yijia Xiao, Xiao Luo, Hua Xu, Wenjin Zheng, and Hongyu Zhao. Geneverse: A collection of open-source multimodal large language models for genomic and proteomic research. In Yaser Al-Onaizan, Mohit Bansal, and Yun-Nung Chen (eds.), *Findings of the Association for Computational Linguistics: EMNLP 2024*, pp. 4819–4836, Miami, Florida, USA, November 2024. Association for Computational Linguistics. doi: 10.18653/v1/2024.findings-emnlp.277. URL <https://aclanthology.org/2024.findings-emnlp.277/>.
- Qing Lyu, Shreya Havaladar, Adam Stein, Li Zhang, Delip Rao, Eric Wong, Marianna Apidianaki, and Chris Callison-Burch. Faithful chain-of-thought reasoning, 2023. URL <https://arxiv.org/abs/2301.13379>.

135 Theo Olausson, Alex Gu, Ben Lipkin, Cedegao Zhang, Armando Solar-Lezama, Joshua Tenenbaum,  
136 and Roger Levy. Linc: A neurosymbolic approach for logical reasoning by combining language  
137 models with first-order logic provers. In *Proceedings of the 2023 Conference on Empirical Methods*  
138 *in Natural Language Processing*, pp. 5153–5176. Association for Computational Linguistics,  
139 2023. doi: 10.18653/v1/2023.emnlp-main.313. URL [http://dx.doi.org/10.18653/v1/2023.](http://dx.doi.org/10.18653/v1/2023.emnlp-main.313)  
140 [emnlp-main.313](http://dx.doi.org/10.18653/v1/2023.emnlp-main.313).

141 OpenAI, Josh Achiam, and Others. Gpt-4 technical report, 2024. URL [https://arxiv.org/abs/](https://arxiv.org/abs/2303.08774)  
142 [2303.08774](https://arxiv.org/abs/2303.08774).

143 Yuval Reif and Roy Schwartz. Beyond performance: Quantifying and mitigating label bias in  
144 llms. In *Proceedings of the 2024 Conference of the North American Chapter of the Association*  
145 *for Computational Linguistics: Human Language Technologies (Volume 1: Long Papers)*, pp.  
146 6784–6798, 2024.

147 Stuart Russell and Peter Norvig. *Artificial Intelligence: A Modern Approach*. Prentice Hall, 3 edition,  
148 2010.

149 Pauli Virtanen, Ralf Gommers, Travis E Oliphant, Matt Haberland, Tyler Reddy, David Cournapeau,  
150 Evgeni Burovski, Pearu Peterson, Warren Weckesser, Jonathan Bright, et al. Scipy 1.0: fundamental  
151 algorithms for scientific computing in python. *Nature methods*, 17(3):261–272, 2020.

152 Xuezhi Wang, Jason Wei, Dale Schuurmans, Quoc Le, Ed Chi, Sharan Narang, Aakanksha Chowdh-  
153 ery, and Denny Zhou. Self-consistency improves chain of thought reasoning in language models,  
154 2023. URL <https://arxiv.org/abs/2203.11171>.

155 Jason Wei, Xuezhi Wang, Dale Schuurmans, Maarten Bosma, Brian Ichter, Fei Xia, Ed Chi, Quoc Le,  
156 and Denny Zhou. Chain-of-thought prompting elicits reasoning in large language models, 2023.  
157 URL <https://arxiv.org/abs/2201.11903>.

158 Fangzhi Xu, Qika Lin, Jiawei Han, Tianzhe Zhao, Jun Liu, and Erik Cambria. Are large language  
159 models really good logical reasoners? a comprehensive evaluation and beyond. *IEEE Transactions*  
160 *on Knowledge and Data Engineering*, 2025.

161 Jundong Xu, Hao Fei, Liangming Pan, Qian Liu, Mong-Li Lee, and Wynne Hsu. Faithful logical  
162 reasoning via symbolic chain-of-thought, 2024. URL <https://arxiv.org/abs/2405.18357>.

163 Kaiyu Yang, Gabriel Poesia, Jingxuan He, Wenda Li, Kristin Lauter, Swarat Chaudhuri, and Dawn  
164 Song. Formal mathematical reasoning: A new frontier in ai, 2024. URL [https://arxiv.org/](https://arxiv.org/abs/2412.16075)  
165 [abs/2412.16075](https://arxiv.org/abs/2412.16075).

166 Xi Ye, Qiaochu Chen, Isil Dillig, and Greg Durrett. SatLM: Satisfiability-aided language models  
167 using declarative prompting. In *Thirty-seventh Conference on Neural Information Processing*  
168 *Systems*, 2023. URL <https://openreview.net/forum?id=TqW5PL1Poi>.

169 Eric Zelikman, Yuhuai Wu, Jesse Mu, and Noah Goodman. Star: Bootstrapping reasoning with  
170 reasoning. *Advances in Neural Information Processing Systems*, 35:15476–15488, 2022.

171 Jiaxin Zhang, Zhongzhi Li, Mingliang Zhang, Fei Yin, Chenglin Liu, and Yashar Moshfeghi. Geoeval:  
172 benchmark for evaluating llms and multi-modal models on geometry problem-solving. *arXiv*  
173 *preprint arXiv:2402.10104*, 2024.

174 Wayne Xin Zhao, Kun Zhou, Junyi Li, Tianyi Tang, Xiaolei Wang, Yupeng Hou, Yingqian Min,  
175 Beichen Zhang, Junjie Zhang, Zican Dong, et al. A survey of large language models. *arXiv*  
176 *preprint arXiv:2303.18223*, 2023.

177 Huaixiu Steven Zheng, Swaroop Mishra, Xinyun Chen, Heng-Tze Cheng, Ed H. Chi, Quoc V Le,  
178 and Denny Zhou. Take a step back: Evoking reasoning via abstraction in large language models,  
179 2024a. URL <https://arxiv.org/abs/2310.06117>.

180 Yaowei Zheng, Richong Zhang, Junhao Zhang, Yanhan Ye, Zheyang Luo, Zhangchi Feng, and  
181 Yongqiang Ma. Llamafactory: Unified efficient fine-tuning of 100+ language models. *arXiv*  
182 *preprint arXiv:2403.13372*, 2024b.

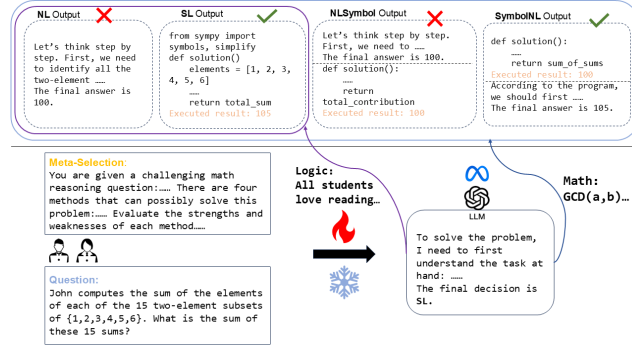


Figure 1: Illustration of HYBRIDMIND. HYBRIDMIND enables the model to analyze the problem and choose the most suitable approach among NL, SL, NLSymbol, and SymbolNL. We have different reasoning paths for logic and mathematical reasoning problems. In this example, HYBRIDMIND selects SL for a mathematical reasoning problem, which leads to the right solution (105).

- 183 Denny Zhou, Nathanael Schärli, Le Hou, Jason Wei, Nathan Scales, Xuezhi Wang, Dale Schuurmans,  
184 Claire Cui, Olivier Bousquet, Quoc V Le, and Ed H. Chi. Least-to-most prompting enables  
185 complex reasoning in large language models. In *The Eleventh International Conference on*  
186 *Learning Representations*, 2023. URL <https://openreview.net/forum?id=WZH7099tgfM>.
- 187 Jin Peng Zhou, Charles Staats, Wenda Li, Christian Szegedy, Kilian Q. Weinberger, and Yuhuai Wu.  
188 Don't trust: Verify – grounding llm quantitative reasoning with autoformalization, 2024a. URL  
189 <https://arxiv.org/abs/2403.18120>.
- 190 Pei Zhou, Jay Pujara, Xiang Ren, Xinyun Chen, Heng-Tze Cheng, Quoc V. Le, Ed H. Chi, Denny  
191 Zhou, Swaroop Mishra, and Huaixiu Steven Zheng. Self-discover: Large language models  
192 self-compose reasoning structures, 2024b. URL <https://arxiv.org/abs/2402.03620>.
- 193 Yukai Zhou, Zhijie Huang, Feiyang Lu, Zhan Qin, and Wenjie Wang. Don't say no: Jailbreaking llm  
194 by suppressing refusal. *arXiv preprint arXiv:2404.16369*, 2024c.

## 195 A Appendix / supplemental material

### 196 A.1 Method illustraion

### 197 A.2 Related Work

### 198 Reasoning with natural language and symbolic language

199 **Reasoning with NL and SL.** Recent research has substantially advanced the logical and mathe-  
200 matical reasoning capabilities of LLMs by incorporating step-by-step reasoning (Kojima et al., 2023;  
201 Wei et al., 2023; Wang et al., 2023; Zhou et al., 2023; Zheng et al., 2024a). These methods encourage  
202 models to produce explicit intermediate reasoning steps in natural language. While these methods  
203 have proven successful across various tasks, they depend entirely on language representations and  
204 can yield incoherent or unreliable solutions when confronted with logically complex problems (Han  
205 et al., 2024; Olausson et al., 2023).

206 Symbolic-based approaches take questions in natural language form and generate symbolic forms  
207 alone such as Python code or First-Order Logic or generate symbolic forms together with natural  
208 language reasoning steps(Gou et al., 2024; Gao et al., 2023; Chen et al., 2023; Imani et al., 2023;  
209 Zhou et al., 2024a). External code executors or symbolic engines are then leveraged for verification or  
210 deriving the final answer These methods guide the model to generate executable symbolic language  
211 forms, enabling direct checks of intermediate or final outputs. Meanwhile, SymbCoT (Xu et al.,  
212 2024) converts premises into a symbolic format and uses an LLM-based solver and verifier to validate  
213 each step. By grounding reasoning in a checkable programmatic or symbolic form, these methods  
214 can offer stronger verifiability than solely natural language-based solutions due to incorporation of  
215 an external execution engine.

