

# Natural Language Embedded Programs for Hybrid Language Symbolic Reasoning

Anonymous ACL submission

## Abstract

How can we perform computations over natural language representations to solve tasks that require symbolic and numeric reasoning? We propose *natural language embedded programs* (NLEP) as a unifying framework for addressing math/symbolic reasoning, natural language understanding, and instruction following tasks. Our approach prompts a language model to generate full Python programs that define functions over data structures which contain natural language representations of structured knowledge. A Python interpreter then executes the generated code and prints the output. Despite using a task-general prompt, we find that this approach can improve upon strong baselines across a range of different tasks including math and symbolic reasoning, text classification, question answering, and instruction following. We further find the generated programs are often interpretable and enable post-hoc verification of the intermediate reasoning steps.

## 1 Introduction

Solving complex language tasks often requires performing computations over natural language representations. For language-based reasoning, chain-of-thought prompting (CoT; Wei et al., 2022) has emerged as a promising approach for surfacing the symbolic reasoning capabilities of large language models (LLMs). However, certain types of computations (e.g., arithmetic) are unnatural to perform in pure language space, and hence present difficulties for LLMs. General-purpose programming languages, on the other hand, provide convenient abstractions as well as predefined libraries and functions for natively implementing many types of symbolic computations, and there has been much recent work on interleaving program calls within CoT-style reasoning to extend the capabilities of LLMs. While promising, existing methods are generally limited to narrow types of

tasks such as math and symbolic reasoning (Chen et al., 2022; Cai et al., 2023; Gao et al., 2023), simple API calling (Schick et al., 2023; Paranjape et al., 2023; Liang et al., 2023a), and database accessing (Cheng et al., 2022). These works moreover rely on task-specific prompts which are hard to generalize across datasets.

This work describes a task-general approach for combining the language-based reasoning capabilities of LLMs with symbolic computations enabled by the use of programs. Specifically, we prompt LLMs to generate *natural language embedded programs* (NLEPs), which are fully executable Python programs containing appropriate package importing, structured natural language representations of knowledge, function definitions for problem solving, and response printing. The generated NLEP is then executed using a Python interpreter that captures the standard output of the program as the response. An example of an NLEP generated by GPT-4 is shown in Figure 1.

NLEPs use code as a scaffold to reason over natural language representations of data and knowledge. This makes our approach different from ToolFormer (Schick et al., 2023) and language model as tool maker (LATM; Cai et al., 2023), which instead use language as the scaffold and interleave API calls within natural language sentences during LLM generation. Compared to program-of-thought (PoT; Chen et al., 2022) and program aided language models (PAL; Gao et al., 2023), which mainly focus on math and symbolic problems, NLEPs handles natural language, math, and symbolic tasks with significantly better generalization ability based on task-general prompts. Compared to existing neuro-symbolic augmentations approaches needing new logical language (Li et al., 2023) and predefined solvers (Pan et al., 2023), NLEP uses a general programming interface thus can handle a significantly wider range of tasks with more flexible programming elements includ-

**Instruction:** Who are the US presidents elected after 1950 and born on Monday?

		NLEP generated by GPT-4
Tool Using		<pre># Step 1: Import necessary libraries from datetime import datetime</pre>
Structured knowledge		<pre># Step 2: Define a dictionary storing the knowledge about US presidents us_presidents = {     "Dwight D. Eisenhower": {"birth_date": "1890-10-14", "term_start": "1953-01-20"},     "John F. Kennedy": {"birth_date": "1917-05-29", "term_start": "1961-01-20"},     "Lyndon B. Johnson": {"birth_date": "1908-08-27", "term_start": "1963-11-22"},     "Richard Nixon": {"birth_date": "1913-01-09", "term_start": "1969-01-20"},     "Gerald Ford": {"birth_date": "1913-07-14", "term_start": "1974-08-09"},     "Jimmy Carter": {"birth_date": "1924-10-01", "term_start": "1977-01-20"},     "Ronald Reagan": {"birth_date": "1911-02-06", "term_start": "1981-01-20"},     "George H. W. Bush": {"birth_date": "1924-06-12", "term_start": "1989-01-20"},     "Bill Clinton": {"birth_date": "1946-08-19", "term_start": "1993-01-20"},     "George W. Bush": {"birth_date": "1946-07-06", "term_start": "2001-01-20"},     "Barack Obama": {"birth_date": "1961-08-04", "term_start": "2009-01-20"},     "Donald Trump": {"birth_date": "1946-06-14", "term_start": "2017-01-20"},     "Joe Biden": {"birth_date": "1942-11-20", "term_start": "2021-01-20"}, }</pre>
Function definition		<pre># Step 3: Define a function that finds the presidents who are elected after a certain year and born on a certain day of the week def find_presidents(presidents, year, day_of_week):     result = []     for president, info in presidents.items():         birth_date = datetime.strptime(info["birth_date"], "%Y-%m-%d")         term_start = datetime.strptime(info["term_start"], "%Y-%m-%d")         if birth_date.strftime("%A") == day_of_week and term_start.year &gt; year:             result.append(president)     return result</pre>
Print answer		<pre># Step 4: Print an answer in natural language using the knowledge and function defined above year = 1950 day_of_week = "Monday" presidents = find_presidents(us_presidents, year, day_of_week) if presidents:     print("The US presidents who were elected after 1950 and born on Mondays are: " + ", ".join(presidents) + ".") else:     print("There are no US presidents who were elected after 1950 and born on Mondays.")</pre>

**Stdout:** The US presidents who were elected after 1950 and born on Mondays are: Gerald Ford, Ronald Reagan, Bill Clinton.

**Figure 1:** A generated NLEP correctly answers the given question while ChatGPT-4 obtains an incorrect answer (link). This NLEP uses the date-weekday conversion tool in the datetime package, constructs structured knowledge about US presidents, implements a selection function, and outputs natural language responses depending on the function output. A more detailed comparison between NLEP and ChatGPT-4 code interpreter is shown in Figure 5.

ing packages, databases, and APIs.

Experiments across math and symbolic reasoning, question answering and instruction following, and text classification tasks demonstrate that (1) NLEPs conducts accurate reasoning on both structured and unstructured tasks and inputs; (2) NLEP’s step-by-step, meta prompting strategy can significantly improve the prompt efficiency across different tasks. As a result, we conclude that programming language prompting with NLEP is more capable and generalizable than existing natural language and neuro-symbolic prompting strategies.

## 2 Approach: NLEP Prompting

An NLEP is a program containing both programming code and natural language. NLEPs use natural language in several different ways. First, it uses natural language comments to guide step-by-step program generation. Second, language is used to represent structured knowledge through Python’s native data structures (e.g., dictionaries and lists). Finally, an NLEP uses language to print fluent responses to the user input by constructing a standard output string containing references to program variables.

The hybrid language-symbolic design of NLEP enables generalized problem solving for natural language, math, symbolic reasoning, and API calling tasks, which have traditionally been tackled by separate mechanisms. This approach combines the benefits of language-based reasoning with program synthesis: comments and knowledge in natural language improve program generation, while the structured/symbolic reasoning powered by program interpreters provides more accurate computations than would have been obtained via direct decoding from LLMs.

An example of an NLEP for answering a question is shown in Figure 5. In the generated program, each section is preceded by comments in natural language, and the defined counting function uses knowledge stored in a key-value dictionary (which itself is generated from GPT-4’s internal knowledge) to find the correct answer. Finally, the answer is printed through a natural language response. In this example, we generated 5 independent NLEPs and found that they achieve 100% accuracy, compared to 60% for ChatGPT-4 and 40% GPT-4 API.

**NLEP structure.** More generally, each NLEP

contains four sections: (1) importing necessary libraries, (2) defining variables containing structured knowledge, (3) implementing problem-solving functions, and (4) printing the response in natural language. Instead of providing direct solutions for each task, we guide the model to arrive at a solution following this four-step process. As show in the example in Figure 1, an NLEP answers the question by constructing a structured knowledge dictionary containing the birthday and start date of the US presidents. To recognize the weekdays, the program utilizes pre-defined functions in the datetime package. The selected answers are stored in a list and then embedded into an output template. The NLEP also handles the situation when no answer is found. The correct answer is then printed by the NLEP.

**Task-general demonstration prompts.** As is standard in chain-of-thought prompting (Nye et al., 2021; Wei et al., 2022), our approach uses demonstration prompts for NLEP generation. However, unlike previous approaches our demonstrations are not task-specific. For example, for all classification tasks we consider we use the *same* demonstration prompt (derived from SST2). Similarly, we use mostly the same prompt for our math and symbolic reasoning tasks. This task-general prompt is similar in spirit to zero-shot chain-of-thought prompting (Kojima et al., 2023) which adds a task-agnostic prompt (“Let’s think step-by-step”) to elicit the reasoning capabilities of LLMs in a task-agnostic way. The prompts used for the various tasks are given in Table 1, and the exact prompts are given in Appendix C. In summary, we use 4 different demonstration prompts across 16 tasks, each of which works well within a task category. Thus, while the proposed method is not fully task-agnostic in the strictest sense of the term, it is still more flexible than previous approaches that combine program synthesis with chain-of-thought prompting (Chen et al., 2022; Gao et al., 2023), which use examples from the dataset to craft prompts.

**Programmatic reasoning for natural language understanding.** Prior works on combining program synthesis with LLM-based reasoning have generally focused on math and symbolic reasoning tasks (Chen et al., 2022; Gao et al., 2023), and it has not been clear how such methods could be extended to address natural language understanding (NLU) tasks. We show that NLEPs can be straightforwardly extended to text classification tasks.

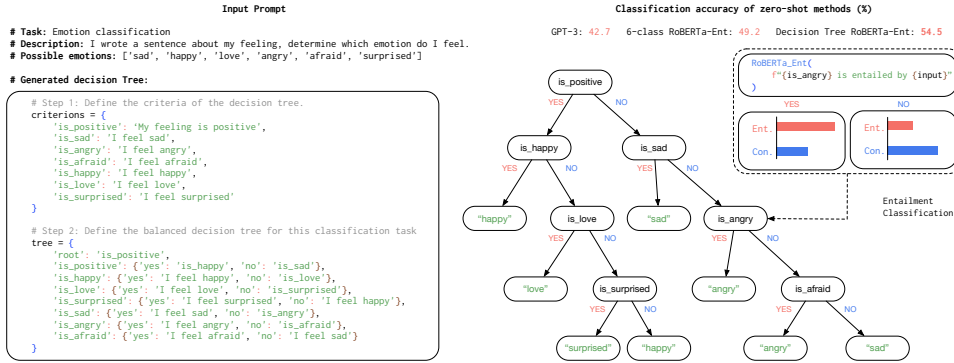
For question answering, we apply NLEP prompting and the target output is constructed by the generated programs. Classification tasks, on the other hand, are handled by a different type of NLEP consisting of a decision tree. Each node of the decision tree is annotated by a simple natural language sentence, and the Yes/No decisions at each node are handled in a zero-shot way by an entailment classifier, which has in general been shown to be an effective approach to zero-shot text classification (Obamuyide and Vlachos, 2018; Condoravdi et al., 2003; Ge et al., 2023). Concretely, given the tree we compute the entailment score between the input and the language description of each node and traverse the decision tree until a leaf node is reached. We emphasize that the topology of the tree and the language description of each node is generated by the prompted LLM. The demonstration prompt for classification tasks is given by a manually constructed example for SST2 (Wang et al., 2018). We find that this prompt can generate NLEPs containing sensible decision trees for various classification tasks without requiring task-specific examples. An example of the generated program and the corresponding decision tree is shown in Figure 2.

### 3 Experiments

We evaluate natural language embedded programs (NLEPs) on 16 tasks across three broad task categories. The tasks and corresponding prompts are summarized in Table 1.

**Math and symbolic reasoning** tasks include Tracking Shuffled Objects (7), Dyck Language, Word Sorting and Chinese Remainder Theorem from BigBench (Srivastava et al., 2023), Scheduling Meeting task from Cai et al. (2023), GSM-Hard benchmark of math word problems from Gao et al. (2023), and Game of 24 (Yao et al., 2023a). We use two examples for all tasks except for Game of 24, for which we applied a word sorting example to elicit stronger game-playing reasoning ability. The exact NLEP prompts we used are given in Appendix C.1 and C.2.

