# ChatLogic: Integrating Logic Programming with Large Language Models for Multi-Step Reasoning

Zhongsheng Wang
*School of Computer Science*
*University of Auckland*
Auckland, New Zealand
zwan516@aucklanduni.ac.nz

Jiamou Liu*
*School of Computer Science*
*University of Auckland*
Auckland, New Zealand
jiamou.liu@auckland.ac.nz

Qiming Bao
*School of Computer Science*
*University of Auckland*
Auckland, New Zealand
qbao775@aucklanduni.ac.nz

Hongfei Rong
*School of Computer Science*
*University of Auckland*
Auckland, New Zealand
hron635@aucklanduni.ac.nz

Jingfeng Zhang
*School of Computer Science*
*University of Auckland*
Auckland, New Zealand
jingfeng.zhang@auckland.ac.nz

*Abstract*—Large language models (LLMs) such as ChatGPT and GPT-4 have demonstrated impressive capabilities in various generative tasks. However, their performance is often hampered by limitations in accessing and leveraging long-term memory, leading to specific vulnerabilities and biases, especially during long interactions. This paper introduces ChatLogic, an innovative framework specifically targeted at LLM reasoning tasks that can enhance the performance of LLMs in multi-step deductive reasoning tasks by integrating logic programming. In ChatLogic, the language model plays a central role, acting as a controller and participating in every system operation stage. We propose a novel method of converting logic problems into symbolic integration with an inference engine. This approach leverages large language models' situational understanding and imitation skills and uses symbolic memory to enhance multi-step deductive reasoning capabilities. Our results show that the ChatLogic framework significantly improves the multi-step reasoning capabilities of LLMs. The source code and data are available at https://github.com/Strong-AI-Lab/ChatLogic.

*Index Terms*—deductive reasoning, multi-step reasoning, large language models

## I. INTRODUCTION

Recent advancements in large language models (LLMs) such as ChatGPT-3.5, GPT-4 [1], and Llama2 [2] have significantly enhanced their capabilities in various industries, proving invaluable in solving complex real-world problems. These models revolutionize sectors like customer service, healthcare, and education through their nuanced contextual comprehension and advanced conversational abilities. However, LLMs face notable challenges in multi-step logic reasoning tasks.

While these models excel in content generation, they consistently struggle to produce coherent responses in tasks requiring multi-step reasoning. Their training methodology, primarily based on the 'next-token prediction' approach, limits their ability to apply logical rules and deep contextual understanding essential for such tasks. For example, Figure 1 shows

how to let LLMs find a reasonable explanation path as the judgment result of the problem in the known randomly disrupted proposition sequence. This represents a critical area for improvement in current LLMs.

Further complicating this issue is the inherent token limitation of LLMs, which becomes apparent in continual dialogues [3]. The token caps in models like GPT-3.5 and GPT-4, while extendable through engineering prompts or technologies like Recursive Model Training [4], still pose a significant constraint. This limitation is particularly pronounced in multi-turn conversations, a common feature in multi-step logic reasoning tasks.
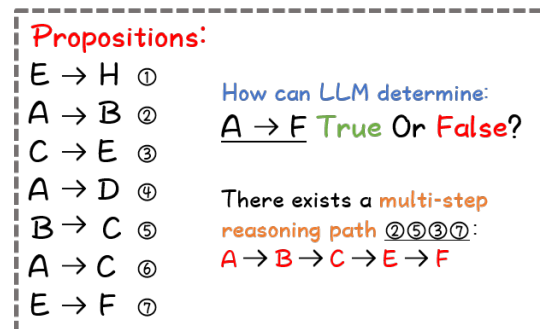


Fig. 1. Demo illustrating how LLMs can effectively identify and follow correct and logical reasoning paths to solve complex multi-step reasoning problems. In this instance, our objective is to let LLMs recognize the presence of an established path, ABCEF, thereby enabling them to accurately deduce that the statement 'A infers F' is true.

To address these limitations, innovative approaches such as external memory augmentation are being explored [5]. This method involves integrating LLMs with extensive databases to enhance their reasoning capabilities [6]. However, this integration brings its own challenges, such as the potential embedding of biases from the retrieval models into the LLMs [7], which could affect their accuracy and stability.

*Corresponding author

Our work introduces ChatLogic, a framework that augments LLMs with a logical reasoning engine to enhance their inferential capabilities. We have innovatively implemented a 'Mixshot Chain of Thought' technique in this framework. This approach significantly enhances the performance of LLMs by combining various prompt engineering methods. Mix-shot CoT efficiently guides the model through logical reasoning steps, achieving improved problem-solving with minimal resource consumption. ChatLogic is designed to be compatible with existing LLMs, significantly increasing their accuracy, especially in high-precision scenarios. The framework orchestrates the functioning of an LLM, enabling it to efficiently generate responses across various tasks.

