# CoT-RAG: Integrating Chain of Thought and Retrieval-Augmented Generation to Enhance Reasoning in Large Language Models

Feiyang Li<sup>1</sup>, Peng Fang<sup>1\*</sup>, Zhan Shi<sup>1</sup>, Arijit Khan<sup>2,3</sup>, Fang Wang<sup>1</sup>, Weihao Wang<sup>4</sup>, Xin Zhang<sup>4</sup>, Yongjian Cui<sup>4</sup>, 

<sup>1</sup>Huazhong University of Science and Technology, Wuhan, China 

<sup>2</sup>Bowling Green State University, Ohio, USA 

<sup>3</sup>Aalborg University, Aalborg, Denmark 

<sup>4</sup>Huawei Technologies Co., Ltd, Shenzhen, China

### **Abstract**

Chain-of-thought (CoT) reasoning boosts large language models' (LLMs) performance on complex tasks but faces two key limitations: a lack of reliability when solely relying on LLM-generated reasoning chains and lower reasoning performance from natural language prompts compared with code prompts. To address these issues, we propose CoT-RAG, a novel reasoning framework with three key designs: (i) Knowledge Graph-driven CoT Generation, featuring knowledge graphs to modulate reasoning chain generation of LLMs, thereby enhancing reasoning credibility; (ii) Learnable Knowledge Case-aware RAG, which incorporates retrieval-augmented generation (RAG) into knowledge graphs to retrieve relevant sub-cases and sub-descriptions, providing LLMs with learnable information; (iii) Pseudo-Program Prompting Execution, which promotes greater logical rigor by guiding LLMs to execute reasoning tasks as pseudo-programs. Evaluations on nine public datasets spanning three reasoning tasks reveal significant accuracy gains—ranging from 4.0% to 44.3%—over state-of-the-art methods. Furthermore, tests on four domain-specific datasets demonstrate exceptional accuracy and efficient execution, underscoring its practical applicability and scalability. Our code and data are available at https: //github.com/hustlfy123/CoT-RAG.

### 1 Introduction

Large language models (LLMs) have garnered significant attention in both academia and industry due to their exceptional performance in natural language processing (NLP) (Kabir et al., 2024; Chiang and Lee, 2023), such as machine translation (Zhu et al., 2024; Wang, 2024), text summarization (Li et al., 2024) and sentiment analysis (Bai et al., 2024). Nonetheless, they exhibit notable limitations in complex tasks requiring arithmetic,

commonsense, and symbolic reasoning (Wei et al., 2022; Gao et al., 2023). To overcome these challenges, chain-of-thought (CoT) reasoning has been introduced, wherein LLMs explicitly generate intermediate reasoning steps prior to reaching a conclusion (Chu et al., 2024). Several CoT variants have emerged, including Manual-CoT (Wei et al., 2022), Zero-shot-CoT (Kojima et al., 2022), PoT (Chen et al., 2023), and PS (Wang et al., 2023a).

Although CoT improves performance on multistep reasoning tasks for LLMs, existing methods still face two main challenges:

(1) The low reliability of relying solely on LLMs to generate reasoning chains. Most existing CoT reasoning methods rely on prompting strategies, such as prompt format (Gao et al., 2023; Chen et al., 2023) and prompt planning (Yao et al., 2023b; Wang et al., 2023a), to guide LLMs in generating reasoning chains. However, due to the inherent black-box nature (Kim et al., 2024) and hallucination issues (Manakul et al., 2023) of LLMs, the generated reasoning steps often contain logical errors or factual inaccuracies (Dubey et al., 2024). For instance, Manual-CoT, Zero-shot-CoT, and PS achieve average accuracies of only 48.4%, 38.9%, and 42.5%, respectively, on the AQuA arithmetic reasoning dataset (Ling et al., 2017), highlighting their generally low performance (Wang et al., 2023a). In vertical domains, e.g., law, medicine, and finance-where errors may compromise human life and critical assets—the inherent unreliability of LLM outputs introduces unquantifiable risks. Consequently, solely depending on LLM-generated reasoning cannot ensure reliable and safe outcomes.

(2) The poorer reasoning performance from natural language prompts compared with code prompts. LLMs typically generate reasoning chains in natural language (NL) to describe intermediate steps. While these NL-based reasoning chains are intuitive for humans, they exhibit lower reasoning accuracies compared with code prompts

<sup>\*</sup>Corresponding author

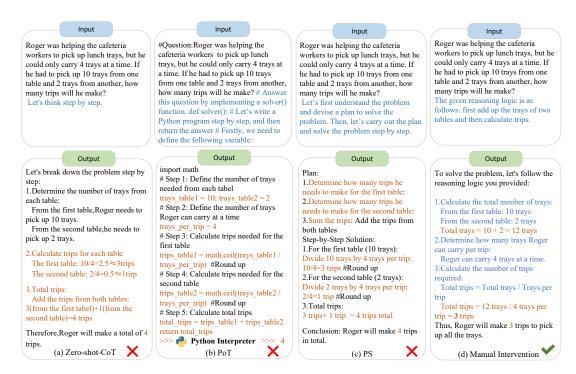


Figure 1: Example inputs and outputs of GPT-40 mini with (a) Zero-shot-CoT (Kojima et al., 2022), (b) PoT (Chen et al., 2023), (c) PS (Wang et al., 2023a) and (d) Manual Intervention on MultiArith (Roy and Roth, 2015).

(Gao et al., 2023; Sun et al., 2023a). Lyu et al. (2023) demonstrate that on the GSM8K arithmetic reasoning dataset (Cobbe et al., 2021), Manual-CoT and LtM (Zhou et al., 2023), which rely on NL-based reasoning chains, achieve average accuracies of 46.6% and 42.3%, respectively, exhibiting inadequate correctness; while Faithful CoT (Lyu et al., 2023) that employs a two-step process where LLMs first generate a symbolic reasoning chain (e.g., in Python) and then execute this chain using an external interpreter (e.g., a Python interpreter), achieves an average accuracy of 64.2%. This underscores that LLMs which use NL prompts achieve inferior reasoning performance compared with code prompts.

To address these challenges, we propose CoT-RAG, an enhanced reasoning framework with three key components. **Knowledge Graph-driven CoT Generation:** To ensure human safety and reduce potential risks in vertical domains, our approach places domain experts at the helm while LLMs augment their capabilities. In particular, experts supply a one-time, coarse-grained decision tree (DT) that encapsulates the underlying reasoning logic for the domain, independent of individual user queries. The LLMs then convert this DT into a detailed knowledge graph (KG) to enhance comprehension. For each user query, the KG is employed to generate reasoning chains by LLMs, thereby

improving the process's controllability, reliability, and adaptability to specific domains. Learnable Knowledge Case-aware RAG: CoT-RAG integrates Retrieval-Augmented Generation into the knowledge graph, retrieving relevant sub-cases and sub-descriptions that supply learnable data to mitigate LLMs' inherent logical errors and factual biases. Moreover, the interactive framework between LLMs and the KG enables dynamic updates to both the graph structure and the case repository. **Pseudo-**Program Prompting Execution: CoT-RAG employs pseudo-program prompting as an alternative to NL prompts, which directs LLMs to execute reasoning tasks via pseudo-programmatic chains, retaining the versatility of NL prompts while improving logical coherence in complex tasks.

We evaluated CoT-RAG on nine public datasets spanning three reasoning tasks. Results indicate that it outperforms existing methods, with accuracy improvements from 4.0% to 44.3%. Additionally, tests on four domain-specific datasets confirm its high accuracy and efficient execution, underscoring its scalable cross-domain performance.

# 2 Motivation

## 2.1 Reasoning Chain Generation

Generally, CoT reasoning relies on prompt strategies to guide LLMs in generating reasoning chains.

However, excessive reliance on LLMs may lead to erroneous reasoning steps, which pose incalculable risks and potential losses in fields related to human life and property due to the tendency of LLMs to make reasoning errors in vertical domains (Huang et al., 2024b; Barile et al., 2024).

To examine this limitation, we compare methods that exclusively depend on LLM-generated reasoning chains, including Zero-shot-CoT (Kojima et al., 2022), PoT (Chen et al., 2023), and PS (Wang et al., 2023a), against manual intervention method (i.e., providing explicit reasoning logic to LLMs). As shown in Figure 1, we input a mathematical problem from the MultiArith dataset (Roy and Roth, 2015) into GPT-40 mini (OpenAI, 2024b), where Zero-shot-CoT, PoT, and PS produce incorrect answers. In contrast, manual intervention in the reasoning logic (Figure 1(d)) allows the LLM to generate correct results. Similarly, as reported by Wang et al. (2023a), Manual-CoT, Zero-shot-CoT, and PS achieve average accuracies of 74.8%, 64.5%, and 68.7% on commonsense reasoning tasks, and 48.4%, 38.9%, and 42.5% on the AQuA arithmetic reasoning dataset (Ling et al., 2017), indicating the generally lower performance of such methods. These findings further support our observations: Solely relying on LLM-generated reasoning chains is insufficient to ensure reliable results, highlighting the necessity of manual intervention in the reasoning process. While some works (Gou et al., 2024; Weng et al., 2023; Paul et al., 2024; Madaan et al., 2023) have integrated verification and refinement into LLMs to reduce errors, they still depend solely on LLMs for evaluation and verification, which results in low accuracy (Chu et al., 2024).

## 2.2 Reasoning Chain Execution

LLMs commonly produce reasoning chains in natural language (NL) to outline intermediate steps. Although these NL-based chains are human-friendly, they exhibit lower reasoning accuracies (Sun et al., 2023a). As stated in Lyu et al. (2023), the average accuracies of Manual-CoT and LtM (Zhou et al., 2023), which execute NL reasoning chains on GSM8K (Cobbe et al., 2021), are 46.6% and 42.3%, respectively, exhibiting their low performance; while Faithful CoT (Lyu et al., 2023) requires the LLM to first generate a symbolic reasoning chain (e.g., in Python) and then execute it using an external interpreter, achieving an average accuracy of 64.2%. Therefore, methods like ProgPrompt (Singh et al., 2022), Code as Policies

(Liang et al., 2023), and AdaPlanner (Sun et al., 2023a) use code prompts instead of NL prompts to reduce ambiguities and improve inference performance. However, Code prompts exhibit three primary limitations. **Complexity:** They use intricate programming symbols and function calls that are often unintelligible to non-programmers. **Scope:** They struggle with general or domain-specific reasoning outside of mathematical contexts. **Language Restriction:** They are confined exclusively to the Python code style. Therefore, it is essential to develop a prompting methodology that combines the broad applicability of natural language prompts with the logical precision of code prompts, all while maintaining clarity.

# 3 Methodology

Overview. To address the aforementioned concerns, we design CoT-RAG, a novel reasoning framework with three stages (Figure 2). First, experts construct and input a coarse-grained decision tree that represents the reasoning logic of problems in a specific domain, where experts only need to build it once, then the LLM transforms it into a knowledge graph for a deeper understanding of its internal logic in the Knowledge Graph-driven CoT Generation phase (§3.1). Next, during the Learnable Knowledge Case-aware RAG phase (§3.2), users input query descriptions related to this domain and the LLM extracts sub-descriptions to update the knowledge graph for accurately generating results in the next stage. Third, the LLM uses Pseudo-Program Prompting Execution (§3.3) to process the updated knowledge graph and produces the ultimate result. Specific algorithm demonstration and time complexity analysis are given in Appendix A, and related notations are shown in Appendix B, which displays significant differences between our decision trees and knowledge graphs w.r.t. their conventional notions.

## 3.1 Knowledge Graph-driven CoT Generation

A decision tree (DT) is a tree-structured algorithm for classification and regression, where the internal nodes "test" a condition, branches represent the test outcome, and the leaf nodes denote the final results (Magerman, 1995). DTs offer robust logical coherence and interpretability, enabling human intervention (Mienye and Jere, 2024; Kalra and Brown, 2024). They are widely used in critical areas such as education (Huang, 2021), finance (Wang, 2021;

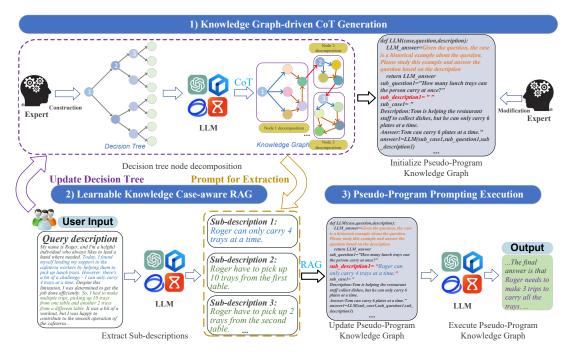


Figure 2: An overview of our CoT-RAG framework.

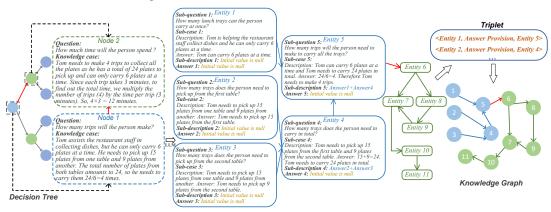


Figure 3: An example of decomposing a decision tree into a knowledge graph.

Chou and Chen, 2024), and medicine (He et al., 2024; Govil et al., 2024), which motivates us to utilize a manually crafted DT that represents the reasoning logic of problems in a specific domain, to modulate the inference processes of LLMs. Distinct from traditional DTs, each node in our DT contains a *Question* and a *Knowledge case*, which represents some descriptions from user input processed by the node (e.g., Figure 3), information supplemented artificially (e.g., Appendix F Table 22), or relevant considerations (e.g., Appendix F Table 28) to assist LLMs in responding the *Question* of the node. Each branch implies that the output of the parent node is fed to the child node.

However, designing a "fine-grained" DT manually consumes a lot of financial and human resources. To alleviate this situation, LLMs can efficiently decompose the "coarse-grained" DT nodes provided by experts, leveraging their proficiency in decomposition (Huang et al., 2024a). Furthermore, Knowledge case can assist LLMs in learning how to decompose each node of the DT, mitigating incorrect results caused by the tendency of LLMs to make errors in vertical fields (Huang et al., 2024b; Barile et al., 2024). Nevertheless, the complex relationships among the new nodes after the decomposition of the DT need to be captured. Considering that knowledge graphs (KGs), featuring clear delineation and interpretability, precisely represent intricate inter-individual relationships and facilitate inference (Sun et al., 2024b), we introduce them to represent decomposed DTs. As depicted in Stage 1 of Figure 2, LLMs transform each node in the coarse-grained DT built by experts into several entities and generate a highly transparent KG, optimizing the decision-making process.

We provide an example in Figure 3 to show the DT node decomposition process. The node 1: "How many trips will the person make?" is decomposed by the LLM into five entities, each corresponding to four attributes: (i) Sub-question, a simplified component of the original complex Question; (ii) Sub-case, which is a concise case corresponding to Sub-question, derived from the LLMdecomposed Knowledge case; (iii) Sub-description, a text description corresponding to Sub-question and Sub-case, its initial value is null and will be assigned by extracting from users' input query descriptions in § 3.2 or obtained from Answer of other entities; and (iv) Answer, which refers to the LLM's output for matching Sub-question, its initial value is null and will be assigned in § 3.3. Edges between entities represent Answer Provision relationships, for example, the triplet *<Entity 1, Answer Provi*sion, Entity 5> indicates that the Answer of Entity 1 is provided to the Sub-description of Entity 5. The Answer of Entity 5 represents the node 1's inference result. Moreover, as node 2 in the DT is the child of node 1, Entity 5 points to Entity 6, which is the entity after *node* 2 is decomposed by the LLM.

### 3.2 Learnable Knowledge Case-aware RAG

§ 3.1 pertains to the introductory phase of the entire CoT-RAG framework. Subsequently, it comes to the application phase targeting users, that is, how to extract brief yet significant information from relatively long query descriptions input by users. Unlike traditional vector-based retrieval in RAG, we utilize LLM-based retrieval (Shen et al., 2024) for long query descriptions, combining *Subquestion* and *Sub-case* of each entity as a prompt generated in the previous stage, to extract relevant descriptions that are assigned to the corresponding *Sub-description*, which has higher accuracy and shorter runtime (refer to § 5 and Appendix C.4).

We present an example to illustrate specific Sub-description extraction process. On the premise that the KG generated from Figure 3 is taken as the output of Stage 1 of CoT-RAG, the user inputs a query description about Roger, as depicted in Stage 2 of Figure 2, and the LLM uses this KG (including Sub-question and Sub-case of Entities 1, 2, and 3) as a prompt to extract the corresponding descriptions as Sub-description. Specifically, Sub-description 1 is "Roger can only carry 4 trays at a time", Sub-description 2 is "Roger has to pick up 10 trays from the first table", and Sub-description 3 is "Roger has to pick up 2 trays from the second

table". These Sub-description and Sub-case will be employed by the LLM to precisely answer the corresponding Sub-question of each entity in § 3.3.

Furthermore, if the query descriptions entered by users are new to the DT, they have the capacity to dynamically update the *Knowledge case* within the DT, which facilitates the LLM in generating a more comprehensive and application-oriented knowledge graph, thereby augmenting the flexibility of the knowledge graph.