**Question answering and instruction following** tasks include the StrategyQA (Geva et al., 2021a), TruthfulQA (Lin et al., 2022), and VicunaQA (Chiang et al., 2023) benchmarks. StrategyQA requires models to answer multi-hop questions with “Yes” or “No”. TruthfulQA and VicunaQA contain questions and instructions requiring free-form responses. VicunaQA also allows us to test how NLEPs perform in the popular instruction-



**Figure 2:** A decision tree structure generated within an NLEP for emotion classification based on task description using an example program for SST2 as the prompt. The branching of each node is decided by a RoBERTa (Liu et al., 2019) text entailment model. This language-based decision tree generated by an NLEP outperforms GPT-3 and entailment-based multi-class prediction (Ge et al., 2023) without needing any task-specific examples (i.e., exemplars specific to the emotion classification dataset).

Domain	Task	Prompt
Math and Symbolic Reasoning	Tracking Shuffled Objects (7)	C.1
	Dyck Language	C.1
	Word Sorting	C.1
	Chinese Remainder Theorem	C.1
	Scheduling Meeting	C.1
	GSM-Hard	C.1
	Game of 24	C.2
Question Answering	StrategyQA	C.1
	TruthfulQA	C.3
	VicunaQA	C.3
Text Classification	SST2	C.4
	Cola	C.4
	Emotion-Classification	C.4
	Amazon Review	C.4
	Hate-Speech	C.4
	Social Bias Frame	C.4

**Table 1:** Summary descriptions of the various tasks considered in this work.

following setting. The evaluation metrics on question answering focus on accuracy, relevance, and factuality of the generated answers. The prompts in Appendix C.1 are used for StrategyQA. For TruthfulQA and VicunaQA, we added an example with a longer response to encourage more detailed response generation.

**Text classification** tasks includes tasks that require understanding of both natural language inputs and labels. We evaluate NLEP on movie-review classification (SST2; Socher et al., 2013), linguistic acceptance (COLA; Warstadt et al., 2019), emotion classification (Saravia et al., 2018), amazon review (Ni et al., 2019), hate speech detection (de Gibert et al., 2018), and stereotypes recognition (Sap et al., 2019). We use the prompts in Appendix C.1 for model-free classification. For decision tree generation, the prompts in Appendix C.4 are applied.

### 3.1 Math and Symbolic Reasoning

We compare NLEP prompting with chain-of-thought (CoT; Wei et al., 2022), program-of-thought (PoT; Chen et al., 2022), and LLMs as tool makers (LATM; Cai et al., 2023). We also compare against tree-of-thought (ToT; Yao et al., 2023a) on the Game of 24 benchmark, where ToT outperforms CoT by a significant margin (but requires many more calls to the LLM). We evaluate CoT and PoT with both task-general and task-specific demonstrations. Since LATM needs in-domain input-output pairs to create tools, we only report the results with task-specific LATM.

**Task-general prompting.** For task-general prompts we use two examples as the in-context demonstration for the math and symbolic reasoning benchmarks (see Table 1 and Appendix C). For CoT, we present two examples with intermediate reasoning represented in natural language rather than as programs. Our task-general PoT implementation takes the math and symbolic reasoning lines similar as (Chen et al., 2022) and (Gao et al., 2023), but without the step-by-step programming scheme in NLEP as an ablation.

**Task-specific prompting baselines.** We report the task-specific prompting performance as an “upper bound” for each task. For CoT, we use the same prompting settings (from 3 to 8-shot) adopted in previous studies (Cobbe et al., 2021; Cai et al., 2023; Fu et al., 2023). For PoT, we use the same in-context examples as in the task-specific CoT examples, but provide intermediate reasoning steps in Python code. On the GSM-Hard benchmark, we adopt the demonstrations (9-shot) for GSM8K used in (Chen et al., 2022). For the Chinese Remainder Theorem and Scheduling Meeting bench-

marks, we construct the in-context examples with the first three successful instances of task-general PoT. For LATM, we evaluate its performance on Tracking Shuffled Objects (7) using the provided tool and cite the results for other tasks from (Cai et al., 2023). Details are shown in Appendix D.

Program synthesis approaches (PoT and NLEP) may sometimes generate non-executable programs if lack task-specific programming demonstration. For both approaches, we select certain benchmarks to resample up to three additional programs if the returned program failed at execution. Since this condition is triggered only if program execution fails, there is no label leakage. We discuss this further in Section 4 and provide results details in Appendix B.

### 3.1.1 Results

We show the main results of NLEP prompting on six math and symbolic reasoning tasks in Table 2. An example of NLEP generated for solving a Dyck language problem is shown in Figure 3(a).

**GPT-4 Results.** Among the three approaches which use task-general prompts, NLEP outperforms both CoT and PoT on 5 of 6 tasks. The large performance gap between NLEP and CoT suggests that programmatic reasoning can enable more accurate answers. Compared to PoT, NLEP achieves significantly higher average accuracy, especially on the Dyck Language (66.4%→91.6%) and the Chinese Remainder Theorem (84.4%→97.2%) tasks. On GSM-Hard, we confirm the same phenomenon discovered by (Gao et al., 2023) where language does not further benefit the calculation accuracy with GPT-4.

NLEP also achieves comparable performance to task-specific, few-shot prompting methods. Notably, our method achieves the best performance on Tracking Shuffled Objects (7) and Dyck Language, and outperforms task-specific CoT on many benchmarks. On the Word Sorting benchmark, NLEP only fails on one instance where the input word sequence contains “steelmake” and GPT-4 automatically corrected it to “steelmaker”. We find that the high scores of task-specific PoT on Word Sorting and Chinese Remainder Theorem come from the generally applicable programming code from the in-context demonstrations.

**GPT-3.5 Results.** We observe significant performance degradation with GPT-3.5, presumably due to its limited programming capabilities. However NLEP still achieves the best average performance,

exhibiting significant improvement on 5 of 6 tasks over all baselines. On the Dyck Language benchmark, program-based strategies (PoT and NLEP with task-general prompts) failed to accomplish the problem without task-specific examples, highlighting the need for strong backbone LLMs.

**Game of 24 results.** Table 3 shows the results on the challenging Game of 24 task from (Yao et al., 2023a). Our approach also surpasses the oracle setup of IO/CoT, which calculates the success rate of IO/CoT by considering the best of 100 samples for each instance. However, unlike ToT which requires in-context demonstrations for each decomposed sub-task, NLEP prompting achieves a significant performance gain over ToT (b=1) without requiring a computationally expensive multi-chain inference procedure.

### 3.2 QA and Instruction Following

**StrategyQA.** Experiment results are presented in Table 4. With GPT-4, NLEP achieves the best performance under the task-general prompt setting and is competitive with the task-specific CoT. With GPT-3.5, although the scores of code-based strategies decrease more than CoT (PoT: 18.4%, NLEP: 20.1%, task-general CoT: 10.5%, task-specific CoT: 10.1%), NLEP still exceeds PoT by a significant margin. An example of output is shown in 3(b).

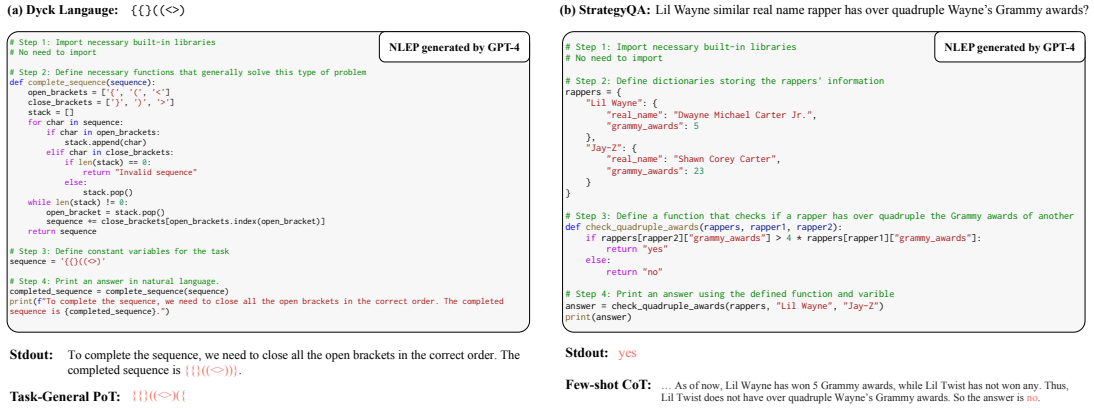
**TruthfulQA.** We also evaluate how NLEP prompting influences the factuality of question answering with the TruthfulQA benchmark (Lin et al., 2022). A fine-tuned GPT-3 model is applied for automatic scoring. In this experiment, we compare the vanilla auto-regressive text generation method against NLEP. Traditionally, such question answering tasks have been solved only with black-box language model without explicit symbolic computations due to the complexity of test questions.

The results are shown in Table 5. With GPT-4, the truth score of NLEP prompting strategy is close to standard LLM-based generation, while the informativeness score is higher. However, performance degrades significant with GPT-3.5-Turbo, indicating a strong dependence on the programming ability of the underlying language model.

**VicunaQA.** Results on the VicunaQA benchmark are shown in Figure 4, where we follow the standard approach and evaluate the answers using GPT-4. We find that GPT-4 prefers its own generations, which are generally more detailed than GPT-3.5-Turbo and NLEP responses. To control for the bias due to response lengths, we also as-

Tasks / Method	GPT-4						GPT-3.5-Turbo					
	(a) Task-Specific			(b) Task-General			(c) Task-Specific			(d) Task-General		
	CoT	PoT	LATM	CoT	PoT	NLEP	CoT	PoT	CoT	PoT	NLEP	
Tracking Shuffled Objects	100.0	100.0	100.0	81.2	98.4	100.0	68.0	6.8	51.2	88.4	74.4	
Dyck Language	63.6 <sup>†</sup>	60.8	87.5 <sup>†</sup>	39.6	66.4	91.6	20.4 <sup>†</sup>	28.4	38.0	4.0	7.2	
Word Sorting	90.9 <sup>†</sup>	100.0	99.1 <sup>†</sup>	84.4	99.6	99.6	59.2 <sup>†</sup>	100.0	75.2	100.0	99.6	
Chinese Remainder Theorem	0.0 <sup>†</sup>	100.0	100.0 <sup>†</sup>	0.0	84.4	97.2	0.0 <sup>†</sup>	100.0	0.0	72.4	96.4	
Scheduling Meeting	55.6 <sup>†</sup>	75.2	100.0 <sup>†</sup>	82.8	85.2	93.2	18.9 <sup>†</sup>	33.6	39.6	49.2	85.6	
GSM-Hard	57.4	74.1	–	54.9	69.3	67.7	45.0	63.4	42.8	52.2	54.1	
<b>Average</b>	61.3	85.0	97.3	57.2	83.9	91.6	35.3	55.4	41.1	61.0	69.6	

**Table 2:** Performance on math and symbolic reasoning tasks with both task-specific and task-general demonstration prompts. <sup>†</sup> stands for results from (Cai et al., 2023). LATM results are not available for GSM-Hard benchmark as it is hard to derive a generally applicable tool function for all test cases.



**Figure 3:** NLEP generated for solving Dyck language and StrategyQA problems. For Dyck, the instruction is “Complete the rest of the sequence, making sure that the parentheses are closed properly.” For StrategyQA, the instruction is “Answer the question with yes or no.”

Prompt	Method	Accuracy (%)
Task-specific	CoT	4
	ToT (b = 1)	45
	ToT (b = 5)	74
Task-general	PoT	52
	NLEP	66

**Table 3:** Performance on the Game of 24 benchmark. CoT and ToT stand for chain-of-thought (Wei et al., 2022) and tree-of-thought (Yao et al., 2023a) prompting respectively. <sup>†</sup> shows the results from (Yao et al., 2023a).

389 sess all responses without the requirement about  
390 details using another evaluation prompt. The eval-  
391 uation prompts with and without the requirement on  
392 details is shown in Appendix E.1 and E.2.

393 As we demonstrate in Figure 4, this assessment  
394 leads to different results on GPT-4. After remov-  
395 ing the detail requirement in the automatic scor-  
396 ing pipeline, NLEP achieves better performance.  
397 This suggests that NLEP can help GPT-4 gener-  
398 ate accurate, factual, and relevant responses. How-  
399 ever, human-generated programs for pretraining the  
400 GPT-4 models usually do not embed long pieces  
401 of natural language. As a result, the responses

generated by NLEP have a limited level of detail. 402

### 3.3 Text Classification 403

404 Finally, we evaluate whether NLEPs can be applied  
405 to solve text classification tasks that have tradition-  
406 ally been difficult for pure program synthesis-based  
407 approaches. As discussed in section 2, we manually  
408 construct a decision tree NLEP for SST2 and use it  
409 as a prompt to guide GPT models to generate deci-  
410 sion trees for other tasks only with task and label  
411 descriptions. An example input and output NLEP  
412 generated by GPT-4 for emotion classification is  
413 shown in Figure 2.