At the heart of ChatLogic is the transformation of natural language into logical symbols, a process executed through pyDatalog. The primary objective of ChatLogic is to reinforce the stability of the reasoning process, ensuring that LLMs can handle intricate reasoning tasks with enhanced reliability and precision. The main characteristics of our framework are summarized below:

- ChatLogic, by combining LLMs with pyDatalog, translates natural language queries into logic programs, enhancing inference accuracy. This improvement is notably evident in multi-step reasoning tasks, as demonstrated on datasets such as PARARULE-Plus[1], CONCEPTRULES V1, and CONCEPTRULES V2.
- ChatLogic mitigates information loss, effectively addressing the long sequence limitation prevalent in adopting LLMs for multi-step reasoning tasks.
- ChatLogic incorporates automated enhancements for logic program execution, including a syntax correction module. This module refines a logic program by learning from previous executions, significantly improving the practical application and effectiveness of the generated code.

## II. RELATED WORK

### A. LLMs Reasoning:

LLMs are considered to have reasoning abilities similar to human cognition [8]. Despite facing challenges in multi-step logical tasks involving contemporary information or complex logical sequences [9], emerging approaches like self-consistency [10] show promise in enhancing performance, particularly in areas such as arithmetic and common sense reasoning. The effectiveness of causal reasoning pathways [11] is also crucial, ensuring that the output of LLMs is accurate but also transparent and verifiable. However, the most impactful method is the Chain of Thought (CoT) [12], which reveals the intermediate reasoning steps used by these models in problem-solving, allowing for continual self-correction of logical thinking, thereby greatly enhancing the rationality of their reasoning capabilities. But they all more or less exposed the weakness of extracting effective content from long and disordered information. We are trying to find ways to improve this shortcoming effectively.

### B. LLMs Code Generation:

LLMs have demonstrated the ability to generate code in various programming languages to meet users' specific needs [13]. However, how to directly apply the generated code to actual environments remains an issue to be resolved. In terms of optimization, the SELF-DEBUGGING approach [14] leads the post-code generation phase. It endows LLMs with the ability to debug their output, reinforcing the theme of continual refinement in the generated code. LOGIC-LM [15] creates a deterministic symbolic solver that expresses reasoning in a specific symbolic format to obtain actual results. The existing flaw we hope to improve is that while LLMs may not fully understand demo samples or generate "fantasies" that deviate from reality due to lack of pyDatalog knowledge in the pretraining data, they can still simply "imitate" Produce high-precision output. Our ultimate goal is to generate code that perfectly meets the requirements, can be deployed directly locally, and get inference results through the user's computer with a basic Python language environment installed.

### C. LLMs Prompt Engineering:

Prompt engineering in LLMs functions akin to psychological suggestions, guiding the model towards specific predictions [16]. Few-shot learning emphasizes training models with the least labelled data for optimized task performance. Notably, models like GPT-3 can handle tasks with a few examples, comparable to fine-tuned models [17]. Enhanced by prompt engineering, their reasoning capabilities are magnified. Zero-sample prompt [18] completely relies on the model's huge intrinsic knowledge and training corpus and assumes full responsibility for solving the problem. This means that we do not need to make any downstream task-related fine-tuning of LLMs; we only provide task content and expect good performance of LLMs. Surprisingly, although guidance is limited, in many cases, it often produces results that exceed expectations. In addition, zero-sample CoT [19] is also considered to be the best inference prompt currently; by using a simple prompt: 'Let's think step by step', the specific prompt and the corresponding two-stage key answer extraction prompt technology has significantly improved multiple inference-related zero-shot tasks. surpassed previous zero-shot learning.

Although not every situation using Zero-shot CoT results in optimal content output, in most cases targeting specific downstream tasks, LLMs have shown potential in the current field to effectively perform tasks using small sample techniques and combining external enhancement symbols [20]. In ChatLogic, we create independent prompt templates for different links in the framework and call them independently. Preliminary results suggest a promising direction by emphasizing LLMs' inherent reasoning capabilities and combining them with basic symbolic rules.

---

[1]https://huggingface.co/datasets/qbao775/PARARULE-Plus

## III. TASK DEFINITION

In our experience with advanced LLMs such as Llama2 and GPT-4, we've noticed their impressive ability to convert text into formal structures like math equations [21] and programming languages [22]. However, these models sometimes struggle with complex, multi-step reasoning tasks. The challenge escalates as the reasoning depth increases, and LLMs often miss key reasoning steps.

Acknowledging these characteristics, our primary goal is to boost the capability of LLMs to effectively represent problems in logic programming languages, particularly pyDatalog. This Python library integrates the logic programming paradigm and is particularly useful for declarative reasoning and complex querying. It allows for sophisticated rule-based logic and inference to be seamlessly incorporated into Python applications, enhancing their capabilities for decision-making processes.

We specify the multi-step deductive reasoning problem as Fact, Rule, and Query. The data sets used in the paper experiments all have this format:

- Fact: A Fact $F = \{f_1, f_2 \cdots, f_n\}$ is a sequence of sentences, each declarative sentence $f_i$ having a subject, a predicate verb, and an object, like 'Bob is poor' and 'Dogs like cats'. Predicates have negative expressions.
- Rule: A Rule $R = \{r_1, r_2 \cdots, r_n\}$ is a sequence of declarative sentences with conditional judgment, including initial features and inferential features, like 'If someone is poor then they are bad.'
- Query: A Query $Q$ is also a declarative sentence consistent with the sentence expression format.

To accomplish this goal, we aim to address the following two subtasks.