## 3.3 Pseudo-Program Prompting Execution

The first two designs enhance the credibility and interpretability of the reasoning chains generated by LLMs. However, there remains a critical issue: how can reasoning chains be represented to ensure the rigorous logical execution of LLMs? Inspired by QDMRPS (Huang et al., 2024a) and AdaPlanner (Sun et al., 2023a), we propose Pseudoprogram Prompting (PsePrompting), which allows the LLM to represent the knowledge graph reasoning chain as pseudo-programs, referred to as the pseudo-program knowledge graph (PKG). Its initial version created by the LLM of Stage 1 is shown in Table 11 of Appendix E, where Entities 1, 2, and 3's Sub-description are filled based on user input in Stage 2 (Table 12 in Appendix E). The final result of the PKG executed by the LLM in Stage 3 is depicted in Table 13 of Appendix E. The LLM processes each entity of the PKG sequentially, it first retrieves the Sub-case and Sub-description for learning, then generates the Answer. We can observe that PsePrompting features a simple logical structure and is easily understandable. Additionally, it demonstrates broad applicability for addressing queries that demand domain-specific syntax or reasoning paradigms (Appendix F), eliminates reliance on external interpreters, and supports extension to programming languages such as C++ and Java (Appendix G).

# 4 Experimental setup

## 4.1 Datasets

General domains: Our CoT-RAG is evaluated on nine benchmark datasets from three categories of reasoning problems: Arithmetic Reasoning: AQUA (Ling et al., 2017), GSM8K (Cobbe et al., 2021), MultiArith (Roy and Roth, 2015), and SingleEq (Koncel-Kedziorski et al., 2015); Commonsense Reasoning: HotpotQA (Wolfson et al., 2020), CSQA (Talmor et al., 2019), and SIQA (Sap

et al., 2019); **Symbolic Reasoning:** Last Letter Concatenation (Wei et al., 2022), and Coin Flip (Wei et al., 2022).

**Vertical domains:** To further demonstrate the scalability of our method across different vertical domains, we evaluate it on four open-source datasets from the legal, financial, and logic fields: LawBench (LaB) (Fei et al., 2024), LegalBench (LeB), CFBenchmark (CFB) (Lei et al., 2023), and AGIEval (AGI) (Zhong et al., 2024).

Additionally, following GraphRAG (Edge et al., 2024) and Graph-CoT (Jin et al., 2024), we employ an LLM to adapt the datasets to satisfy our testing needs, where the specific details and descriptions of datasets are shown in Appendix H.

### 4.2 Baselines

**General domains:** We compare our CoT-RAG with following methods that focus on CoT: (1) Manual-CoT (Wei et al., 2022) uses a "thought chain" prompt to guide LLMs through a step-bystep solution process, leading to a detailed answer. (2) Zero-shot-CoT (Kojima et al., 2022) encourages LLMs to generate reasoning steps automatically by appending "Let's think step by step" to the question. (3) **Complex-CoT** (Fu et al., 2023) represents a simple and effective example selection scheme for multi-step reasoning. (4) Auto-CoT (Zhang et al., 2023) automates the generation of high-quality prompts, improving efficiency and accuracy in reasoning tasks. (5) PS (Wang et al., 2023a) first generates a coarse task plan, followed by a fine-grained solution process. (6) KD-CoT (Wang et al., 2023) modifies reasoning traces in CoT via a retriever to interact with external knowledge stored in an unstructured knowledge base. (7) IRCoT (Trivedi et al., 2023) interleaves retrieval with steps in a CoT, in turn using retrieved results to improve CoT. Notice that both KD-CoT and IRCoT retrieve external knowledge from an unstructured knowledge base, which is different from our highly-structured and logical DTs and KGs. (8) QDMRPS (Huang et al., 2024a) decomposes problems into QDMR-based directed acyclic graphs and reasons step-by-step based on dependencies. (9) **Iter-CoT** (Sun et al., 2024a) prompts LLMs to self-correct their errors in reasoning chains by leveraging iterative bootstrapping. (10) **KG-CoT** augments LLMs (Zhao et al., 2024) via a graph reasoning model that generates explicit reasoning paths over factual KGs, e.g., Freebase, which is different from our specialized DTs and KGs, as

well as our KG-focused logical reasoning. (11) **ZEUS** (Kumar et al., 2025) improves CoT prompting by utilizing uncertainty estimates to select effective demonstrations without needing access to model parameters. (12) **Pattern-CoT** (Zhang et al., 2025) employs reasoning patterns to enhance CoT prompting effectiveness. (13) In addition, we evaluate the intrinsic capabilities of the LLM, where we only deal with the original problem input to the LLM without using any additional methods, namely **Zero-shot**. Methods like PoT (Chen et al., 2023) and Faithful CoT (Lyu et al., 2023), which cannot handle complex non-mathematics reasoning tasks, especially in the field of commonsense reasoning, are excluded from this comparison.

**Vertical domains:** We utilize the Faiss (Pinecone, 2024) vector database to replace the LLM-based retrieval in CoT-RAG with vector-based retrieval, where six variants are set up according to different indexes (Meta-Research, 2024), namely CoT-RAG (IndexFlatL2), CoT-RAG (IndexFlatIP), CoT-RAG (IndexIVFFlat), CoT-RAG (IndexLSH), CoT-RAG (IndexPQ), and CoT-RAG (IndexIVFPQ). We also compare the scalability of CoT-RAG under the condition of zero-expert, that is, CoT-RAG (Zeroexpert), where LLMs replace experts to generate decision trees. Additionally, besides KG-CoT (§ 4.2), we also incorporate eight state-of-the-art graphform LLM-based RAG methods: RoG (LUO et al., 2024) first generates KG-grounded relation paths as plans, then uses them to retrieve reasoning paths from KGs for LLMs to reason reliably; Graph-CoT (Jin et al., 2024) enhances LLMs by encouraging them to conduct iterative reasoning on the graph; ToG (Sun et al., 2024b) employs LLMs for iterative beam search on KGs, finding the best paths and returning top results.; ToG-2 (Ma et al., 2025) utilizes KGs to connect documents via entities, which facilitates deep and knowledge-guided context retrieval, thereby enhancing the reasoning ability of LLMs; RRKG (Ji et al., 2024) combines explainable knowledge graphs with LLMs to enhance complex reasoning capabilities; **AtomR** (Xin et al., 2024) is a framework that enables LLMs to conduct accurate heterogeneous knowledge reasoning at the atomic level; GraphRAG (Edge et al., 2024) combines the advantages of RAG and queryfocused summarization, and it uses knowledge graphs to store source document data for efficient global question answering; PoG (Tan et al., 2025) enhances LLM reasoning by integrating knowledge reasoning paths from KGs.

Method	AQuA	GSM8K	MultiArith	SingEq	HotpotQA	CSQA	SIQA	Letter	Coin	Average	
			Without exte	rnal knowle	edge and exem	plars					
Zero-shot	42.6	77.8	95.9	86.8	80.1	74.5	77.3	35.8	76.7	71.9	
Zero-shot-CoT	43.4	78.3	96.7	88.5	81.4	75.6	78.0	34.5	75.6	72.4	
PS	50.1	82.8	96.9	89.4	83.0	74.2	76.3	44.7	79.5	75.2	
QDMRPS	47.3	83.8	95.2	90.7	86.7	76.6	77.8	36.5	76.3	74.5	
With exemplars											
Manual-CoT	54.3	85.8	97.2	92.3	85.7	79.6	82.4	39.6	79.2	77.3	
Auto-CoT	47.8	82.4	97.5	91.6	86.1	76.4	80.6	41.0	81.2	76.1	
Complex-CoT	51.7	83.5	96.6	92.8	82.8	76.9	78.5	37.7	81.9	75.8	
Iter-CoT	51.6	80.0	97.8	93.4	64.8	76.9	77.3	41.8	77.5	73.5	
ZEUS	51.9	88.4	97.3	92.8	84.9	77.4	81.7	42.8	82.5	77.7	
Pattern-CoT	52.8	85.3	97.7	91.3	82.5	76.3	78.9	41.4	83.7	76.7	
			Wii	th external l	knowledge						
KD-CoT	22.3	68.4	76.0	62.3	79.9	85.6	90.8	24.6	58.3	63.1	
IRCoT	20.7	65.6	78.3	65.1	87.5	82.8	87.9	28.2	54.5	63.4	
KG-CoT	12.3	66.8	78.6	61.5	<u>73.5</u>	88.9	92.1	26.7	53.2	<u>61.5</u>	
			With both ext	ernal know	ledge and exen	ıplars					
CoT-RAG (ours)	65.7	94.7	98.5	98.7	98.4	97.9	98.7	54.6	94.7	89.1	

Table 1: Accuracy on nine datasets from three categories of reasoning tasks using ERNIE-Speed-128K. Throughout the tables in this paper, the best results are highlighted in bold, and the poorest results are underlined.

Method	LaB	LeB	CFB	AGI	Average
Graph-for	m LLM-b	ased RA	G method	ds	
GraphRAG	94.8	97.5	73.1	54.6	80.0
KG-CoT	89.3	93.6	72.8	32.6	72.1
ToG	86.7	90.2	68.3	64.2	77.4
PoG	93.8	91.7	89.5	45.3	80.1
RRKG	91.3	92.4	74.7	27.4	71.5
RoG	90.4	88.1	88.7	67.5	83.7
Graph-CoT	54.7	68.2	63.1	24.7	52.7
ToG-2	92.5	93.7	76.6	70.8	83.4
AtomR	83.6	82.5	77.3	37.5	70.2
Ve	ariants oj	CoT-RA	.G		
CoT-RAG (IndexFlatL2)	92.7	93.7	85.2	67.8	84.9
CoT-RAG (IndexFlatIP)	91.5	91.9	87.1	72.4	85.7
CoT-RAG (IndexIVFFlat)	93.6	92.1	86.8	75.3	87.0
CoT-RAG (IndexLSH)	90.8	92.5	88.1	74.6	86.5
CoT-RAG (IndexPQ)	93.1	92.9	87.4	73.9	86.8
CoT-RAG (IndexIVFPQ)	93.6	93.1	87.8	76.2	87.7
CoT-RAG (Zero-expert)	93.6	94.7	86.3	74.8	87.4
CoT-RAG (ours)	99.3	98.6	94.7	88.3	95.2

Table 2: Accuracy on four datasets from vertical domains using GPT-40 mini

Moreover, we evaluate these methods across five NL-based LLMs: ERNIE-Speed-128K (Baidu, 2025), ERNIE-3.5-128K (Baidu, 2025), GLM-4-flash (Zhipuai, 2025), GPT-40 mini (OpenAI, 2024b), and GPT-4o (OpenAI, 2024a). Following Wei et al. (2022); Kojima et al. (2022), we invoke LLMs that are not fine-tuned via API and set the temperature to 0 to ensure deterministic outputs. However, there is one difference, that is, we set the max tokens to 1000 to accommodate the need for a longer context. Moreover, in line with Wang et al. (2023a), for Manual-CoT and Auto-CoT, we typically use 8 demonstration examples across most tasks, 4 examples for AQuA and Last Letter tasks, and 7 examples for CSQA. Complex-CoT selects exemplars with most complex rationales as demonstrations. The remaining baselines use default settings. In our proposed CoT-RAG, each DT node's Knowledge case contains only one demonstration example (e.g., Node 1 of Figure 3).

# 5 Experimental Results

Table 1 presents our results on ERNIE-Speed-128K, with additional results from other LLMs available in Appendix D. The experimental data demonstrate that CoT-RAG significantly improves reasoning accuracy across all datasets compared to existing CoT techniques, ranging from 4.0% to 44.3%. Specifically illustrated by representative results, compared to Manual-CoT, Zero-shot-CoT, Auto-CoT, PS, QDMRPS, KD-CoT, Iter-CoT and KG-CoT, the average accuracy across datasets increases by 4.0%-15.2%, 7.3%-23.0%, 5.0%-17.1%, 5.7%-18.5%, 4.8%-19.5%, 13.7%-41.1%, 4.2%-21.3%, and 13.8%-44.3%, respectively. The overall improvement in accuracy can be attributed to the comprehensive optimizations of our method in terms of credible CoT generation and rigorous instruction execution.

Table 2 indicates that, compared with the baselines, our CoT-RAG has increased accuracy by 8.9% to 80.6%. Specifically, compared to CoT-RAG methods based on different vector indexes, the average accuracy of our method has increased by 8.6%-12.1%. In particular, we observe that information loss and high running time may occur when using vectors to store and retrieve knowledge. Furthermore, the average accuracy of CoT-RAG (Zero-expert) is 7.8% lower than our CoT-RAG with expert involvement, which further highlights the indispensable role of experts in vertical domain applications. Other graph-form LLMbased RAG methods, including KG-CoT, ToG, PoG, RRKG, etc., target factual knowledge graph question answering (KGQA) and do not cover logical reasoning in vertical domains, showing a relatively low accuracy. In particular, these methods

have a very poor accuracy on the AGIEval dataset with high logical requirements. Additionally, it is clearly observable that the average accuracy of CoTRAG variants significantly outperforms graph-form LLM-based RAG methods. The above experimental analysis further justifies the strong scalability and sufficient novelty of our CoT-RAG. The relevant analysis regarding ablation results, robustness study, the relationship between accuracy and reasoning complexity, the runtime and resource usage can be found in Appendix C.

## 6 Related Work

## 6.1 Chain of Thought Reasoning

To utilize LLMs' reasoning capabilities, Wei et al. (2022) propose Chain-of-Thought prompting, adding reasoning steps before the answer to improve performance. Subsequent works have enhanced CoT reasoning in areas including prompt format (Gao et al., 2023; Chen et al., 2023; Lyu et al., 2023), selection (Yao et al., 2023a; Yu et al., 2024; Kumar et al., 2025; Zhang et al., 2025), ensemble (Weng et al., 2023; Li et al., 2023; Hu et al., 2024; Fu et al., 2023), decomposition (Zhou et al., 2023; Khot et al., 2023; Press et al., 2023; Huang et al., 2024a), model fine-tuning (Liao et al., 2025; Yeo et al., 2025), and planning (Wang et al., 2023a; Sun et al., 2024c; Wang et al., 2023b). Chen et al. (2023) design PoT, separating computation and reasoning with code-trained LLMs. Lyu et al. (2023) introduce Faithful CoT, combining natural and symbolic languages to enhance interpretability. Kojima et al. (2022) propose Zero-shot-CoT, automating reasoning steps with minimal human input. Wang et al. (2023a) address errors in Zero-shot-CoT by introducing PS prompting, which refines coarse task plans. Huang et al. (2024a) use QDMR-based graphs for structured problem decomposition. Kumar et al. (2025) design ZEUS to improve CoT prompting by utilizing uncertainty estimates to select effective demonstrations. Zhang et al. (2025) introduce Pattern-CoT, which employs reasoning patterns to enhance CoT prompting effectiveness. However, these methods solely rely on the reasoning and evaluation of the LLM itself, leading to low reliability. We refer readers to the survey (Chu et al., 2024) for more related works.

# 6.2 RAG Reasoning with Knowledge Graph

Retrieval-Augmented Generation (RAG) enhances LLMs by using relevant content retrieved from

knowledge sources, aiming to mitigate the blackbox nature and hallucination issues and improve text generation quality (Mao et al., 2024; Liu et al., 2024). Previous works have explored different ways in which LLMs leverage retrieved or generated text as external knowledge to boost reasoning (Lewis et al., 2020; Izacard et al., 2023; Sun et al., 2023b; Hagström et al., 2023). Recent research has further enhanced RAG reasoning by incorporating knowledge graphs (LUO et al., 2024; Ji et al., 2024; Ma et al., 2025). Sun et al. (2024b) introduce ToG, an innovative framework that enhances LLMs through interaction with knowledge graphs and expanding reasoning paths via beam search. Edge et al. (2024) propose GraphRAG, which combines the strengths of RAG and queryfocused summarization, using knowledge graphs to store source document data for efficient global question answering. Saleh et al. (2024) design SG-RAG, a zero-shot method, utilizing structured information from knowledge graphs to enable LLMs to accurately answer multi-hop questions. Ma et al. (2025) propose ToG-2, which utilizes KGs to connect documents via entities, facilitating deep and knowledge-guided context retrieval and enhancing the reasoning ability of LLMs. Our work incorporates RAG to enhance the reasoning of LLMs based on the generated knowledge graphs.

### 7 Conclusion

To address the issues of low reliability in reasoning chains generated solely by LLMs and poorer reasoning performance from natural language prompts compared with code prompts, we propose CoT-RAG, a novel reasoning framework with three key designs: (i) Knowledge Graph-driven CoT Generation, which introduces knowledge graphs to regulate LLMs' reasoning chain generation, thereby enhancing reasoning reliability; (ii) Learnable Knowledge Case-aware RAG, intergrating RAG into knowledge graphs to retrieve relevant subcases and sub-descriptions for learnable information; (iii) Pseudo-Program Prompting Execution, which inspires LLMs to execute logically reasoning tasks in pseudo-programs. Evaluations across nine public datasets demonstrate that CoT-RAG outperforms existing methods, particularly in four domain-specific datasets, validating its powerful cross-domain scalability.

# 8 Limitations

Although the CoT-RAG framework demonstrates effectiveness, it faces two key limitations. First, its implementation relies on LLMs with advanced program understanding and execution capabilities, thus publicly available LLMs, especially those with a smaller scale (e.g., 7B or 13B parameters), fall short on these requirements, necessitating our use of proprietary LLMs. This constraint limits the framework's generalizability and precludes the evaluation of how LLMs' parameter sizes impact CoT-RAG performance. Second, the construction of decision trees is influenced by the expert's domain-specific knowledge and background. In future, we plan to explore more automated methods in decision tree design. For example, in the literature of vertical domains (e.g., court proceeding records, medical case history, etc.), there are already descriptions of knowledge cases. We shall explore automated coarse-grained DT construction from such domain knowledge available in the form of text. Moreover, we will explore more widely applicable prompting techniques to adapt to small language models.