414 We compare NLEP against two baseline meth-  
415 ods. Our first baseline uses the zero-shot classifica-  
416 tion method proposed in (Ge et al., 2023) (“multi-  
417 class prompting”). This method uses the same  
418 entailment models but makes the prediction with-  
419 out the tree structure. Our second baseline asks  
420 a human expert to design a decision tree for each  
421 task also based on the SST-2 example. The re-  
422 sults shown in Table 6 show that NLEP generated  
423 by GPT-4 outperforms multi-class prompting and

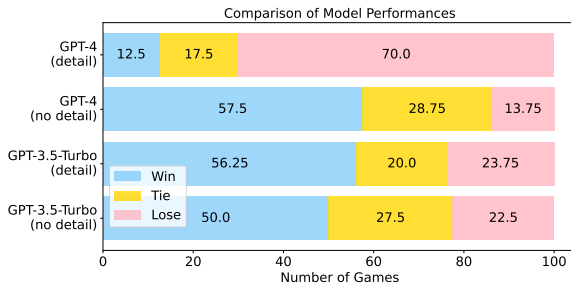
setting	GPT-4				GPT-3.5-Turbo			
	Task-specific	Task-general			Task-specific	Task-general		
	CoT	CoT	PoT	NLEP (ours)	CoT	CoT	PoT	NLEP (ours)
StrategyQA	<b>81.7</b>	78.6	68.6	81.2	71.6	68.1	50.2	61.1

**Table 4:** Performance on the StrategyQA benchmark. The experimental setup is the same as in Table 2. Note that LLMs do not always generate “Yes” or “No”, and we only predict the “Yes” label if the “Yes” string is generated explicitly. See Appendices C.1 and D for implementation details.

Foundation Model	Mode	True	Info	True * Info
GPT-4	Text	<b>76.01</b>	97.55	73.56
	NLEP	75.76	<b>99.63</b>	<b>75.40</b>
GPT-3.5-Turbo	Text	68.91	98.90	67.93
	NLEP	61.69	97.18	59.00

**Table 5:** Performance of GPT-4 and GPT-3.5-Turbo on the TruthfulQA benchmark.

Model	# NLEP >Text	Detail	% Score	% Length Bias
GPT-4	23.75	yes	93.08	72.72
		no	<b>105.06</b>	26.67
GPT-3.5-Turbo	38.75	yes	101.22	<b>3.13</b>
		no	102.50	10.34



**Figure 4:** Automatic evaluations of NLEP against standard LLM-based generation with different models. # NLEP >Text means that the % of NLEP responses containing more tokens than the baseline. **Detail** means if the evaluation metric considers details and response lengths. **Score** stands for the scores received by NLEP divided by the baseline scores (>100 means NLEP is better). **Win**, **tie**, and **lose** stand for the % of evaluation cases resulting in each category. **Length Bias** shows how much the evaluation pipeline prefers longer or shorter answers (lower means fairer, introduced in Appendix E.3).

human-generated tree baselines on most datasets.

**Model-free NLEP.** We also tried using the task-general prompt shown in C.1 to generate NLEPs that directly use programs to solve these tasks. These programs do not need any neural models and are hence very efficient (e.g., finishing the entire validation set in about 2 seconds on CPUs). The results can be found in Table 6 (“Model-free NLEP”). While not achieving the performance of entailment-based methods, the generated NLEP significantly outperforms random baselines, suggesting that this may be a viable approach for quickly extracting

simple and interpretable classifiers from LLMs.

## 4 Discussion

**Execution failures and retries.** While the task-general PoT and NLEP lack programming demonstrations for the target task, GPT-4 in general is able to generate bug-free programs as presented in Appendix B Table 10. Notably, both PoT and NLEP obtain execution error rate of 0 on Tracking Shuffled Objects (7) and Word Sort tasks. One advantage of the program synthesis approaches such as PoT and NLEP is that non-executable programs can be identified and filtered. This gives LLMs the chance to “self-correct” and generate new answers, and we take advantage of this in our math and symbolic reasoning tasks by generating up to three programs if there is an execution failure on certain benchmarks. (For fair comparison we apply this reattempting scheme to PoT as well). We ablate on this mechanism in Appendix B, Tables 7, 9 and 10. Besides effectively reducing the execution error as presented in Table 10, these retries greatly enhance the reasoning accuracy. In particular, 12% and 15.6% improvement is observed on the Chinese Remainder Theorem and the Scheduling Meeting tasks in Table 7(b). In this work we only experiment extra retries with larger temperatures for diverse sampling and leave more advanced “self-correction” algorithms (e.g., those that make use of error messages (Cai et al., 2023; Hu et al., 2023)) for future work.

**Different foundation LLMs for NLEP.** The large performance gaps of CodeLlama-7b-Instruct (Rozière et al., 2023), GPT-3.5-Turbo, and GPT-4 in Table 2 and 8 indicate that strong programming ability of underlying LLMs is vital to generate accurate responses and achieve performance improvements with NLEP. For example, on the Dyck Language task, GPT-3.5-Turbo only achieves 7.2% accuracy while GPT-4 achieves 91.6% accuracy. TruthfulQA experiments also show that NLEP could *hurt* the factuality of GPT-3.5-Turbo. Surprisingly, zero-shot CodeLlama-7b (Rozière

Model	Method	Performance (Num. Classes)					
		cola (2)	emotion (6)	amazon (5)	hsd (2)	sbic (3)	Average
RoBERTa	Multi-class Prompting	65.87	49.2	33.31	67.78	52.99	53.83
	Human-Generated Tree	<b>69.03</b>	22.20	26.88	64.85	<b>58.37</b>	48.27
	NLEP w/ GPT-3.5	56.66	35.1	33.46	67.36	38.25	46.17
	NLEP w/ GPT-4	68.94	<b>54.5</b>	<b>38.88</b>	<b>70.92</b>	55.95	<b>57.65</b>
DeBERTa	Multi-class Prompting	53.50	51.93	37.01	67.78	59.08	53.86
	Human-Generated Tree	<b>69.22</b>	32.15	33.00	<b>72.18</b>	55.02	52.31
	NLEP w/ GPT-3.5	49.66	39.00	36.18	70.29	52.49	49.52
	NLEP w/ GPT-4	68.36	<b>55.4</b>	<b>40.2</b>	70.08	<b>59.68</b>	<b>58.74</b>
None	Model-free NLEP w/o Tree	69.13	40.55	25.76	59.62	37.63	46.54

**Table 6:** Zero-shot performance of different human and LLM-generated classification schemes. The GPT-4 generated decision trees consistently exhibit significant improvement. For model-free NLEP, generated code can be executed on the entire validation set in 2 seconds and notably surpasses the random baseline, with cola notably matching the state-of-the-art performance.

et al., 2023) trained using NLEP-style data (without in-domain examples) demonstrates superiority on Tracking Shuffled Objects (7) benchmark over NLEP prompted GPT-3.5 and Word Sorting benchmark over task-general CoT prompted GPT-3.5, even with significantly fewer parameters (see details in Appendix B). It shows the potential for effective training of compact large language models, enabling them to achieve performance levels comparable to those of extremely large models.

## 5 Related Work

**Large language models for reasoning.** State-of-the-art LLMs (OpenAI, 2022, 2023; Touvron et al., 2023; Zeng et al., 2022) have shown very strong performance on complicated reasoning tasks, including commonsense (Geva et al., 2021b), math (Cobbe et al., 2021), symbolic reasoning (Suzgun et al., 2022), and programming (Austin et al., 2021; Chen et al., 2021). Tackling such tasks with LLMs often requires prompting them with demonstrations that elicit their reasoning capabilities. (Wei et al., 2022) proposed chain-of-thought prompting technique that encourages language to generate answers step-by-step. (Wang et al., 2022) found that self-consistency can further improve the performance of chain of thoughts reasoning ability. (Kojima et al., 2023) discovered that LLMs can perform reasoning without any demonstrations through adding the incantation “Let’s think step-by-step”. Tree of thoughts (Yao et al., 2023a) and graph of thoughts (Yao et al., 2023b; Besta et al., 2023) were proposed to tackle tasks that require more complicated reasoning processes. These improved reasoning methods apply chain of thoughts as the atomic reasoning step but organize reasoning

“chains” through more advanced mechanisms.

**Programs and tools.** Previous studies have found that some limitations of LLMs can be overcome by combining program synthesis techniques with prompt-based learning. Program of thoughts (Chen et al., 2022) and program aided language model (Gao et al., 2023) both translate mathematical questions to equations and use the python interpreter to ensure the correctness of the calculations. Another line of related work for enabling LLMs to use tools is through interleaving API calls during LLM generation (Qin et al., 2023; Liang et al., 2023b; Mialon et al., 2023; Tang et al., 2023). APIs can aid many tasks that are challenging for LLMs by providing tailored tools (e.g., calculators, search) that can solve specific tasks. Toolformer (Schick et al., 2023) addresses reasoning tasks by using predefined tools, and LLMs as tool makers (LATM) can implement functions solving a class of tasks based on few-shot examples (Cai et al., 2023). With these solutions, the correctness of a prediction can be ensured if correct API is called and correct inputs are selected. Existing works on combining program synthesis and tool usage with LLMs generally rely on task-specific prompts, in contrast to the more task-general prompt explored in the present work.

## 6 Conclusion

This work describes natural language embedded programs (NLEP), which flexibly combine natural language reasoning with program synthesis within prompt-based learning to tackle a variety of tasks. Our experiments demonstrate that NLEPs expand the scope of applications that can be addressed by program synthesis by more closely incorporating natural language during code generation.



## 548 Limitation

549 We found that the NLEP prompts are not suit-  
550 able for generating long-form natural language re-  
551 sponses. Experimental results on VicunaQA show  
552 that most responses generated by NLEP prompting  
553 have fewer tokens than responses obtained from  
554 usual LLM generation. This feature is expected, be-  
555 cause most naturally-occurring programs (on which  
556 the LLMs were pretrained) do not contain large  
557 chunks of natural language. Future work could con-  
558 sider incorporating (potentially synthetically gener-  
559 ated) programs with longer-form natural language  
560 within the pretraining set to enable the application  
561 of NLEP to more involved NLG tasks.

## 562 Ethical Statement

563 This work intends to design a accurate and inter-  
564 pretable reasoning framework for language that en-  
565 tails more transparent and responsible LLM appli-  
566 cations. However, the program generation method  
567 is more capable to handle different tasks in both  
568 areas of natural and program languages, infecting  
569 both humans and computing systems. As a re-  
570 sult, we believe program generation models need  
571 stronger alignment and careful management.

## 572 References

573 Jacob Austin, Augustus Odena, Maxwell Nye, Maarten  
574 Bosma, Henryk Michalewski, David Dohan, Ellen  
575 Jiang, Carrie Cai, Michael Terry, Quoc Le, et al. 2021.  
576 Program synthesis with large language models. *arXiv*  
577 *preprint arXiv:2108.07732*.

578 Maciej Besta, Nils Blach, Ales Kubicek, Robert Ger-  
579 stenberger, Lukas Gianinazzi, Joanna Gajda, Tomasz  
580 Lehmann, Michal Podstawski, Hubert Niewiadomski,  
581 Piotr Nyczyk, et al. 2023. Graph of thoughts: Solv-  
582 ing elaborate problems with large language models.  
583 *arXiv preprint arXiv:2308.09687*.

584 Tianle Cai, Xuezhi Wang, Tengyu Ma, Xinyun Chen,  
585 and Denny Zhou. 2023. Large language models as  
586 tool makers. *arXiv preprint arXiv:2305.17126*.