**Meta selection with LLMs** Zhou et al. (2024b); Gao et al. (2024) use meta prompts to enable LLMs to choose multiple fundamental reasoning modules, such as critical thinking and step-by-step reasoning, and organize them into a clear reasoning structure that guides the model during decoding. However, these methods only adopt natural language reasoning while HYBRIDMIND exploits the of natural language reasoning and symbolic language reasoning, taking into consideration the characteristics of logical and mathematical reasoning problems

### A.3 Example figures

Table 3: HYBRIDMIND output example

Input	
<b>Premises</b> All runners who ranked in the top 1000 in the Boston Marathon finished the marathon within 3 hours. All runners from the University Running Club ranked in the top 1000 in the Boston Marathon. Some of the runners wearing T-shirts with the bulldog logo are majoring in Computer Science. Tim is a runner and if Tim finishes the Boston Marathon within 3 hours, then he is not both from University Running Club and ranked in the top 1000 in the Boston Marathon. Tim is a runner and if Tim majors in computer science, then he is either wearing a T-shirt with the bulldog logo or is a runner from the University Running Club. <b>Conclusion</b> Tim majors in computer science.	<b>Input</b> <b>Question</b> Six-digit integers will be written using each of the digits 1 through 6 exactly once per six-digit integer. How many different positive integers can be written such that all pairs of consecutive digits of each integer are relatively prime? (Note: 1 is relatively prime to all integers.) <b>Rationale</b> ... given the limited number of permutations and the straightforward nature of checking the pairwise condition—the SL (Symbolic Language) method is the most effective. ... <b>Outputs</b> Selected method: SL (Symbolic Language).
<b>Output</b> ... While the premises are not overly complex, the relationships between the entities involved are intertwined and require careful consideration. Given the structured nature of the premises and the need for precise logical deductions ... Selected method: SL (Symbolic Language).	

(a) HYBRIDMIND output on FOLIO.

(b) HYBRIDMIND output on MATH.

### A.4 Method details

The SymbolNL approach emphasizes the importance of analyzing the code’s execution through natural language reasoning, allowing for corrections and insights even if the initial code has errors. In contrast, NLSymbol starts with a natural language outline to clarify the problem before writing the code, which can enhance the likelihood of successful implementation. These two formats leverage both reasoning forms to improve problem-solving effectiveness. The optimal choice among these approaches varies significantly across problem types — some problems benefit from the precision of symbolic reasoning, while others require the flexibility of natural language decomposition (Olausson et al., 2023; Chen et al., 2023; Zhou et al., 2024a).

This observation motivates HYBRIDMIND, a meta-selection framework that analyzes each problem’s characteristics to determine the most suitable reasoning approach. As illustrated in Figure 1, before using a specific approach to derive the final answer, HYBRIDMIND uses a meta-selection module to first analyze the problem and decide which reasoning approach to apply among NL, SL, SymbolNL and NLSymbol, before using the selected approach to generate a solution.

The core intuition behind HYBRIDMIND is that not all problems benefit equally from the same reasoning approach (Zhou et al., 2024c). Some problems may require step-by-step natural language reasoning, while others are better suited for symbolic solutions. HYBRIDMIND empowers the model to dynamically adapt its strategy by selecting the most appropriate method, ensuring flexibility and maximizing performance across different problem types.

### A.5 Choice of symbolic language

In mathematical reasoning and logical reasoning, the choice between symbolic languages such as Python and formal logical systems like first-order logic (FOL) arises from their fundamental differences in execution and purpose (Russell & Norvig, 2010). Python is a procedural language designed for computation, making it well-suited for mathematical reasoning tasks that involve

numerical calculations, symbolic algebra, and algorithmic problem-solving. We include examples where NL reasoning is better and where SL reasoning is better in the appendix. With built-in support for arithmetic operations, iterative processes, and specialized libraries such as NumPy and SymPy, Python efficiently performs both exact and approximate computations (Gao et al., 2023; Zhou et al., 2024a). In contrast, FOL provides a declarative framework for logical reasoning, allowing for the formal specification of knowledge, relationships, and inferential rules (Han et al., 2024; Olausson et al., 2023; Lyu et al., 2023). Unlike Python, which executes a sequence of computational steps, FOL is primarily used for theorem proving and rule-based inference, relying on formal logical operators, quantifiers, and axioms to establish truths within a domain. While Python can simulate logical reasoning through symbolic computation, its execution semantics is step-wise and does not allow general deduction, which may result in unnecessary steps to reach the final value (Russell & Norvig, 2010). Consequently, we select Python to be the symbolic language for math reasoning, and FOL to be the symbolic language for logical reasoning.

## A.6 Training data generation

We present a systematic procedure for generating labeled training data for our meta-selector. We use a state-of-the-art LLM, GPT-4o to solve a logical or math reasoning problem based on NL, SL, NLSymbol, and SymbolNL. Each problem is solved in four aforementioned ways using the same LLM. Labels are automatically created by measuring answer correctness with exact match: whichever approach yields the correct solution is treated as the “label” for that problem. If multiple approaches work, one is chosen at random<sup>2</sup>. These “correctness-labeled” examples are then used to train a meta-selector model that can predict which solution strategy is most likely to be correct when given a new problem.

This training data generation method has two major advantages. 1) The data generation is fully automatic with the label directly generated with the original logical or math reasoning dataset answer labels. 2) Different prompting strategies can complement one another; the meta-selector can learn to pick the more suitable strategy for a given problem.

## A.7 Result details

Our experiments are designed to test whether large language models can reason more effectively by dynamically choosing between natural language reasoning and symbolic language reasoning (or mixing both) rather than relying on a single approach. We explore whether the model can analyze each problem first and then pick the most suitable method — natural language, symbolic language, or a combination to arrive at a solution. By doing so, we aim to see if different types of problems in math and logical reasoning benefits from different styles of reasoning, and whether a “meta-selection” strategy can yield higher accuracy. Implementation details can be found in Appendix B.4.

## A.8 Datasets

In our experiments, we consider two challenging datasets: MATH (Hendrycks et al., 2021) for evaluating mathematical reasoning and FOLIO (Han et al., 2024) for evaluating logical reasoning. These two datasets are widely adopted (OpenAI et al., 2024; Xu et al., 2025) and offer a range of categories and difficulty levels, enabling us to conduct various analyses. The MATH dataset comprises 7,498 training samples and 5,000 testing samples in total across 7 categories and 5 difficulty levels; and the FOLIO dataset contains 1,001 training samples and 226 testing samples spanning 2 categories. Our finetuning methods are trained on the entire training set and evaluated based on the entire testing set. Since we reported model performance by category, our experimental setting consists of a total of 9 subsets. Detailed sample counts for each category are provided in Appendix B.5.

## A.9 Model choice

We selected several advanced LLMs as the base models for our meta-selector. For *open-source* LLMs, we employ Llama-3.1-8B-Instruct; for *proprietary* LLMs, we utilize GPT-3.5-turbo, GPT-4o, and

<sup>2</sup>Notably, having 4 labels for whether each of the method solves the problem correctly produces a more rigorous setting, but it leads to severe label imbalance which is hard to be mitigated in our initial experiments



o3-mini.<sup>3</sup> Finetuning with our training data reduces the bias of keeping selecting NL reasoning, but hurts the final performance. We then transform these base models into meta-selectors by applying either finetuning or prompting. To ensure a fair evaluation, we use 8 shots for both NL and SL, whereas SymbolNL and NLSymbol both use 4 shots at each stage. Details of the complete prompt texts are in the Appendix B.6.

For open-source models, we consider both zero-shot prompting approach and fine-tuning approaches including Supervised Fine-Tuning (SFT) (Zhao et al., 2023; Zheng et al., 2024b; Liu et al., 2024) and iterative finetuning strategy, STaR (Zelikman et al., 2022), to train a meta-selector. Our fine-tuning experiments are based on LLaMA-factory (Zheng et al., 2024b), which is widely employed in LLM development (Liu et al., 2024; Zhao et al., 2023). Proprietary models are based on both the zero-shot prompting and few-shots prompting approaches to act as a meta-selector.

## A.10 NL, SL, NLSymbol and SymbolNL

We perform meta selection over NL, SL, NLSymbol, and SymbolNL for MATH since all of the four methods achieve high performance on MATH. Furthermore, the oracle accuracy of at least one of the methods is at least 10.3% higher than that of only one other methods, showing the potential of these methods to be complementary. However, NLSymbol and SymbolNL both achieve very low performance on FOLIO. The first-order logic formulas generated by NLSymbol in the second stage tend to be unexecutable and SymbolNL and the output generated by SymbolNL in the first stage would distract the model from generating a reasoning path. Therefore, for FOLIO we adopt two-class meta-selection over NL and SL without considering NLSymbol and SymbolNL.

## A.11 Analysis on MATH

**HYBRIDMIND can mitigate the bias of choosing correct methods.** The bias of LLMs is shown to have negative effects towards generating correct outputs in reasoning tasks (Gupta et al., 2023; Reif & Schwartz, 2024). We noticed that small-scale models before fine-tuning, tended always to choose NL as their final solution, and thus, it is important to investigate the output distribution by choices under both mathematical reasoning and logical reasoning questions. We then investigated the frequency of choice to solve problems in the MATH and FOLIO datasets, shown in Table 7. Details and case studies are included in the Appendix.

**Case study.** Table 5a provides an example where HYBRIDMIND accurately discovers complex logical relationships between entities in different premises and identifies “the need for precise logical deductions to reach the conclusion.” It successfully selects SL. This underscores its capacity to dynamically select the most suitable reasoning approach.

# B Analysis

## B.1 Analysis on FOLIO

The FOLIO dataset consists of two subsets: WikiLogic and HybLogic. WikiLogic presents simpler logical reasoning tasks (1 – 5 reasoning steps), whereas HybLogic includes more logically complex problems (5 – 8 reasoning steps). HYBRIDMIND improves model performance on HybLogic set by 13.04% (Table 1).

**Meta-selector statistics.** Initially, the meta selector utilized only a simple instruction prompt, providing limited performance improvement (Vanilla in Table 4a). Subsequent prompt engineering informs the model of the weaknesses of NL: NL reasoning can be prone to errors if the problem is especially logically challenging and its chain-of-thought becomes long or convoluted. It also informs the model of the weakness of SL reasoning (weakness info): SL is not ideal if the premises and conclusion rely heavily on nuanced interpretations of language, or involve concepts that are hard

<sup>3</sup>In our initial experiments, we have also tested other models with 3-shot and 8-shot prompting, such as Qwen2.5-14B-Instruct, Qwen-32B-Coder-Instruct, and CodeLlama for meta-selection, however, these models only selected reasoning with NL for all of the tested examples, even when they are given multiple shots, which rendered them a non-optimal choice for testing with our method.

Table 4: Analysis on FOLIO

Strategy	WikiLogic	HybLogic	Avg.
NL	80.18	60.87	70.35
Vanilla	80.18	62.61	71.24
+ wn info	80.18	59.13	69.47
+ wn info + 2-shots	79.28	73.91	76.55

(a) Comparison of different prompting strategies on the FOLIO *test* set using GPT-3.5-Turbo.