### A. Augmenting the Inferential Abilities of LLMs

Our preliminary objective is to exploit LLMs' one-shot generalization and zero-shot thinking capabilities $C$ (the collective name for both technologies). To achieve this, we aim to familiarize the model with the intricacies of symbolic language, specifically utilizing pyDatalog. We introduce them to this language through meticulously crafted examples that cover all potential edge cases. The structured pyDatalog syntax and these detailed prompts example set $S$ are provided to improve skill C and are essential to guide LLMs in understanding and handling multi-step inference.

$$pyDatalogCode = C(LLM(F, R, Q), S)$$

The code generated with high precision is executed by the local compiler to obtain accurate multi-step inference results.

$$Result = LocalExecution(pyDatalogCode)$$

### B. Amplifying the Executability of Automated Code Generation Processes

The translation of LLMs from text to code [23] is often imperfect on the first try and may contain errors. We aim to design a specific module $M$ in ChatLogic that ensures high-accuracy alignment of natural language and translation codes. The generated code should be easily executable locally, producing the desired results directly. Essentially, it further improves code generation quality so that it can be performed accurately and reasonably.

$$pyDatalogCode = M(C(LLM(F, R, Q), S))$$

## IV. CHATLOGIC

This section provides a detailed overview of the ChatLogic framework, particularly emphasising the finer details of its constituent parts and our innovative strategy regarding Mix-shot CoT.

### A. Framework Overview

The ChatLogic framework comprises four primary phases: input processing, Semantic Correction, Syntax Correction, and local execution response, as meticulously illustrated in Figure 2. The entire process, from problem input to result output, is depicted in the image as a demonstration. It is worth noting that the initial version of the logic code generated by LLMs at the beginning, after continuous revision through two or more iterations within two modules (semantic and syntactic correction), as a direct code solution for multi-step reasoning Executability is significantly improved, resulting in more precise results. Moreover, the enhanced refinement of this code substantially bolsters its executability, contributing to improved system performance and accuracy. The blue and green boxes are the intermediate self-correction processes that appear in the semantic correction and syntax correction modules and are generated correspondingly by the ChatLogic framework during the experiment.

Algorithm 1 delves deeper into the comprehensive algorithmic process for querying response data in ChatLogic. Apart from the locally executed part, all sub-tasks within ChatLogic are controlled and driven by LLMs acting as components. It consists of two loops, each corresponding to the content of two correction phases. We observe that LLMs excel at semantic corrections, and with limited modifications, correct text translations can be achieved. In lines 5 and 6 of the code, we employed zero-shot CoT to assist in determining the textual similarity of two propositions. Based on the judgment, we update the 'DifferentFlag' label, which influences the progression of the loop process of 'Semantic Correction'. However, syntax corrections are unreliable; they may get stuck in an infinite loop, repeatedly performing meaningless tasks. To address this issue, we consider introducing an upper loop limit. Although this somewhat diminishes ChatLogic's inferential capabilities, it significantly enhances the framework's robustness, making it better suited for multi-step deductive reasoning tasks.

### B. Mix-shot CoT

Our innovative mix-shot CoT (Chain of Thought) approach represents a groundbreaking hybrid methodology, blending the strengths of zero-shot CoT and one-shot learning to create a more versatile and effective learning paradigm for language
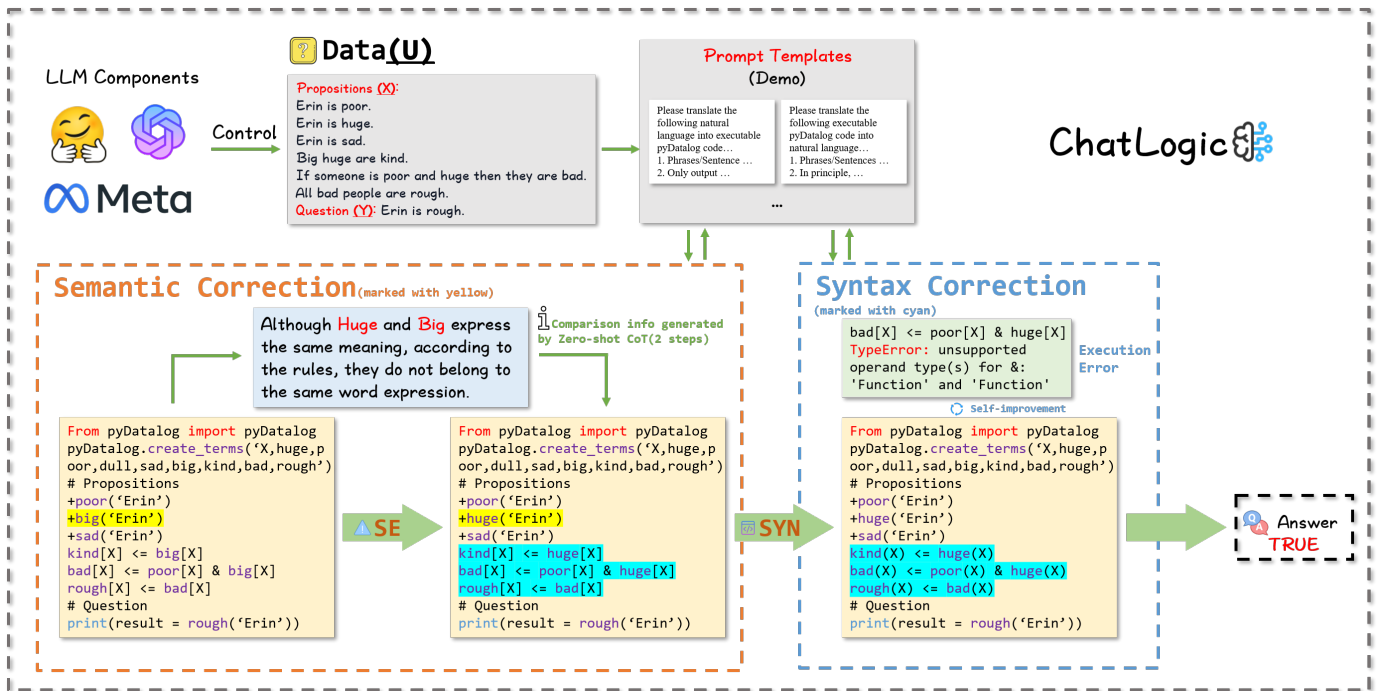
Fig. 2. ChatLogic containing more details uses LLMs as controllers, calls appropriate demonstration examples from Prompt Templates, guides the two modules of semantic correction(SE) and syntax correction(SYN) to output correct code, and produces execution results. This excerpts a specific question in PARARULE-Plus and the code generation process. The yellow portion represents the achievements of SE, and the cyan portion represents SYN.