587 Mark Chen, Jerry Tworek, Heewoo Jun, Qiming  
588 Yuan, Henrique Ponde de Oliveira Pinto, Jared Kap-  
589 lan, Harri Edwards, Yuri Burda, Nicholas Joseph,  
590 Greg Brockman, Alex Ray, Raul Puri, Gretchen  
591 Krueger, Michael Petrov, Heidy Khlaaf, Girish Sas-  
592 try, Pamela Mishkin, Brooke Chan, Scott Gray,  
593 Nick Ryder, Mikhail Pavlov, Alethea Power, Lukasz  
594 Kaiser, Mohammad Bavarian, Clemens Winter,  
595 Philippe Tillet, Felipe Petroski Such, Dave Cum-  
596 mings, Matthias Plappert, Fotios Chantzis, Eliza-  
597 beth Barnes, Ariel Herbert-Voss, William Hebgen  
598 Guss, Alex Nichol, Alex Paino, Nikolas Tezak, Jie

Tang, Igor Babuschkin, Suchir Balaji, Shantanu Jain, 599  
William Saunders, Christopher Hesse, Andrew N. 600  
Carr, Jan Leike, Josh Achiam, Vedant Misra, Evan 601  
Morikawa, Alec Radford, Matthew Knight, Miles 602  
Brundage, Mira Murati, Katie Mayer, Peter Welinder, 603  
Bob McGrew, Dario Amodei, Sam McCandlish, Ilya 604  
Sutskever, and Wojciech Zaremba. 2021. [Evaluating 605  
large language models trained on code](#). 606

Wenhu Chen, Xueguang Ma, Xinyi Wang, and 607  
William W Cohen. 2022. Program of thoughts 608  
prompting: Disentangling computation from reason- 609  
ing for numerical reasoning tasks. *arXiv preprint 610  
arXiv:2211.12588*. 611

Zhoujun Cheng, Tianbao Xie, Peng Shi, Chengzu 612  
Li, Rahul Nadkarni, Yushi Hu, Caiming Xiong, 613  
Dragomir Radev, Mari Ostendorf, Luke Zettlemoyer, 614  
et al. 2022. Binding language models in symbolic 615  
languages. In *The Eleventh International Conference 616  
on Learning Representations*. 617

Wei-Lin Chiang, Zhuohan Li, Zi Lin, Ying Sheng, 618  
Zhanghao Wu, Hao Zhang, Lianmin Zheng, Siyuan 619  
Zhuang, Yonghao Zhuang, Joseph E. Gonzalez, Ion 620  
Stoica, and Eric P. Xing. 2023. [Vicuna: An open- 621  
source chatbot impressing gpt-4 with 90%\\* chatgpt 622  
quality](#). 623

Karl Cobbe, Vineet Kosaraju, Mohammad Bavarian, 624  
Mark Chen, Heewoo Jun, Lukasz Kaiser, Matthias 625  
Plappert, Jerry Tworek, Jacob Hilton, Reiichiro 626  
Nakano, et al. 2021. Training verifiers to solve math 627  
word problems. *arXiv preprint arXiv:2110.14168*. 628

Cleo Condoravdi, Dick Crouch, Valeria de Paiva, Rein- 629  
hard Stolle, and Daniel G. Bobrow. 2003. [Entailment, 630  
intensionality and text understanding](#). In *Proceedings 631  
of the HLT-NAACL 2003 Workshop on Text Meaning,*  
632 pages 38–45. 633

Ona de Gibert, Naiara Perez, Aitor García-Pablos, and 634  
Montse Cuadros. 2018. [Hate speech dataset from 635  
a white supremacy forum](#). In *Proceedings of the 636  
2nd Workshop on Abusive Language Online (ALW2),*  
637 pages 11–20, Brussels, Belgium. Association for  
638 Computational Linguistics. 639

Yao Fu, Litu Ou, Mingyu Chen, Yuhao Wan, Hao Peng, 640  
and Tushar Khot. 2023. [Chain-of-thought hub: A 641  
continuous effort to measure large language models’ 642  
reasoning performance](#). 643

Luyu Gao, Aman Madaan, Shuyan Zhou, Uri Alon, 644  
Pengfei Liu, Yiming Yang, Jamie Callan, and Gra- 645  
ham Neubig. 2023. Pal: Program-aided language  
646 models. In *International Conference on Machine 647  
Learning*, pages 10764–10799. PMLR. 648

Jiaxin Ge, Hongyin Luo, Yoon Kim, and James Glass. 649  
2023. Entailment as robust self-learner. *arXiv 650  
preprint arXiv:2305.17197*. 651

Mor Geva, Daniel Khashabi, Elad Segal, Tushar Khot, 652  
Dan Roth, and Jonathan Berant. 2021a. [Did aristotle 653  
use a laptop? a question answering benchmark with 654  
implicit reasoning strategies](#). 655

656	Mor Geva, Daniel Khashabi, Elad Segal, Tushar Khot, Dan Roth, and Jonathan Berant. 2021b. Did Aristotle Use a Laptop? A Question Answering Benchmark with Implicit Reasoning Strategies. <i>Transactions of the Association for Computational Linguistics (ACL)</i> .	711
657		712
658		713
659		714
660		715
661		716
662	Yi Hu, Haotong Yang, Zhouchen Lin, and Muhan Zhang. 2023. Code prompting: a neural symbolic method for complex reasoning in large language models.	717
663		718
664		719
665	Takeshi Kojima, Shixiang Shane Gu, Machel Reid, Yutaka Matsuo, and Yusuke Iwasawa. 2023. Large language models are zero-shot reasoners.	720
666		721
667		722
668	Ziyang Li, Jiani Huang, and Mayur Naik. 2023. Scallop: A language for neurosymbolic programming. <i>Proceedings of the ACM on Programming Languages</i> , 7(PLDI):1463–1487.	723
669		724
670		725
671		726
672	Yaobo Liang, Chenfei Wu, Ting Song, Wenshan Wu, Yan Xia, Yu Liu, Yang Ou, Shuai Lu, Lei Ji, Shaoguang Mao, Yun Wang, Linjun Shou, Ming Gong, and Nan Duan. 2023a. TaskMatrix.AI: Completing Tasks by Connecting Foundation Models with Millions of APIs. <i>arXiv:2303.16434</i> .	727
673		728
674		729
675		730
676		731
677		732
678	Yaobo Liang, Chenfei Wu, Ting Song, Wenshan Wu, Yan Xia, Yu Liu, Yang Ou, Shuai Lu, Lei Ji, Shaoguang Mao, Yun Wang, Linjun Shou, Ming Gong, and Nan Duan. 2023b. Taskmatrix.ai: Completing tasks by connecting foundation models with millions of apis.	733
679		734
680		735
681		736
682		737
683		738
684	Stephanie Lin, Jacob Hilton, and Owain Evans. 2022. Truthfulqa: Measuring how models mimic human falsehoods. In <i>Proceedings of the 60th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)</i> , pages 3214–3252.	739
685		740
686		741
687		742
688		743
689	Yinhan Liu, Myle Ott, Naman Goyal, Jingfei Du, Mandar Joshi, Danqi Chen, Omer Levy, Mike Lewis, Luke Zettlemoyer, and Veselin Stoyanov. 2019. Roberta: A robustly optimized bert pretraining approach. <i>arXiv preprint arXiv:1907.11692</i> .	744
690		745
691		746
692		747
693		748
694	Qing Lyu, Shreya Havaldar, Adam Stein, Li Zhang, Delip Rao, Eric Wong, Marianna Apidianaki, and Chris Callison-Burch. 2023. Faithful chain-of-thought reasoning.	749
695		750
696		751
697		752
698	Grégoire Mialon, Roberto Dessi, Maria Lomeli, Christoforos Nalmpantis, Ram Pasunuru, Roberta Raileanu, Baptiste Rozière, Timo Schick, Jane Dwivedi-Yu, Asli Celikyilmaz, Edouard Grave, Yann LeCun, and Thomas Scialom. 2023. Augmented language models: a survey.	753
699		754
700		755
701		756
702		757
703		758
704	Jianmo Ni, Jiacheng Li, and Julian McAuley. 2019. Justifying recommendations using distantly-labeled reviews and fine-grained aspects. In <i>Proceedings of the 2019 Conference on Empirical Methods in Natural Language Processing and the 9th International Joint Conference on Natural Language Processing (EMNLP-IJCNLP)</i> , pages 188–197.	759
705		760
706		761
707		762
708		763
709		764
710		765
	Maxwell Nye, Anders Johan Andreassen, Guy Gur-Ari, Henryk Michalewski, Jacob Austin, David Bieber, David Dohan, Aitor Lewkowycz, Maarten Bosma, David Luan, Charles Sutton, and Augustus Odena. 2021. Show your work: Scratchpads for intermediate computation with language models.	716
	Abiola Obamuyide and Andreas Vlachos. 2018. Zero-shot relation classification as textual entailment. In <i>Proceedings of the First Workshop on Fact Extraction and VERification (FEVER)</i> , pages 72–78, Brussels, Belgium. Association for Computational Linguistics.	717
	OpenAI. 2022. Introducing chatgpt.	722
	OpenAI. 2023. Gpt-4 technical report.	723
	Liangming Pan, Alon Albalak, Xinyi Wang, and William Yang Wang. 2023. Logic-lm: Empowering large language models with symbolic solvers for faithful logical reasoning. <i>arXiv preprint arXiv:2305.12295</i> .	724
		725
		726
		727
		728
	Bhargavi Paranjape, Scott Lundberg and Sameer Singh, Hannaneh Hajishirzi, Luke Zettlemoyer, and Marco Tulio Ribeiro. 2023. ART: Automatic multi-step reasoning and tool-use for large language models. <i>arXiv:2303.09014</i> .	729
		730
		731
		732
		733
	Yujia Qin, Shengding Hu, Yankai Lin, Weize Chen, Ning Ding, Ganqu Cui, Zheni Zeng, Yufei Huang, Chaojun Xiao, Chi Han, Yi Ren Fung, Yusheng Su, Huadong Wang, Cheng Qian, Runchu Tian, Kunlun Zhu, Shihao Liang, Xingyu Shen, Bokai Xu, Zhen Zhang, Yining Ye, Bowen Li, Ziwei Tang, Jing Yi, Yuzhang Zhu, Zhenning Dai, Lan Yan, Xin Cong, Yaxi Lu, Weilin Zhao, Yuxiang Huang, Junxi Yan, Xu Han, Xian Sun, Dahai Li, Jason Phang, Cheng Yang, Tongshuang Wu, Heng Ji, Zhiyuan Liu, and Maosong Sun. 2023. Tool learning with foundation models.	734
		735
		736
		737
		738
		739
		740
		741
		742
		743
		744
		745
	Baptiste Rozière, Jonas Gehring, Fabian Gloeckle, Sten Sootla, Itai Gat, Xiaoqing Ellen Tan, Yossi Adi, Jingyu Liu, Tal Remez, Jérémy Rapin, Artyom Kozhevnikov, Ivan Evtimov, Joanna Bitton, Manish Bhatt, Cristian Canton Ferrer, Aaron Grattafiori, Wenhan Xiong, Alexandre Défossez, Jade Copet, Faisal Azhar, Hugo Touvron, Louis Martin, Nicolas Usunier, Thomas Scialom, and Gabriel Synnaeve. 2023. Code llama: Open foundation models for code.	746
		747
		748
		749
		750
		751
		752
		753
		754
	Maarten Sap, Saadia Gabriel, Lianhui Qin, Dan Jurafsky, Noah A Smith, and Yejin Choi. 2019. Social bias frames: Reasoning about social and power implications of language. <i>arXiv preprint arXiv:1911.03891</i> .	755
		756
		757
		758
	Elvis Saravia, Hsien-Chi Toby Liu, Yen-Hao Huang, Junlin Wu, and Yi-Shin Chen. 2018. CARER: Contextualized affect representations for emotion recognition. In <i>Proceedings of the 2018 Conference on Empirical Methods in Natural Language Processing</i> , pages 3687–3697, Brussels, Belgium. Association for Computational Linguistics.	759
		760
		761
		762
		763
		764
		765