Subset	w/o wn info		w/ wn info	
	NL	SL	NL	SL
WikiLogic	111	0	102	9
HybLogic	110	5	31	84

(b) Distribution of NL vs. SL selected by HYBRIDMIND on WikiLogic and HybLogic w/o weakness info and w/ weakness info.

338 to formalize in FOL. Details and case studies are included in the Appendix. We also show detailed  
 339 analysis on MATH in the Appendix.

## 340 B.2 Analysis on MATH

Table 5: HYBRIDMIND output example

Input
<b>Premises</b> All runners who ranked in the top 1000 in the Boston Marathon finished the marathon within 3 hours. All runners from the University Running Club ranked in the top 1000 in the Boston Marathon. Some of the runners wearing T-shirts with the bulldog logo are majoring in Computer Science. Tim is a runner and if Tim finishes the Boston Marathon within 3 hours, then he is not both from University Running Club and ranked in the top 1000 in the Boston Marathon. Tim is a runner and if Tim majors in computer science, then he is either wearing a T-shirt with the bulldog logo or is a runner from the University Running Club. <b>Conclusion</b> Tim majors in computer science.
Output
... While the premises are not overly complex, the relationships between the entities involved are intertwined and require careful consideration. Given the structured nature of the premises and the need for precise logical deductions ... <b>Selected method: SL (Symbolic Language).</b>

(a) HYBRIDMIND output on FOLIO.

Input
<b>Question</b> Six-digit integers will be written using each of the digits 1 through 6 exactly once per six-digit integer. How many different positive integers can be written such that all pairs of consecutive digits of each integer are relatively prime? (Note: 1 is relatively prime to all integers.)
Rationale
... given the limited number of permutations and the straightforward nature of checking the pairwise condition—the SL (Symbolic Language) method is the most effective. ...
Outputs
<b>Selected method: SL (Symbolic Language).</b>

(b) HYBRIDMIND output on MATH.

341 **HYBRIDMIND can mitigate the bias of choosing correct methods.** The bias of LLMs is shown  
 342 to have negative effects towards generating correct outputs in reasoning tasks (Gupta et al., 2023; Reif  
 343 & Schwartz, 2024). We noticed that small-scale models before fine-tuning, tended always to choose  
 344 NL as their final solution, and thus, it is important to investigate the output distribution by choices  
 345 under both mathematical reasoning and logical reasoning questions. Therefore, we investigated  
 346 the frequency of choice to solve problems in the MATH and FOLIO datasets, shown in Table 7.  
 347 Considering the case of random choice, the proportion of each choice will be 25% for MATH and 50  
 348 % for FOLIO, while the outputs from our method do not follow this proportion, and our model did  
 349 not choose the solution randomly but performs rigorous inference. Moreover, base models before  
 350 fine-tuning had 100% NL choices in the outputs, where our generated outputs only contained 94.10%  
 351 NL choices for MATH and 58.84% NL choices for FOLIO. Considering the high performance of  
 352 NL in solving mathematical reasoning problems, the outputs from our method should have a higher  
 353 proportion of NL choices while it can also select correct choices for other questions that NL cannot  
 354 resolve, and this improvement also happened in our testing results based on the FOLIO dataset. We  
 355 also tried to downsample NL in our training data and increase the proportion of SL for both problems  
 356 as a preliminary attempt to mitigate existing label bias, but the final performance was reduced by  
 357 5.3% in this setting. Therefore, the improvement of HYBRIDMIND can also be interpreted by the  
 358 contribution of reducing selection bias compared with other baselines.

359 **Case study** Table 5b provides an example from the number theory category, where HYBRIDMIND  
 360 selects the correct reasoning path in the problem that NL cannot address. HYBRIDMIND makes

highest percentage of improvement in the problems from this category, and HYBRIDMIND will utilize the shots containing the description of problem statistics (e.g., length) as well as the most effective solution and then make a final decision. This specific property allows us to handle complicated metathetical reasoning problems with a stronger meta-selector. We include more examples where NL reasoning is better and where SL reasoning is better in the appendix.

### B.3 Comparing FOLIO with MATH

HYBRIDMIND produces more substantial improvement on FOLIO than on MATH, as manifested by the smaller p-value. We hypothesize this could be because that FOL formulas are designed to capture the original FOL reasoning question. In contrast, reasoning with Python code like SL for MATH follows a procedural approach that more closely resembles natural language (NL) reasoning than generating FOL formulas for FOLIO. The smaller difference between best-of-2 performance and the NL or SL approach alone indicates that these two classes are more alike, which ultimately makes the classifier more difficult to train because distinguishing between more similar classes is more challenging. For future work, we will explore auto-formalization using Isabelle-form formal specifications as SL for MATH reasoning (Yang et al., 2024; Zhou et al., 2024a).

**Prompting results.** In the experiment of zero-shot prompting Llama-3.1-8B-Instruct model, the model performance falls closer to the average of the four approaches. This suggests that more advanced models, like o3-mini, are better at selecting the optimal approach based on the problem. Moreover, prompting the Llama-3.1-8B model to act as a meta-selector cannot surpass NL reasoning with GPT-4o, which shows the limitation of the model’s capacity for selecting the correct strategy for solving mathematical reasoning problems. On the other hand, prompting advanced models as a meta selector cannot lead to better performance either, which also demonstrates the constraint of prompting in this type of question.

**Finetuning results.** Methods incorporating NL reasoning (NL, NLSymbol, and SymbolNL) outperform SL, emphasizing the importance of NL in mathematical reasoning. Overall, HYBRIDMIND demonstrates higher selection accuracy, effectively identifying the most appropriate approach for each problem. HYBRIDMIND also surpasses selectors with the same scale or larger scales in different categories by 0.3% at least and 2.0% at most. Our improvement is statistically significant, as validated by Wilcoxon Rank Sum test (Virtanen et al., 2020) ( $p = 0.002$ ) shown in Appendix B.7.

HYBRIDMIND particularly works well in questions belonging to geometry, number theory, precalculus, and prealgebra, and thus it shows the strong capacity of a fine-tuned model for solving complicated geometric or numerical problems. Questions with these types are discussed to be more challenging than other categories for general LLMs (Zhang et al., 2024; Ahn et al., 2024) to solve. A prompt and paired radiational generated by HYBRIDMIND can be found in Table 5b, which matches well with regular reasoning paths as a meta selector. Moreover, for finetuning experiments, we have tried different types of meta-selectors, which are summarized in Appendix B.4. The best performer for MATH is a meta-selector trained based on labels created with four methods. Our final choice also outperforms STaR with fewer shots. This could be because that training the meta-selector for math reasoning requires more examples.

Finally, while majority voting performs well on problems with higher accuracy, it struggles with more challenging problems, such as those in the precalculus category. We note that both majority voting (MV)<sup>4</sup> and random choosing cannot surpass our method as well as NL, especially in challenging problems in the precalculus category. This observation demonstrates the necessity of fine-tuning for improving LLMs’ mathematical reasoning.

Overall, in this section, we demonstrate that a smaller-scale model can be finetuned as a meta-selector in selecting the most suitable strategy for a given mathematical reasoning problem, leading to better performance than state-of-the-art baselines.

<sup>4</sup>When comparing majority voting, it’s important to consider the trade-off between effectiveness and cost. Majority voting requires calling the language model to execute each of the four methods, whereas HYBRIDMIND only necessitates executing a single selected method.

## B.4 Implementation details

**Prompting-based HYBRIDMIND.** In our experiment, we considered prompting different LLMs with various instructions to work as a meta selector for generating correct solution based on the given question. Regarding the format of prompting, we tried two different zero-shot instruction design and made comparison. We also tried few-shot prompting design, whose number of examples (shots) in the instruction is at most 2. The shot contains question, answer, and the rationale generated by GPT-o3-mini.

**Fine-tuning-based HYBRIDMIND.** We also performed experiments to fine-tune Llama 3.1-8B to perform complicated reasoning tasks. We prepared the training datasets by sampling the choices and paired one correct choice with one question. The implementation of SFT is based on Llama-factory and the hyper-parameter setting is same as the example configuration template. The implementation of STaR is modified from the original code base. STaR utilizes models to generate rationale and then select the correct samples as well as the generated rationales to perform fine-tuning. To fine-tune our models, we utilize both NVIDIA H100 GPU and Together AI API. Details of our training codes can be found in the attached zip file to this submission.

**Information of baseline models.** In our experiment, we considered prompting closed-source LLMs, including GPT-4o, GPT-3.5-turbo, GPT-o3-mini, and fine-tuning open-sourced LLM including Llama 3.1-8b. The prompting experiment is performed with different strategies (NL, SL, SymbolNL, NLSymbol) for solving the question directly. MV represents performing majority voting based on the methods’ outputs. We also considered fine-tuning the base model with different strategies, including making binary choices (NL or SL), or making 4-class choices (NL, SL, SymbolNL, NLSymbol). The fine-tuning process is for making a meta selector.

## B.5 Statistics of testing data

Please check Tables 6, 7, and 8.

Category	Algebra	Counting & Probability	Geometry	Number & Theory	Intermediate & Algebra	Precalculus	Prealgebra	Total
Number	1187	474	479	540	903	546	871	5000

Table 6: Number of samples in the testing set by categories for MATH.

Method	MATH		FOLIO	
	Freq.	Prop.	Freq.	Prop.
NL	4705	94.10%	133	58.84%
SL	124	2.48%	93	41.16%
SymbolNL	163	3.26%	-	-
NLSymbol	8	0.16%	-	-

Table 7: The frequency and proportion of selected choices by HYBRIDMIND on MATH and FOLIO test set. We implemented 4-choice selector for MATH dataset and 2-choice selector for FOLIO.

Category	WikiLogic	HybLogic	Total
Number	111	115	226

Table 8: Number of samples in the testing set by categories for FOLIO.

## B.6 Prompts and shots

Below are the prompts and shots used for the four approaches: NL, SL, SymbolNL, and NLSymbol.