---

**Algorithm 1** The Algorithm of ChatLogic

**Input**: U ← Rules supplemented based on the close-world assumption
X ← Proposition group (contains facts and rules)
Y ← Question
**Output**: TRUE/FALSE (Answer to Y given U, X)

1: DifferentFlag = TRUE
                        ▷ Semantic Correction
2: **while** DifferentFlag **do**
3:   Code ← PropositionTransformation($X, Y, U$)   ▷ Generate logic program based on close-world assumption
4:   RevProposition ← ReverseTransformation(Code, U)     ▷ Convert code back to natural language
5:   DifferentInfo ← TextComparison($(X, Y)$, RevProposition)
6:   DifferentFlag ← JudgeInfo(DifferentInfo)   ▷ Determining semantic similarity status with zero-shot CoT in 2 steps
7: **end while**
                         ▷ Syntax Correction
8: ExecutionError = NULL                  ▷ Execution result record
9: **while** Code cannot be executed **do**
10:   Code ← CodeImprovement(Code, ExecutionError)       ▷ Improve code based on error info
11:   **if** Running Time Overflow **then**
12:     Terminate WHILE Loop
13:   **end if**
14: **end while**
15: **return** CodeExecution(Code)        ▷ Get results by executing pyDatalog code locally

---

models. As mentioned before, zero-shot CoT uses LLMs' generative capabilities and chain thinking patterns to complete natural language tasks without targeted training or fine-tuning. One-shot learning provides a well-written task completion sample manually to guide LLMs to imitate the same workflow to achieve high-precision task completion. At its core, mix-shot CoT leverages the language model's innate ability for autonomous sub-task identification; a character follows established patterns accurately es it with the precision of one-shot learning through strategically chosen demonstration examples.

This dual approach allows for dynamic adaptation to various tasks' complexity and specific requirements. Mix-shot CoT guides the model using high-quality demonstration examples as templates, enhancing its precision and contextual depth for tasks demanding high accuracy and nuanced understanding, like in the method mentioned in this paper to convert natural language questions into logical codes. In addition, in scenarios that require more extensive analysis, such as text similarity comparison between generated propositions and original propositions in ChatLogic, LLMs are given greater autonomy to leverage their analytical capabilities and generate innovative solutions, thereby developing their ability to navigate massive amounts of information and generate unique insights. In addition, using the two-step calling process of zero-shot cot, during the second call to LLMs, we will extract valid key information from the first step as a status label to guide the normal operation of the entire framework.

Moreover, mix-shot CoT is designed to cultivate an adaptable learning process in language models, harmoniously combining structured guidance with the freedom of exploration. This flexibility is crucial in enabling the model to follow established patterns, innovate, and adapt to diverse tasks.

In our comparative analysis, as depicted in Table I, our mix-shot CoT methodology showcases a considerable leap in performance by striking an optimal balance between the grounded precision of one-shot learning and the generative flexibility of zero-shot CoT. Our approach reduces hallucination significantly and boasts the highest task-specific accuracy due to its judicious use of high-quality demonstration examples. While it does not entirely eliminate the need for demonstrations like zero-shot CoT, nor does it match the minimal hallucination levels of one-shot learning, the mix-shot CoT's enhanced adaptability and efficiency make it a powerful tool in the realm of prompt engineering. By acknowledging the limitations of requiring some demonstrations and not being as inherently scalable as zero-shot CoT, our mix-shot CoT nonetheless stands out for its pragmatic effectiveness in real-world applications where precision and adaptability are paramount.

## V. EVALUATION

In this section, we conduct experiments to evaluate the effectiveness of ChatLogic-augmented LLMs. Our experimental results show that **ChatLogic+LLMs** significantly outperform baseline LLMs, highlighting the advantages of using logical symbols to enhance the multi-step reasoning capabilities of LLMs.