## 9 Acknowledgements

This work was supported by the National Key Research and Development Program of China (No. 2023YFB4502701), the National Natural Science Foundation of China (Grant Nos.U22A2027, 82090044 and 62402187), the China Postdoctoral Science Foundation (Nos. GZB20240243 and 2024M751009), and the Postdoctoral Project of Hubei Province (No.2024HBBHCXA024). Arijit Khan acknowledges support from the Novo Nordisk Foundation Grant NNF22OC0072415.

### References

Vahan Arsenyan, Spartak Bughdaryan, Fadi Shaya, Kent Wilson Small, and Davit Shahnazaryan. 2024. Large language models for biomedical knowledge graph construction: Information extraction from EMR notes. In *Proceedings of the 23rd Workshop on Biomedical Natural Language Processing*, pages 295–317, Bangkok, Thailand. Association for Computational Linguistics.

Yinhao Bai, Zhixin Han, Yuhua Zhao, Hang Gao, Zhuowei Zhang, Xunzhi Wang, and Mengting Hu. 2024. Is compound aspect-based sentiment analysis addressed by LLMs? In Findings of the Association for Computational Linguistics: EMNLP 2024, pages

7836–7861, Miami, Florida, USA. Association for Computational Linguistics.

Baidu. 2025. Ernie-speed-128k. https://cloud.baidu.com/product-s/qianfan\_home.

Joseph Barile, Alex Margolis, Grace Cason, Rachel Kim, Saia Kalash, Alexis Tchaconas, and Ruth Milanaik. 2024. Diagnostic accuracy of a large language model in pediatric case studies. *JAMA Pediatrics*, 178.

Wenhu Chen, Xueguang Ma, Xinyi Wang, and William W. Cohen. 2023. Program of thoughts prompting: Disentangling computation from reasoning for numerical reasoning tasks. *Transactions on Machine Learning Research*.

Cheng-Han Chiang and Hung-yi Lee. 2023. Can large language models be an alternative to human evaluations? In *Proceedings of the 61st Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 15607–15631, Toronto, Canada. Association for Computational Linguistics.

Jui-Sheng Chou and Ke-En Chen. 2024. Optimizing investment portfolios with a sequential ensemble of decision tree-based models and the fbi algorithm for efficient financial analysis. *Applied Soft Computing*, 158:111550.

Zheng Chu, Jingchang Chen, Qianglong Chen, Weijiang Yu, Tao He, Haotian Wang, Weihua Peng, Ming Liu, Bing Qin, and Ting Liu. 2024. Navigate through enigmatic labyrinth a survey of chain of thought reasoning: Advances, frontiers and future. In *Proceedings of the 62nd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 1173–1203, Bangkok, Thailand. Association for Computational Linguistics.

Karl Cobbe, Vineet Kosaraju, Mohammad Bavarian, Mark Chen, Heewoo Jun, Lukasz Kaiser, Matthias Plappert, Jerry Tworek, Jacob Hilton, Reiichiro Nakano, Christopher Hesse, and John Schulman. 2021. Training Verifiers to Solve Math Word Problems. *arXiv e-prints*, arXiv:2110.14168.

Avinava Dubey, Zhe Feng, Rahul Kidambi, Aranyak Mehta, and Di Wang. 2024. Auctions with llm summaries. In *Proceedings of the 30th ACM SIGKDD Conference on Knowledge Discovery and Data Mining*, KDD '24, page 713–722, New York, NY, USA. Association for Computing Machinery.

Darren Edge, Ha Trinh, Newman Cheng, Joshua Bradley, Alex Chao, Apurva Mody, Steven Truitt, and Jonathan Larson. 2024. From Local to Global: A Graph RAG Approach to Query-Focused Summarization. *arXiv e-prints*, arXiv:2404.16130.

Sarah Elhammadi, Laks V.S. Lakshmanan, Raymond Ng, Michael Simpson, Baoxing Huai, Zhefeng Wang, and Lanjun Wang. 2020. A high precision pipeline for financial knowledge graph construction.

- In *Proceedings of the 28th International Conference on Computational Linguistics*, pages 967–977, Barcelona, Spain (Online). International Committee on Computational Linguistics.
- Zhiwei Fei, Xiaoyu Shen, Dawei Zhu, Fengzhe Zhou, Zhuo Han, Alan Huang, Songyang Zhang, Kai Chen, Zhixin Yin, Zongwen Shen, Jidong Ge, and Vincent Ng. 2024. LawBench: Benchmarking legal knowledge of large language models. In *Proceedings of the 2024 Conference on Empirical Methods in Natural Language Processing*, pages 7933–7962, Miami, Florida, USA. Association for Computational Linguistics.
- Yao Fu, Hao Peng, Ashish Sabharwal, Peter Clark, and Tushar Khot. 2023. Complexity-based prompting for multi-step reasoning. In *The Eleventh International Conference on Learning Representations*.
- Luyu Gao, Aman Madaan, Shuyan Zhou, Uri Alon, Pengfei Liu, Yiming Yang, Jamie Callan, and Graham Neubig. 2023. Pal: Program-aided language models. In *International Conference on Machine Learning*, pages 10764–10799. PMLR.
- Zhibin Gou, Zhihong Shao, Yeyun Gong, yelong shen, Yujiu Yang, Nan Duan, and Weizhu Chen. 2024. CRITIC: Large language models can self-correct with tool-interactive critiquing. In *The Twelfth International Conference on Learning Representations*.
- Tanishk Govil, Tarek Hamid, Kimberly Helm, Elliot Greenberg, Kevin Landrum, J. Todd R. Lawrence, Theodore J. Ganley, and Amanda Watson. 2024. Using decision tree classifier to increase screening test sensitivity for the prediction of acl retear. In *Companion of the 2024 on ACM International Joint Conference on Pervasive and Ubiquitous Computing*, Ubi-Comp '24, page 171–175, New York, NY, USA. Association for Computing Machinery.
- Neel Guha, Julian Nyarko, Daniel E. Ho, Christopher Re, Adam Chilton, Aditya Narayana, Alex Chohlas-Wood, Austin Peters, Brandon Waldon, Daniel Rockmore, Diego Zambrano, Dmitry Talisman, Enam Hoque, Faiz Surani, Frank Fagan, Galit Sarfaty, Gregory M. Dickinson, Haggai Porat, Jason Hegland, Jessica Wu, Joe Nudell, Joel Niklaus, John J Nay, Jonathan H. Choi, Kevin Tobia, Margaret Hagan, Megan Ma, Michael Livermore, Nikon Rasumov-Rahe, Nils Holzenberger, Noam Kolt, Peter Henderson, Sean Rehaag, Sharad Goel, Shang Gao, Spencer Williams, Sunny Gandhi, Tom Zur, Varun Iyer, and Zehua Li. 2023. Legalbench: A collaboratively built benchmark for measuring legal reasoning in large language models. In Thirty-seventh Conference on Neural Information Processing Systems Datasets and Benchmarks Track.
- Shailja Gupta, Rajesh Ranjan, and Surya Narayan Singh. 2024. A Comprehensive Survey of Retrieval-Augmented Generation (RAG): Evolution, Current Landscape and Future Directions. *arXiv e-prints*, arXiv:2410.12837.

- Lovisa Hagström, Denitsa Saynova, Tobias Norlund, Moa Johansson, and Richard Johansson. 2023. The effect of scaling, retrieval augmentation and form on the factual consistency of language models. In *Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing*, pages 5457–5476, Singapore. Association for Computational Linguistics.
- Yuxin He, Buzhou Tang, and Xiaoling Wang. 2024. Generative models for automatic medical decision rule extraction from text. In *Proceedings of the 2024 Conference on Empirical Methods in Natural Language Processing*, pages 7034–7048, Miami, Florida, USA. Association for Computational Linguistics.
- Chi Hu, Yuan Ge, Xiangnan Ma, Hang Cao, Qiang Li, Yonghua Yang, Tong Xiao, and Jingbo Zhu. 2024. RankPrompt: Step-by-step comparisons make language models better reasoners. In *Proceedings of the 2024 Joint International Conference on Computational Linguistics, Language Resources and Evaluation (LREC-COLING 2024)*, pages 13524–13536, Torino, Italia. ELRA and ICCL.
- Jinfeng Huang, Qiaoqiao She, Wenbin Jiang, Hua Wu, Yang Hao, Tong Xu, and Feng Wu. 2024a. QDMR-based planning-and-solving prompting for complex reasoning tasks. In *Proceedings of the 2024 Joint International Conference on Computational Linguistics, Language Resources and Evaluation (LREC-COLING 2024)*, pages 13395–13406, Torino, Italia. ELRA and ICCL.
- Xia Huang. 2021. Analysis of public physical education teaching and quality evaluation in colleges and universities based on decision tree algorithm. In 2020 2nd International Workshop on Artificial Intelligence and Education, WAIE 2020, page 69–73, New York, NY, USA. Association for Computing Machinery.
- Yudong Huang, Hongyang Du, Xinyuan Zhang, Dusit Niyato, Jiawen Kang, Zehui Xiong, Shuo Wang, and Tao Huang. 2024b. Large language models for networking: Applications, enabling techniques, and challenges. *Netwrk. Mag. of Global Internetwkg.*, 39(1):235–242.
- Gautier Izacard, Patrick Lewis, Maria Lomeli, Lucas Hosseini, Fabio Petroni, Timo Schick, Jane Dwivedi-Yu, Armand Joulin, Sebastian Riedel, and Edouard Grave. 2023. Atlas: Few-shot learning with retrieval augmented language models. *Journal of Machine Learning Research*, 24(251):1–43.
- Yixin Ji, Kaixin Wu, Juntao Li, Wei Chen, Mingjie Zhong, Xu Jia, and Min Zhang. 2024. Retrieval and reasoning on KGs: Integrate knowledge graphs into large language models for complex question answering. In *Findings of the Association for Computational Linguistics: EMNLP 2024*, pages 7598–7610, Miami, Florida, USA. Association for Computational Linguistics.
- Bowen Jin, Chulin Xie, Jiawei Zhang, Kashob Kumar Roy, Yu Zhang, Zheng Li, Ruirui Li, Xianfeng

- Tang, Suhang Wang, Yu Meng, and Jiawei Han. 2024. Graph chain-of-thought: Augmenting large language models by reasoning on graphs. In *Findings of the Association for Computational Linguistics: ACL 2024*, pages 163–184, Bangkok, Thailand. Association for Computational Linguistics.
- Mohsinul Kabir, Mohammed Saidul Islam, Md Tahmid Rahman Laskar, Mir Tafseer Nayeem, M Saiful Bari, and Enamul Hoque. 2024. BenLLM-eval: A comprehensive evaluation into the potentials and pitfalls of large language models on Bengali NLP. In Proceedings of the 2024 Joint International Conference on Computational Linguistics, Language Resources and Evaluation (LREC-COLING 2024), pages 2238–2252, Torino, Italia. ELRA and ICCL.
- Akansha Kalra and Daniel S. Brown. 2024. Can differentiable decision trees learn interpretable reward functions?
- Tushar Khot, Harsh Trivedi, Matthew Finlayson, Yao Fu, Kyle Richardson, Peter Clark, and Ashish Sabharwal. 2023. Decomposed prompting: A modular approach for solving complex tasks. In *The Eleventh International Conference on Learning Representations*.
- Jaehyung Kim, Dongyoung Kim, and Yiming Yang. 2024. Learning to correct for QA reasoning with black-box LLMs. In *Proceedings of the 2024 Conference on Empirical Methods in Natural Language Processing*, pages 8916–8937, Miami, Florida, USA. Association for Computational Linguistics.
- Takeshi Kojima, Shixiang Shane Gu, Machel Reid, Yutaka Matsuo, and Yusuke Iwasawa. 2022. Large language models are zero-shot reasoners. *Advances in neural information processing systems*, 35:22199–22213.
- Rik Koncel-Kedziorski, Hannaneh Hajishirzi, Ashish Sabharwal, Oren Etzioni, and Siena Dumas Ang. 2015. Parsing algebraic word problems into equations. *Transactions of the Association for Computational Linguistics*, 3:585–597.
- Shanu Kumar, Saish Mendke, Karody Lubna Abdul Rahman, Santosh Kurasa, Parag Agrawal, and Sandipan Dandapat. 2025. Enhancing zero-shot chain of thought prompting via uncertainty-guided strategy selection. In *Proceedings of the 31st International Conference on Computational Linguistics*, pages 2003–2025, Abu Dhabi, UAE. Association for Computational Linguistics.
- Yang Lei, Jiangtong Li, Dawei Cheng, Zhijun Ding, and Changjun Jiang. 2023. CFBenchmark: Chinese Financial Assistant Benchmark for Large Language Model. *arXiv e-prints*, arXiv:2311.05812.
- Patrick Lewis, Ethan Perez, Aleksandra Piktus, Fabio Petroni, Vladimir Karpukhin, Naman Goyal, Heinrich Küttler, Mike Lewis, Wen-tau Yih, Tim Rocktäschel, Sebastian Riedel, and Douwe Kiela. 2020.

- Retrieval-augmented generation for knowledgeintensive nlp tasks. In *Advances in Neural Information Processing Systems*, volume 33, pages 9459– 9474. Curran Associates, Inc.
- Dongyuan Li, Ying Zhang, Zhen Wang, Shiyin Tan, Satoshi Kosugi, and Manabu Okumura. 2024. Active learning for abstractive text summarization via LLM-determined curriculum and certainty gain maximization. In *Findings of the Association for Computational Linguistics: EMNLP 2024*, pages 8959–8971, Miami, Florida, USA. Association for Computational Linguistics.
- Yifei Li, Zeqi Lin, Shizhuo Zhang, Qiang Fu, Bei Chen, Jian-Guang Lou, and Weizhu Chen. 2023. Making language models better reasoners with step-aware verifier. In *Proceedings of the 61st Annual Meeting of the Association for Computational Linguistics* (Volume 1: Long Papers), pages 5315–5333, Toronto, Canada. Association for Computational Linguistics.
- Jacky Liang, Wenlong Huang, Fei Xia, Peng Xu, Karol Hausman, Brian Ichter, Pete Florence, and Andy Zeng. 2023. Code as policies: Language model programs for embodied control. In 2023 IEEE International Conference on Robotics and Automation (ICRA), pages 9493–9500.
- Huanxuan Liao, Shizhu He, Yupu Hao, Xiang Li, Yuanzhe Zhang, Jun Zhao, and Kang Liu. 2025. SKIntern: Internalizing symbolic knowledge for distilling better CoT capabilities into small language models. In *Proceedings of the 31st International Conference on Computational Linguistics*, pages 3203–3221, Abu Dhabi, UAE. Association for Computational Linguistics.
- Wang Ling, Dani Yogatama, Chris Dyer, and Phil Blunsom. 2017. Program induction by rationale generation: Learning to solve and explain algebraic word problems. In *Proceedings of the 55th Annual Meeting of the Association for Computational Linguistics* (*Volume 1: Long Papers*), pages 158–167, Vancouver, Canada. Association for Computational Linguistics.
- Aiwei Liu, Qiang Sheng, and Xuming Hu. 2024. Preventing and detecting misinformation generated by large language models. In *Proceedings of the 47th International ACM SIGIR Conference on Research and Development in Information Retrieval*, SIGIR '24, page 3001–3004, New York, NY, USA. Association for Computing Machinery.
- LINHAO LUO, Yuan-Fang Li, Gholamreza Haffari, and Shirui Pan. 2024. Reasoning on graphs: Faithful and interpretable large language model reasoning. In *The Twelfth International Conference on Learning Representations*.
- Qing Lyu, Shreya Havaldar, Adam Stein, Li Zhang, Delip Rao, Eric Wong, Marianna Apidianaki, and Chris Callison-Burch. 2023. Faithful chain-ofthought reasoning. In *Proceedings of the 13th In*ternational Joint Conference on Natural Language

- Processing and the 3rd Conference of the Asia-Pacific Chapter of the Association for Computational Linguistics (Volume 1: Long Papers), pages 305–329, Nusa Dua, Bali. Association for Computational Linguistics.
- Shengjie Ma, Chengjin Xu, Xuhui Jiang, Muzhi Li, Huaren Qu, Cehao Yang, Jiaxin Mao, and Jian Guo. 2025. Think-on-graph 2.0: Deep and faithful large language model reasoning with knowledge-guided retrieval augmented generation. In *The Thirteenth International Conference on Learning Representations*.
- Aman Madaan, Niket Tandon, Prakhar Gupta, Skyler Hallinan, Luyu Gao, Sarah Wiegreffe, Uri Alon, Nouha Dziri, Shrimai Prabhumoye, Yiming Yang, Shashank Gupta, Bodhisattwa Prasad Majumder, Katherine Hermann, Sean Welleck, Amir Yazdanbakhsh, and Peter Clark. 2023. Self-refine: Iterative refinement with self-feedback. In *Thirty-seventh Conference on Neural Information Processing Systems*
- David M. Magerman. 1995. Statistical decision-tree models for parsing. In 33rd Annual Meeting of the Association for Computational Linguistics, pages 276–283, Cambridge, Massachusetts, USA. Association for Computational Linguistics.
- Potsawee Manakul, Adian Liusie, and Mark Gales. 2023. SelfCheckGPT: Zero-resource black-box hallucination detection for generative large language models. In *Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing*, pages 9004–9017, Singapore. Association for Computational Linguistics.
- Yuren Mao, Xuemei Dong, Wenyi Xu, Yunjun Gao, Bin Wei, and Ying Zhang. 2024. Fit-rag: Black-box rag with factual information and token reduction. *ACM Trans. Inf. Syst.* Just Accepted.
- Renê Mendes, Dimas Oliveira, and Victor Garcia. 2024. Application of generative AI as an enterprise wikibase knowledge graph Q&A system. In *Proceedings of the 1st Workshop on Knowledge Graphs and Large Language Models (KaLLM 2024)*, pages 35–42, Bangkok, Thailand. Association for Computational Linguistics.
- Meta-Research. 2024. Faiss indexes. https://github.com/facebookresearch/faiss/wiki/Faiss-indexes.
- Ibomoiye Domor Mienye and Nobert Jere. 2024. A survey of decision trees: Concepts, algorithms, and applications. *IEEE Access*, 12:86716–86727.
- OpenAI. 2024a. Gpt-4o. https://openai.com/index/hello-gpt-4o.
- OpenAI. 2024b. Gpt-4o-mini. https://openai.com/index/gpt-4o-mini-advancing-cost-efficient-intelligence.