## 434 B.6.1 NL

```

435 nl_system_prompt = '''
436 You are a helpful assistant who is good at solving math problems. You
437     should follow the guidelines below:
438 - Present the final result in LaTeX using a '\boxed{}' without any units.
439 - Utilize the 'pi' symbol, and simplify all fractions and square roots
440     without converting them to decimal values.
441 '''
442
443 nl_instruction_prompt = "Please_think_step_by_step._"
444
445 nl_math_shots= [
446     '''Question: Kevin Kangaroo begins hopping on a number line at 0. He
447         wants to get to 1, but he can hop only  $\frac{1}{3}$  of the distance.
448         Each hop tires him out so that he continues to hop  $\frac{1}{3}$  of
449         the remaining distance. How far has he hopped after five hops?
450         Express your answer as a common fraction.
451 Answer: Let's think step by step
452 Kevin hops  $\frac{1}{3}$  of the remaining distance with every hop.
453 His first hop takes  $\frac{1}{3}$  closer.
454 For his second hop, he has  $\frac{2}{3}$  left to travel, so he hops forward  $\frac{2}{3} \cdot \frac{1}{3}$ .
455 For his third hop, he has  $(\frac{2}{3})^2$  left to travel, so he hops forward  $(\frac{2}{3})^2 \cdot \frac{1}{3}$ .
456 In general, Kevin hops forward  $(\frac{2}{3})^{k-1} \cdot \frac{1}{3}$  on his  $k$ th hop.
457 We want to find how far he has hopped after five hops.
458 This is a finite geometric series with first term  $\frac{1}{3}$ , common ratio  $\frac{2}{3}$ ,
459     and five terms.
460 Thus, Kevin has hopped  $\frac{\frac{1}{3}(1-(\frac{2}{3})^5)}{1-\frac{2}{3}} = \boxed{\frac{211}{243}}$ .
461 The answer is  $\frac{211}{243}$ ''',
462
463     '''Question: What is the area of the region defined by the equation  $x^2 + y^2 - 7 = 4y - 14x + 3$ ?
464 Answer: Let's think step by step
465 We rewrite the equation as  $x^2 + 14x + y^2 - 4y = 10$  and then complete
466     the square,
467     resulting in  $(x+7)^2 - 49 + (y-2)^2 - 4 = 10$ ,
468     or  $(x+7)^2 + (y-2)^2 = 63$ .
469 This is the equation of a circle with center  $(-7, 2)$  and radius  $\sqrt{63}$ ,
470     so the area of this region is  $\pi r^2 = \boxed{63\pi}$ .
471 The answer is  $63\pi$ ''',
472
473     '''Question: If  $x^2 + y^2 = 1$ , what is the largest possible value of  $|x| + |y|$ ?
474 Answer: Let's think step by step
475 If  $(x, y)$  lies on the circle,
476     so does  $(x, -y)$ ,  $(-x, -y)$ , and  $(-x, y)$ , (which all give the same
477     value of  $|x| + |y|$ ),
478     so we can assume that  $x \geq 0$  and  $y \geq 0$ .
479     Then  $|x| + |y| = x + y$ . Squaring, we get
480      $(x + y)^2 = x^2 + 2xy + y^2 = 1 + 2xy$ .
481     Note that  $(x - y)^2 \geq 0$ .
482     Expanding, we get  $x^2 - 2xy + y^2 \geq 0$ , so  $2xy \leq x^2 + y^2 = 1$ .
483     Hence,  $1 + 2xy \leq 2$ , which means  $x + y \leq \sqrt{2}$ .
484     Equality occurs when  $x = y = \frac{1}{\sqrt{2}}$ ,
485     so the maximum value of  $|x| + |y|$  is  $\boxed{\sqrt{2}}$ .
486     The answer is  $\sqrt{2}$ '''
487 ]
488
489
490
491
492
493
494
495
496

```

497 '''Question: If  $f(x) = \frac{ax+b}{cx+d}$ ,  $abcd \neq 0$  and  $f(f(x)) = x$  for  
 498 all  $x$  in the domain of  $f$ , what is the value of  $a+d$ ?  
 499 Answer: Let's think step by step  
 500 The condition  $f(f(x)) = x$  means that  $f$  is the inverse of itself,  
 501 so its graph is symmetrical about the line  $y = x$ .  
 502 With a rational function of this form, we will have two asymptotes:  
 503 a vertical one at  $x = -d/c$  if  $cx+d$  does not divide  $ax+b$ ,  
 504 and a horizontal one at  $y = a/c$ ,  
 505 if we take the limit of  $f(x)$  as  $x$  goes to  $\pm\infty$ .  
 506 In order for  $f$  to be its own inverse, the intersection of the  
 507 asymptotes must lie on the line  $y = x$   
 508 so that it and its asymptotes reflect onto themselves.  
 509 This means that  $-d/c = a/c$ ,  
 510 and therefore  $-d = a$  and  $a+d = \boxed{0}$ .  
 511 The answer is 0''',  
 512  
 513  
 514 '''Question: A math teacher requires Noelle to do one homework assignment  
 515 for each of the first five homework points she wants to earn; for  
 516 each of the next five homework points, she needs to do two homework  
 517 assignments; and so on, so that to earn the  $n^{\text{th}}$  homework  
 518 point, she has to do  $n \div 5$  (rounded up) homework assignments. For  
 519 example, when she has 11 points, it will take  $12 \div 5 = 2.4 \rightarrow 3$   
 520 homework assignments to earn her  $12^{\text{th}}$  point.  
 521 What is the smallest number of homework assignments necessary to  
 522 earn a total of 25 homework points?  
 523 Answer: Let's think step by step  
 524 Noelle only has to do 1 homework assignment to earn her first point,  
 525 and the same is true for each of her first five points.  
 526 She must then do 2 homework assignments to earn her sixth point, seventh  
 527 point, and so on, up to her tenth point.  
 528 Continuing, we see that Noelle must do a total of  $1+1+1+1+1+2+2+2+2+2+5+5+5+5+5$   
 529 homework assignments to earn 25 points.  
 530 This sum may be rewritten as  $5(1+2+3+4+5) = 5(15) = \boxed{75}$ .  
 531 The answer is 75''',  
 532  
 533  
 534 '''Question: The quadratic equation  $x^2 + mx + n = 0$  has roots that are twice  
 535 those of  $x^2 + px + m = 0$ , and none of  $m$ ,  $n$ , and  $p$  is zero. What  
 536 is the value of  $n/p$ ?  
 537 Answer: Let's think step by step  
 538 Let  $r_1$  and  $r_2$  be the roots of  $x^2 + px + m = 0$ .  
 539 Since the roots of  $x^2 + mx + n = 0$  are  $2r_1$  and  $2r_2$ , we have the  
 540 following relationships:  $\begin{cases} m = r_1 + r_2 \\ n = 4r_1 r_2 \\ p = -(r_1 + r_2) \end{cases}$   
 541  $m = r_1 + r_2$ ,  $n = 4r_1 r_2$ ,  $p = -(r_1 + r_2)$ , and  
 542  $m = -2(r_1 + r_2)$ .  
 543  $\therefore$  So  $\begin{cases} n = 4m \\ p = -\frac{1}{2}m \end{cases}$ ,  
 544  $\therefore \frac{n}{p} = \frac{4m}{-\frac{1}{2}m} = \boxed{-8}$ .  
 545  
 546 Alternatively, the roots of  $\left(\frac{x}{2}\right)^2 + p\left(\frac{x}{2}\right) + m = 0$   
 547 are twice those of  $x^2 + px + m = 0$ .  
 548 Since the first equation is equivalent to  $x^2 + 2px + 4m = 0$ ,  
 549 we have  $\begin{cases} m = 2p \\ n = 4m \end{cases}$ , so  $\frac{n}{p} = \boxed{8}$ .  
 550 The answer is 8''',  
 551  
 552  
 553  
 554 '''Question: Expand  $(2z^2 + 5z - 6)(3z^3 - 2z + 1)$ .  
 555 Answer: Let's think step by step  
 556 
$$\begin{array}{ccccccc} & & & & & & \\ & & & & & & \\ & & & & & & \\ & & & & & & \\ & & & & & & \\ & & & & & & \\ & & & & & & \end{array}$$
  
 557 
$$\begin{array}{ccccccc} & & & & & & \\ & & & & & & \\ & & & & & & \\ & & & & & & \\ & & & & & & \\ & & & & & & \\ & & & & & & \end{array}$$
  
 558 
$$\begin{array}{ccccccc} & & & & & & \\ & & & & & & \\ & & & & & & \\ & & & & & & \\ & & & & & & \\ & & & & & & \\ & & & & & & \end{array}$$
  
 559 
$$\begin{array}{ccccccc} & & & & & & \\ & & & & & & \\ & & & & & & \\ & & & & & & \\ & & & & & & \\ & & & & & & \\ & & & & & & \end{array}$$
  
 560 
$$\begin{array}{ccccccc} & & & & & & \\ & & & & & & \\ & & & & & & \\ & & & & & & \\ & & & & & & \\ & & & & & & \\ & & & & & & \end{array}$$
  
 561 
$$\begin{array}{ccccccc} & & & & & & \\ & & & & & & \\ & & & & & & \\ & & & & & & \\ & & & & & & \\ & & & & & & \\ & & & & & & \end{array}$$

```

562 & & & -18z^3 & & +12z & -6 & \\
563 & & +15z^4 & & -10z^2 & +5z & & \\
564 + & 6z^5 & & -4z^3 & +2z^2 & & & \\
565 \cline{1-7}\rule{0pt}{0.17in}
566 & 6z^5 & +15z^4 & -22z^3 & - 8z^2 & +17z & -6 &
567 \end{array}$$
568 The answer is  $6z^5+15z^4-22z^3-8z^2+17z-6$ '''',
569
570
571 '''Question: Find the mean of all solutions for  $x$  when  $x^3 + 3x^2 - 10$ 
572  $x = 0$ .
573 Answer: Let's think step by step
574 First, we factor the equation as  $x(x^2 + 3x - 10) = 0$ .
575 So, one solution is  $x=0$  and the other two solutions are the solutions
576 to  $x^2 + 3x-10=0$ .
577 We could either factor the quadratic, or note that the sum of the
578 solutions to this quadratic is  $-(3/1)=-3$ ,
579 so the mean of the three solutions to the original equation is  $-3/3=\$ 
580  $\boxed{-1}$ .
581 The answer is -1
582 ''',
583 ]

```

## 584 B.6.2 SL

```

585
586
587 sl_system_prompt = '''
588 You are a helpful assistant who is good at sloving math problems and
589 writing code. You should should follow the guidelines below:
590 - Utilize the `pi` symbol and `Rational` from Sympy for  $\pi$  and
591 fractions, and simplify all fractions and square roots without
592 converting them to decimal values
593 - You should only write code blocks and the function name should be `
594 solution` and the returned value should be the final answer.
595 '''
596
597 sl_instruction_prompt = "Let's use python to solve the math problem."
598
599
600 sl_math_shots= [
601 '''Question: Find the coefficient of  $x^3$  when  $3(x^2 - x^3+x) + 3(x + 2x$ 
602  $^2- 3x^2 + 3x^5+x^3) - 5(1+x-4x^3 - x^2)$  is simplifie.
603
604 ```python
605 from sympy import symbols, simplify
606
607 def solution():
608     """Find the coefficient of  $x^3$  when  $3(x^2 - x^3+x) + 3(x + 2x^3- 3x$ 
609  $^2 + 3x^5+x^3) - 5(1+x-4x^3 - x^2)$  is simplified."""
610     x = symbols('x')
611     expr = 3*(x**2 - x**3 + x) + 3*(x + 2*x**3 - 3*x**2 + 3*x**5 + x**3)
612     - 5*(1 + x - 4*x**3 - x**2)
613     simplified_expr = simplify(expr)
614
615     x3_coefficient = simplified_expr.as_coefficients_dict()[x**3]
616     result = x3_coefficient
617     return result
618     ''',
619
620 '''Question: The surface area of a sphere with radius  $r$  is  $4\pi r^2$ .
621 Including the area of its circular base, what is the total surface
622 area of a hemisphere with radius 6 cm? Express your answer in terms
623 of  $\pi$ .