### A. Datasets and Metrics

All reasoning questions in PARARULE-Plus adhere to the closed-world assumption, totalling approximately 400,000 samples. It features linguistic information in two contextual scenarios: People and Animal. For this dataset, we conducted our experiment by randomly selecting 50 instances from each depth level in both the Animal and People categories,

combining them to form a set of 100 test cases for each depth level, ranging from *Depth=2* to *Depth=5*.

In addition to PARARULE-Plus, we also incorporate the CONCEPTRULES V1[2] and CONCEPTRULES V2[3] datasets in our study. These datasets contain samples that require multi-step reasoning, with depths up to 3, making them suitable for evaluating models' capabilities in complex reasoning tasks. They are available in both simplified and full versions. For each version of the CONCEPTRULES datasets, we initially consolidated all data from the train, test, and dev sets into a single pool. We randomly sampled 100 instances from this unified dataset for our tests.

In our experiments with ChatGPT, GPT-4, and Llama 2-7B, we aimed to establish a baseline for the reasoning capabilities of LLMs as documented in the literature [24]. A significant part of our study was the implementation of ChatLogic, a framework designed to enhance the accuracy of these models' reasoning. This involved testing configurations like ChatGPT vs. ChatLogic (ChatGPT) across uniform scenarios using instances from the PARARULE-Plus and CONCEPTRULES datasets. The crux of our hypothesis is that if models augmented with ChatLogic demonstrate improved reasoning performance across various difficulty levels compared to their baseline, it would be a strong indicator of ChatLogic's effectiveness. Such results would suggest that ChatLogic could be a valuable addition to the field of artificial intelligence and natural language processing, affirming its utility in advancing the reasoning capabilities of LLMs.

### B. LLMs Configuration

In the ChatLogic framework invocation, ensuring the controllability of the text is paramount. For ChatGPT and GPT-4, the model invocation versions are respectively "**gpt-3.5-turbo**" and "**gpt-4**", and both have the hyperparameter temperature set to 0 so that we can have more control over what the LLM outputs. For Llama 2-7B, using an NVIDIA GeForce RTX 3090 with 24GB of memory, we utilized the transformer-based version **meta-llama/Llama-2-7b-chat-hf** [4] provided by Huggingface.

### C. Intermediate Process

Figure 3 provides an illustrative comparison between two approaches for handling the PARARULE-Plus dataset. On the left, we can observe the inference process generated directly by ChatGPT. While it appears initially logical, upon closer examination, it becomes apparent that this inference process contains numerous logical inconsistencies.

In contrast, on the right side of the figure, we showcase the enhanced performance of ChatGPT when coupled with the ChatLogic framework. This presentation includes only a selected excerpt of the final generated code. When executed locally, this integration produces remarkably precise inference

[2]https://bit.ly/3uVemXG
[3]https://bit.ly/3PApIIB
[4]https://huggingface.co/meta-llama/Llama-2-7b-chat-hf

| Attributes | Prompt Engineering Methods | | |
|---|---|---|---|
| | One-shot Learning [17] | Zero-shot CoT [19] | Mix-shot CoT (Ours) |
| Hallucination | Least hallucination | High hallucination | Less hallucination |
| Accuracy | Poor performance | Relatively high accuracy | Highest accurate response to task |
| Expansibility | Different demos needed for different tasks | No demo needed, only simple guidance required | Fewer demos, and improved performance than fewer shots |

**Propositions:**
**Facts:**
Anne is huge. Anne is strong. Fiona is thin.
Fiona is little. Alan is clever. Dave is rough. Dave is bad.
**Rules:**
If someone is not big then they are rough.
(Continue on the right)

If someone is not sad then they are kind.
If someone is clever then they are quiet.
If someone is quiet and not poor then they are smart.
If someone is rough and not big then they are dull.
If someone is dull then they are small.
**Question:**
**Anne is tiny?**

**ChatGPT:**
To determine whether "Anne is tiny" is correct, let's examine the given propositions step-by-step and follow the chain of reasoning.
From: Anne is huge.
Based on Proposition 8: "If someone is not big then they are rough. **"Since Anne is huge, she is big."** Hence, she is not rough.
We need to determine if Anne can be categorized as small and tiny.
(WRONG! " Huge " is not the same as " big ". We should not infer downward based on " Anne is big ", but based on the closed-world assumption: The proposition: "Anne is big" is not clearly defined, we should think that Anne is not big.)
Therefore, we don't have enough evidence to say that "Anne is tiny" based on the given propositions.
**Answer: 0 (The statement "Anne is tiny" is wrong based on the provided propositions.)**

**ChatLogic:**
```python
# Declare the pyDatalog variables
pyDatalog.create_terms('X,huge,strong,little,
rough,clever,sad,kind,quiet,tiny,big […]')
# Define the facts
+huge('Anne')
+strong('Anne')
+little('Fiona') […]
# Define the rules
rough(X) <= ~big(X)
kind(X) <= ~sad(X)
quiet(X) <= clever(X) […]
# Query the knowledge base
result = tiny('Anne')
if result:
    print(1)
else:
    print(0)
```

```
result: 1
Anne is tiny.
```

Fig. 3. Comparison based on the PARARULE-Plus dataset shows that while ChatGPT, even with CoT reasoning, often leads to incorrect inferences, the ChatLogic framework (also driven by ChatGPT) in most cases accurately generates pyDatalog code, highlighting its more reliable reasoning proficiency.

results, addressing the logical inconsistencies encountered in the unassisted ChatGPT inference process.