766	Timo Schick, Jane Dwivedi-Yu, Roberto Dessì, Roberta	Yao Yao, Zuchao Li, and Hai Zhao. 2023b. Be-	823
767	Raileanu, Maria Lomeli, Luke Zettlemoyer, Nicola	yond chain-of-thought, effective graph-of-thought	824
768	Cancedda, and Thomas Scialom. 2023. Toolformer:	reasoning in large language models. <i>arXiv preprint</i>	825
769	Language models can teach themselves to use tools.	<i>arXiv:2305.16582</i> .	826
770	<i>arXiv preprint arXiv:2302.04761</i> .		
771	Richard Socher, Alex Perelygin, Jean Wu, Jason	Aohan Zeng, Xiao Liu, Zhengxiao Du, Zihan Wang,	827
772	Chuang, Christopher D. Manning, Andrew Ng, and	Hanyu Lai, Ming Ding, Zhuoyi Yang, Yifan Xu,	828
773	Christopher Potts. 2013. <a href="#">Recursive deep models for</a>	Wendi Zheng, Xiao Xia, et al. 2022. Glm-130b:	829
774	<a href="#">semantic compositionality over a sentiment treebank</a> .	An open bilingual pre-trained model. <i>arXiv preprint</i>	830
775	In <i>Proceedings of the 2013 Conference on Empiri-</i>	<i>arXiv:2210.02414</i> .	831
776	<i>cal Methods in Natural Language Processing</i> , pages		
777	1631–1642, Seattle, Washington, USA. Association		
778	for Computational Linguistics.		
779	Aarohi Srivastava, Abhinav Rastogi, Abhishek Rao,		
780	et al. 2023. <a href="#">Beyond the imitation game: Quanti-</a>		
781	<a href="#">fying and extrapolating the capabilities of language</a>		
782	<a href="#">models</a> .		
783	Mirac Suzgun, Nathan Scales, Nathanael Schärli, Se-		
784	bastian Gehrmann, Yi Tay, Hyung Won Chung,		
785	Aakanksha Chowdhery, Quoc V Le, Ed H Chi, Denny		
786	Zhou, , and Jason Wei. 2022. Challenging big-bench		
787	tasks and whether chain-of-thought can solve them.		
788	<i>arXiv preprint arXiv:2210.09261</i> .		
789	Qiaoyu Tang, Ziliang Deng, Hongyu Lin, Xianpei Han,		
790	Qiao Liang, Boxi Cao, and Le Sun. 2023. <a href="#">Toolalpa-</a>		
791	<a href="#">Generalized tool learning for language models with</a>		
792	<a href="#">3000 simulated cases</a> .		
793	Hugo Touvron, Thibaut Lavril, Gautier Izacard, Xavier		
794	Martinet, Marie-Anne Lachaux, Timothée Lacroix,		
795	Baptiste Rozière, Naman Goyal, Eric Hambro,		
796	Faisal Azhar, et al. 2023. Llama: Open and effi-		
797	cient foundation language models. <i>arXiv preprint</i>		
798	<i>arXiv:2302.13971</i> .		
799	Alex Wang, Amanpreet Singh, Julian Michael, Felix		
800	Hill, Omer Levy, and Samuel Bowman. 2018. Glue:		
801	A multi-task benchmark and analysis platform for		
802	natural language understanding. In <i>Proceedings of</i>		
803	<i>the 2018 EMNLP Workshop BlackboxNLP: Analyz-</i>		
804	<i>ing and Interpreting Neural Networks for NLP</i> , pages		
805	353–355.		
806	Xuezhi Wang, Jason Wei, Dale Schuurmans, Quoc Le,		
807	Ed Chi, and Denny Zhou. 2022. Self-consistency im-		
808	proves chain of thought reasoning in language mod-		
809	els. <i>arXiv preprint arXiv:2203.11171</i> .		
810	Alex Warstadt, Amanpreet Singh, and Samuel R. Bow-		
811	man. 2019. <a href="#">Neural network acceptability judgments</a> .		
812	<i>Transactions of the Association for Computational</i>		
813	<i>Linguistics</i> , 7:625–641.		
814	Jason Wei, Xuezhi Wang, Dale Schuurmans, Maarten		
815	Bosma, Ed Chi, Quoc Le, and Denny Zhou. 2022.		
816	Chain of thought prompting elicits reasoning in large		
817	language models. <i>arXiv preprint arXiv:2201.11903</i> .		
818	Shunyu Yao, Dian Yu, Jeffrey Zhao, Izhak Shafran,		
819	Thomas L Griffiths, Yuan Cao, and Karthik		
820	Narasimhan. 2023a. Tree of thoughts: Deliberate		
821	problem solving with large language models. <i>arXiv</i>		
822	<i>preprint arXiv:2305.10601</i> .		

## A Additional Example Comparing with ChatGPT-4

## B Additional results and analysis on math and symbolic reasoning

We present the detailed experimental results of math and symbolic reasoning tasks in Tables 7 to 9, with execution failure analysis in Table 10.

**GPT Results.** We report the results of Table 2 with more details in Table 7. The effect of extra retries described in Section 4 is highlighted with (→).

**CodeLlama-7b-Instruct Results.** To investigate the effect of NLEP on compact large language models, we report the results with CodeLlama-7b-Instruct (Rozière et al., 2023) in Table 8. Following the guidance of the instruction-following models<sup>1</sup>, we employ a chat session to include task-specific and task-general prompts as previous turns by interleaving the “user” and “assistant” messages with a system message “Provide answers in Python” at the beginning. Hence, we only treat bug-free Python programs that have the desired results after execution as correct answers, regardless of natural language outputs since we explicitly prompt CodeLlama to generate the answer in Python Unlike the prominent performance of GPT-4, the positive impact of NLEP with CodeLlama-7b-Instruct is diminished due to the much smaller model size and greatly reduced programming capability. Although NLEP prompting outperforms the task-general PoT by a large margin on Chinese Remainder Theorem and Scheduling Meeting benchmarks, a non-negligible performance gap is observed between NLEP and task-specific PoT on most tasks.

To further investigate the benefits of NLEP, we fine-tune a CodeLlama-7b (Rozière et al., 2023) model using NLEP-style instances, resulting in a variant that we term NLEP-CodeLlama. Note that our training corpus does not include specific evaluation tasks. During the evaluation phase, we adopt zero-shot prompting strategy, where the model is provided with only test instances without in-context demonstrations. As presented in Table 8(c), zero-shot NLEP-CodeLlama exhibits consistent performance improvement on 5 of 6 tasks. The only exception is the Chinese Remainder Theorem benchmark, which is notably more complex in nature. Remarkably, zero-shot NLEP-CodeLlama demonstrates superior performance on Word Sort-

ing benchmark when compared to task-general CoT prompted GPT-3.5-Turbo, and outperforms NLEP prompted GPT-3.5-Turbo on Tracking Shuffled Objects (7) benchmark, despite a considerably lower parameter size.

**Game of 24 Results.** We present the detailed experimental results on the Game of 24 benchmark in Table 9. The effect of extra retries described in Section 4 is highlighted with (→).

**Execution Failure Analysis.** We present the execution failure statistics of code-based reasoning strategies in Table 10. The effect of extra retries described in Section 4 is highlighted with (→). Note that different from task-specific PoT with demonstrations showing how to return the desired outputs in Python program, e.g.,

```
# Python code, return ans
Alice = "striker"
Bob = "right winger"
Claire = "left winger"
Dave = "benchwarmer"
Eve = "goalkeeper"
Fred = "center midfielder"
Gertrude = "cheerleader"
Eve, Claire = Claire, Eve
Gertrude, Alice = Alice, Gertrude
Fred, Bob = Bob, Fred
Dave, Fred = Fred, Dave
Fred, Bob = Bob, Fred
Bob, Eve = Eve, Bob
Claire, Alice = Alice, Claire
ans = Gertrude
```

we need to design rules to extract the target answers from the execution results of task-general PoT and NLEP since they are allowed to generate free-form outputs. For example, given the generated programs,

```
# Step 1: Import necessary built-in libraries
# No need to import

# Step 2: Define necessary functions that generally
solve this type of problem
def swap_positions(positions, swaps):
    for swap in swaps:
        positions[swap[0]], positions[swap[1]] =
            positions[swap[1]], positions[swap[0]]
    return positions

# Step 3: Define constant variables for the task
positions = {
    "Alice": "striker",
    "Bob": "right winger",
    "Claire": "left winger",
    "Dave": "benchwarmer",
    "Eve": "goalkeeper",
    "Fred": "center midfielder",
    "Gertrude": "cheerleader"
}

swaps = [
    ("Eve", "Claire"),
    ("Gertrude", "Alice"),
    ("Fred", "Bob"),
    ("Dave", "Fred"),
    ("Fred", "Bob"),
    ("Bob", "Eve"),
    ("Claire", "Alice")
]

# Step 4: Print an answer in natural language.
final_positions = swap_positions(positions, swaps)
```

<sup>1</sup><https://github.com/facebookresearch/codellama>

	Tracking Shuffled Objects (7)	Dyck Language	Word Sorting	Chinese Remainder Theorem	Scheduling Meeting	GSM-Hard
(a) Task-Specific Prompting: GPT-4						
CoT	<u>100.0</u>	63.6 <sup>†</sup>	90.9 <sup>†</sup>	0.0 <sup>†</sup>	55.6 <sup>†</sup>	57.4
PoT	<u>100.0</u>	60.8	<u>100.0</u>	<u>100.0</u>	75.2	<u>74.1</u>
LATM	<u>100.0</u>	87.5 <sup>†</sup>	99.1 <sup>†</sup>	<u>100.0</u> <sup>†</sup>	<u>100.0</u> <sup>†</sup>	-
(b) Task-General Prompting: GPT-4						
CoT	81.2	39.6	84.4	0.0	82.8	54.9
PoT	98.4	66.4	<u>99.6</u>	76.4 (→84.4)	<u>84.4</u> (→85.2)	<u>69.3</u>
NLEP (Ours)	<u>100.0</u>	<u>91.6</u>	<u>99.6</u>	<u>85.2</u> (→97.2)	<u>77.6</u> (→93.2)	<u>67.7</u>
(c) Task-Specific Prompting: GPT-3.5-Turbo						
CoT	<u>68.0</u>	20.4 <sup>†</sup>	59.2 <sup>†</sup>	0.0 <sup>†</sup>	18.9 <sup>†</sup>	45.0
PoT	6.8	<u>28.4</u>	<u>100.0</u>	<u>100.0</u>	<u>33.6</u>	<u>63.4</u>
(d) Task-General Prompting: GPT-3.5-Turbo						
CoT	51.2	<u>38.0</u>	75.2	0.0	39.6	42.8
PoT	<u>88.4</u>	4.0	<u>100.0</u>	58.4 (→72.4)	46.8 (→49.2)	39.0 (→52.2)
NLEP (Ours)	74.4	7.2	99.6	<u>94.8</u> (→96.4)	<u>75.2</u> (→85.6)	<u>50.9</u> (→54.1)

**Table 7:** Performance on six reasoning tasks. <sup>†</sup> stands for results from LATM (Cai et al., 2023). Results with <sup>†</sup> and of LATM are reported on the last 240 instances with the first 10 instances as training and validation sets for tool making according to LATM’s design. LATM is not appropriate for GSM-Hard benchmark as it is hard to derive a generally applicable tool function for all test cases. We mainly experiment LATM with GPT-4 as the tool maker since (Cai et al., 2023) found that GPT-3.5 fails in all 5 trials on hard tasks like Tracking Shuffled Objects (5). If the generated tool is not general enough or only suitable for training samples, the tool using phase will fail. We perform experiments using GPT-4 and GPT-3.5-Turbo with a sampling temperature of 0 for all settings except PoT and NLEP on GSM-Hard in (b) which use a temperature of 0.5 to increase the sampling diversity. Since task-general PoT and NLEP lack task-specific programming instruction, they may generate non-executable Python programs. We select some settings and give each instance failed at execution up to three additional trials with temperature=0.4 to diversify the possible outputs. No label leakage is involved in this process as only the success of execution is used as a judgement. We report the results with these extra retries on execution failures in (→). The highest score among each sub-table (a), (b), (c) and (d) is underlined and the best result for each task is in **bold**.

```
952 print(f"At the end of the match, Gertrude is playing
953       {final_positions['Gertrude']}".)
```

954 we need to extract “striker”, the target answer,  
955 from the execution results “At the end of the match,  
956 Gertrude is playing striker.”.

957 Although task-specific PoT explicitly instructs  
958 the model to generate a required line of code  
959 ans={desired\_result} to capture the final result  
960 by adding

```
961 # Python code, return ans
```

962 in the instruction and explicitly provides task-  
963 specific examples in the prompt, e.g.,

```
964 # Python code, return ans
965 ... ..
966 ans = Gertrude
```

967 we found that GPT models sometimes fail to  
968 accomplish the requirement. On Tracking Shuf-  
969 fled Objects (7) benchmark, only 17 programs  
970 out of 250 test cases generated by GPT-3.5 suc-  
971 cessfully contain the required line of code, i.e.,  
972 ans={desired\_result}, which explains the huge  
973 number of failures (233). In addition to already  
974 having task-specific programming examples, the

975 inability to distinguish between the erroneous pro-  
976 grams and lack of required line of code is another  
977 reason why we do not apply the error retries on  
978 task-specific PoT.

	Tracking Shuffled Objects (7)	Dyck Language	Word Sorting	Chinese Remainder Theorem	Scheduling Meeting	GSM-Hard
(a) Task-Specific Prompting: CodeLlama-7b-Instruct						
PoT	95.6	15.2	78.0	100.0	32.0	23.9
(b) Task-General Prompting: CodeLlama-7b-Instruct						
PoT	21.2	0.8	98.0	0.0	4.0	22.9
NLEP (Ours)	30.0	0.8	93.2	18.8	24.8	15.2
(c) Zero-shot Prompting: NLEP Trained CodeLlama-7b						
Zero-shot	84.4	1.2	98.4	0.0	34.4	16.8

**Table 8:** Performance on six math and symbolic reasoning tasks. We directly prompt CodeLlama-7b-Instruct (Rozière et al., 2023) with task-specific and task-general demonstrations in (a) and (b). The corresponding experimental setup remains consistent with these outlined in Table 2 except we employ a chat session. In this instance, we incorporate the in-context demonstrations as previous turns by interleaving the “user” and “assistant” messages. We further train CodeLlama-7b (Rozière et al., 2023) with NLEP examples and report the zero-shot performance in (c). We adhere to the configuration employed in the GPT-series experiments, wherein we prepend the in-context demonstrations before each test instance.