```

```

624
625 ```python
626 import math
627
628 def solution():
629     """The surface area of a sphere with radius $r$ is $4\pi r^2$.
630     Including the area of its circular base, what is the total
631     surface area of a hemisphere with radius 6 cm? Express your
632     answer in terms of $\pi$"""
633     radius = 6
634
635     # Surface area of the hemisphere
636     hemisphere_area = 2 * math.pi * radius**2
637
638     # Area of the circular base
639     base_area = math.pi * radius**2
640
641     # Total surface area
642     total_surface_area = hemisphere_area + base_area
643
644     # Formatting the result in LaTeX
645     result = r'\pi'.format(total_surface_area / math.pi)
646     return result
647     """',
648
649 '''Question: Monica tosses a fair 6-sided die. If the roll is a prime
650 number, then she wins that amount of dollars (so that, for example,
651 if she rolls 3, then she wins 3 dollars). If the roll is composite,
652 she wins nothing. Otherwise, she loses 3 dollars. What is the
653 expected value of her winnings on one die toss? Express your answer
654 as a dollar value to the nearest cent.
655
656 ```python
657 def solution():
658     """Monica tosses a fair 6-sided die. If the roll is a prime number,
659     then she wins that amount of dollars (so that, for example, if
660     she rolls 3, then she wins 3 dollars). If the roll is composite,
661     she wins nothing. Otherwise, she loses 3 dollars. What is the
662     expected value of her winnings on one die toss? Express your
663     answer as a dollar value to the nearest cent."""
664     # Probabilities of each outcome
665     prime_prob = 1 / 6
666     composite_prob = 1 / 3
667     otherwise_prob = 1 / 6
668
669     # Expected value of each outcome
670     prime_expected_value = (2 * prime_prob) + (3 * prime_prob) + (5 *
671     prime_prob)
672     composite_expected_value = 0 * composite_prob
673     otherwise_expected_value = -3 * otherwise_prob
674
675     # Total expected value
676     total_expected_value = prime_expected_value +
677     composite_expected_value + otherwise_expected_value
678
679     # Dollar value to the nearest cent
680     result = "{:.2f}".format(total_expected_value)
681     return result
682     """',
683
684 '''Question: Given $\mathbf{a} = \begin{pmatrix} -7 \\ 0 \\ 1 \end{pmatrix}$
685 and $\mathbf{b} = \begin{pmatrix} 4 \\ 2 \\ -1 \end{pmatrix}$, find $\mathbf{a} - 3 \mathbf{b}$.
686
687
688 Solution:

```



```

689 ```python
690 import numpy as np
691
692 def solution()
693     """Given  $\mathbf{a} = \begin{pmatrix} -7 \\ 0 \\ 1 \end{pmatrix}$ 
694         and  $\mathbf{b} = \begin{pmatrix} 4 \\ 2 \\ -1 \end{pmatrix}$ ,
695         find  $\mathbf{a} - 3 \mathbf{b}$ ."""
696     a = np.array([-7, 0, 1])
697     b = np.array([4, 2, -1])
698
699     result = a - 3 * b
700
701     result = r'\begin{{pmatrix}} {} \\\ {} \\\ {} \end{{pmatrix}}'.format(
702         result[0], result[1], result[2])
703     return result
704     '''',
705
706 '''Question: The endpoints of a diameter of circle  $M$  are  $(-1,-4)$  and
707  $(-7,6)$ . What are the coordinates of the center of circle  $M$ ?
708 Express your answer as an ordered pair.
709 ```python
710 def solution():
711     """The endpoints of a diameter of circle  $M$  are  $(-1,-4)$  and
712          $(-7,6)$ . Find the coordinates of the center of circle  $M$ ."""
713     x1, y1 = -1, -4
714     x2, y2 = -7, 6
715
716     # Midpoint formula
717     center_x = (x1 + x2) / 2
718     center_y = (y1 + y2) / 2
719
720     # Result as an ordered pair
721     result = (center_x, center_y)
722     return result
723     '''',
724
725 '''Question: Find the remainder when  $2x^6 - x^4 + 4x^2 - 7$  is divided by  $x^2 + 4x + 3$ .
726
727 ```python
728 from sympy import symbols, div
729
730 def solution():
731     """Find the remainder when  $2x^6 - x^4 + 4x^2 - 7$  is divided by  $x^2 + 4x + 3$ 
732         ."""
733     x = symbols('x')
734     numerator = 2*x**6 - x**4 + 4*x**2 - 7
735     denominator = x**2 + 4*x + 3
736
737     quotient, remainder = div(numerator, denominator)
738
739     return remainder
740     '''',
741
742 '''Question: Find  $x$  so that the vectors  $\begin{pmatrix} 2 \\ 5 \end{pmatrix}$ 
743 and  $\begin{pmatrix} x \\ -3 \end{pmatrix}$  are orthogonal.
744 ```python
745 from sympy import symbols, Eq, solve
746
747 def solution():
748     """Find  $x$  so that the vectors  $(2, 5)$  and  $(x, -3)$  are orthogonal."""
749     x = symbols('x')
750     # Dot product of the vectors should be zero for orthogonality
751     equation = Eq(2*x + 5*(-3), 0)
752     result = solve(equation, x)
753     return result[0]

```

```

754     """
755     '''Question: If  $x$  is positive and  $x^2 = 729$ , what is the value of  $x$ 
756     ?
757     ```python
758     from sympy import symbols, Eq, solve
759
760     def solution():
761         """If  $x$  is positive and  $x^2 = 729$ , what is the value of  $x$ ?"""
762         x = symbols('x')
763         equation = Eq(x**2, 729)
764         solutions = solve(equation, x)
765
766         # Since x is positive, we select the positive solution
767         result = [sol for sol in solutions if sol > 0][0]
768         return result
769     """
770 ]

```

### 771 B.6.3 Second stage of SymbolNL

```

772 symbolnl_instructio_prompt = "Please_think_step_by_step_about_the_
773     question_based_on_the_provided_code(it_may_be_wrong)_and_the_executed
774     _output._"
775
776 symbolnl_math_shots = [
777     '''Question: How many distinct positive factors does 32 have?
778     Code: ```python
779     def solution():
780         """How many distinct positive factors does 32 have?"""
781         number = 32
782         factors = set()
783
784         for i in range(1, int(number**0.5) + 1):
785             if number % i == 0:
786                 factors.add(i)
787                 factors.add(number // i)
788
789         result = len(factors)
790         return result
791     ```
792     Output: 6
793     Answer: To determine the number of distinct positive factors of 32, we
794         can follow these steps:
795
796     1. Prime Factorization:
797         - First, we find the prime factorization of 32.
798         -  $(32 = 2^5)$ 
799
800     2. Number of Factors:
801         - The formula to find the number of distinct positive factors of a
802             number based on its prime factorization is to add 1 to each of the
803             exponents in the factorization and then multiply the results.
804         - For  $(32 = 2^5)$ , the exponent is 5.
805         - Adding 1 to the exponent:  $(5 + 1 = 6)$ .
806
807     3. Conclusion:
808         - Therefore, the number of distinct positive factors of 32 is 6.
809
810     Thus, the final result is:
811      $\boxed{6}$ '''
812
813     '''Question: In parallelogram  $ABCD$ , angle  $B$  measures  $110^\circ$ .
814         What is the number of degrees in the measure of angle  $C$ ?
815     Code: ```python

```

```

816 def solution():
817     """In parallelogram $ABCD$, angle $B$ measures $110^\circ$. What is
818         the number of degrees in the measure of angle $C$?"""
819     angle_B = 110
820
821     # In a parallelogram, opposite angles are equal and adjacent angles
822     are supplementary
823     angle_C = 180 - angle_B
824
825     return angle_C
826
827 Output: 70
828 Answer: In a parallelogram, opposite angles are equal and adjacent angles
829 are supplementary. Given that angle  $\angle B$  measures  $110^\circ$ 
830 , we need to find the measure of angle  $\angle C$ .
831
832 Step-by-step reasoning:
833 1. In a parallelogram, adjacent angles are supplementary. This means that
834 the sum of the measures of angle  $\angle B$  and angle  $\angle C$  is  $180^\circ$ .
835
836 2. Therefore, we can write the equation:  $\angle B + \angle C = 180^\circ$ .
837
838 3. Substituting the given measure of angle  $\angle B$  into the equation:  $110^\circ + \angle C = 180^\circ$ .
839
840 4. Solving for angle  $\angle C$ :  $\angle C = 180^\circ - 110^\circ = 70^\circ$ .
841
842 Thus, the measure of angle  $\angle C$  is  $\boxed{70}$ .
843
844
845 '''Question: In triangle $ABC$, we have that $E$ and $F$ are midpoints of
846 sides  $\overline{AC}$  and  $\overline{AB}$ , respectively. The area of
847  $\triangle ABC$  is 24 square units. How many square units are in the
848 area of  $\triangle CEF$ ?
849 Code: ```python
850 from sympy import Rational
851
852 def solution():
853     """Calculate the area of triangle CEF given that E and F are
854         midpoints of sides AC and AB respectively, and the area of
855         triangle ABC is 24 square units."""
856     # Area of triangle ABC
857     area_ABC = 24
858
859     # Triangle CEF is a median triangle, which is 1/4 the area of
860     triangle ABC
861     area_CEF = Rational(1, 4) * area_ABC
862
863     return area_CEF
864
865 Output: 6
866 Answer: The problem involves finding the area of triangle $CEF$ where $E$
867 and $F$ are midpoints of sides  $\overline{AC}$  and  $\overline{AB}$ 
868 respectively in triangle $ABC$. Given that the area of triangle $ABC$
869 is 24 square units, we need to determine the area of triangle $CEF$.
870
871 Since $E$ and $F$ are midpoints of their respective sides, line segment  $\overline{EF}$ 
872 is parallel to side  $\overline{BC}$  and is half its
873 length. This makes triangle $CEF$ a median triangle of triangle $ABC$.
874 A median triangle is one-fourth the area of the original triangle
875 because it is formed by drawing lines from each vertex of the
876 original triangle to the midpoints of the opposite sides.
877
878 Thus, the area of triangle $CEF$ is calculated as:
879 
$$[\text{Area of } \triangle CEF] = \frac{1}{4} \times [\text{Area of } \triangle ABC]$$

880 
$$= \frac{1}{4} \times 24$$