It is worth noting that after experiments, we found that before calling the text data in the PARARULE-Plus data set, LLMs can not complete the missing information for the relevant text content based on the closed-world assumption, resulting in a significant decrease in reasoning efficiency. For example, the beginning of a multi-step deductive reasoning chain process must be the existing attributes of an object, which we call "meta attributes". Meta attributes must be attributes the object possesses rather than attributes derived through inference. This is missing in the data set based on the closed-world assumption mentioned in the paper, so it needs to be supplemented manually. To avoid this situation, we built a script to extract the subject and attribute affiliations in the text through Named Entity Recognition (NER) and complete the missing information in the proposition in the artificially synthesized natural language format. Then, use the ChatLogic framework to enhance text reasoning on LLMs. For example, in the **'Rules'** in Figure 3, through observation, we find that no sentence can deduce that someone has the attribute of **'big'**,

so this attribute must appear in the definition of character attributes. If we do not see tasks with **'big'** attributes in **'Facts'**, then we default that they do not have **'big'** attributes, which requires us to manually add them to **'Facts'**. This part also belongs to the specific content represented by **U** in Algorithm 1 and Figure 2. The same approach was not utilized for the two distinct versions (simplified & full) of each of the CONCEPTRULES datasets (V1 & V2), as these versions did not require the supplementation of additional information.

## VI. RESULT

The experimental findings, as illustrated in Tables II and III, unequivocally showcase the ChatLogic framework's significant enhancement of LLMs' performance, surpassing the Baseline with considerably higher accuracy in most scenarios. While native models demonstrate competency in answering straightforward questions, they exhibit limitations in complex multi-step reasoning tasks, reducing accuracy in more challenging questions. In stark contrast, the amalgamation of ChatLogic with LLMs consistently manifests superior accuracy across various levels of question difficulty. This highlights the critical role of augmenting LLMs with logical symbolic

| Model | Method | Depth=2 | Depth=3 | Depth=4 | Depth=5 | Total |
|-------|--------|---------|---------|---------|---------|-------|
| GPT-3.5 | Base | 0.4 | 0.34 | 0.32 | 0.3 | 0.344 |
| | Zero-shot CoT | 0.42 | 0.42 | 0.41 | 0.3 | 0.3875 |
| | ChatLogic | **0.49** | **0.56** | **0.65** | **0.41** | **0.5275** |
| GPT-4 | Base | 0.65 | 0.75 | 0.42 | 0.4 | 0.555 |
| | Zero-shot CoT | **0.72** | 0.72 | 0.62 | **0.7** | 0.69 |
| | ChatLogic | **0.72** | **0.8** | **0.7** | **0.7** | **0.73** |
| Llama 2-7B | Base | 0.11 | 0.06 | 0.01 | 0.01 | 0.0475 |
| | Zero-shot CoT | 0.15 | **0.13** | 0.08 | 0.06 | 0.105 |
| | ChatLogic | **0.2** | **0.13** | **0.22** | **0.18** | **0.1825** |

*Note: This table evaluates accuracy on the PARARULE-Plus dataset, with a score of 1 indicating perfect accuracy. It includes comparisons with the 'Base' and 'Zero-shot CoT' methods. The 'ChatLogic' framework generally surpasses other models, highlighting its superior effectiveness with LLMs. For each depth level, the best performances are highlighted in bold.*

| Model | Method | CONCEPTRULES V1 | | CONCEPTRULES V2 | |
|-------|--------|-----------------|------|-----------------|------|
| | | simplified | full | simplified | full |
| GPT-3.5 | Base | 0.57 | 0.55 | 0.5 | 0.51 |
| | Zero-shot CoT | 0.63 | 0.51 | 0.7 | 0.67 |
| | ChatLogic | **0.69** | **0.67** | **0.79** | **0.74** |
| GPT-4 | Base | 0.95 | 0.94 | 0.89 | 0.86 |
| | Zero-shot CoT | **0.96** | **0.97** | **0.95** | **0.94** |
| | ChatLogic | **0.96** | 0.96 | 0.94 | **0.94** |
| Llama 2-7B | Base | 0.32 | 0.29 | 0.31 | 0.24 |
| | Zero-shot CoT | 0.42 | 0.41 | 0.33 | 0.3 |
| | ChatLogic | **0.48** | **0.49** | **0.37** | **0.36** |

*Note: This table shows accuracy comparisons on the CONCEPTRULES V1 and V2 datasets, where a score of 1 indicates perfect accuracy. The comparison includes both simplified and full versions of the datasets. Notably, GPT-4 with 'Zero-shot CoT' closely matches or occasionally surpasses our 'ChatLogic' framework in performance. The best performance results for each version of the datasets are highlighted in bold.*

operations in multi-step reasoning. By adopting this methodology, we ensure the retention of comprehensive information in natural language, effectively preventing omissions and the accumulation of errors that could compromise reasoning outcomes. Moreover, this approach enhances the transparency of the reasoning process, thereby elevating the credibility and traceability of the results.