- Debjit Paul, Mete Ismayilzada, Maxime Peyrard, Beatriz Borges, Antoine Bosselut, Robert West, and Boi Faltings. 2024. REFINER: Reasoning feedback on intermediate representations. In *Proceedings of the 18th Conference of the European Chapter of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 1100–1126, St. Julian's, Malta. Association for Computational Linguistics.
- Boci Peng, Yun Zhu, Yongchao Liu, Xiaohe Bo, Haizhou Shi, Chuntao Hong, Yan Zhang, and Siliang Tang. 2024. Graph Retrieval-Augmented Generation: A Survey. *arXiv e-prints*, arXiv:2408.08921.
- Pinecone. 2024. Faiss. https://www.pinecone.io/learn/series/faiss/faiss-tutorial/.
- Ofir Press, Muru Zhang, Sewon Min, Ludwig Schmidt, Noah Smith, and Mike Lewis. 2023. Measuring and narrowing the compositionality gap in language models. In *Findings of the Association for Computational Linguistics: EMNLP 2023*, pages 5687–5711, Singapore. Association for Computational Linguistics.
- Subhro Roy and Dan Roth. 2015. Solving general arithmetic word problems. In *Proceedings of the 2015 Conference on Empirical Methods in Natural Language Processing*, pages 1743–1752, Lisbon, Portugal. Association for Computational Linguistics.
- Ahmmad O. M. Saleh, Gokhan Tur, and Yucel Saygin. 2024. SG-RAG: Multi-hop question answering with large language models through knowledge graphs. In *Proceedings of the 7th International Conference on Natural Language and Speech Processing (IC-NLSP 2024)*, pages 439–448, Trento. Association for Computational Linguistics.
- Maarten Sap, Hannah Rashkin, Derek Chen, Ronan Le Bras, and Yejin Choi. 2019. Social IQa: Commonsense reasoning about social interactions. In Proceedings of the 2019 Conference on Empirical Methods in Natural Language Processing and the 9th International Joint Conference on Natural Language Processing (EMNLP-IJCNLP), pages 4463–4473, Hong Kong, China. Association for Computational Linguistics.
- Phillip Schneider, Tim Schopf, Juraj Vladika, Mikhail Galkin, Elena Simperl, and Florian Matthes. 2022. A decade of knowledge graphs in natural language processing: A survey. In Proceedings of the 2nd Conference of the Asia-Pacific Chapter of the Association for Computational Linguistics and the 12th International Joint Conference on Natural Language Processing (Volume 1: Long Papers), pages 601–614, Online only. Association for Computational Linguistics.
- Tao Shen, Guodong Long, Xiubo Geng, Chongyang Tao,Yibin Lei, Tianyi Zhou, Michael Blumenstein, andDaxin Jiang. 2024. Retrieval-augmented retrieval:Large language models are strong zero-shot retriever.

- In Findings of the Association for Computational Linguistics: ACL 2024, pages 15933–15946, Bangkok, Thailand. Association for Computational Linguistics.
- Ishika Singh, Valts Blukis, Arsalan Mousavian, Ankit Goyal, Danfei Xu, Jonathan Tremblay, Dieter Fox, Jesse Thomason, and Animesh Garg. 2022. Progprompt: Generating situated robot task plans using large language models. In Second Workshop on Language and Reinforcement Learning.
- Haotian Sun, Yuchen Zhuang, Lingkai Kong, Bo Dai, and Chao Zhang. 2023a. Adaplanner: Adaptive planning from feedback with language models. In Advances in Neural Information Processing Systems, volume 36, pages 58202–58245. Curran Associates, Inc.
- Jiashuo Sun, Yi Luo, Yeyun Gong, Chen Lin, Yelong Shen, Jian Guo, and Nan Duan. 2024a. Enhancing chain-of-thoughts prompting with iterative bootstrapping in large language models. In *Findings of the Association for Computational Linguistics: NAACL 2024*, pages 4074–4101, Mexico City, Mexico. Association for Computational Linguistics.
- Jiashuo Sun, Chengjin Xu, Lumingyuan Tang, Saizhuo Wang, Chen Lin, Yeyun Gong, Lionel Ni, Heung-Yeung Shum, and Jian Guo. 2024b. Think-on-graph: Deep and responsible reasoning of large language model on knowledge graph. In *The Twelfth International Conference on Learning Representations*.
- Simeng Sun, Yang Liu, Shuohang Wang, Dan Iter, Chenguang Zhu, and Mohit Iyyer. 2024c. PEARL: Prompting large language models to plan and execute actions over long documents. In *Proceedings of the 18th Conference of the European Chapter of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 469–486, St. Julian's, Malta. Association for Computational Linguistics.
- Zhiqing Sun, Xuezhi Wang, Yi Tay, Yiming Yang, and Denny Zhou. 2023b. Recitation-augmented language models. In *The Eleventh International Conference on Learning Representations*.
- Alon Talmor, Jonathan Herzig, Nicholas Lourie, and Jonathan Berant. 2019. CommonsenseQA: A question answering challenge targeting commonsense knowledge. In *Proceedings of the 2019 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long and Short Papers)*, pages 4149–4158, Minneapolis, Minnesota. Association for Computational Linguistics.
- Xingyu Tan, Xiaoyang Wang, Qing Liu, Xiwei Xu, Xin Yuan, and Wenjie Zhang. 2025. Paths-over-graph: Knowledge graph empowered large language model reasoning. In *Proceedings of the ACM on Web Conference 2025*, WWW '25, page 3505–3522, New York, NY, USA. Association for Computing Machinery.

- Harsh Trivedi, Niranjan Balasubramanian, Tushar Khot, and Ashish Sabharwal. 2023. Interleaving retrieval with chain-of-thought reasoning for knowledge-intensive multi-step questions. In *Proceedings of the 61st Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 10014–10037, Toronto, Canada. Association for Computational Linguistics.
- Keheng Wang, Feiyu Duan, Sirui Wang, Peiguang Li, Yunsen Xian, Chuantao Yin, Wenge Rong, and Zhang Xiong. 2023. Knowledge-Driven CoT: Exploring Faithful Reasoning in LLMs for Knowledge-intensive Question Answering. *arXiv e-prints*, arXiv:2308.13259.
- Lei Wang, Wanyu Xu, Yihuai Lan, Zhiqiang Hu, Yunshi Lan, Roy Ka-Wei Lee, and Ee-Peng Lim. 2023a. Plan-and-solve prompting: Improving zero-shot chain-of-thought reasoning by large language models. In *Proceedings of the 61st Annual Meeting of the Association for Computational Linguistics* (Volume 1: Long Papers), pages 2609–2634, Toronto, Canada. Association for Computational Linguistics.
- Xiaoyi Wang. 2021. Research on risk and supervision of financial big data application based on decision tree algorithm. In 2021 4th International Conference on Information Systems and Computer Aided Education, ICISCAE 2021, page 478–481, New York, NY, USA. Association for Computing Machinery.
- Yiheng Wang. 2024. Large language models evaluate machine translation via polishing. In *Proceedings* of the 2023 6th International Conference on Algorithms, Computing and Artificial Intelligence, ACAI '23, page 158–163, New York, NY, USA. Association for Computing Machinery.
- Zihao Wang, Shaofei Cai, Guanzhou Chen, Anji Liu, Xiaojian Ma, and Yitao Liang. 2023b. Describe, explain, plan and select: Interactive planning with LLMs enables open-world multi-task agents. In *Thirty-seventh Conference on Neural Information Processing Systems*.
- Jason Wei, Xuezhi Wang, Dale Schuurmans, Maarten Bosma, brian ichter, Fei Xia, Ed H. Chi, Quoc V Le, and Denny Zhou. 2022. Chain of thought prompting elicits reasoning in large language models. In *Advances in Neural Information Processing Systems*.
- Yixuan Weng, Minjun Zhu, Fei Xia, Bin Li, Shizhu He, Shengping Liu, Bin Sun, Kang Liu, and Jun Zhao. 2023. Large language models are better reasoners with self-verification. In *Findings of the Association for Computational Linguistics: EMNLP 2023*, pages 2550–2575, Singapore. Association for Computational Linguistics.
- Tomer Wolfson, Mor Geva, Ankit Gupta, Matt Gardner, Yoav Goldberg, Daniel Deutch, and Jonathan Berant. 2020. Break it down: A question understanding benchmark. *Transactions of the Association for Computational Linguistics*, 8:183–198.

Lianwei Wu, Yuan Rao, Yongqiang Zhao, Hao Liang, and Ambreen Nazir. 2020. DTCA: Decision tree-based co-attention networks for explainable claim verification. In *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*, pages 1024–1035, Online. Association for Computational Linguistics.

Amy Xin, Jinxin Liu, Zijun Yao, Zhicheng Lee, Shulin Cao, Lei Hou, and Juanzi Li. 2024. AtomR: Atomic Operator-Empowered Large Language Models for Heterogeneous Knowledge Reasoning. *arXiv e-prints*, arXiv:2411.16495.

Shunyu Yao, Dian Yu, Jeffrey Zhao, Izhak Shafran, Thomas L. Griffiths, Yuan Cao, and Karthik R Narasimhan. 2023a. Tree of thoughts: Deliberate problem solving with large language models. In *Thirty-seventh Conference on Neural Information Processing Systems*.

Shunyu Yao, Jeffrey Zhao, Dian Yu, Nan Du, Izhak Shafran, Karthik R Narasimhan, and Yuan Cao. 2023b. React: Synergizing reasoning and acting in language models. In *The Eleventh International Conference on Learning Representations*.

Edward Yeo, Yuxuan Tong, Morry Niu, Graham Neubig, and Xiang Yue. 2025. Demystifying Long Chain-of-Thought Reasoning in LLMs. *arXiv e-prints*, arXiv:2502.03373.

Junchi Yu, Ran He, and Zhitao Ying. 2024. THOUGHT PROPAGATION: AN ANALOGICAL APPROACH TO COMPLEX REASONING WITH LARGE LANGUAGE MODELS. In *The Twelfth International Conference on Learning Representations*.

Yufeng Zhang, Xuepeng Wang, Lingxiang Wu, and Jinqiao Wang. 2025. Enhancing chain of thought prompting in large language models via reasoning patterns. *Proceedings of the AAAI Conference on Artificial Intelligence*, 39(24):25985–25993.

Zhuosheng Zhang, Aston Zhang, Mu Li, and Alex Smola. 2023. Automatic chain of thought prompting in large language models. In *The Eleventh International Conference on Learning Representations*.

Penghao Zhao, Hailin Zhang, Qinhan Yu, Zhengren Wang, Yunteng Geng, Fangcheng Fu, Ling Yang, Wentao Zhang, Jie Jiang, and Bin Cui. 2024. Retrieval-Augmented Generation for Al-Generated Content: A Survey. *arXiv e-prints*, arXiv:2402.19473.

Ruilin Zhao, Feng Zhao, Long Wang, Xianzhi Wang, and Guandong Xu. 2024. Kg-cot: Chain-of-thought prompting of large language models over knowledge graphs for knowledge-aware question answering. In *Proceedings of the Thirty-Third International Joint Conference on Artificial Intelligence, IJCAI-24*, pages 6642–6650. International Joint Conferences on Artificial Intelligence Organization. Main Track.

Zhipuai. 2025. glm-4-flash. https://open.bigmodel.cn/dev/api/normal-model/glm-4.

Wanjun Zhong, Ruixiang Cui, Yiduo Guo, Yaobo Liang, Shuai Lu, Yanlin Wang, Amin Saied, Weizhu Chen, and Nan Duan. 2024. AGIEval: A human-centric benchmark for evaluating foundation models. In *Findings of the Association for Computational Linguistics: NAACL 2024*, pages 2299–2314, Mexico City, Mexico. Association for Computational Linguistics.

Denny Zhou, Nathanael Schärli, Le Hou, Jason Wei, Nathan Scales, Xuezhi Wang, Dale Schuurmans, Claire Cui, Olivier Bousquet, Quoc V Le, and Ed H. Chi. 2023. Least-to-most prompting enables complex reasoning in large language models. In *The Eleventh International Conference on Learning Representations*.

Wenhao Zhu, Hongyi Liu, Qingxiu Dong, Jingjing Xu, Shujian Huang, Lingpeng Kong, Jiajun Chen, and Lei Li. 2024. Multilingual machine translation with large language models: Empirical results and analysis. In *Findings of the Association for Computational Linguistics: NAACL 2024*, pages 2765–2781, Mexico City, Mexico. Association for Computational Linguistics.

# A CoT-RAG Algorithm and Time Complexity Analysis

As shown in Algorithm 1, the CoT-RAG framework involves three stages: 1) Knowledge Graphdriven CoT Generation (Lines 1-9), 2) Learnable Knowledge Case-aware RAG (Lines 11-18), and 3) Pseudo-Program Prompting Execution (Lines 20-23). First, experts construct and input a coarsegrained decision tree related to a field (Line 3), and then ask the LLM to decompose the decision tree into a fine-grained, highly structured knowledge graph (Line 4), which is represented by the Pseudo-Program Prompting we proposed (Line 5), referred to as the initialized pseudo-program knowledge graph (initial\_PKG). This can further be modified by the experts (Lines 6-8). Next, users enter query descriptions related to this field (Line 13), which can help update the decision tree benefiting from their usefulness (Lines 14-16). According to the initial\_PKG, the LLM subsequently extracts the corresponding sub-descriptions from the query description to update the initial\_PKG, namely updated\_PKG (Lines 17). Finally, the LLM executes updated\_PKG to output the final result (Lines 20-23).

The time complexity of the CoT-RAG algorithm is determined by three aforementioned stages.

# Algorithm 1 CoT-RAG: Enhanced Reasoning Framework Based on CoT and RAG

```
1: Function KnowledgeGraphDrivenCoTGeneration(field, ExpertModify):
         // field depends on the actual scenario of CoT-RAG application
 2:
 3:
         decision\_tree \leftarrow ExpertBuildDecisionTree(field)
 4:
         knowledge\_graph \leftarrow LLMDecomposeDecisionTree(decision\_tree)
         initial_PKG \leftarrow LLMGeneratePKG(knowledge\_graph)
 5:
         if ExpertModify is true then
 6:
              initial_PKG \leftarrow ExpertModifyKnowledgeGraph(initial_PKG)
 7:
 8:
         end
         return initial PKG
 9:
10:
11: Function LearnableKnowledgeCaseAwareRAG(initial_PKG, decision_tree, is_useful):
12:
         // Get the query description entered by the user
         query\_description \leftarrow GetUserQueryDescription()
13:
         if is_useful is true then
14:
              updated_decision_tree \leftarrow Update(query_description, decision_tree)
15:
         end
16:
         updated_PKG \leftarrow LLMExtractSubDescriptions(initial_PKG, query_description)
17:
18:
         return updated_PKG
19:
20: Function PseudoProgramPromptingExecution(updated_PKG):
         // Generate final results
21:
22:
         final\_result \leftarrow LLMGenerateFinalResult(updated\_PKG)
23:
         return final result
```

In Stage 1, the algorithm decomposes a coarsegrained decision tree into a fine-grained knowledge graph. The time complexity of this process is O(N), where N represents the number of nodes in the decision tree. In Stage 2, the algorithm extracts subdescriptions from user input and updates the knowledge graph. The time complexity of this stage is O(M), with M denoting the number of entities in the knowledge graph. In Stage 3, the algorithm inputs the updated knowledge graph into LLMs to generate the final result, with a time complexity of O(M) as well. Overall, because the number of entities M in the knowledge graph has a linear relationship with decision tree nodes' number N in our experiments, the total time complexity of the algorithm is O(N).

# **B** Frequently-Used Notations

Table 3 shows the frequently used notations in this paper.

# C Analysis

We also perform an extensive analysis of CoT-RAG for a deeper insight into the function of each component (§C.1), the robustness of exemplars (§C.2),

the relationship between accuracy and reasoning complexity (§C.3), and the runtime and resource usage (§C.4). GPT-40 mini is chosen as the benchmark model for all ensuing analyses. To ensure transparency and reproducibility, we have released the codebase, datasets, and manually designed decision trees at https://github.com/hustlfy123/CoT-RAG, enabling the research community to replicate and extend our findings.