```

```

881
882 Simplifying this, we get:
883 
$$\text{Area of } \triangle CEF = 6$$

884
885 Therefore, the area of triangle $CEF$ is  $6$  square units.'',
886
887
888 '''Question: What is the greatest two-digit multiple of 13?
889 Code: ```python
890 def solution():
891     """Find the greatest two-digit multiple of 13."""
892     # Start from the greatest two-digit number and go downwards
893     for number in range(99, 12, -1):
894         if number % 13 == 0:
895             return number
896     ...
897 Output: 91
898 Answer: To find the greatest two-digit multiple of 13, we start by
899     considering the largest two-digit number, which is 99. We need to
900     find the largest number less than or equal to 99 that is divisible by
901     13.
902
903 We perform the division of 99 by 13 to find the largest integer quotient:
904 
$$\frac{99}{13} \approx 7.615$$

905 The integer part of this division is 7, indicating that the largest
906     multiple of 13 within the two-digit range is  $13 \times 7$ .
907
908 Calculating this product:
909 
$$13 \times 7 = 91$$

910
911 Thus, the greatest two-digit multiple of 13 is 91. The Python code
912     provided confirms this by iterating from 99 downwards and checking
913     divisibility by 13, stopping at the first number that satisfies this
914     condition, which is 91.
915
916 Therefore, the final answer is:
917 
$$91$$

918 ]

```

#### 919 B.6.4 Second stage of NLSymbol

```

920 nlsymbol_instruction_prompt = "Write a Python code that translates a
921     natural language (NL) reasoning path into executable code to answer a
922     given question. The output of the generated code should be the final
923     answer to the question."
924
925 nlsymbol_math_shots = [
926     '''Question: The area of the floor in a square room is 225 square feet.
927     The homeowners plan to cover the floor with rows of 6-inch by 6-inch
928     tiles. How many tiles will be in each row?
929 Reasoning Path: Let's think step by step:
930
931 The area of the floor is given as 225 square feet. Since the room is
932     square, each side of the room is the square root of the area. Thus,
933     each side of the room is  $\sqrt{225} = 15$  feet.
934
935 Each tile is 6 inches by 6 inches, which is equivalent to  $0.5$  feet by
936      $0.5$  feet since there are 12 inches in a foot.
937
938 To find out how many tiles fit along one side of the room, we divide the
939     length of the room by the length of one tile:
940 
$$\frac{15 \text{ feet}}{0.5 \text{ feet/tile}} = 30 \text{ tiles}$$

941
942 ]

```

```

943
944 Therefore, there will be  $\boxed{30}$  tiles in each row.
945 Code: ```python
946 from sympy import sqrt
947
948 def solution():
949     # Area of the floor in square feet
950     area = 225
951
952     # Since the room is square, calculate the side length of the room
953     side_length = sqrt(area)
954
955     # Each tile's side length in feet (6 inches = 0.5 feet)
956     tile_length = 0.5
957
958     # Calculate the number of tiles in each row
959     tiles_per_row = side_length / tile_length
960
961     return tiles_per_row
962
963 """
964 '''Question: When the set of natural numbers is listed in ascending order
965     , what is the smallest prime number that occurs after a sequence of
966     five consecutive positive integers all of which are nonprime?
967 Reasoning Path: Let's think step by step:
968
969 To find the smallest prime number that occurs after a sequence of five
970 consecutive nonprime positive integers, we need to identify such a
971 sequence first.
972
973 1. Start by checking small numbers. The first few nonprime numbers are 4,
974    6, 8, 9, and 10. However, these do not form a sequence of five
975    consecutive nonprime numbers.
976
977 2. Continue checking further:
978    - 8, 9, 10, 12, 14 are nonprime, but they are not consecutive.
979    - 24, 25, 26, 27, 28 are nonprime and consecutive.
980
981 3. Check the number following 28:
982    - 29 is the next number, and it is a prime number.
983
984 Thus, the smallest prime number that occurs after a sequence of five
985 consecutive nonprime positive integers (24, 25, 26, 27, 28) is 29.
986
987 The answer is  $\boxed{29}$ .
988 Code: ```python
989 from sympy import isprime
990
991 def solution():
992     # Start checking from a reasonable number
993     n = 1
994     while True:
995         # Check if the next five numbers are all nonprime
996         if all(not isprime(n + i) for i in range(5)):
997             # Check the number right after the sequence of five nonprimes
998             if isprime(n + 5):
999                 return n + 5
1000     n += 1
1001 """
1002
1003 '''Question: What is the slope of the line that is tangent to a circle at
1004 point (5,5) if the center of the circle is (3,2)? Express your
1005 answer as a common fraction.
1006 Reasoning Path: Let's think step by step:
1007

```

```

1008 The slope of the line connecting the center of the circle at point
1009  $((3,2))$  to the point on the circle  $((5,5))$  is given by the
1010 formula for the slope between two points:
1011  $m = \frac{y_2 - y_1}{x_2 - x_1} = \frac{5 - 2}{5 - 3} = \frac{3}{2}$ 
1012
1013 The line that is tangent to the circle at the point  $((5,5))$  will be
1014 perpendicular to the line connecting the center of the circle to this
1015 point. The slope of a line perpendicular to another line with slope
1016  $m$  is the negative reciprocal of  $m$ . Therefore, the slope of
1017 the tangent line is:
1018  $m_{\text{tangent}} = -\frac{1}{\frac{3}{2}} = -\frac{2}{3}$ 
1019
1020 Thus, the slope of the line that is tangent to the circle at point
1021  $((5,5))$  is  $(-\frac{2}{3})$ .
1022 Code: ```python
1023 from sympy import Rational
1024
1025 def solution():
1026     # Slope of the line connecting the center of the circle to the point
1027     (5,5)
1028     slope_radius = Rational(5 - 2, 5 - 3)
1029
1030     # Slope of the tangent line, which is the negative reciprocal of the
1031     slope of the radius
1032     slope_tangent = -1 / slope_radius
1033
1034     return slope_tangent
1035
1036 ```',
1037
1038 '''Question: What is the greatest integer less than 100 for which the
1039 greatest common factor of that integer and 18 is 3?
1040 Reasoning Path: Let's think step by step
1041
1042 We are looking for the greatest integer less than 100 for which the
1043 greatest common factor (GCF) with 18 is 3.
1044
1045 First, note that 18 can be factored into prime factors as  $(18 = 2 \times 3^2)$ .
1046
1047 For the GCF of a number  $n$  and 18 to be 3,  $n$  must:
1048 1. Be divisible by 3 (to include the factor of 3).
1049 2. Not be divisible by 2 or 9 (to avoid increasing the GCF beyond 3).
1050
1051 We need to find the largest integer less than 100 that meets these
1052 criteria. We look for numbers that are multiples of 3 but not
1053 multiples of 2 or 9.
1054
1055 The largest multiple of 3 under 100 is 99. We check if it is divisible by
1056 2 or 9:
1057 - 99 is not divisible by 2 (since it is odd).
1058 - 99 is divisible by 9 (since  $(9 + 9 = 18)$ , and 18 is divisible by 9).
1059
1060 Since 99 does not work (as it is divisible by 9), we check the next
1061 largest multiple of 3, which is 96.
1062 - 96 is divisible by 2 (even number), so it does not work.
1063
1064 Next, we check 93:
1065 - 93 is not divisible by 2 (odd number).
1066 - 93 is not divisible by 9 (since  $(9 + 3 = 12)$ , and 12 is not divisible
1067 by 9).
1068
1069 Thus, 93 meets the criteria of being divisible by 3 but not by 2 or 9.
1070 Therefore, the greatest integer less than 100 for which the GCF with
1071 18 is 3 is  $(93)$ .
1072

```

```

1073 Code ```python
1074 def solution():
1075     from math import gcd
1076     # Start from the largest number less than 100 and check downwards
1077     for n in range(99, 0, -1):
1078         if gcd(n, 18) == 3:
1079             return n
1080
1081 # The function will return the greatest integer less than 100 for which
1082 the GCD with 18 is 3
1083 ```',
1084 ]

```

## 1085 B.7 Significance test

Method	Significance Level
GPT-3.5-Turbo compare prompt (less)	0.05
GPT-4o compare prompt (less)	0.05
o3-mini prompt (less)	0.257
HYBRIDMIND (more)	0.05

Table 9: Testing of significance levels of MATH dataset between HYBRIDMIND and NL.

Method	Significance Level
GPT-4o compare prompt (less)	0.080
o3-mini prompt (less)	0.05
HYBRIDMIND (more)	0.05

Table 10: Testing of significance levels of FOLIO dataset between HYBRIDMIND and NL.