In our analysis of the PARARULE-Plus dataset, the Chat-Logiour work's significant value and relevant model ('Base') and 'Zero-shot CoT' in most scenarios. Notably, GPT-4, in conjunction with ChatLogic, exhibits exceptional performance on questions of higher complexity (Depth=4 and Depth=5), underscoring ChatLogic's robust capability in handling intricate problems. Regarding Llama 2-7B, despite its weaker baseline performance, it shows significant improvement at all depth levels when assisted by ChatLogic. This indicates the framework's versatility in enhancing multi-step reasoning across different models.

Observing the CONCEPTRULES V1 and V2 datasets makes a notable shift in GPT-4's performance evident. With

Zero-shot CoT, GPT-4 either parallels or slightly surpasses ChatLogic in many cases, particularly in the full version of the CONCEPTRULES V2 dataset. The performance difference between GPT-4 and ChatLogic on these datasets is more nuanced than their performance on the PARARULE-Plus dataset. The inherent robustness of the GPT-4 model, likely due to its larger parameter count, already demonstrates formidable capabilities. This finding underscores the future need for more sophisticated datasets to challenge the upper-performance limits of advanced LLMs. Additionally, it's observed that ChatLogic primarily enhances models with a smaller parameter count by providing appropriate guidance, thereby boosting their performance on the task. This reaffirms our work's significant value and relevance, especially in optimizing models not inherently equipped with extensive computational resources.

| Dataset | Model | Base | SE | SE+SYN |
|---------|-------|------|-----|--------|
| CONCEPTRULES V1 | GPT-3.5 | 0.63 | 0.68 | 0.7 |
| | GPT-4 | 0.92 | 0.96 | 0.96 |
| | Llama 2-7B | 0.31 | 0.60 | 0.62 |
| CONCEPTRULES V2 | GPT-3.5 | 0.6 | 0.73 | 0.8 |
| | GPT-4 | 0.92 | 0.93 | 0.95 |
| | Llama 2-7B | 0.33 | 0.52 | 0.53 |
| PARARULE-Plus | GPT-3.5 | 0.26 | 0.5 | 0.62 |
| | GPT-4 | 0.54 | 0.64 | 0.7 |
| | Llama 2-7B | 0.1 | 0.16 | 0.16 |

*Note: Two modules respectively improve the executability of code, Semantic Correction (SE) and Semantic Correction + Syntax Correction (SE+SYN).*

## VII. ABLATION STUDY

We have introduced two modules to improve code execution from semantic and syntax perspectives. To demonstrate their role in aiding LLMs' multi-step reasoning, we will separately assess how each module affects the code's executability. We anticipate a gradual increase in successful executions, which would validate the effectiveness of our approach.

For the PARARULE-Plus dataset, we increased our data sampling to 100 random samples from the entire dataset, ensuring that the selection was shuffled entirely, with no specific order applied to parameters such as 'Depth (2-5)' and 'Pattern (Animal & People)'. For the CONCEPTRULES V1 and V2 datasets, we thoroughly shuffled all data from the 'simplified' and 'full' versions of each dataset and randomly selected 100 samples. To assess the executability rate of the generated code, we deployed it locally, evaluating solely based on the absence of error messages and the correctness of the output content.

The results are presented in Table IV. In our comparison, we examined the enhancements in code execution rate achieved by LLMs with the gradual integration of different modules. Relative to the baseline, both modules demonstrated incremental improvements in execution rates. The Syntax Correction module still proved valuable despite Llama 2 not utilizing as much code text for pre-training as GPT-3.5 and GPT-4. While it didn't significantly increase the execution rate, its contribution to refining code quality is noteworthy. Furthermore, it's important to highlight that GPT-4, due to its extremely advanced capabilities in some subtasks, has seemingly reached a 'performance ceiling' on the current datasets. This suggests that the dataset's limitations may somewhat constrain its potential.

## VIII. LIMITATION AND FUTURE WORK

Through experimental evaluations on multiple mainstream LLMs, we observed that ChatLogic+LLMs outperformed native LLMs in performance. The impressive performance demonstrates the effectiveness of our work. However, some issues have also been exposed. PARARULE-Plus is based on the closed-world assumption in question-answering data. Additionally, datasets like CONCEPTRULES V1 and V2, which are also artificially constructed, lack natural linguistic expression, which may not fully represent real-world complexities. When confronted with more complex sentences in the context of an open-world assumption, importing, integrating, and inferring external knowledge information that is expressed differently still poses challenges. Despite the valuable results from our experiments in enhancing code reliability, it's essential to acknowledge that the optimization module's applicability is currently limited to specific datasets. The carefully designed prompt sample cases are also optimized for specific data sets and are not a universal prompt template. Future developments should focus on creating adaptable optimization components [25] to address a wider array of scenarios and data sources.

## REFERENCES

[1] OpenAI, "Gpt-4 technical report," *ArXiv*, vol. abs/2303.08774, 2023. [Online]. Available: https://api.semanticscholar.org/CorpusID: 257532815

[2] H. Touvron, L. Martin, K. Stone, P. Albert, A. Almahairi, Y. Babaei, N. Bashlykov, S. Batra, P. Bhargava, S. Bhosale *et al.*, "Llama 2: Open foundation and fine-tuned chat models," *arXiv preprint arXiv:2307.09288*, 2023.

[3] A. J. Thirunavukarasu, D. S. J. Ting, K. Elangovan, L. Gutierrez, T. F. Tan, and D. S. W. Ting, "Large language models in medicine," *Nature medicine*, pp. 1–11, 2023.