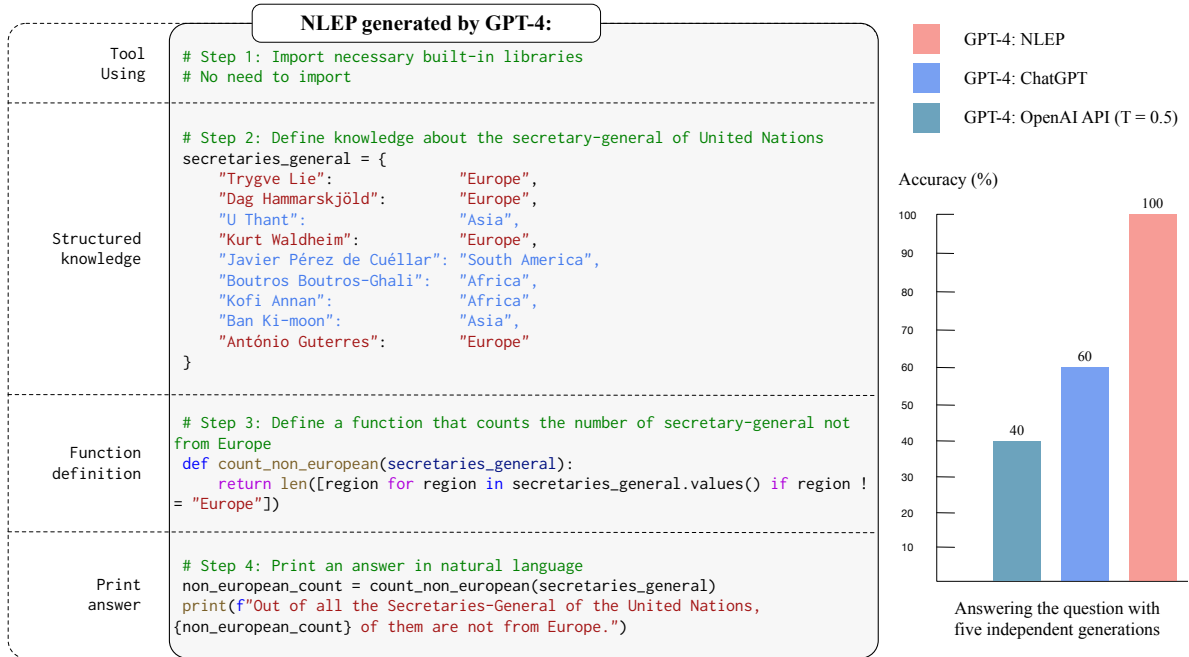
setting	Task-Specific						Task-General	
	IO	CoT	IO (best of 100)	CoT (best of 100)	ToT (b=1)	ToT (b=5)	PoT	NLEP (ours)
Game of 24	7.3 <sup>†</sup>	4.0 <sup>†</sup>	33.0 <sup>†</sup>	49.0 <sup>†</sup>	45.0 <sup>†</sup>	<b>74.0<sup>†</sup></b>	52 (→52)	63 (→66)

**Table 9:** Performance on Game of 24 benchmark. <sup>†</sup> stands for results from (Yao et al., 2023a).

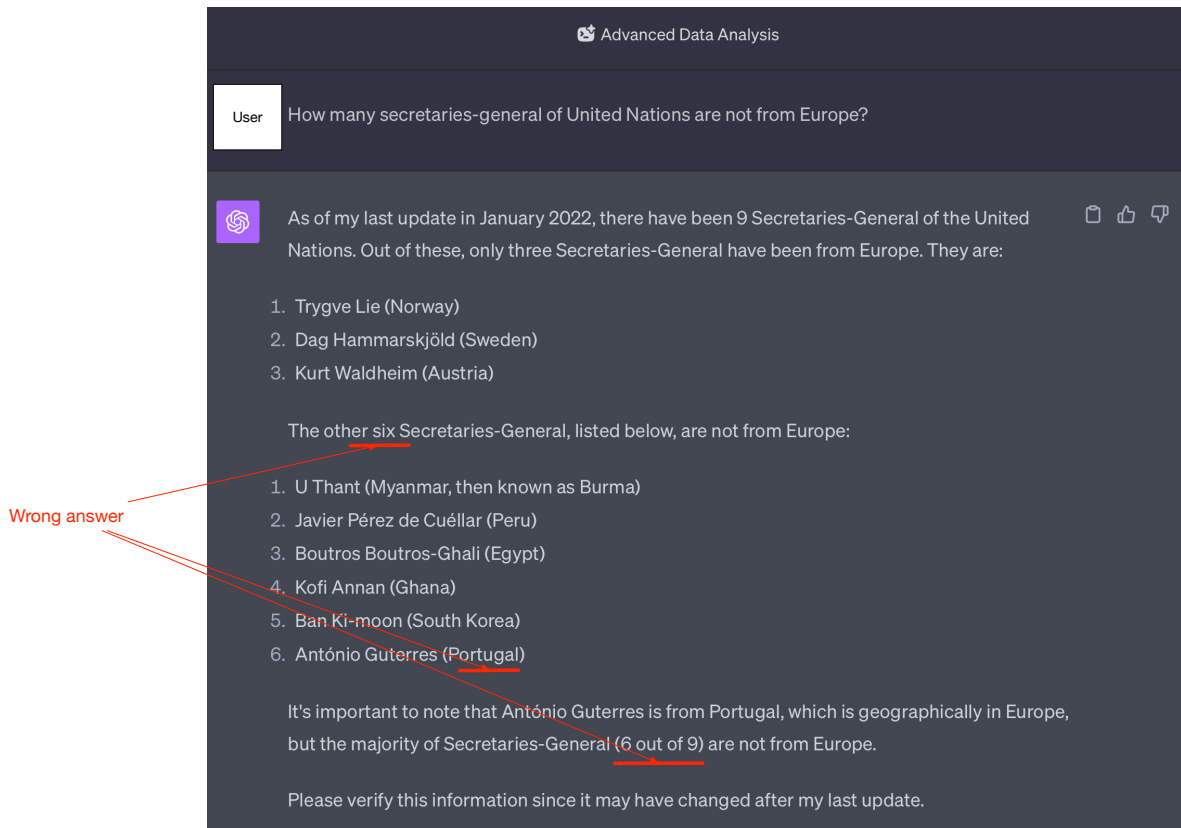
	GPT-4			GPT-3.5-Turbo		
	Task-Specific	Task-General		Task-Specific	Task-General	
	PoT	PoT	NLEP	PoT	PoT	NLEP
Track Shuffled Objects (7)	0	0	0	233	26	24
Dyck Language	16	24	10	32	81	26
Word Sort	0	0	0	0	0	0
Chinese Remainder Theorem	0	32 (→0)	37 (→6)	0	46 (→7)	4 (→0)
Scheduling Meeting	0	3 (→0)	43 (→0)	2	15 (→2)	36 (→0)
GSM-Hard	17	6	8	31	464 (→145)	95 (→13)

**Table 10:** Execution failure statistics on six math and symbolic reasoning tasks. Results with extra retries are reported in (→). For task-specific PoT, we report the execution error statistics with None as the return value of safe\_execute() function following the source code of (Chen et al., 2022): <https://github.com/wenhuchen/Program-of-Thoughts/blob/main/tool.py>. It includes instances where the generated programs do not contain the required line of code: ans={desired\_result}, which are explicitly required in the instruction and given in the prompt demonstration. Under this scenario, we cannot capture the execution results of task-specific PoT.

**Instruction:** How many secretaries-general of United Nations are **not from Europe**?



**Stdout:** Out of all the Secretaries-General of the United Nations, **5 of them are not from Europe**.



**Figure 5:** NLEP answering a question which requires numeric reasoning of structured knowledge. ChatGPT-4 code interpreter (currently the advanced data analysis option) constantly prefers to answer this question with plain natural language.

979 **C Prompts for Task-General Strategies**

980 **C.1 Prompts for Table 2 and 4**

981 We list the prompts for the task-general chain-of-  
982 thought (CoT), our implementation of program-of-  
983 thoughts (PoT) and the proposed natural language  
984 embedded programs (NLEP) strategies in the fol-  
985 lowing code segments. They share the same two  
986 NLEP examples (one is for natural language rea-  
987 soning and the other is for mathematical reasoning)  
988 but with different forms of intermediate reasoning  
989 steps (e.g., code, text etc.) to evaluate the average  
990 performance of different strategies.

991 **Prompt for task-general chain-of-thought**  
992 **(CoT) in Table 2 and 4.** The detailed intermediate  
993 natural language reasoning chains are generated by  
994 prompting GPT-4 given the input and target output.

```
995 '''  
996 Answer the problem based on the given instruction  
997 and input.  
998  
999 ### Instruction: Identify the odd one out.  
1000 ### Input: Twitter, Instagram, Telegram  
1001 ### Answer:  
1002 Let's think step by step.  
1003 1. Start by understanding the task instruction. The  
1004 task is to identify  
1005 the odd one out from a given list.  
1006 2. Look at the input. The input consists of three  
1007 items: Twitter,  
1008 Instagram, and Telegram.  
1009 3. Identify what these items are. Twitter and  
1010 Instagram are social media  
1011 platforms, while Telegram is a messaging app.  
1012 4. Compare the items to find the odd one out.  
1013 Twitter and Instagram are
```

```
primarily used for sharing information, images, and  
videos. On the other  
1015 hand, Telegram is mainly used for instant messaging  
1016 and voice-over-IP  
1017 service.  
1018 5. Determine the odd one out based on the comparison  
1019 . In this case,  
1020 Telegram is the odd one out because it serves a  
1021 different primary  
1022 function compared to Twitter and Instagram.  
1023 6. Formulate the target output. The target output  
1024 should clearly state  
1025 that Telegram is the odd one out and provide the  
1026 reason why it is so. The  
1027 reason being that Twitter and Instagram are social  
1028 media platforms mainly  
1029 for sharing information, images, and videos while  
1030 Telegram is a cloud-  
1031 based instant messaging and voice-over-IP service.  
1032 The correct answer is Telegram.  
1033
```

```
### Instruction: Use the given data to calculate the  
median.  
### Input: [2, 3, 7, 8, 10]  
### Answer:  
Let's think step by step.  
1. Start by understanding the task, which is to  
calculate the median of a  
given data set. The median is the middle value in a  
sorted, ascending or  
descending, list of numbers.  
2. Look at the given input, which is a list of  
numbers: [2, 3, 7, 8, 10].  
3. Notice that the list is already sorted in  
ascending order. If it wasn'  
t, the first step would be to sort it.  
4. Understand that to find the median, we need to  
find the middle value.  
If the list has an odd number of observations, the  
median is the middle  
number. If the list has an even number of  
observations, the median is the  
average of the two middle numbers.  
5. Count the number of values in the list. There are  
5 values, which is  
an odd number, so the median will be the middle  
value.  
6. Identify the middle value. Since there are 5  
values, the middle value  
is the 3rd value when counting from either end.  
7. Find the 3rd value in the list, which is 7.  
8. Conclude that the median of the given data set is  
7.  
The correct answer is 7.  
...  
1037  
1038  
1039  
1040  
1041  
1042  
1043  
1044  
1045  
1046  
1047  
1048  
1049  
1050  
1051  
1052  
1053  
1054  
1055  
1056  
1057  
1058  
1059  
1060  
1061  
1062  
1063  
1064  
1065  
1066  
1067  
1068  
1069
```