## 1086 B.8 Analysis between SL and NL

1087 Please check Tables 11, 12, 13, and 14.

## 1088 B.9 Meta-Selection prompt

1089 Here is 2-class meta-selection prompt for MATH dataset.

```

1090 two_class_shots = [
1091     '''
1092     Question: Four standard, six-sided dice are to be rolled. If the
1093     product of their values turns out to be an even number, what is
1094     the probability their sum is odd? Express your answer as a common
1095     fraction.
1096
1097     Response: Below is an evaluation of the two methods for this specific
1098     four-dice probability problem:
1099
1100     Natural Language (Language Only):
1101     Strengths: This method lets you break the problem down step by step.
1102     You can explain your reasoning about conditionals (the even
1103     product and odd sum) and manually handle the combinatorial cases.
1104     This approach can offer insight into why the answer is what it
1105     is.
1106     Weaknesses: It requires careful bookkeeping to avoid mistakes in
1107     counting cases and may become tedious when managing several
1108     conditions simultaneously. There is more risk of a small
1109     oversight that leads to an error in the final probability.
1110
1111     Symbolic Language (Coding Only):
1112     Strengths: By writing a short Python program to iterate over all 6^4
1113     possible outcomes, you can systematically and accurately count
1114     the desired cases. It is very effective for this problem because

```

<b>Question</b>
For how many integers $a$ satisfying $1 \leq a \leq 23$ is it true that $a^{-1} \equiv a \pmod{24}$ ?
<b>NL</b>
<p>Let's think step by step. We need to find integers <math>a</math> such that <math>1 \leq a \leq 23</math> and <math>a^{-1} \equiv a \pmod{24}</math>. This condition implies that <math>a \cdot a \equiv 1 \pmod{24}</math>, or equivalently, <math>a^2 \equiv 1 \pmod{24}</math>. This congruence can be rewritten as <math>a^2 - 1 \equiv 0 \pmod{24}</math>, which factors as <math>(a - 1)(a + 1) \equiv 0 \pmod{24}</math>. Therefore, <math>(a - 1)(a + 1)</math> must be divisible by 24. The prime factorization of 24 is <math>2^3 \times 3</math>. Thus, <math>(a - 1)(a + 1)</math> must be divisible by <math>2^3 \times 3</math>. Since <math>a - 1</math> and <math>a + 1</math> are consecutive integers, one of them is even, and the other is odd. Therefore, one of them must be divisible by 4, and the other by 2, to ensure divisibility by <math>2^3 = 8</math>. Additionally, one of <math>a - 1</math> or <math>a + 1</math> must be divisible by 3.</p> <p>We solve these systems of congruences using the Chinese Remainder Theorem:</p> <p><b>Case 1:</b> <math>a \equiv 17 \pmod{24}</math></p> <p><b>Case 2:</b> <math>a \equiv 7 \pmod{24}</math></p> <p>Thus, there are two integers <math>a</math> that satisfy the condition: <math>a = 7</math> and <math>a = 17</math>. Therefore, the number of integers <math>a</math> is 2.</p>
<b>SL</b>
<p><b>Python</b></p> <pre> from sympy import symbols, Eq, solve, mod_inverse def solution():     count = 0     for a in range(1, 24):         try:             # Calculate the modular inverse of a             modulo 24             a_inv = mod_inverse(a, 24)             # Check if the modular inverse is equal             to a             if a_inv == a:                 count += 1             except ValueError:                 # If a modular inverse does not exist,                 continue         continue     return count solution() </pre> <p><b>Executed Result</b></p> <p>8</p>

Table 11: An example where SL outperforms NL. While NL reasoning leads to logical errors in analysis and overlooks some possible scenarios, resulting in incorrect conclusions, writing Python code is more elegant and intuitive.

```

1115         the search space is small enough for exhaustive enumeration,
1116         thereby reducing human error.
1117     Weaknesses: It might seem like overkill for a problem that is
1118         solvable by combinatorial reasoning, and you must be careful to
1119         correctly implement the conditions in code. However, given the
1120         manageable sample space, this method is straightforward and
1121         reliable.
1122

```



<b>Question</b>
How many days in a year have a 3 in the tens digit of the day of the month?
<b>NL</b>
<p>Let's think step by step. We need to count the number of days in a year where the tens digit of the day of the month is 3. This means we are looking for days numbered 30 through 39 in each month.</p> <p>1. Identify the months with days 30-39:</p> <ul style="list-style-type: none"> <li>- The months with at least 30 days are January, March, April, May, June, July, August, September, October, November, and December.</li> <li>- February has only 28 or 29 days, so it does not have any days with a tens digit of 3.</li> </ul> <p>2. Count the days for each applicable month:</p> <ul style="list-style-type: none"> <li>- For each of the months with at least 30 days, the days 30 and 31 have a tens digit of 3.</li> <li>- Therefore, each of these months contributes 2 days (30 and 31) to our count.</li> </ul> <p>3. Calculate the total number of such days:</p> <ul style="list-style-type: none"> <li>- There are 11 months with at least 30 days.</li> <li>- Each of these months contributes 2 days.</li> <li>- Total number of days = <math>11 \times 2 = 22</math>. Thus, the number of days in a year that have a 3 in the tens digit of the day of the month is 22.</li> </ul>
<b>SL</b>
<p><b>Python</b></p> <pre>def solution():     # Not all months have 31 days, so we need     # to count carefully     # Months with 31 days: January, March, May,     # July, August, October, December     months_with_31_days = 7     # Months with 30 days: April, June,     # September, November     months_with_30_days = 4     # February has 28 or 29 days, so it never     # has a 30th or 31st     # Days with a 3 in the tens digit:     # - 30th: occurs in every month except     # February     # - 31st: occurs only in months with 31 days     days_with_3_in_tens_digit = (months_with_30_days + months_with_31_days) + months_with_31_days     return days_with_3_in_tens_digit solution()</pre> <p><b>Executed Result</b></p> <p>18</p>

Table 12: An example where SL outperforms NL. Python code offers a more accurate and structured approach. While NL reasoning overlooks some edge cases and leads to incorrect conclusions, SL systematically handles the details and produces the correct result.

1123 For this specific question, because the number of outcomes is small  
1124 and the conditions can be programmed explicitly, the program-  
1125 aided approach minimizes error risk and offers a quick,  
1126 systematic solution.  
1127  
1128 Selected method: SL (Symbolic Language, Coding Only)

---

**Question**

---

Find the product of all positive integer values of  $c$  such that  $3x^2 + 7x + c = 0$  has two real roots.

---

**NL**

---

To determine the product of all positive integer values of  $c$  such that the quadratic equation  $3x^2 + 7x + c = 0$  has two real roots, we need to ensure that the discriminant of the quadratic is positive. The discriminant  $\Delta$  of a quadratic equation  $ax^2 + bx + c = 0$  is given by:

$$\Delta = b^2 - 4ac$$

.....

Solving this inequality for  $c$ , we get:

$$49 > 12c \Rightarrow c < \frac{49}{12}$$

Calculating  $\frac{49}{12}$ , we find:

$$\frac{49}{12} \approx 4.0833$$

Since  $c$  must be a positive integer, the possible values for  $c$  are 1, 2, 3, and 4. Now, we find the product of these values:

$$1 \times 2 \times 3 \times 4 = 24$$

Thus, the product of all positive integer values of  $c$  such that the quadratic equation has two real roots is 24.

---

**SL**

---

**Python**

```
from sympy import symbols, solve, Rational
def solution():
    x, c = symbols('x c')
    # For the quadratic equation to have
    # two real roots, the discriminant must be
    # non-negative
    # Discriminant of  $3x^2 + 7x + c = 0$  is  $7^2 - 4 \cdot 3 \cdot c$ 
    discriminant = 7**2 - 4*3*c
    # Solve for c when discriminant is greater
    # than 0
    c_values = solve(discriminant > 0, c)
    # c must be a positive integer
    positive_integer_c_values = [i for i in
    range(1, int(c_values[0]) + 1)]
    # Calculate the product of all positive
    # integer values of c
    product = 1
    for value in positive_integer_c_values:
        product *= value
    return product
solution()
```

**Executed Result**

**TypeError: 'And' object is not subscriptable**

Table 13: An example where NL outperforms SL. When writing python code to solve the problem, another key factor is that the language model has learned this symbolic language well. Otherwise, it may produce incorrect code, even if it understands the right concepts.

```
1129
1130 ' ',
1131 ' ',
```

---

**Question**

---

How many vertical asymptotes does the graph of  $y = \frac{x-3}{x^2+7x-30}$  have?

---

**NL**

---

Let's think step by step. To find the vertical asymptotes of the function  $y = \frac{x-3}{x^2+7x-30}$ , we need to determine where the denominator is equal to zero, as these are the points where the function is undefined and may have vertical asymptotes. First, we set the denominator equal to zero:

$$x^2 + 7x - 30 = 0.$$

Next, we factor the quadratic equation:

$$x^2 + 7x - 30 = (x + 10)(x - 3) = 0.$$

Setting each factor equal to zero gives us the solutions:

$$x + 10 = 0 \Rightarrow x = -10,$$

$$x - 3 = 0 \Rightarrow x = 3.$$

These are the potential vertical asymptotes. However, we must check if these points are also zeros of the numerator, which would indicate a hole instead of an asymptote. The numerator is  $x - 3$ . Setting it equal to zero gives:

$$x - 3 = 0 \Rightarrow x = 3.$$

Since  $x = 3$  is a zero of both the numerator and the denominator, it is a removable discontinuity (a hole) rather than a vertical asymptote.

Therefore, the only vertical asymptote occurs at  $x = -10$ .

Thus, the graph of  $y = \frac{x-3}{x^2+7x-30}$  has 1 vertical asymptote.

---

**SL**

---

**Python**

```
from sympy import symbols, solve
def solution():
    x = symbols('x')
    denominator = x**2 + 7*x - 30

    # Find the roots of the denominator
    roots = solve(denominator, x)

    # The number of vertical asymptotes is the
    # number of distinct real roots
    num_vertical_asymptotes = len(roots)

    return num_vertical_asymptotes
solution()
```

**Executed Result**

2

---

Table 14: An example where NL outperforms SL. SL doesn't consider the case where the numerator equals zero, which highlights the weakness of SL reasoning when facing certain edge cases or special conditions.

```
1132 Question: A piece of cheese is located at $(12,10)$ in a coordinate
1133 plane. A mouse is at $(4,-2)$ and is running up the line $y=-5
1134 x+18$. At the point $(a,b)$ the mouse starts getting farther
1135 from the cheese rather than closer to it. What is $a + b$?
```

1136 Response: Below is an evaluation of each approach for this particular  
 1137 problem:  
 1138

1139 Natural Language (NL):  
 1140 Strengths:  
 1141 - It is especially useful for problems that rely on geometric insight  
 1142 and visual reasoning. For instance, one can break down the  
 1143 problem by first finding the point on the line that is closest to  
 1144 the cheese (using ideas like the perpendicular from the cheese  
 1145 to the mouse's path) and then compute the coordinates, which  
 1146 directly leads to a conceptual understanding of why that point  
 1147 minimizes the distance.  
 1148 - The natural language breakdown allows the solver to easily explain,  
 1149 verify, and understand each step in the geometric reasoning  
 1150 process.  
 1151 Weaknesses:  
 1152 - It relies on the solver's ability to carefully reason through each  
 1153 step without skipping important details. Errors might occur if a  
 1154 step is overlooked or misinterpreted, but these issues are often  
 1155 easier to spot and correct when the reasoning is laid out in  
 1156 natural language.  
 1157