# C.1 Ablation Results w.r.t. CoT-RAG Components

The CoT-RAG framework enhances LLMs' reasoning accuracy. To quantify the contribution of each component to accuracy improvement, we conduct an ablation analysis by removing different parts of the framework. We test four variants on arithmetic reasoning datasets (AQuA, GSM8K) and commonsense reasoning datasets (HotpotQA, CSQA): No node decomposition: Decision tree nodes are not decomposed into knowledge graphs; No RAG: Removal of entity *Sub-case* in the knowledge graph; No PsePrompting: Replace Pseudoprogram Prompting with NL prompts; No expert inspection: The initial PKG is not inspected by an expert.

Notation	Description
KG(s)	A Knowledge Graph(s) is a semantically structured network, which uses triplets to describe cross-domain entity relationships and attribute values (Mendes et al., 2024; Elhammadi et al., 2020). It integrates general facts, domain-specific knowledge (such as in mathematics, law, finance, and biomedical), and common-sense knowledge, providing a knowledge base and reasoning foundation for complex tasks like cross-domain intelligent decision-making and comprehensive knowledge-based question-answering (Schneider et al., 2022; Arsenyan et al., 2024). Distinct from traditional KGs, the nodes in our KGs include attributes such as Sub-question/Sub-case. The edges represent reasoning dependencies (i.e., "Answer Provision" relationships) rather than factual relationships as in the classic KGs and place a pronounced emphasis on the "structuring of reasoning logic" As depicted in Figure 3, the triple <entity 1,="" 5="" answer="" entity="" provision,=""> indicates that the answer of Entity 1 is provided to the sub-description of Entity 5, which aids in the reasoning and answer-generation processes of Entity 5. This relationship clarifies the information flow and dependency among different entities within the KG. It enables LLMs to leverage relevant information during the reasoning process, thereby enhancing the controllability and reliability of reasoning.</entity>
DT(s)	A Decision Tree(s) is a tree-structured machine-learning algorithm that finds extensive application in classification and regression tasks (Wu et al., 2020; Huang, 2021). Its fundamental principle entails recursively partitioning a dataset into progressively smaller subsets, thereby giving rise to a dendritic structure. In this construct, each internal node signifies a test of a characteristic or attribute, each branch represents a particular value of that characteristic, and each leaf node represents a classification result or a regression value (He et al., 2024; Magerman, 1995). Different from the concept described above, the decision tree in this paper refers to a tree-shaped structure constructed by experts, where there are well-defined logical relationships between parent and child nodes. The parent nodes represent more general or antecedent questions. The information and reasoning results from parent nodes are passed on to child nodes, facilitating the resolution of problems at the child node. For example, in the decision tree in Figure 3 about transporting dishes in a restaurant, the parent node could be "Calculate the number of times restaurant staff transport dishes," while a child node might be "Calculate the total time for transporting dishes." Branching conditions are determined by domain-specific logic defined by experts. Parent nodes point to child nodes to construct a complete reasoning path, enabling LLMs to reason step-by-step following this structured logic.
RAG	Retrieval-Augmented Generation incorporates a retrieval mechanism within the text-generation process (Gupta et al., 2024). It empowers LLMs to retrieve pertinent information from external knowledge repositories, and integrate it into the generated outputs, endowing the generated content with enhanced accuracy, pertinence, and factual fidelity. RAG can be generally categorized into three types: Documen-based Retrieval RAG, Graph-based Retrieval RAG, and Multi-modal Retrieval RAG (Zhao et al., 2024; Peng et al., 2024). Distinct from the above-mentioned vector-based retrieval RAG, the RAG in this paper harnesses the formidable retrieval capabilities inherent in LLMs. It retrieves relevant content from the knowledge cases of the decision tree and the query descriptions input by users, reducing runtime and enhancing accuracy (Appendix C).

Table 3: Frequently-Used Notations.

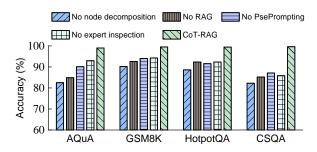


Figure 4: Ablation study results: accuracy when we remove different parts of CoT-RAG.

Figure 4 presents the experimental results, showing that accuracy decreases for all variants across datasets. Notably, removing node decomposition leads to the most significant drop in accuracy, highlighting its crucial role in the framework. In GSM8K, HotpotQA, and CSQA, RAG and expert inspection have similar contributions. However,

in AQuA, RAG's contribution is evidently higher, probably due to the difficulty of the dataset's mathematical tasks. In Table 1, AQuA's lower accuracy suggests its more challenging tasks, which may require knowledge retrieval. PsePrompting's contribution is moderate, probably because the average complexity of these datasets is low and thus inadequate to manifest its advantages.

### C.2 Robustness Results

Considering that robustness is a key performance indicator, we further conduct a comprehensive robustness analysis on CoT-RAG. In addition to the original CoT-RAG (Table 4), we test several variants to assess the impact of different factors. Specifically, we evaluate the effect of PsePrompting language by using C++ and Java to generate PKG (Appendix G). We also investigate the influence of knowledge cases by manually replacing all *Knowl-*

Method	AQuA	GSM8K	MultiArith	SingEq	HotpotQA	CSQA	SIQA	Letter	Coin	Average		
			Without exte	rnal knowle	edge and exemp	olars						
Zero-shot	82.4	92.6	96.8	97.7	87.2	83.3	87.1	92.2	98.7	90.9		
Zero-shot-CoT	81.8	93.3	96.3	97.3	87.7	82.7	87.8	92.7	98.0	90.8		
PS	82.1	92.8	97.7	98.1	89.9	84.9	86.5	97.3	99.9	92.1		
QDMRPS	83.5	94.3	98.5	98.8	90.8	83.9	85.9	96.7	98.3	92.3		
With exemplars												
Manual-CoT	90.2	95.3	98.1	98.0	92.3	89.0	93.9	97.7	100	94.9		
Auto-CoT	84.3	94.6	97.3	98.6	90.0	85.3	90.1	97.0	99.8	93.0		
Complex-CoT	90.8	93.9	97.2	98.0	91.3	86.8	89.3	96.8	98.8	93.7		
Iter-CoT	92.6	96.3	97.6	98.4	84.8	89.8	92.3	94.5	99.5	94.0		
ZEUS	91.8	95.0	97.9	98.3	91.8	88.6	89.3	98.3	98.7	94.4		
Pattern-CoT	86.7	94.6	97.9	98.3	88.3	85.8	88.6	94.3	99.5	92.7		
			Wii	th external l	knowledge							
KD-CoT	72.4	82.8	85.3	86.6	84.2	94.7	96.2	87.3	90.2	86.6		
IRCoT	76.6	83.5	82.9	88.2	86.5	92.8	95.9	85.8	93.7	87.3		
KG-CoT	<u>67.2</u>	84.6	88.0	84.9	84.7	93.6	96.8	89.5	92.3	86.8		
			With both ext	ernal know	ledge and exen	iplars						
CoT-RAG (ours)	99.0	99.5	99.9	99.6	99.4	99.6	98.3	99.9	100	99.5		

Table 4: Accuracy on nine datasets from three categories of reasoning tasks using GPT-40 mini.

Method	AQuA	GSM8K	MultiArith	SingEq	HotpotQA	CSQA	SIQA	Letter	Coin	Average	
			Without exte	rnal knowle	edge and exemp	olars					
Zero-shot	78.8	91.3	96.5	93.3	84.8	80.3	87.1	91.8	81.4	87.3	
Zero-shot-CoT	80.3	92.0	96.7	92.9	84.3	80.8	86.6	92.2	81.9	87.5	
PS	83.4	93.6	97.3	95.6	86.6	84.9	87.0	96.8	94.6	91.1	
QDMRPS	85.7	94.9	98.2	94.3	89.3	86.3	89.5	97.3	92.7	92.0	
With exemplars											
Manual-CoT	87.9	94.5	98.6	96.4	90.5	89.3	92.1	98.3	95.7	93.7	
Auto-CoT	82.6	95.3	98.2	95.5	91.4	87.3	91.5	96.4	96.8	92.8	
Complex-CoT	84.3	94.2	97.7	96.1	87.8	88.2	90.6	96.6	95.2	92.3	
Iter-CoT	85.8	95.4	98.5	97.2	87.8	86.7	92.6	97.8	96.9	93.2	
ZEUS	85.8	94.9	98.8	96.1	88.8	85.2	90.4	96.6	93.7	92.3	
Pattern-CoT	83.2	94.7	98.2	95.8	87.6	85.4	89.3	95.8	96.2	91.8	
			Wi	th external i	knowledge						
KD-CoT	64.3	84.3	82.9	82.2	81.8	93.5	92.5	79.8	88.6	83.3	
IRCoT	57.0	81.5	84.5	84.7	82.6	93.7	94.2	75.4	91.0	82.7	
KG-CoT	61.7	85.7	84.3	82.0	83.1	93.8	92.4	76.7	90.5	83.4	
			With both ext	ernal know	ledge and exen	iplars					
CoT-RAG (ours)	96.3	99.8	99.4	99.1	98.3	99.8	99.7	99.9	99.7	99.1	

Table 5: Accuracy on nine datasets from three categories of reasoning tasks using ERNIE-3.5-128K.

edge case in the decision tree twice and comparing inference accuracy. Following Wei et al. (2022), two other co-authors of this paper (A and B) assume expert roles to contribute to the decision tree design, allowing us to compare the impact of different experts on our method.

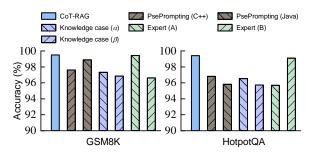


Figure 5: Robustness of CoT-RAG.

Figure 5 presents the experimental results across the GSM8K and HotpotQA datasets for all vari-

ants. Analysis shows that compared to the original CoT-RAG, the accuracy error for all variants does not exceed 4%, demonstrating the strong robustness of CoT-RAG. Furthermore, the participation of Expert A and Expert B caused noticeable accuracy fluctuations, indicating that the alignment between an expert's domain knowledge and the CoT-RAG framework can influence the quality of decision trees and final inference outcomes. This observation highlights the synergy between expert knowledge and the framework.

# C.3 Accuracy Results w.r.t. Varying Reasoning Complexity

To further showcase the superiority of our method, we contrast the average accuracies of diverse reasoning complexities on GSM8K (arithmetic reasoning) and HotpotQA (commonsense reasoning). Reasoning complexity is gauged by the entity

Method	AQuA	GSM8K	MultiArith	SingEq	HotpotQA	CSQA	SIQA	Letter	Coin	Average		
			Without exte	rnal knowle	edge and exem	plars						
Zero-shot	57.4	86.3	97.5	97.2	87.7	84.2	86.6	55.7	79.3	81.3		
Zero-shot-CoT	59.0	86.8	97.8	97.0	87.3	83.3	86.4	54.8	79.8	81.4		
PS	61.9	94.7	98.2	97.8	89.9	85.3	87.9	60.4	83.3	84.4		
QDMRPS	57.1	93.9	98.0	98.3	91.2	85.0	86.7	55.1	84.5	83.3		
With exemplars												
Manual-CoT	62.3	93.8	97.9	98.1	90.9	87.1	91.6	60.8	89.9	85.8		
Auto-CoT	57.8	92.0	98.4	97.9	92.4	84.2	89.3	58.3	88.7	84.3		
Complex-CoT	60.6	92.8	97.9	98.5	88.7	86.4	90.4	55.3	83.6	83.8		
Iter-CoT	64.6	94.7	97.7	97.6	82.4	86.8	92.8	57.5	92.2	85.1		
ZEUS	61.5	90.7	98.2	98.6	88.9	84.8	89.6	62.7	91.3	85.1		
Pattern-CoT	59.7	92.8	98.1	97.5	86.2	84.8	86.9	55.8	88.4	83.4		
			Wi	th external i	knowledge							
KD-CoT	38.8	76.0	83.5	80.8	74.8	89.7	92.5	29.8	68.4	70.5		
IRCoT	<u>35.1</u>	78.6	85.8	78.3	93.5	91.6	93.8	35.6	76.2	74.3		
KG-CoT	35.4	<u>72.2</u>	<u>78.6</u>	<u>75.6</u>	82.5	91.5	92.2	<u>25.4</u>	64.0	<u>68.6</u>		
			With both ext	ernal know	ledge and exen	ıplars						
CoT-RAG (ours)	72.3	99.2	99.5	99.8	98.6	98.8	98.1	64.3	98.7	92.2		

Table 6: Accuracy on nine datasets from three categories of reasoning tasks using GLM-4-flash.

Method	AQuA	GSM8K	MultiArith	SingEq	HotpotQA	CSQA	SIQA	Letter	Coin	Average	
			Without exte	rnal knowl	edge and exemp	plars					
Zero-shot	86.3	94.2	97.7	97.4	89.6	83.2	89.4	94.5	98.6	92.3	
Zero-shot-CoT	86.2	94.5	97.3	97.8	89.3	83.9	89.4	94.3	98.4	92.3	
PS	87.4	94.8	98.7	98.5	89.4	86.8	91.9	97.1	98.9	93.7	
QDMRPS	89.0	95.3	99.2	99.5	90.4	89.0	92.3	97.5	98.8	94.6	
With exemplars											
Manual-CoT	91.4	96.5	99.0	99.3	91.5	88.9	92.6	97.6	100	95.2	
Auto-CoT	88.7	95.6	98.8	98.8	91.8	87.3	91.7	96.8	99.2	94.3	
Complex-CoT	89.0	95.3	97.7	98.2	89.7	85.8	90.9	96.8	98.6	93.6	
Iter-CoT	93.2	96.1	99.2	98.7	89.5	90.5	91.7	97.3	99.4	95.1	
ZEUS	92.7	95.8	98.2	98.6	92.0	86.7	91.4	98.6	99.2	94.8	
Pattern-CoT	88.3	95.8	98.6	98.3	86.4	84.7	90.4	96.0	99.4	93.1	
			Wi	th external	knowledge						
KD-CoT	75.3	81.9	87.4	88.3	84.8	93.6	96.7	85.2	91.5	87.2	
IRCoT	75.9	84.2	83.6	86.4	83.8	93.2	94.9	89.7	92.8	87.2	
KG-CoT	77.6	83.2	86.8	83.6	85.9	93.2	97.1	82.6	93.6	87.1	
			With both ext	ernal know	ledge and exen	ıplars					
CoT-RAG (ours)	96.5	98.7	99.7	99.0	98.8	99.5	99.6	99.9	100	99.1	

Table 7: Accuracy on nine datasets from three categories of reasoning tasks using GPT-4o.

counts post-original decision tree decomposition. Considering that the maximum number of entities split from the two datasets is 9, we compared the accuracies with four levels of complexity: < 5, 6,7, and  $\geq$  8. According to Table 1, we select representative baselines for comparison: Zero-shot-CoT, PS, Manual-CoT, Auto-CoT, and Iter-CoT. As shown in Figure 6, for problems with over 8 entities, on GSM8K, CoT-RAG's average accuracy outperforms Manual-CoT by 3.4%, Zero-shot-CoT by 15.3%, Auto-CoT by 4.5%, PS by 6.5%, and Iter-CoT by 5.1%. On HotPotQA, the gains are 8.7%, 15.8%, 7.1%, 11.5%, and 10.3% respectively, attesting to our method's efficacy in complex reasoning tasks. Moreover, as the number of the entity increases, the superiority of CoT-RAG over the competitors enlarges progressively.

# C.4 Runtime and Resource Usage Results

We conduct more comparisons with CoT approaches on average runtime for answering each question in each dataset using GPT-40 mini, where our running time only includes *Stage 2* and *Stage 3* of our CoT-RAG, that is, CoT inference time to answer a user's question. Table 8 reveals that our method generally outperforms PS, QDMRPS, KD-CoT, IRCoT, KG-CoT, and ZEUS in terms of running time. However, it shows slightly higher runtime compared to Manual-CoT, Zero-shot, Zero-shot-CoT, and Auto-CoT, due to longer text processing workflow. Overall, it ranks in the middle range regarding the running time.

Additionally, we also compare with the average running time and token resource consumption of the graph-form LLM-based RAG methods. Table 9 indicates that, compared with GraphRAG, CoT-

Method	AQuA	GSM8K	MultiArith	SingEq	HotpotQA	CSQA	SIQA	Letter	Coin	Average	
			Without exte	rnal knowle	edge and exemp	olars					
Zero-shot	2.58	2.29	1.53	1.42	2.21	1.25	1.42	1.64	1.38	1.75	
Zero-shot-CoT	3.36	3.03	1.87	1.69	2.43	1.34	1.50	1.78	1.66	2.07	
PS	5.27	4.76	4.43	3.61	6.24	4.52	4.57	2.94	3.24	4.40	
QDMRPS	5.84	5.73	4.86	4.17	6.78	5.83	4.96	3.37	4.43	5.11	
With exemplars											
Manual-CoT	2.63	2.42	2.09	1.73	3.14	2.67	1.97	2.02	1.98	2.29	
Auto-CoT	4.04	3.87	3.27	2.08	3.68	3.18	2.14	2.60	2.27	3.01	
Complex-CoT	3.86	3.31	2.93	2.65	4.18	2.74	2.43	2.18	2.09	2.93	
Iter-CoT	4.05	3.79	3.11	2.85	3.94	3.71	3.65	2.67	2.48	3.36	
ZEUS	7.23	6.74	6.52	6.81	7.17	5.61	5.95	4.43	4.61	6.12	
Pattern-CoT	3.64	3.37	3.16	2.94	3.88	2.91	2.61	2.26	2.14	2.99	
			Wi	th external l	knowledge						
KD-CoT	7.34	7.89	6.96	7.61	8.83	8.54	8.37	6.32	5.86	7.52	
IRCoT	8.67	8.12	8.37	7.94	9.06	8.32	9.43	7.16	7.53	8.29	
KG-CoT	8.03	7.44	7.11	7.63	8.46	8.26	7.81	5.94	5.81	7.39	
			With both ext	ernal know	ledge and exen	iplars					
CoT-RAG (ours)	4.25	4.21	4.12	3.22	4.97	4.18	3.84	3.24	2.64	3.85	

Table 8: Runtime (sec.) on nine datasets from three categories of reasoning tasks using GPT-40 mini.