**Prompt for task-general NLEP in Table 2 and**  
4. The intermediate program reasoning chains are  
generated by prompting GPT-4 given the input and  
target output.

```
Write a bug-free Python program that can generate  
the answer to the given instruction when  
correctly executed.  
1074  
1075  
1076  
1077  
### Instruction: Identify the odd one out.  
### Input: Twitter, Instagram, Telegram  
### Python program:  
...  
# Step 1: Import necessary built-in libraries  
from collections import OrderedDict  
1081  
1082  
# Step 2: Define necessary functions that generally  
solve this type of problem  
1083  
1084  
def find_odd_one_out(services, input_services):  
1085  
1086  
1087  
1088  
1089  
1090  
1091  
1092  
1093  
1094  
1095  
1096  
# Step 3: Define constant variables for the task
```



```

1097 services = OrderedDict([
1098     ("Twitter", "a social media platform mainly for
1099     sharing information, images and videos"),
1100     ("Instagram", "a social media platform mainly
1101     for sharing information, images and videos
1102     "),
1103     ("Telegram", "a cloud-based instant messaging
1104     and voice-over-IP service"),
1105 ])
1106
1107 input_services = ["Twitter", "Instagram", "Telegram
1108 "]
1109
1110 # Step 4: Print an answer in natural language.
1111 odd_one_out = find_odd_one_out(services,
1112     input_services)
1113 if odd_one_out:
1114     other_services = [service for service in
1115         input_services if service != odd_one_out]
1116     print(f"The odd one out is {odd_one_out}. {
1117         other_services[0]} and {other_services[1]}
1118         are {services[other_services[0]]} while {
1119         odd_one_out} is {services[odd_one_out]}.\n
1120         The correct answer is {odd_one_out}.")
1121 ...
1122
1123 ### Instruction: Use the given data to calculate the
1124 median.
1125 ### Input: [2, 3, 7, 8, 10]
1126 ### Python program:
1127 ...
1128 # Step 1: Import necessary built-in libraries
1129 # No need to import
1130
1131 # Step 2: Define necessary functions that generally
1132 solve this type of problem
1133 def calculate_median(data):
1134     data.sort()
1135     length = len(data)
1136     if length % 2 == 0:
1137         return (data[length//2] + data[length//2 -
1138             1]) / 2
1139     else:
1140         return data[length//2]
1141
1142 # Step 3: Define constant variables for the task
1143 data = [2, 3, 7, 8, 10]
1144
1145 # Step 4: Print an answer in natural language.
1146 median = calculate_median(data)
1147 print(f"To find the median of a data set, we need to
1148 arrange the data in ascending order and then
1149 find the middle value. In this case, the given
1150 data is already arranged in ascending order.
1151 Since there are {len(data)} values in the data
1152 set, the median will be the middle value, which
1153 is the {len(data)//2 + 1}rd value. Hence, the
1154 median of the given data set is {median}.\n
1155 The correct answer is {median}.")
1156 ...

```

**Prompt for our task-general program-of-thoughts (PoT) implementation in Table 2 and 4.** The intermediate program reasoning chains are obtained from NLEP. We keep the key implementation steps of Python programs but eliminate the step-by-step programming instructions and detailed natural language elaboration as an ablation.

```

1164 Write a bug-free Python program that can generate
1165 the answer to the given instruction when
1166 correctly executed.
1167
1168 ### Instruction: Identify the odd one out.
1169 ### Input: Twitter, Instagram, Telegram
1170 ### Python program:
1171 ...
1172 softwares = {
1173     'social media': ['twitter', 'instagram'],
1174     'communication': ['telegram']
1175 }
1176 for genre, apps in softwares.items():
1177     if len(apps) == 1:

```

```

1178         print(apps[0])
1179 ...
1180 ### Instruction: Use the given data to calculate the
1181 median.
1182 ### Input: [2, 3, 7, 8, 10]
1183 ### Python program:
1184 ...
1185 data = [2, 3, 7, 8, 10]
1186 data.sort()
1187 length = len(data)
1188 if length % 2 == 0:
1189     print((data[length//2] + data[length//2 - 1]) /
1190         2)
1191 else:
1192     print(data[length//2])
1193 ...
1194

```

## C.2 Prompts for Table 3

The Game of 24 task is much more challenging and we replace the first example in Appendix C.1 with a word sorting example to elicit stronger reasoning ability.

**Prompt for task-general NLEP in Table 3.** The intermediate program reasoning chains are generated by prompting GPT-4 given the input and target output.

```

1204 Write a bug-free Python program that can generate
1205 the answer to the given instruction when
1206 correctly executed.
1207
1208 ### Instruction: Arrange the following words to make
1209 the longest possible word.
1210 ### Input: the, had, not, been
1211 ### Python program:
1212 ...
1213 # Section 1: Define necessary functions and
1214 calculate intermediate variables
1215 def longest_word(words):
1216     from itertools import permutations
1217     all_words = [''.join(p) for p in permutations
1218         (''.join(words))]
1219     all_words.sort(key=len, reverse=True)
1220     with open('english_words.txt') as word_file: #
1221         Assuming you have a list of english words
1222         english_words = set(word.strip().lower() for
1223             word in word_file)
1224     for word in all_words:
1225         if word.lower() in english_words:
1226             return word
1227     return None
1228
1229 # Section 2: Define constant variables
1230 words = ["the", "had", "not", "been"]
1231
1232 # Section 3: Insert variables in text outputs using
1233 f-strings.
1234 longest = longest_word(words)
1235 if longest:
1236     print(f"The longest word that can be made from
1237         the letters in the words '{', '.join(words
1238         )}' is '{longest}'.")
1239 ...
1240
1241 ### Instruction: Use the given data to calculate the
1242 median.
1243 ### Input: [2, 3, 7, 8, 10]
1244 ### Python program:
1245 ...
1246 # Step 1: Import necessary built-in libraries
1247 # No need to import
1248
1249 # Step 2: Define necessary functions that generally
1250 solve this type of problem
1251 def calculate_median(data):
1252     data.sort()
1253     length = len(data)
1254     if length % 2 == 0:

```

```

1255         return (data[length//2] + data[length//2 -
1256             1]) / 2
1257     else:
1258         return data[length//2]
1259
1260 # Step 3: Define constant variables for the task
1261 data = [2, 3, 7, 8, 10]
1262
1263 # Step 4: Print an answer in natural language.
1264 median = calculate_median(data)
1265 print(f"To find the median of a data set, we need to
1266     arrange the data in ascending order and then
1267     find the middle value. In this case, the given
1268     data is already arranged in ascending order.
1269     Since there are {len(data)} values in the data
1270     set, the median will be the middle value, which
1271     is the {len(data)//2 + 1}rd value. Hence, the
1272     median of the given data set is {median}.")
1273 ...

```

1274 **Prompt for our task-general program-of-**  
1275 **thoughts (PoT) implementation in Table 3.** The  
1276 intermediate program reasoning chains are ob-  
1277 tained from NLEP. We keep the key implemen-  
1278 tation steps of Python programs but eliminate the  
1279 step-by-step programming instructions and detailed  
1280 natural language elaboration as an ablation.

```

1281 Write a bug-free Python program that can generate
1282 the answer to the given instruction when
1283 correctly executed.
1284
1285 ### Instruction: Arrange the following words to make
1286 the longest possible word.
1287 ### Input: the, had, not, been
1288 ### Python program:
1289 ...
1290 def longest_word(words):
1291     from itertools import permutations
1292     all_words = [''.join(p) for p in permutations
1293                 (''.join(words))]
1294     all_words.sort(key=len, reverse=True)
1295     with open('english_words.txt') as word_file: #
1296         Assuming you have a list of english words
1297         english_words = set(word.strip().lower() for
1298                             word in word_file)
1299     for word in all_words:
1300         if word.lower() in english_words:
1301             return word
1302     return None
1303
1304 words = ["the", "had", "not", "been"]
1305
1306 longest = longest_word(words)
1307 if longest:
1308     print(longest)
1309 ...
1310
1311 ### Instruction: Use the given data to calculate the
1312 median.
1313 ### Input: [2, 3, 7, 8, 10]
1314 ### Python program:
1315 ...
1316 data = [2, 3, 7, 8, 10]
1317 data.sort()
1318 length = len(data)
1319 if length % 2 == 0:
1320     print((data[length//2] + data[length//2 - 1]) /
1321           2)
1322 else:
1323     print(data[length//2])
1324 ...

```

### 1325 C.3 Prompts for NLEP in Table 5 and Figure 1326 4

1327 For experiments in TruthfulQA and VicunaQA,  
1328 we added the following example into the NLEP  
1329 prompts shown in Appendix C.1 to encourage gen-

```

erating more detailed responses: 1330
# Write a bug-free Python program that can generate 1331
the answer to the given instruction when 1332
correctly executed. Do not ask for user input. 1333
For reasoning tasks, define functions first and 1334
then define variables. For knowledge intensive 1335
tasks, define variables before defining 1336
functions. Do not define any variable that 1337
directly stores the final answer. If there is a 1338
knowledge definition step, use dictionaries to 1339
store both the knowledge and detailed 1340
explanation. 1341
### Instruction: Discuss the causes of the Great 1342
Depression 1343
### Input: None 1344
### Python program: 1345
... 1346
# Step 1: Import necessary built-in libraries 1347
# No need to import 1348
# Step 2: Define dictionaries storing detailed 1349
knowledge about the great depression 1350
depression_name = "The Great Depression" 1351
depression_period = "1929-1939" 1352
depression_countries = "the United States and 1353
countries around the world" 1354
depression_causes = { 1355
    "Stock Market Crash of 1929": "In October of 1356
1929, the stock market experienced a 1357
significant fall that wiped out millions of 1358
investors. This event is considered by 1359
many to be the initial trigger of the Great 1360
Depression.", 1361
    "Overproduction": "During the 1920s, many 1362
industries produced more goods than 1363
consumers wanted or could afford. This 1364
ultimately led to a decline in demand for 1365
goods, causing job loss, lower wages, and 1366
business failure.", 1367
    "High Tariffs and War Debts": "Protectionist 1368
trade policies in the form of high tariffs 1369
led to a decline in global trade, as other 1370
countries retaliated with tariffs of their 1371
own. Additionally, many countries were 1372
struggling to repay war debts, which led to 1373
economic instability.", 1374
    "Bank Failures": "As demand for goods declined, 1375
many banks began to fail, causing a loss of 1376
confidence in the banking system. This led 1377
to a massive withdrawal of money from 1378
banks, causing even more banks to fail.", 1379
    "Drought Conditions": "The Dust Bowl was a 1380
severe drought and dust storm that hit the 1381
Great Plains region of the United States in 1382
the 1930s. This had a significant impact 1383
on agriculture, causing many farmers to 1384
lose their land and livelihoods which 1385
worsened the effects of the depression." 1386
} 1387
# Step 3: Define necessary functions that generally 1388
solve this type of problem 1389
# Do not need to define functions 1390
# Step 4: Print the answer and explain in natural 1391
language by calling the information in the 1392
defined knowledge dictionary 'depression_causes' 1393
print(f"{depression_name} was a period of economic 1394
decline that lasted from {depression_period}, 1395
making it the longest-lasting depression in 1396
modern history. It affected not only { 1397
depression_countries}, causing substantial 1398
social and economic upheaval.\n") 1399
print(f"There were several major causes of { 1400
depression_name}, which include:\n") 1401
# List causes and explanations in 'depression_causes' 1402
with a for-loop. 1403
for i, (cause, description) in enumerate( 1404
    depression_causes.items(), 1): 1405
    print(f"{i}. {cause} - {description}\n") 1406
print(f"Overall, {depression_name} was caused by a 1407
combination of factors, including economic, 1408
environmental, and political factors. Its 1409
impact was widespread, affecting millions of 1410
1411
1412
1413
1414
1415
1416

```

```

1417     ... people around the world.")
1418     ...
1419
1420     ### Instruction: Identify the odd one out.
1421     ### Input: Twitter, Instagram, Telegram
1422     ### Python program:
1423     ...
1424     # Step 1: Import necessary built-in libraries
1425     from collections import OrderedDict
1426
1427     # Step 2: Define dictionaries storing detailed
1428     knowledge about the main function of each
1429     application
1430     services = {
1431         "Twitter": "a social media platform mainly for
1432         sharing information, images and videos",
1433         "Instagram": "a social media platform mainly for
1434         sharing information, images and videos",
1435         "Telegram": "a cloud-based instant messaging and
1436         voice-over-IP service",
1437     }
1438
1439     # Step 3: Define a function that finds the different
1440     application
1441     def find_odd_one_out(services, input_services):
1442         descriptions = [services[service] for service in
1443         input_services]
1444         for description in descriptions:
1445             if descriptions.count(description) == 1:
1446                 return input_services[descriptions.index
1447                 (description)]
1448             return None
1449
1450     # Step 4: Print the answer in natural language by
1451     calling the information stored in 'services'
1452     and the defined function 'find_odd_one_out'
1453     input_services = ["Twitter", "Instagram", "Telegram
1454     "]
1455     odd_one_out = find_odd_one_out(services,
1456     input_services)
1457     if odd_one_out:
1458         other_services = [service for service in
1459         input_services if service != odd_one_out]
1460         print(f"The odd one out is {odd_one_out}. {
1461         other_services[0]} and {other_services[1]}
1462         are {services[other_services[0]]} while {
1463         odd_one_out} is {services[odd_one_out]}.")
1464     ...
1465
1466     ### Instruction: Calculate the total surface area of
1467     a cube with a side length of 5 cm.
1468     ### Input: None
1469     ### Python program:
1470     ...
1471     # Step 1: Import necessary built-in libraries
1472     # No need to import
1473
1474     # Step 2: Define a function that calculate the
1475     surface area of cubes
1476     def calculate_surface_area(side_length):
1477         return 6 * (side_length ** 2)
1478
1479     # Step 3: Define dictionaries storing the cube
1480     information
1481     cube = {
1482         "side_length": 5 # Side length of the cube
1483     }
1484
1485     # Step 4: Print a step-by-step calculation answer in
1486     natural language using the defined function
1487     and variable
1488     side_length = cube["side_length"]
1489     surface_area = calculate_surface_area(side_length)
1490     print(f"The surface area of a cube is found by
1491     calculating the area of one of its faces and
1492     multiplying it by six (since a cube has six
1493     faces). The area of a cube face is simply its
1494     side length squared.\n")
1495     print(f"Thus for this particular cube:")
1496     print(f"Surface Area = 6 x (Side Length)\^2")
1497     print(f"                = 6 x ({side_length} cm)\^2")
1498     print(f"                = 6 x {side_length**2} cm\^2")
1499     print(f"                = {surface_area} cm\^2")
1500     print(f"The total surface area of this cube is {
1501     surface_area} square centimeters.")