1158 Symbolic Language (SL):  
 1159 Strengths:  
 1160 - A Python program can be very effective for computing an answer if  
 1161 one sets up the right system (for example, defining the distance  
 1162 function and finding its minimum via analytic methods or  
 1163 optimization). It is useful if one has a routine set of tools (like  
 1164 computational geometry libraries or symbolic manipulators) that  
 1165 can automate the calculation accurately.  
 1166 - It minimizes the potential for simple arithmetic mistakes, as the  
 1167 computer does the heavy lifting once the input is correctly  
 1168 specified.  
 1169 Weaknesses:  
 1170 - Setting up the correct algorithm in code (especially for a geometry  
 1171 problem that involves minimizing a distance) can be nontrivial.  
 1172 One must translate the geometric intuition correctly into code,  
 1173 and a mistake there might lead to an incorrect answer without an  
 1174 immediately clear explanation.  
 1175 - The program might properly produce a number but without conveying  
 1176 the underlying geometric reasoning, which can reduce the  
 1177 educational value of the solution process for problems like this.  
 1178

1179 For this specific problem—a geometry problem involving distances and  
 1180 perpendicular intersections—the Chain-of-Thought method is likely  
 1181 the more effective approach. It promotes understanding of the  
 1182 underlying geometric reasoning and highlights the process behind  
 1183 finding the point where the mouse's distance from the cheese is  
 1184 minimized.  
 1185

1186 Selected method: NL (Natural Language, Language Only)'''  
 1187 '''

1188 Question: Find the minimum of the function  
 1189  $\sqrt{\frac{xy}{x^2 + y^2}}$  in the domain  $\frac{2}{5} \leq x \leq \frac{1}{2}$   
 1190  $\frac{1}{2}$  and  $\frac{1}{3} \leq y \leq \frac{3}{8}$ .  
 1191 Response: Let's review what each approach offers for this problem.  
 1192

1193 Natural Language (NL) method:  
 1194 Strengths: It allows step-by-step reasoning, considering the  
 1195 structure of the algebraic expression and boundary behavior,  
 1196 which may reveal insights such as where the minimum might occur (for  
 1197 example, at the endpoints or at a critical point) and can be  
 1198 done entirely by hand.  
 1199 Weaknesses: The algebra can become quite messy given the rational  
 1200 function and the two-dimensional domain. Tracking all possible

cases (including boundary and potential critical points) increases the risk of error or oversight. The method relies on human intuition and careful logical decomposition, which can be challenging under time constraints.

Symbolic Language (SL) method:

Strengths: A Python program can brute-force or use well-tested optimization routines to accurately sample the function over the given domain. Because the domain is compact and small in extent, a numerical search (or even more refined optimization) is straightforward and likely to yield a correct answer with less chance of human error.

Weaknesses: One must be careful to cover the entire domain correctly and account for potential issues with numerical precision. Additionally, writing the optimization code correctly also requires some programming skill, though for a two-variable bounded problem this is well within reach.

For this specific question-finding the minimum of a rational function over a small, compact domain-the program-aided approach is more directly effective. It reduces the intellectual overhead of dealing with all the subtle boundary and interior cases manually and tends to be more reliable given the tight constraints.

Selected method: SL (Symbolic Language, Coding Only)

'''

Question: Someone observed that  $6! = 8 \cdot 9 \cdot 10$ . Find the largest positive integer  $n$  for which  $n!$  can be expressed as the product of  $n - 3$  consecutive positive integers.

Response: For this particular problem the trade-offs are as follows:

Natural Language (NL) Strengths and Weaknesses:

- Strength: A step-by-step reasoning process can provide insight into the structure of the equation (relating factorials to products of sequences) and might lead to an elegant, human-readable argument.
- Weakness: The reasoning becomes quite intricate because one needs to relate rapidly growing factorials to a product of a long chain of consecutive integers. Handling large numbers and justifying bounds without resorting to heavy computations is tricky and error-prone for this specific task.

Symbolic Language (SL) Strengths and Weaknesses:

- Strength: A computer program can iterate through candidate values of  $n$  (and compute factorials and products of consecutive integers) with precision. This method offloads the heavy number crunching and lets one quickly test many cases, which is especially valuable given the enormous growth of factorials.
- Weakness: A coding solution requires careful handling of very large numbers (or the use of logarithms/approximations) to avoid computational inefficiencies or overflow issues. Still, with modern programming languages and libraries this issue can be managed.

For this specific question, where the answer hinges on a precise numerical property and where direct computational verification seems feasible (once the search boundaries are known or estimated), the Symbolic Language (SL) approach is more effective. It allows for a systematic exploration of candidates with the precision that the problem demands without running into lengthy manual derivations.

Selected method: SL (Symbolic Language, Coding Only)'''

1265 '''Question: Evaluate  $(2-w)(2-w^2)\cdots(2-w^{10})$  where  $w=e^{2\pi i/11}$ .\$

1266

1267 Response: Below is an evaluation of both approaches for this problem:

1268

1269 Natural Language (nl, Language Only):

1270 - Strengths: This method allows one to recognize the structure of the

1271 product. In this particular problem the factors are  $(2 - w^k)$

1272 for the nontrivial 11th roots of unity, and one can quickly

1273 connect this to the cyclotomic polynomial formula  $(x^{11} - 1)/(x -$

1274  $1)$ . This insight leads directly to an exact, elegant solution

1275 without resorting to numerical approximation.

1276 - Weaknesses: Using pure logical reasoning requires a strong

1277 familiarity with algebraic properties of roots of unity and

1278 cyclotomic polynomials. If one's background is less strong in

1279 these areas, there is a risk of overlooking key steps or making

1280 errors in derivation.

1281

1282 Symbolic Language (SL, Coding Only):

1283 - Strengths: A Python program (or similar) can compute the product

1284 numerically, providing a quick verification of the result. It

1285 automatically handles the multiplication of many complex numbers,

1286 which is useful when manual calculation is cumbersome.

1287 - Weaknesses: The numerical approach may introduce floating-point

1288 rounding errors, especially when working with complex

1289 exponentials. It also may hide the underlying structure of the

1290 problem and does not provide the exact algebraic insight that

1291 leads to a neat closed-form answer. Additionally, setting up the

1292 program correctly requires careful implementation of complex

1293 arithmetic.

1294

1295 For this specific question, since the underlying algebraic structure

1296 (via cyclotomic polynomials) offers an elegant, exact answer, the

1297 chain-of-thought method is more effective and insightful.

1298

1299 Selected method: NL (Natural Language, Language Only)

1300 '''

1301 '''

1302 Question: In trapezoid \$ABCD\$ the lengths of the bases \$AB\$ and \$CD\$

1303 are 8 and 17 respectively. The legs of the trapezoid are extended

1304 beyond \$A\$ and \$B\$ to meet at point \$E\$. What is the ratio of

1305 the area of triangle \$EAB\$ to the area of trapezoid \$ABCD\$?

1306 Express your answer as a common fraction.

1307 Response: For this trapezoid geometry problem, the NL method allows

1308 you to gradually set up the problem by assigning variables,

1309 drawing relationships (such as the similarity of triangles), and

1310 then carefully deriving the necessary formulas. This process can

1311 enhance understanding and make it easier to check each individual

1312 step, especially when dealing with geometric properties and

1313 ratios.

1314

1315 On the other hand, the SL method would involve writing a Python

1316 program to symbolically manipulate or numerically approach the

1317 problem. While Python can be very powerful, setting up the

1318 symbolic algebra needed for a clear derivation of an algebraic

1319 ratio may require additional packages or careful coding. This

1320 approach might risk making the solution less transparent and

1321 could introduce potential programming errors for a problem that

1322 is more elegantly handled through traditional reasoning.

1323

1324 Given that this problem benefits from step-by-step geometric

1325 reasoning and clear algebraic manipulation, the chain-of-thought

1326 method is more effective and straightforward for this specific

1327 question.

1328

1329 Selected method: NL (Natural Language, Language Only)

```

1330     ''',
1331     '''
1332     Question: A 4-inch by 6-inch picture is enlarged for framing by
1333             tripling its dimensions. A 2-inch-wide border is then placed
1334             around each side of the enlarged picture, as shown. Thin metal
1335             framing is sold only in increments of one foot. What is the
1336             minimum number of linear feet of framing that must be purchased
1337             to go around the perimeter of the border?
1338
1339 [asy]
1340
1341 draw((0,0)--(14,0)--(14,20)--(0,20)--cycle,linewidth(2));
1342
1343 draw((4,4)--(10,4)--(10,16)--(4,16)--cycle);
1344
1345 label("border",(7,17),N);
1346
1347 label("picture",(7,8),N);
1348
1349 label("frame",(14,5),E);
1350
1351 draw((17.5,7.5)--(14.5,7.5),Arrow);
1352 draw((10.5,7.5)--(13.5,7.5),Arrow);
1353
1354 [/asy]
1355 Response: Below is an evaluation of the two methods for this particular
1356         problem, followed by the selected method.
1357
1358 Natural Language (NL) - Language Only:
1359 Strengths:
1360 - It lets the solver reason through each step, clarifying the enlargement
1361   process, how the border is added to both dimensions, and the
1362   conversion from inches to feet.
1363 - The explanation can detail the arithmetic (tripling dimensions, adding
1364   2 inches on each side for the border, computing the new perimeter,
1365   and then converting inches to feet while accounting for the need to
1366   round up) so that every piece of the problem is understood.
1367 - This method can easily communicate the reasoning to others.
1368 Weaknesses:
1369 - If not careful, details might be overlooked or arithmetic errors made.
1370 - The reliance on manual computations (even if done carefully) might lead
1371   to mistakes if any step is not double-checked.
1372
1373 Symbolic Language (SL) - Coding Only:
1374 Strengths:
1375 - A short Python program can accurately perform the arithmetic once the
1376   proper logic is implemented.
1377 - Code can handle unit conversion and rounding (for example, using math.
1378   ceil) without any miscalculation due to human error.
1379 - It is especially useful when the operations are straightforward and
1380   mostly numerical, as in this case.
1381 Weaknesses:
1382 - The problem requires careful interpretation (tripling the dimensions,
1383   adding a border to both sides) which must be encoded correctly; any
1384   misinterpretation could lead to a coding error.
1385 - Writing code for such a simple arithmetic problem might be overkill and
1386   could miss conveying the geometric understanding behind the steps.
1387 - There can also be distraction in the form of programming syntax rather
1388   than focusing on the underlying mathematics.
1389
1390 For THIS SPECIFIC QUESTION, the arithmetic and geometry are simple and
1391 direct. It is crucial to make sure that each conceptual step (
1392 enlargement, addition of borders, perimeter computation, and
1393 conversion unit) is clearly understood and communicated. The chain-of
1394 -thought method naturally lends itself to this clarity without the

```

```
1395     overhead of translating the problem into code. While coding could
1396     give the answer reliably, it may hide the reasoning steps that are
1397     important to verify each part of the computation.
1398
1399     Selected method: NL (Natural Language, Language Only)
1400     '''
1401 ]
```