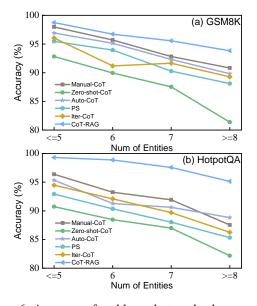


Figure 6: Accuracy of problems that can be decomposed into different numbers of entities.

RAG achieves an overall improvement, with 29.2% reduction in average runtime, and 33.4% decrease in average token consumption. Compared with the rest of the baselines, although CoT-RAG consumes more tokens, it shows an improvement in the average running time, in the range 13.9%-63.0%. Furthermore, Table 2 indicates that, compared with the baselines, our CoT-RAG has increased accuracy by 8.9% to 80.6%, which evidently demonstrates that the benefits far outweigh the costs.

However, the running time mentioned above does not incorporate *Stage 1* of our CoT-RAG, which consists of two components: the construction of DTs by experts and the decomposition of DTs into KGs by LLMs. Therefore, we conduct

a further exploration of the time it takes Expert A and Expert B to design DTs, and subsequently, the running time of LLMs to convert these DTs into detailed KGs, across four vertical domain datasets as presented in Table 10. Analysis demonstrates that the time cost of DT design is significantly influenced by an expert's expertise, while LLM-based KG conversion time is lower. Although *Stage 1* consumes a great deal of time, it only needs to be constructed once in an offline manner. When subsequently faced with thousands of user questions in real time, the running time of our CoT-RAG (*Stage 2 + Stage 3*) is lower than that of all the baseline methods (Table 9). The more questions from users there are, the more obvious our advantages will be.

Furthermore, in practical applications, experts can conduct desensitization preprocessing (i.e., deidentifying, masking, or replacing sensitive data) when constructing decision trees. Therefore, there is **no risk of data leakage** during the process in which LLMs decompose the decision trees to construct knowledge graphs.

## D Accuracy Results with Varying LLMs

Table 4, Table 5, Table 6, and Table 7 are the experimental results of GPT-40 mini, ERNIE-3.5-128K, GLM-4-flash, and GPT-40, respectively.

# E Example Outputs of the three stages of CoT-RAG

Tabel 11, Table 12, and Table 13 are example outputs for the three stages of CoT-RAG, respectively.

Method	LawBer	nch	LegalBe	nch	CFBenchi	nark	AGIEv	al
Method	Runtime (s)	Token	Runtime (s)	Token	Runtime (s)	Token	Runtime (s)	Token
		Graph-fo	rm LLM-based	RAG met	hods			
GraphRAG	5.85	2883	4.72	1850	5.08	2026	4.89	2817
KG-CoT	6.56	265	6.93	235	6.84	286	7.27	316
ToG	<u>10.4</u>	843	<u>11.7</u>	876	<u>9.66</u>	728	<u>13.5</u>	774
PoG	8.70	547	11.2	529	9.43	486	8.92	516
RRKG	6.84	583	6.31	542	7.24	647	7.67	628
RoG	6.27	628	5.92	583	5.38	516	6.64	674
Graph-CoT	7.24	1368	7.82	1237	6.73	1297	8.13	1312
ToG-2	9.53	1097	10.9	987	8.62	1128	10.2	1035
AtomR	6.34	1762	6.92	1531	7.08	1621	7.86	1687
		I	Variants of CoT	-RAG				
CoT-RAG (IndexFlatL2)	4.67	896	4.89	875	5.12	1028	5.23	1257
CoT-RAG (IndexFlatIP)	4.52	923	4.96	815	4.73	985	5.32	1316
CoT-RAG (IndexIVFFlat)	4.04	908	4.23	892	4.09	1028	4.53	1342
CoT-RAG (IndexLSH)	4.49	988	4.37	842	4.66	979	4.67	1268
CoT-RAG (IndexPQ)	4.32	1002	4.46	878	4.86	996	4.89	1225
CoT-RAG (IndexIVFPQ)	4.12	975	4.25	866	4.42	1005	4.62	1243
CoT-RAG (Zero-expert)	3.14	1532	3.27	1346	3.98	1653	4.31	1826
CoT-RAG (ours)	3.02	1557	3.46	1298	3.95	1643	4.11	1861

Table 9: Runtime and token consumption on four datasets from vertical domains using GPT-40 mini.

Expert	LawBench		LegalBench		CFBen	chmark	AGIEval	
Expert	DT design	DT to KG	DT design	DT to KG	DT design	DT to KG	DT design	DT to KG
A	262.3	5.7	232.7	4.3	179.2	3.8	454.8	7.2
В	130.8	6.2	469.1	6.9	58.8	3.2	574.3	6.4

Table 10: Average running time (sec.) for Experts A and B in decision tree design of various vertical domains and converting DTs to KGs using GPT-40 mini.

# F Examples of pseudo-program knowledge graphs from different datasets

Tables 14 to 37 are examples of pseudo-program knowledge graphs and their output results of GPT-40 mini for each dataset, where we can evidently observe that the pseudo-program format exhibits generality and scalability, enabling it to meet the requirements of various domains to great extent. In this study, the pseudo-program is designed based on reasoning chains within knowledge graphs. Its core logic involves guiding LLMs to conduct reasoning through well-defined steps. Although it is not designed with domain-specific syntax, due to its simple logical structure and comprehensibility, it can be flexibly adapted to different domain-specific reasoning paradigms. This is achieved by adjusting elements such as sub-questions, sub-cases, and sub-descriptions of nodes.

Consider the two pseudo-program prompts in Table 32 and Table 34 as examples. They belong to the law and finance, respectively. We ob-

serve that the variable names of these two pseudo-program prompt words are all "sub\_questionX", "sub\_caseX", "sub\_descriptionX", and "answerX", where "X" represents a specific number, and the logical connections are all made around these variable names. Therefore, when the application field shifts from law to finance, one only needs to modify the corresponding variable contents, the number of variables, and the logical relationships in Table 32, and then it can be transformed into Table 34, enabling fully automatic domain adaptation.

# G Pseudo-program knowledge graph in different programming language forms

Tables 38 to 41 are examples of pseudo-program knowledge graphs in both C++ and Java languages, along with their output results.

## **H** Descriptions of Datasets

**General domains: Arithmetic Reasoning:** (1) the AQUA (Ling et al., 2017) dataset of algebraic word problems with NL rationales, (2) the GSM8K

(Cobbe et al., 2021) dataset of high quality linguistically diverse grade school math word problems created by human problem writers, (3) the MultiArith (Roy and Roth, 2015) dataset of mathematical problems that necessitates numerous inference steps for resolution, (4) the SingleEq (Koncel-Kedziorski et al., 2015) dataset of algebraic problems calls for the solution of equations; Commonsense Reasoning: (5) the HotpotQA (Wolfson et al., 2020) dataset of commonsense questions based on Wikipedia that requires reasoning using multiple supporting problems, (6) the CSQA (Talmor et al., 2019) benchmark dataset of multiplechoice questions that require different types of commonsense knowledge to obtain the correct answers, (7) the SIQA (Sap et al., 2019) dataset of questions which focus on inferring people's behavior and its social impact; Symbolic Reasoning: (8) the Last Letter Concatenation (Wei et al., 2022) dataset of questions requiring the last letters of words in a name to be concatenated (e.g., "Donald Trump"→"dp"), and (9) the Coin Flip (Wei et al., 2022) dataset of questions on whether a coin is still heads up after it is flipped or not flipped based on steps given in the questions.

Vertical domains: LawBench (LaB) (Fei et al., 2024): A Chinese legal benchmark including tasks such as entity recognition, reading comprehension, and crime amount calculation; LegalBench (LeB) (Guha et al., 2023): A American legal benchmark featuring 162 legal reasoning tasks; CFBenchmark (CFB) (Lei et al., 2023): A Chinese financial benchmark assessing the performance of LLMs in finance; AGIEval (AGI) (Zhong et al., 2024): A benchmark for evaluating human cognitive task performance, where we select the dataset with a focus on logic-based question answering.

However, in our experiments, the problem categories in existing open-source question-answering datasets are too diverse to directly meet our testing requirements. Thus, following GraphRAG (Edge et al., 2024) and Graph-CoT (Jin et al., 2024), we utilize five LLMs in § 4.2 to generate datasets that suit our testing needs. Specifically, we first select 200 questions with distinct reasoning logics from each open-source dataset. Then, we prompt each LLM to generate four new questions for each selected question. These new questions had the same reasoning logic but different content. As a result, each domain-specific dataset yields a collection of 200 question sets, each containing 21 questions (one original question plus twenty new

questions generated by the five LLMs) with the same reasoning logic, namely a total of 4200 questions for each dataset, whose scale is far larger than the 125 questions per dataset of GraphRAG (Edge et al., 2024), the 2579 questions per dataset on average of KD-CoT (Wang et al., 2023), the 947 questions per dataset on average of Iter-CoT (Sun et al., 2024a), the 618 questions per dataset on average of Pattern-CoT (Zhang et al., 2025), the 500 questions per dataset on average of IRCoT (Trivedi et al., 2023), the 174 questions per domain on average of Graph-CoT (Jin et al., 2024), and the 1334 questions per dataset on average of PoG (Tan et al., 2025), evincing the substantial magnitude and adequacy of our datasets. The specific prompt is shown in Table 42 and the generated datasets are available at https://github.com/hustlfy123/CoT-RAG.

```
def LLM(case, question, description):
  LLM_answer=Given the question, the case is an example about the question. Please study this example and
answer the question based on the description
  return LLM_answer
sub_question1="How many lunch trays can the person carry at once?"
sub_case1='
Description: Tom is helping the restaurant staff collect dishes and he can only carry 6 plates at a time.
Answer: Tom can carry 6 plates at a time."
sub_description1=" |
answer1=LLM(sub_case1, sub_question1,sub_description1)
sub_question2="How many trays does the person need to pick up from the first table?"
sub_case2='
Description: Tom needs to pick up 15 plates from one table and 9 plates from another.
Answer: Tom needs to pick up 15 plates from the first table."
answer2=LLM(sub_case2,sub_question2,sub_description2)
sub_question3="How many trays does the person need to pick up from the second table?"
sub case3='
Description: Tom needs to pick up 15 plates from one table and 9 plates from another.
Answer: Tom needs to pick up 9 plates from the second table."
sub_description3=" '
answer3=LLM(sub_case3,sub_question3,sub_description3)
sub_question4="How many trays does the person need to carry in total?"
sub_case4=
Description: Tom needs to pick up 15 plates from the first table and 9 plates from the second table.
Answer: 15+9=24. Therefore Tom needs to carry 24 plates in total.'
#Answering sub_question4 requires relying on the answers of sub_question2 and sub_question3, namely answer2
and answer3
sub_description4=answer2+answer3
answer4=LLM(sub_case4,sub_question4,sub_description4)
sub_question5="How many trips will the person need to make to carry all the trays?"
Description: Tom can carry 6 plates at a time and Tom needs to carry 24 plates in total.
Answer: 24/6=4. Therefore Tom needs to carry 24 plates in total."
#Answering sub_question5 requires relying on the answers of sub_question1 and sub_question4, namely answer1
and answer4
sub_description5=answer1+answer4
answer5=LLM(sub_case5,sub_question5,sub_description5)
#The final answer
final_answer=answer5
print(final_answer)
```

Table 11: The initialized pseudo-program knowledge graph on MultiArith, where entity 1, 2, and 3's *Sub-description* are **null** values. (Question: Roger was helping the cafeteria workers pick up lunch trays, but he could only carry 4 trays at a time. If he had to pick up 10 trays from one table and 2 trays from another, how many trips will he make?)

times, which means you will answer corresponding questions based on different cases and descriptions. Please output the final result of this program text in natural language. def LLM(case,question,description): LLM answer=Given the question, the case is an example about the question. Please study this example and answer the question based on the description return LLM answer sub\_question1="How many lunch trays can the person carry at once?" sub\_case1=" Description: Tom is helping the restaurant staff collect dishes, but he can only carry 6 plates at a time. Answer: Tom can carry 6 plates at a time." sub\_description1=" Roger can only carry 4 trays at a time. " answer1=LLM(sub\_case1, sub\_question1,sub\_description1) sub\_question2="How many trays does the person need to pick up from the first table?" sub\_case2=" Description: Tom needs to pick up 15 plates from one table and 9 plates from another. Answer: Tom needs to pick up 15 plates from the first table." sub\_description2="Roger have to pick up 10 trays from the first table." answer2=LLM(sub\_case2,sub\_question2,sub\_description2) sub\_question3="How many trays does the person need to pick up from the second table?" sub\_case3=' Description: Tom needs to pick up 15 plates from one table and 9 plates from another. Answer: Tom needs to pick up 9 plates from the second table." sub\_description3="Roger have to pick up 2 trays from the second table." answer3=LLM(sub\_case3,sub\_question3,sub\_description3) sub\_question4="How many trays does the person need to carry in total?" sub case4=' Description: Tom needs to pick up 15 plates from the first table and 9 plates from the second table. Answer: 15+9=24. Therefore Tom needs to carry 24 plates in total." #Answering sub\_question4 requires relying on the answers of sub\_question2 and sub\_question3, namely answer2 and answer3 sub\_description4=answer2+answer3 answer4=LLM(sub\_case4,sub\_question4,sub\_description4) sub\_question5="How many trips will the person need to make to carry all the trays?" sub\_case5=1 Description: Tom can carry 6 plates at a time and Tom needs to carry 24 plates in total. Answer: 24/6=4. Therefore Tom needs to carry 24 plates in total." #Answering sub\_question5 requires relying on the answers of sub\_question1 and sub\_question4, namely answer1 and answer4 sub\_description5=answer1+answer4 answer5=LLM(sub\_case5,sub\_question5,sub\_description5) #The final answer final answer=answer5 print(final\_answer)

#Given a program text, your role is the LLM function, which has three parameters: case, question, and description. The description is textual, the question requires you to answer based on the description, and the case involves processed scenarios or supplementary information or related considerations to assist you in answering the question. You need to strictly follow the program logic to execute. During execution, the LLM function will be called multiple

Table 12: The updated pseudo-program knowledge graph on MultiArith, where entity 1, 2, and 3's *Sub-description* are **assigned** values. (Question: Roger was helping the cafeteria workers pick up lunch trays, but he could only carry 4 trays at a time. If he had to pick up 10 trays from one table and 2 trays from another, how many trips will he make?)

To derive the final answer step by step, we will process the given sub-questions, descriptions, cases, and collect the responses according to the logic of the program.

### 1. Subquestion 1:

- Question: How many lunch trays can the person carry at once?
- Description: Roger can only carry 4 trays at a time.
- Case: The example describes Tom, who can carry 6 plates at a time.
- Answer: Roger can carry 4 trays at a time.

#### 2. Subquestion 2:

- Question: How many trays does the person need to pick up from the first table?
- Description: Roger have to pick up 10 trays from the first table.
- Case: The example describes Tom needing to pick up 15 plates from one table and 9 plates from another.
- Answer: Roger needs to pick up 10 trays from the first table.

### 3. Subquestion 3:

- Question: How many trays does the person need to pick up from the second table?
- Description: Roger have to pick up 2 trays from the second table.
- Case: The example describes Tom needing to pick up 15 plates from one table and 9 plates from another.
- Answer: Roger needs to pick up 2 trays from the second table.

### 4. Subquestion 4:

- Question: How many trays does the person need to carry in total?
- Description: Based on answers from subquestions 2 and 3, Roger's total tray count is 10 (from the first table) + 2 (from the second table).
- Case: The example explains how Tom calculates his total number of plates.
- Answer: 10 + 2 = 12. Therefore, Roger needs to carry 12 trays in total.

## 5. Subquestion 5:

- Question: How many trips will the person need to make to carry all the trays?
- Description: Based on answers from subquestions 1 and 4, Roger can carry 4 trays at a time and he needs to carry 12 trays in total.
- Case: The example outlines how Tom calculates the number of trips needed.
- Answer: 12 (total trays) / 4 (trays per trip) = 3 trips. Therefore, Roger needs to make 3 trips to carry all the trays.

Finally, by compiling all the answers, we find that:

Final Answer: Roger needs to make 3 trips to carry all the trays.

Table 13: The output of executing the pseudo-program knowledge graph using GPT-40 mini in Table 12.