[4] A. Bulatov, Y. Kuratov, and M. S. Burtsev, "Scaling transformer to 1m tokens and beyond with rmt," *arXiv preprint arXiv:2304.11062*, 2023.

[5] Z. Zhong, T. Lei, and D. Chen, "Training language models with memory augmentation," *arXiv preprint arXiv:2205.12674*, 2022.

[6] S. Borgeaud, A. Mensch, J. Hoffmann, T. Cai, E. Rutherford, K. Millican, G. B. Van Den Driessche, J.-B. Lespiau, B. Damoc, A. Clark *et al.*, "Improving language models by retrieving from trillions of tokens," in *International conference on machine learning*. PMLR, 2022, pp. 2206–2240.

[7] O. Khattab, K. Santhanam, X. L. Li, D. Hall, P. Liang, C. Potts, and M. Zaharia, "Demonstrate-search-predict: Composing retrieval and language models for knowledge-intensive nlp," *arXiv preprint arXiv:2212.14024*, 2022.

[8] J. Huang and K. C.-C. Chang, "Towards reasoning in large language models: A survey," *arXiv preprint arXiv:2212.10403*, 2022.

[9] A. Creswell, M. Shanahan, and I. Higgins, "Selection-inference: Exploiting large language models for interpretable logical reasoning," *arXiv preprint arXiv:2205.09712*, 2022.

[10] X. Wang, J. Wei, D. Schuurmans, Q. Le, E. Chi, S. Narang, A. Chowdhery, and D. Zhou, "Self-consistency improves chain of thought reasoning in language models," *arXiv preprint arXiv:2203.11171*, 2022.

[11] A. Creswell and M. Shanahan, "Faithful reasoning using large language models," *arXiv preprint arXiv:2208.14271*, 2022.

[12] J. Wei, X. Wang, D. Schuurmans, M. Bosma, F. Xia, E. Chi, Q. V. Le, D. Zhou *et al.*, "Chain-of-thought prompting elicits reasoning in large language models," *Advances in Neural Information Processing Systems*, vol. 35, pp. 24 824–24 837, 2022.

[13] Z. Yang, A. Ishay, and J. Lee, "Coupling large language models with logic programming for robust and general reasoning from text," *arXiv preprint arXiv:2307.07696*, 2023.

[14] X. Chen, M. Lin, N. Schärli, and D. Zhou, "Teaching large language models to self-debug," *arXiv preprint arXiv:2304.05128*, 2023.

[15] L. Pan, A. Albalak, X. Wang, and W. Y. Wang, "Logic-lm: Empowering large language models with symbolic solvers for faithful logical reasoning," *arXiv preprint arXiv:2305.12295*, 2023.

[16] Y. Wang, Q. Yao, J. T. Kwok, and L. M. Ni, "Generalizing from a few examples: A survey on few-shot learning," *ACM computing surveys (csur)*, vol. 53, no. 3, pp. 1–34, 2020.

[17] T. Brown, B. Mann, N. Ryder, M. Subbiah, J. D. Kaplan, P. Dhariwal, A. Neelakantan, P. Shyam, G. Sastry, A. Askell *et al.*, "Language models are few-shot learners," *Advances in neural information processing systems*, vol. 33, pp. 1877–1901, 2020.

[18] L. Reynolds and K. McDonell, "Prompt programming for large language models: Beyond the few-shot paradigm," in *Extended Abstracts of the 2021 CHI Conference on Human Factors in Computing Systems*, 2021, pp. 1–7.

[19] T. Kojima, S. S. Gu, M. Reid, Y. Matsuo, and Y. Iwasawa, "Large language models are zero-shot reasoners," *Advances in neural information processing systems*, vol. 35, pp. 22 199–22 213, 2022.

[20] C. H. Song, J. Wu, C. Washington, B. M. Sadler, W.-L. Chao, and Y. Su, "Llm-planner: Few-shot grounded planning for embodied agents with large language models," *arXiv preprint arXiv:2212.04088*, 2022.

[21] J. He-Yueya, G. Poesia, R. E. Wang, and N. D. Goodman, "Solving math word problems by combining language models with symbolic solvers," *arXiv preprint arXiv:2304.09102*, 2023.

[22] P. Vaithilingam, T. Zhang, and E. L. Glassman, "Expectation vs. experience: Evaluating the usability of code generation tools powered by large language models," in *Extended Abstracts of the 2022 CHI Conference on Human Factors in Computing Systems*, ser. CHI EA '22. New York, NY, USA: Association for Computing Machinery, 2022. [Online]. Available: https://doi.org/10.1145/3491101.3519665

[23] F. J. Budinsky, M. A. Finnie, J. M. Vlissides, and P. S. Yu, "Automatic code generation from design patterns," *IBM Systems Journal*, vol. 35, no. 2, pp. 151–171, 1996.

[24] C. Hu, J. Fu, C. Du, S. Luo, J. Zhao, and H. Zhao, "Chatdb: Augmenting llms with databases as their symbolic memory," *arXiv preprint arXiv:2306.03901*, 2023.

[25] R. Marvie, "Picolo: A simple python framework for introducing component principles," in *Euro Python Conference*. Citeseer, 2005.