```

## C.4 Prompts for Table 6 and ??

We use the following prompt for the entailment-based NLEP results in Table 6. The model-free result uses the NLEP prompt shown in C.1.

```

"""Write a Python function that constructs a
decision tree according to the given examples
that can generate the correct label of the
given classification task."""

### Available functions (shared for all tasks):

# Returns whether the hypothesis is entailed by the
premise.
def entailment(hypothesis, premise, model, tokenizer):
    proposition = f'{{hypothesis}} is entailed by {{
premise}}.'
    inputs = tokenizer(proposition, return_tensors
    ="pt", truncation=True, padding=True,
    max_length=128)
    outputs = model(**inputs)['logits'][0]
    ent_label = int(outputs[0] > outputs[2])
    if ent_label == 1:
        return 'yes'
    else:
        return 'no'

# Use the constructed decision tree to predict the
label of the sentence.
def tree_predict(sentence, criterions, tree, model,
tokenizer):
    node = tree['root']
    while node not in POSSIBLE_CLASSES:
        ent_label = entailment(criterions[node],
sentence, model, tokenizer)
        node = tree[node][ent_label]
    return node

### Task: Movie review classification
### Possible classes: [positive, negative]
### Examples:
"""
- contains no wit, only labored gags
  - [The movie is wise|The movie is not wise|], [
the story is fun|the story is not boring
|], [the review is positive|the review is
negative|]
- that loves its characters and communicates
something rather beautiful about human nature
  - [The characters are lovely|The characters are
awful|], [the script is touching|the
script is dry|], [the review is positive|
the review is negative|]
- on the worst revenge-of-the-nerds cliches the
filmmakers could dredge up
  - [The movie is novel|The movie is mostly
platitudes|], [the review is negative|]
- are more deeply thought through than in most right
-thinking films
  - [The takeaway of the movie is profound|The
idea of the movie is shallow|], [the
review is positive|the review is negative
|]
"""

### Define possible classes
POSSIBLE_CLASSES = ['positive', 'negative']

### Python program:
...
def get_decision_tree(sentence, model, tokenizer):
    # Step 1: define criterions of the decision tree
    criterions = [
        'This movie is interesting',
        'The movie has a good script',
        'The characters are awesome',
        'This movie is wise'
    ]

```

```

1586 # Step 2: define the Decision Tree for
1587 classification
1588 tree = {
1589     'root': 0,
1590     0: {'yes': 1, 'no': 3},
1591     1: {'yes': 'positive', 'no': 2},
1592     2: {'yes': 'positive', 'no': 'negative'},
1593     3: {'yes': 'positive', 'no': 'negative'}
1594 }
1595 return criterions, tree
1596 ...
1597

```

1598 **When we test the SST-2 performance based on**  
1599 **a generated Cola decision tree in Table ??, we use**  
1600 **the following prompt:**

```

1601 Write a Python function that constructs a decision
1602 tree according to the given examples that can
1603 generate the correct label of the given
1604 classification task.
1605
1606 ### Available APIs(shared for all tasks):
1607
1608 # Returns whether the hypothesis is entailed by the
1609 premise.
1610 def entailment(hypothesis, premise, model, tokenizer
1611 ):
1612     proposition = f'{hypothesis} is entailed by {
1613     premise}.'
1614     inputs = tokenizer(proposition, return_tensors
1615     ="pt", truncation=True, padding=True,
1616     max_length=128)
1617     outputs = model(**inputs)['logits'][0]
1618     ent_label = int(outputs[0] > outputs[2])
1619     if ent_label == 1:
1620         return 'yes'
1621     else:
1622         return 'no'
1623
1624 # Use the constructed decision tree to predict the
1625 label of the sentence.
1626 def tree_predict(sentence, criterions, tree, model,
1627 tokenizer):
1628     node = tree['root']
1629     while node not in POSSIBLE_CLASSES:
1630         ent_label = entailment(criterions[node],
1631 sentence, model, tokenizer)
1632         node = tree[node][ent_label]
1633     return node
1634
1635 ### Task: Grammar correctness classification
1636 ### Possible classes: ['acceptable', 'unacceptable']
1637
1638 ### Define possible classes
1639 POSSIBLE_CLASSES = ['acceptable', 'unacceptable']
1640
1641 ### Decision Tree Logic:
1642 - If verbs are not correctly constructed, the
1643 sentence is immediately labeled as unacceptable
1644 .
1645 - If verbs are correct:
1646 The tree then checks if the sentence has correct
1647 punctuation
1648 - If incorrect, label the sentence as
1649 unacceptable
1650 - If correct:
1651 The next criterion to be assessed is the
1652 subject-verb agreement.
1653 - If subject and verb disagree, label the
1654 sentence as unacceptable.
1655 - If they agree, check for sentence
1656 fragments.
1657 - If the sentence is a fragment, label
1658 it as unacceptable.
1659 - If it is not a sentence fragment,
1660 label the sentence as acceptable.
1661
1662 ### Python code for the decision tree:
1663
1664 ```python
1665 def get_decision_tree(sentence, model, tokenizer):
1666     # Step 1: define criterions of the decision tree
1667     criterions = {
1668         'correct_verbs': 'The verbs are correctly
1669         constructed in the sentence',

```

```

'correct_punctuation': 'The sentence is
1670 punctuated correctly',
1671 'subject_verb_agreement': 'The subject and
1672 verb agree in the sentence',
1673 'no_sentence_fragments': 'The sentence is
1674 not a fragment',
1675 }
1676
1677 # Step 2: define the balanced decision tree for
1678 this classification task
1679 tree = {
1680     'root': 'correct_verbs',
1681     'correct_verbs': {'yes': '
1682     correct_punctuation', 'no': '
1683     unacceptable'},
1684     'correct_punctuation': {'yes': '
1685     subject_verb_agreement', 'no': '
1686     unacceptable'},
1687     'subject_verb_agreement': {'yes': '
1688     no_sentence_fragments', 'no': '
1689     unacceptable'},
1690     'no_sentence_fragments': {'yes': 'acceptable'
1691     ', 'no': 'unacceptable'}
1692 }
1693
1694 return criterions, tree
1695 ...
1696

```

The input format of target tasks is 1697

```

1698 ### Task: Grammar correctness classification
1699 ### Possible classes: [acceptable, unacceptable]

```

## D Implementation Details for Task-Specific Strategies

We detail the few-shot chain-of-thought (CoT) and program-of-thought (PoT) prompting under the task-specific setting in Tables 2 and 4:

- **Tracking Shuffled Objects (7).** We use the same 3-shot examples as used by previous work for both task-specific CoT and PoT. The three examples are related to Tracking Shuffled Objects (3) and the models need to learn from demonstrations and generalize to seven objects test cases. The difference between CoT and PoT lies on the format of intermediate reasoning: CoT adopts natural language as the reasoning chains while we transform the thought process into concise Python code for PoT.
- **Dyck Language.** We cite the results of CoT from LATM (Cai et al., 2023) and transform the reasoning steps of the 3-shot examples used by previous chain-of-thought work into Python code for PoT. In order to evaluate the generalization ability of program-of-thought prompting, we try to avoid directly giving generally applicable Python program that can be used for all test instances.
- **Word Sorting.** We cite the results of CoT from LATM (Cai et al., 2023) and transform the reasoning steps of the 3-shot examples used by previous chain-of-thought work into Python code for PoT. Since the task can be effectively resolved by just few lines of code, i.e., read in the given input and use `sorted()` function, e.g.,

```
# Python code, return ans
words = ['oven', 'costume', 'counterpart']
ans = " ".join(sorted(words))
```

it can be regarded that the generally applicable tool is already given in the input prompt.

- **Chinese Remainder Theorem.** We cite the results of CoT from LATM (Cai et al., 2023). We build the in-context examples (3-shot) with the first three successful instances of task-general PoT as we construct the Python reasoning chains from the generated programs of task-general PoT with GPT-4. Indeed, for this complicated task, the provided program in the demonstration can also be regarded as a generally applicable tool. That’s a main reason why

task-specific PoT can obtain 100% accuracy on this benchmark.

- **Scheduling Meeting.** We cite the results of CoT from LATM (Cai et al., 2023). We build the in-context examples (3-shot) with the first three successful instances of task-general PoT as we construct the Python reasoning chains from the generated programs of task-general PoT with GPT-4. However, unlike giving the “ground-truth” problem solving tool for Chinese Remainder Theorem, the provided Python reasoning chains can only derive the correct answer for each specific demonstration question but can not be directly applied to all scenarios due to the complexity of the problem. We hope to compare this setup with Chinese Remainder Theorem and evaluate the performance of task-specific PoT on complicated tasks through different in-context learning demonstrations.
- **GSM-Hard.** We use the same 8-shot examples as used by previous work on GSM8K dataset for CoT GSM-Hard. For PoT, we adopt the 9-shot examples on GSM8K dataset from program-of-thought (Chen et al., 2022) containing concise Python code as reasoning chains.
- **StrategyQA.** We remove 1 example that appears in the development set from the 6-shot demonstration of previous work (Lyu et al., 2023) for CoT. As PoT is not designed and applied for natural language question answering task, we did not reproduce task-specific PoT results on StrategyQA benchmark.

## E Evaluation Prompts for VicunaQA

We have two metrics for VicunaQA. The first metric assesses the level of details and biases to long responses generated by GPT-4, while the other metric does not ask for details.

### E.1 Evaluation prompt asking for details.

```
prompt = f'''[Question]\n{ques_str}
[The Start of Assistant 1's Answer]\n{gpt4_res}
\n[The End of Assistant 1's Answer]
[The Start of Assistant 2's Answer]\n{target_res}
\n[The End of Assistant 2's Answer]
[System]
We would like to request your feedback on the
performance of two AI assistants in response to
the user question displayed above.\nPlease
rate the helpfulness, relevance, accuracy,
level of details of their responses. Each
assistant receives an overall score on a scale
of 1 to 10, where a higher score indicates
better overall performance.\nPlease first
output a single line containing only two values
indicating the scores for Assistant 1 and 2,
respectively. The two scores are separated by a
space. In the subsequent line, please provide
a comprehensive explanation of your evaluation,
avoiding any potential bias and ensuring that
the order in which the responses were presented
does not affect your judgment.'''
```

### E.2 Evaluation prompt not asking for details.

```
prompt = f'''[Question]\n{ques_str}
[The Start of Assistant 1's Answer]\n{gpt4_res}
\n[The End of Assistant 1's Answer]
[The Start of Assistant 2's Answer]\n{target_res}
\n[The End of Assistant 2's Answer]
[System]
We would like to request your feedback on the
performance of two AI assistants in response to
the user question displayed above.\nPlease
rate the relevance and accuracy of their
responses. Each assistant receives an overall
score on a scale of 1 to 10, where a higher
score indicates better overall performance.\n
Please first output a single line containing
only two values indicating the scores for
Assistant 1 and 2, respectively. The two scores
are separated by a space. In the subsequent
line, please provide a comprehensive
explanation of your evaluation, avoiding any
potential bias and ensuring that the order in
which the responses were presented does not
affect your judgment. Do not bias on either
longer or shorter answers.'''
```

### E.3 Calculation of Length Bias

Suppose we have  $N$  evaluation cases, each receiving 2 candidate responses. A GPT-4 scorer decides the winner between the candidates.  $a$  stands for the number of cases where a candidate response with more tokens wins. The length bias is calculated by

$$lb = \left| \frac{a}{N} - 0.5 \right| * 2 \quad (1)$$