#Given a program text, your role is the LLM function, which has three parameters: case, question, and description. The description is textual, the question requires you to answer based on the description, and the case involves processed scenarios or supplementary information or related considerations to assist you in answering the question. You need to strictly follow the program logic to execute. During execution, the LLM function will be called multiple times, which means you will answer corresponding questions based on different cases and descriptions. Please output the final result of this program text in natural language.

```
def LLM(case, question, description):
```

LLM\_answer=Given the question, the case is an example about the question. Please study this example and answer the question based on the description.

return LLM answer

sub\_question1="Assuming the speed of the first person is v kilometers per hour, what is the speed of the second person in kilometers per hour?"

sub\_case1=" Description: Cyclist A's rate is 20% faster than Cyclist B's.

Answer: Given that Cyclist A has a speed 20% faster than Cyclist B, the speed of Cyclone A is v+0.2v=1.2v kilometers per hour."

sub\_description1="Friend P's rate is 15% faster than Friend Q's."
answer1=LLM(sub\_case1,sub\_question1,sub\_description1)

sub\_question2="How to express the time t (expressed in an equation containing v) when two people start from both ends and walk towards each other at the same time?"
sub\_case2="

Description:Two cyclists plan to ride along a 60-km route, starting at opposite ends of the route at the same time. The speed of Cyclist A is 1.2v, and the speed of Cyclist B is v.

Answer: Two people are walking towards each other, with a total distance of 60 kilometers. According to the formula for the encounter problem, "Encounter time=total distance  $\div$  sum of two people's speeds", the sum of two people's speeds is v+1.2v=2.2v kilometers per hour, so the encounter time t=60/2.2v=300/11v hours"

#Answering sub\_question2 requires relying on the answer of sub\_question1, namely answer1 sub\_description2="Two friends plan to walk along a 43-km trail, starting at opposite ends of the trail at the same

time."+answer1
answer2=LLM(sub\_case2,sub\_question2,sub\_description2)

sub\_question3="How many kilometers will the first person walk when they meet?" sub\_case3="

Description: The speed of Cyclist A is 1.2v and the encounter time t=60/2.2v=300/11v hours.

Answer: The speed of Cyclone A is 1.2v kilometers per hour, and the encounter time is 300/11v hours. According to the distance=speed x time, the distance traveled by Cyclone A is S=1.2v x 300/11v=360/11=32.73 kilometers." #Answering sub\_question3 requires relying on the answers of sub\_question1 and sub\_question2, namely answer1 and answer2

sub\_description3=answer1+answer2

answer3=LLM(sub\_case3,sub\_question3,sub\_description3)

#The final answer final\_answer=answer3 print(final\_answer)

Table 14: The updated pseudo-program knowledge graph on AQuA. (Question: Two friends plan to walk along a 43-km trail, starting at opposite ends of the trail at the same time. If Friend P's rate is 15% faster than Friend Q's, how many kilometers will Friend P have walked when they pass each other?)

Let's break down the program logic step by step and determine the final answer.

Step-by-Step Breakdown:

- 1. First Call to LLM (sub\_question1):
- Question: Assuming the speed of the first person is v kilometers per hour, what is the speed of the second person in kilometers per hour?
- Case: Cyclist A's rate is 20% faster than Cyclist B's.
- Description: Friend P's rate is 15% faster than Friend Q's.

#### Angwer

- Since Friend P's rate is 15% faster than Friend Q's, the speed of Friend P is v + 0.15v = 1.15v kilometers per hour.

answer1 = "1.15v"

- 2. Second Call to LLM (sub\_question2):
- Question: How to express the time t (expressed in an equation containing v) when two people start from both ends and walk towards each other at the same time?
- Case: Two cyclists plan to ride along a 60-km route, starting at opposite ends of the route at the same time. The speed of Cyclist A is 1.2v, and the speed of Cyclist B is v.
- Description: Two friends plan to walk along a 43-km trail, starting at opposite ends of the trail at the same time.
- The description from 'answer1' tells us that Friend P's speed is 1.15v (instead of 1.2v), and Friend Q's speed is v.

#### Answer

- The total distance is 43 kilometers.
- The sum of the speeds is v + 1.15v = 2.15v.
- The encounter time is t = Total distance / Sum of the speeds = 43/2.15v.

answer2 = "43 / 2.15v"

- 3. Third Call to LLM (sub\_question3):
- Question: How many kilometers will the first person walk when they meet?
- Case: The speed of Cyclist A is 1.2v and the encounter time t = 300/11v hours.
- Description: The description from 'answer1' tells us the speed of Friend P is 1.15v, and 'answer2' gives the encounter time equation as t = 43/2.15v.

### Answer:

- The speed of Friend P is 1.15v km/h.
- The encounter time is t = 43/2.15v hours.
- The distance walked by Friend P is:

Distance = Speed  $\times$ Time = 1.15v  $\times$  43/2.15v = 49.45  $\times$  2.15 = 23 kilometers.

answer3 = "23 kilometers"

### Final Answer:

The first person (Friend P) will walk 23 kilometers when they meet.

Final Output: 23 kilometers

Table 15: The output of executing the pseudo-program knowledge graph in Table 14.

The description is textual, the question requires you to answer based on the description, and the case involves processed scenarios or supplementary information or related considerations to assist you in answering the question. You need to strictly follow the program logic to execute. During execution, the LLM function will be called multiple times, which means you will answer corresponding questions based on different cases and descriptions. Please output the final result of this program text in natural language. def LLM(case, question, description): LLM\_answer=Given the question, the case is an example about the question. Please study this example and answer the question based on the description. return LLM\_answer sub\_question1="How many times does the person write letters to each friend per week?" sub\_case1= Description: Tom writes a letter to 4 different friends once a week. Answer: Tom writes to each friend once a sub\_description1="James writes a letter to 2 different friends twice a week." answer1=LLM(sub\_case1,sub\_question1,sub\_description1) sub\_question2="How many letters does the person write per week?" sub\_case2=" Description: Tom writes a letter to 4 different friends once a week and he writes to each friend once a week. Answer: Because there are 4 different friends and I write letters to each friend once a week, the total number of letters written per week is 4 x 1=4 letters.' #Answering sub\_question2 requires relying on the answer of sub\_question1, namely answer1 sub\_description2="James writes a letter to 2 different friends."+answer1 answer2=LLM(sub\_case2,sub\_question2,sub\_description2) sub\_question3="What is the total number of pages written by the person per week?" sub\_case3=" Description: The total number of letters written by Tom per week is 4 letters and each letter is 5 pages long. Answer: Given that each letter has 5 pages, a total of 4 letters are written per week. Based on the total number of pages=the number of pages per letter x the total number of letters, the total number of pages written per week is 5 x 4=20 pages." #Answering sub question3 requires relying on the answers of sub question2, namely answer2 sub\_description3="Each letter is 3 pages long."+answer2 answer3=LLM(sub\_case3,sub\_question3,sub\_description3) sub\_question4="What is the total number of pages written by the person a year?" sub\_case4=" Description: One year is about 52 weeks and the total number of pages written by Tom per week is 20 pages. Answer: Writing 20 pages per week, and knowing that there are approximately 52 weeks in a year, according to the total number of pages per year=pages per week x weeks per year, the total number of pages written in a year is 1040 pages." #Answering sub\_question4 requires relying on the answers of sub\_question3, namely answer3 sub description4=answer3 answer4=LLM(sub\_case4,sub\_question4,sub\_description4) #The final answer

#Given a program text, your role is the LLM function, which has three parameters: case, question, and description.

Table 16: The updated pseudo-program knowledge graph on GSM8K. (Question: James writes a 3-page letter to 2 different friends twice a week. How many pages does he write a year?)

final\_answer=answer4
print(final\_answer)

Let's break down the program step by step and follow the logic to find the final answer.

Step-by-Step Breakdown:

- 1. First Call to LLM (sub\_question1):
- Question: How many times does the person write letters to each friend per week?
- Case: The case provides an example where Tom writes a letter to 4 different friends once a week.
- Description: James writes a letter to 2 different friends twice a week.

Answer: Based on the description, James writes a letter to each friend twice a week.

answer1 = "twice"

- 2. Second Call to LLM (sub\_question2):
- Question: How many letters does the person write per week?
- Case: The case tells us that Tom writes a letter to 4 different friends once a week.
- Description: The description tells us James writes a letter to 2 different friends, and \*\*twice\*\* a week (from 'answer1').

### Calculation:

- James writes 2 letters (one to each friend) twice a week, so the total number of letters James writes per week is:

 $2 \times 2 = 4$  letters

answer2 = "4 letters"

- 3. Third Call to LLM (sub\_question3):
- Question: What is the total number of pages written by the person per week?
- Case: This case mentions that Tom writes 4 letters per week, and each letter is 5 pages long.
- Description: The description tells us that each letter is 3 pages long (adjusted for James based on 'answer2').

### Calculation:

- James writes 4 letters per week, and each letter is 3 pages long.
- Total pages per week:

```
4 \times 3 = 12 pages per week
answer3 = "12 pages"
```

- 4. Fourth Call to LLM (sub\_question4):
- Question: What is the total number of pages written by the person in a year?
- Case: The case tells us that one year is approximately 52 weeks, and the total number of pages written per week is 12 pages (from 'answer3').
- Description: We use the same number of pages per week (12), and multiply by the number of weeks in a year.

### Calculation:

- Total pages per year:

```
12 \times 52 = 624 pages per year
```

answer4 = "624 pages"

### Final Answer:

The total number of pages James writes in a year is 624 pages.

Final Output: 624 pages

Table 17: The output of executing the pseudo-program knowledge graph in Table 16.

processed scenarios or supplementary information or related considerations to assist you in answering the question. You need to strictly follow the program logic to execute. During execution, the LLM function will be called multiple times, which means you will answer corresponding questions based on different cases and descriptions. Please output the final result of this program text in natural language. def LLM(case, question, description): LLM\_answer=Given the question, the case is an example about the question. Please study this example and answer the question based on the description. return LLM\_answer sub\_question1="What was the total number of shells initially found by the person?" sub\_case1=' Description: Lily picked up 80 seashells by the seaside. Answer: The total number of shells initially found by Lily initially was 80." sub\_description1="Joan found 70 seashells on the beach." answer1=LLM(sub\_case1,sub\_question1,sub\_description1) sub\_question2="How many shells does this person have left?" sub\_case2=" Description: Lily has 32 seashells left. Answer: Lily has 32 seashells left." sub\_description2="Joan has 27 seashell left." sub\_question3="How many shells did this person give to others?" sub\_case3=" Description: The total number of shells initially found by Lily initially was 80 and she has 32 seashells left. Answer: Subtracting the remaining number of shells from Lily's initial total number of shells gives her the number

#Given a program text, your role is the LLM function, which has three parameters: case, question, and description. The description is textual, the question requires you to answer based on the description, and the case involves

Table 18: The updated pseudo-program knowledge graph on SingleEq. (Question: Joan found 70 seashells on the beach. she gave Sam some of her seashells. She has 27 seashell left. How many seashells did she give to Sam?)

of shells she gives to others. That is to say, the number of shells given=the initial total number of shells - the

#Answering sub\_question3 requires relying on the answers of sub\_question1 and sub\_question2, namely answer1

remaining number of shells, that is, 80-32=48 (pieces)."

answer3=LLM(sub\_case3,sub\_question3,sub\_description3)

sub\_description3=answer1+answer2

and answer2

#The final answer final\_answer=answer3 print(final\_answer) Let's walk through the program step by step to determine the final answer.

Step-by-Step Breakdown:

- 1. First Call to LLM (sub\_question1):
- Question: What was the total number of shells initially found by the person?
- Case: The case provides an example where Lily initially found 80 seashells.
- Description: Joan found 70 seashells on the beach.

Answer: The total number of shells initially found by Joan was 70, based on the description provided.

answer1 = "70"

- 2. Second Call to LLM (sub\_question2):
- Question: How many shells does this person have left?
- Case: The case provides an example where Lily has 32 seashells left. Description: Joan has 27 seashells left.

Answer: According to the description, Joan has 27 seashells left.

answer2 = "27"

- 3. Third Call to LLM (sub\_question3):
- Question: How many shells did this person give to others?
- Case: This case involves calculating how many shells were given by Lily based on the initial number of shells and how many are left. The formula provided is:

Shells given = Initial shells - Remaining shells

- Description: We are given the answers from previous questions: Joan initially found 70 shells and has 27 left.

### Calculation:

- Initially, Joan found 70 seashells.
- She has 27 seashells left.
- Therefore, Joan gave away 70 27 = 43 seashells.

answer3 = "43"

4. Final Answer:

The final answer is "43," indicating how many shells Joan gave to others.

Final Output: 43

Table 19: The output of executing the pseudo-program knowledge graph in Table 18.

#Given a program text, your role is the LLM function, which has three parameters: case, question, and description. The description is textual, the question requires you to answer based on the description, and the case involves processed scenarios or supplementary information or related considerations to assist you in answering the question. You need to strictly follow the program logic to execute. During execution, the LLM function will be called multiple times, which means you will answer corresponding questions based on different cases and descriptions. Please output the final result of this program text in natural language.

### def LLM(case, question, description):

LLM\_answer=Given the question, the case is an example about the question. Please study this example and answer the question based on the description.

return LLM\_answer

sub\_question1="Who was the member of the 'Mother Love Bone' band who passed away before the release of the album 'Apple'?"

sub\_case1=" Please focus on the correlation between the death of band members and the upcoming release of the album."

sub\_description1="Frontman Andrew Wood's personality and compositions helped to catapult the group to the top of the burgeoning late 1980s/early 1990s Seattle music scene. Wood died only days before the scheduled release of the band's debut album, "Apple", thus ending the group's hopes of success. The album was finally released a few months later."

answer1=LLM(sub\_case1,sub\_question1,sub\_description1)

sub\_question2="What band has Andrew Wood joined (except for 'Mother Love Bone')?"
sub\_case2="Please pay attention to the description of the person's band experience."

sub\_description2="Return to Olympus is the only album by the alterna tive rock band Malfunkshun. It was released after the band had broken up and after lead singer Andrew Wood (later of Mother Love Bone) had died of a drug overdose in 1990."

answer2=LLM(sub\_case2,sub\_question2,sub\_description2)

sub\_question3="What band was Andrew Wood in before joining Mother Love Bone?"

sub\_case3="Please pay attention to the order in which the person joined the band."

#Answering sub\_question3 requires relying on the answers of sub\_question1 and sub\_question2, namely answer1 and answer2

sub\_description3=answer1+answer2

answer3=LLM(sub\_case3,sub\_question3,sub\_description3)

#The final answer final\_answer=answer3 print(final\_answer)

Table 20: The updated pseudo-program knowledge graph on HotpotQA. (Question: What was the former band of the member of Mother Love Bone who died just before the release of "Apple"? Paragraph A, Return to Olympus: Return to Olympus is the only album by the alternative rock band Malfunkshun. It was released after the band had broken up and after lead singer Andrew Wood (later of Mother Love Bone) had died of a drug overdose in 1990. Stone Gossard, of Pearl Jam, had compiled the songs and released the album on his label, Loosegroove Records. Paragraph B, Mother Love Bone: Mother Love Bone was an American rock band that formed in Seattle, Washington in 1987. The band was active from 1987 to 1990. Frontman Andrew Wood's personality and compositions helped to catapult the group to the top of the burgeoning late 1980s/early 1990s Seattle music scene. Wood died only days before the scheduled release of the band's debut album, "Apple", thus ending the group's hopes of success. The album was finally released a few months later.)

Let's break down the program and follow the steps logically to determine the final answer.

Step-by-Step Breakdown:

- 1. First Call to LLM (sub\_question1):
- Question: Who was the member of the "Mother Love Bone" band who passed away before the release of the album "Apple"?
- Case: The case asks us to focus on the correlation between the death of a band member and the upcoming release of the album.
- Description: The description mentions that the frontman, Andrew Wood, passed away just days before the release of "Apple," which was the band's debut album.

Answer: The band member who passed away before the release of the album "Apple" is Andrew Wood.

answer1 = "Andrew Wood"

- 2. Second Call to LLM (sub\_question2): Question: What band has Andrew Wood joined (except for "Mother Love Bone")?
- Case: The case tells us to pay attention to the description of the person's band experience.
- Description: The description talks about Andrew Wood being the lead singer of a band called Malfunkshun, which released the album "Return to Olympus" after the band broke up. This band existed before Andrew Wood joined "Mother Love Bone."

Answer: Before joining "Mother Love Bone," Andrew Wood was part of the band Malfunkshun.

answer2 = "Malfunkshun"

- 3. Third Call to LLM (sub\_question3):
- Question: What band was Andrew Wood in before joining Mother Love Bone?
- Case: This call asks to pay attention to the order in which Andrew Wood joined bands.
- Description: The description includes the answers from previous questions. From answer1, we know Andrew Wood was in "Mother Love Bone." From answer2, we know he was in "Malfunkshun" before joining "Mother Love Bone."

Answer: Before joining "Mother Love Bone," Andrew Wood was in "Malfunkshun."

answer3 = "Malfunkshun"

4. Final Answer:

The final answer is "Malfunkshun."

Final Output: Malfunkshun

Table 21: The output of executing the pseudo-program knowledge graph in Table 20.

#Given a program text, your role is the LLM function, which has three parameters: case, question, and description. The description is textual, the question requires you to answer based on the description, and the case involves processed scenarios or supplementary information or related considerations to assist you in answering the question. You need to strictly follow the program logic to execute. During execution, the LLM function will be called multiple times, which means you will answer corresponding questions based on different cases and descriptions. Please output the final result of this program text in natural language.

### def LLM(case,question,description):

LLM\_answer=Given the question, the case is an example about the question. Please study this example and answer the question based on the description.

return LLM\_answer

sub\_question1="What are the general characteristics of the area represented by each option in terms of housing distribution and isolation?"

sub\_case1="

Neighborhood: The houses in the neighborhood are relatively compact, and there may be a mix of different types of houses. The boundaries are usually not particularly closed, and many houses in the neighborhood may not have obvious isolation facilities between them. They are closely connected to the surrounding streets and other buildings. Subdivision: Zoning is usually a planned area, and some zones may have a unified architectural style and management regulations. Some high-end partitions may have access control facilities to ensure residents' privacy and safety, isolating themselves from the outside world.

City: A city is a large area that includes various functional areas and a wide variety of housing types. Most houses in the city are open plan, closely connected to various commercial facilities, public transportation, etc., and few houses can be completely isolated from other places.

Suburbs: Houses in the suburbs are relatively scattered, but generally do not have strong isolation. The suburbs are mainly for residential purposes, although there may be some independent residences, they are not primarily characterized by isolation.

Street: The street is mainly a transportation channel, and the houses on both sides of the street are mainly distributed along the road, usually without overall isolation measures."

```
sub_description1="Answer Choices: (A) neighborhood (B) subdivision (C) city (D) suburbs (E) street"
answer1=LLM(sub_case1,sub_question1,sub_description1)
```

sub\_question2="Which area is most likely to provide houses with access control facilities and isolation from other areas?"

sub\_case2="Pay attention to comparing the characteristics of each region"

#Answering sub\_question2 requires relying on the answers of sub\_question1, namely answer1

sub\_description2=answer1

answer2=LLM(sub\_case2,sub\_question2,sub\_description2)

#The final answer final\_answer=answer2 print(final\_answer)

Table 22: The updated pseudo-program knowledge graph on CSQA, where *sub-case* is derived from the LLM-decomposed *Knowledge case* of the decision tree built by experts. (Question: He wanted a house that was gated off from other places, where should he start looking? Answer Choices: (A) neighborhood (B) subdivision (C) city (D) suburbs (E) street)

Let's walk through the steps of the program and how it works logically to produce the final result.

Step-by-Step Breakdown:

- 1. First Call to LLM (sub\_question1):
- Question: What are the general characteristics of the area represented by each option in terms of housing distribution and isolation?
- Case: The case provides descriptions of different areas: neighborhood, subdivision, city, suburbs, and street. Each description highlights the characteristics of housing distribution and isolation.
- Description: The answer choices include (A) neighborhood, (B) subdivision, (C) city, (D) suburbs, and (E) street.

Analysis: The focus of the question is on the characteristics of each option in terms of housing distribution and isolation.

- Neighborhood: Houses are relatively compact with no clear isolation between them.
- Subdivision: Has unified architectural styles and may have access control facilities and isolation, especially in high-end areas.
- City: Houses are open plan and closely connected to other buildings and commercial areas; few houses have isolation.
- Suburbs: Houses are scattered but generally lack strong isolation.
- Street: Houses are distributed along the road and generally have no isolation.

Answer: Based on the description, the "subdivision" is most likely to provide houses with access control facilities and isolation from other areas, especially in high-end zones.

answer1 = "B" (subdivision)

- 2. Second Call to LLM (sub\_question2):
- Question: Which area is most likely to provide houses with access control facilities and isolation from other areas?
- Case: The case asks to compare the characteristics of each area to determine which is most likely to have access control and isolation.
- Description: The description relies on the answer from sub\_question1, which identified the subdivision as the area most likely to have access control facilities and isolation.

Answer: The subdivision is the most likely area to provide houses with access control and isolation based on the characteristics described earlier.

answer2 = "B" (subdivision)

3. Final Answer:

The final answer is "B", corresponding to \*\*subdivision\*\*.

Final Output: B

Table 23: The output of executing the pseudo-program knowledge graph in Table 22.

The description is textual, the question requires you to answer based on the description, and the case involves processed scenarios or supplementary information or related considerations to assist you in answering the question. You need to strictly follow the program logic to execute. During execution, the LLM function will be called multiple times, which means you will answer corresponding questions based on different cases and descriptions. Please output the final result of this program text in natural language. def LLM(case, question, description): LLM answer=Given the question, the case is an example about the question. Please study this example and answer the question based on the description. return LLM answer sub\_question1="Is this behavior directly related to understanding and answering the teacher's questions?" sub\_case1=" Asking for a gold star is a behavior where one may expect to receive a reward after answering a question, and it is not something that needs to be done before understanding and answering the question." sub\_description1="Behavior: ask for a gold star" answer1=LLM(sub\_case1,sub\_question1,sub\_description1) sub\_question2="Is this behavior helpful for understanding and answering the teacher's questions?" sub\_case2="Skipping class means not being in the classroom and not being able to hear the teacher's questions" sub description2="Behavior: skip her class answer2=LLM(sub\_case2,sub\_question2,sub\_description2) sub\_question3="Does this behavior help to understand and answer the teacher's questions?" sub\_case3="Only when a person knows the information related to the problem, can they understand the content of the teacher's question well and accurately answer the question based on the information they have mastered" sub\_description3="Behavior: know the information" answer3=LLM(sub\_case3,sub\_question3,sub\_description3) sub\_question4="Which behavior is most likely to occur?" sub\_case4="Please choose the behavior that best fits your understanding and answers the teacher's question in this scenario." #Answering sub\_question4 requires relying on the answers of sub\_question1, sub\_question2 and sub\_question3, namely answer1, answer2 and answer3 sub description4=answer1+answer2+answer3 answer4=LLM(sub\_case3,sub\_question3,sub\_description3) #The final answer final\_answer=answer4

#Given a program text, your role is the LLM function, which has three parameters: case, question, and description.

Table 24: The updated pseudo-program knowledge graph on SIQA. (Question: The teacher asked the class a question and they seemed puzzled. Aubrey understood the question well and answered. What does Aubrey need to do before this? Answer Choices: (A) ask for a gold star (B) skip her class (C) know the information)

print(final\_answer)

Let's break down the logic and the steps for this program:

Step-by-Step Breakdown:

- 1. First Call to LLM (sub\_question1):
- Question: Is this behavior directly related to understanding and answering the teacher's questions?
- Case: Asking for a gold star is a behavior that does not directly contribute to understanding the question. It's more about receiving a reward after the question is answered.
- Description: The behavior is to "ask for a gold star."

Answer: Since asking for a gold star doesn't directly relate to understanding or answering the question, the answer is "No."

- 2. Second Call to LLM (sub\_question2):
- Question: Is this behavior helpful for understanding and answering the teacher's questions?
- Case: Skipping class means not being present to hear the teacher's questions, which would make it impossible to understand or answer the question.
- Description: The behavior is "skip her class."

Answer: Since skipping class means missing the teacher's questions, the answer is "No."

- 3. Third Call to LLM (sub\_question3):
- Question: Does this behavior help to understand and answer the teacher's questions?
- Case: Knowing the information related to the problem is essential to understanding the teacher's question and answering it accurately.
- Description: The behavior is "know the information."

Answer: Knowing the information is crucial for understanding and answering the question, so the answer is "Yes."

- 4. Fourth Call to LLM (sub\_question4):
- Question: Which behavior is most likely to occur?
- Case: This question asks to select the behavior that best fits the context of understanding and answering the teacher's question.
- Description: The description includes the answers from the previous questions: "No," "No," and "Yes."

Answer: The behavior most likely to occur is the one that helps understand and answer the teacher's question, which is knowing the information ("Yes").

Final Answer: know the information

Table 25: The output of executing the pseudo-program knowledge graph in Table 24.

processed scenarios or supplementary information or related considerations to assist you in answering the question. You need to strictly follow the program logic to execute. During execution, the LLM function will be called multiple times, which means you will answer corresponding questions based on different cases and descriptions. Please output the final result of this program text in natural language. def LLM(case,question,description): LLM\_answer=Given the question, the case is an example about the question. Please study this example and answer the question based on the description. return LLM\_answer sub\_question1="What is the last letter of this word?" sub case1=' Description: Pretty Answer: The last letter of 'Pretty' is' y '." sub\_description1="Gavin" answer1=LLM(sub\_case1,sub\_question1,sub\_description1) sub\_question2="What is the last letter of this word?" sub\_case2=' Description: Jada Answer: The last letter of 'Jada' is' a '." sub\_description2="Neha" answer2=LLM(sub\_case2,sub\_question2,sub\_description2) sub\_question3="What is the last letter of this word?" sub\_case3=" Description: Sarita Answer: The last letter of 'Sarita' is' a '." sub\_description3="Asha" answer3=LLM(sub\_case3,sub\_question3,sub\_description3) sub\_question4="What is the last letter of this word?" sub\_case4=" Description: Allen Answer: The last letter of 'Allen' is' n'." sub description4="Baltazar" answer4=LLM(sub\_case4,sub\_question4,sub\_description4) sub\_question5="What is the result of concatenating the last letter obtained from each word in order?" sub\_case5=" Description: The last letter of 'Pretty' is' y', the last letter of 'Jada' is' a', the last letter of 'Sarita' is' a' and the last letter of 'Allen' is' n'. Answer: The last letter of the four words is' y',' a',' n', and when they are concatenated in order, the result is' yaan'.' #Answering sub\_question5 requires relying on the answers of sub\_question1,sub\_question2,sub\_question3 and sub\_question4, namely answer1, answer2, answer3 and answer4 sub\_description5=answer1+answer2+answer3+answer4 answer5=LLM(sub\_case5,sub\_question5,sub\_description5)

#Given a program text, your role is the LLM function, which has three parameters: case, question, and description. The description is textual, the question requires you to answer based on the description, and the case involves

Table 26: The updated pseudo-program knowledge graph on Last Letter Concatenation. (Question: Take the last letters of each words in Gavin Neha Asha Baltazaränd concatenate them.)

#The final answer final\_answer=answer5 print(final\_answer)

```
Step-by-Step Breakdown:
1. First Call to LLM (sub_question1):
- Question: What is the last letter of this word?
- Case: The description indicates that the last letter of "Pretty" is "y".
- Description: The word is "Gavin".
Answer: The last letter of "Gavin" is "n".
answer1 = "n"
2. Second Call to LLM (sub_question2):
- Question: What is the last letter of this word?
- Case: The description indicates that the last letter of "Jada" is "a".
- Description: The word is "Neha".
Answer: The last letter of "Neha" is "a".
answer2 = "a"
3. Third Call to LLM (sub_question3):
- Question: What is the last letter of this word?
- Case: The description indicates that the last letter of "Sarita" is "a".
- Description: The word is "Asha".
Answer: The last letter of "Asha" is "a".
answer3 = "a"
4. Fourth Call to LLM (sub_question4):
- Question: What is the last letter of this word?
- Case: The description indicates that the last letter of "Allen" is "n".
- Description: The word is "Baltazar".
Answer: The last letter of "Baltazar" is "r".
answer4 = "r"
5. Fifth Call to LLM (sub_question5):
- Question: What is the result of concatenating the last letter obtained from each word in order?
- Case: The case gives the last letters of the words: "n" (from Gavin), "a" (from Neha), "a" (from Asha), and "r"
(from Baltazar).
- Description: The concatenation of the last letters "n", "a", "a", and "r" results in "naar".
Answer: The concatenated result of the last letters is "naar".
answer5 = "naar"
Final Answer: naar
```

Let's break down the program step by step to determine the final output:

Table 27: The output of executing the pseudo-program knowledge graph in Table 26.

```
def LLM(case, question, description):
  LLM answer=Given the question, the case is an example about the question. Please study this example and
answer the question based on the description.
  return LLM answer
sub_question1="What is the initial state of the coin?"
sub_case1="Pay attention to the initial state of the coin at this moment."
sub_description1="A coin is heads up"
answer1=LLM(sub_case1,sub_question1,sub_description1)
sub_question2="What is the current state of the coin?"
sub_case2="
Pay attention to whether this operation changes the original state of the coin and outputs the state after this operation."
#Answering sub_question2 requires relying on the answers of sub_question1, namely answer1
sub_description2=answer1+"Sager does not flip the coin"
answer2=LLM(sub_case2,sub_question2,sub_description2)
sub_question3="What is the current state of the coin?"
sub_case3="
Pay attention to whether this operation changes the original state of the coin and outputs the state after this operation"
#Answering sub_question3 requires relying on the answers of sub_question2, namely answer2
sub_description3=answer2+"Zyheir flips the coin"
answer3=LLM(sub_case3,sub_question3,sub_description3)
sub_question4="Is the coin still heads up?"
sub_case4="Pay attention to the coin's transition at this moment"
#Answering sub_question4 requires relying on the answers of sub_question3, namely answer3
sub_description4=answer3
answer4=LLM(sub_case4,sub_question4,sub_description4)
#The final answer
final_answer=answer4
print(final_answer)
```

Table 28: The updated pseudo-program knowledge graph on Coin Flip, where *sub-case* is derived from the LLM-decomposed *Knowledge case* of the decision tree built by experts. (Question: A coin is heads up. sager does not flip the coin. zyheir flips the coin. Is the coin still heads up?)

Let's analyze the program step-by-step and follow the execution logic to determine the final output:

Step-by-Step Breakdown:

- 1. First Call to LLM (sub\_question1):
- Question: What is the initial state of the coin?
- Case: Pay attention to the initial state of the coin at this moment.
- Description: The coin is heads up.

Answer: Based on the description, the initial state of the coin is "heads up."

answer1 = "heads up"

- 2. Second Call to LLM (sub\_question2):
- Question: What is the current state of the coin?
- Case: Pay attention to whether this operation changes the original state of the coin and output the state after the operation.
- Description: The description includes the answer from 'sub\_question1' ("heads up") and adds the information that Sager does not flip the coin.

Answer: Since Sager does not flip the coin, the state remains the same as "heads up."

answer2 = "heads up"

- 3. Third Call to LLM (sub\_question3):
- Question: What is the current state of the coin?
- Case: Pay attention to whether this operation changes the original state of the coin and output the state after the operation.
- Description: The description includes the answer from 'sub\_question2' ("heads up") and adds the information that Zyheir flips the coin.

Answer: Since Zyheir flips the coin, the state changes from "heads up" to "tails up."

answer3 = "tails up"

- 4. Fourth Call to LLM (sub\_question4):
- Question: Is the coin still heads up?
- Case: Pay attention to the coin's transition at this moment.
- Description: The description includes the answer from 'sub\_question3' ("tails up"), so we check if the coin is still heads up

Answer: Since the coin is "tails up" (after Zyheir flipped it), it is no longer heads up.

answer4 = "No"

5. Final Answer:

The final answer is "No" because after the flip, the coin is not heads up.

Table 29: The output of executing the pseudo-program knowledge graph in Table 28.

# def LLM(case,question,description):

LLM\_answer=Given the question, the case is an example about the question. Please study this example and answer the question based on the description.

return LLM\_answer

sub\_question1="What is the amount of cash stolen in this theft?"

sub\_case1="Pay attention to accurately finding sentences that describe the theft and the corresponding amount of stolen cash, to avoid confusion with related descriptions of other theft events."

sub\_description1="At about 11:00 on January 2, 2016, the defendant Liu XX, Song XX, and Hu Mou took the Ouyang special bus and passed a certain section of a certain town in Jianyang District, Nanping City. Hu Mou and Liu XX watched, and Song XX used a blade to cut the right pocket of the victim Liu Mou's coat and stole the small yellow envelope containing 1,500 yuan in cash."

answer1=LLM(sub\_case1,sub\_question1,sub\_description1)

sub\_question2="What is the amount of cash stolen in this theft?"

sub\_case2="Pay attention to accurately finding sentences that describe the theft and the corresponding amount of stolen cash, to avoid confusion with related descriptions of other theft events."

sub\_description2="At about 11:00 on January 20, 2016, the defendants Liu XX, Song XX, and Hu Mou took the Ouyang special bus to a certain section of the road. They saw the victim Zhang Mou 1 sleeping, so Hu Mou and Liu XX watched, and Song XX cut open Zhang Mou 1's coat pocket with a blade, and stole 7,000 yuan in cash tied with cow tendons in the pocket."

answer2=LLM(sub\_case2,sub\_question2,sub\_description2)

sub\_question3="What is the total amount involved in all theft crimes?"

sub\_case3="Pay attention to performing correct addition operations on all previous results to ensure that there are no errors in the calculation process, and do not omit or double calculate the amount."

#Answering sub\_question3 requires relying on the answers of sub\_question1 and sub\_question2, namely answer1 and answer2

sub\_description3=answer1+answer2

answer3=LLM(sub\_case3,sub\_question3,sub\_description3)

#The final answer final\_answer=answer3 print(final\_answer)

Table 30: The updated pseudo-program knowledge graph on LawBench. (Question: What is the total amount of crime involved in the document? Document: After investigation, it was found that during January 2016, the defendant Liu XX, Song XX, and Hu (both of whom have been sentenced) agreed to commit theft on the Ouyang special bus, and agreed that whoever did it would get half of the stolen property, and the other two would share the other half. The specific criminal facts are as follows: 1. At about 11:00 on January 2, 2016, the defendant Liu XX, Song XX, and Hu Mou took the Ouyang special bus and passed a certain section of a certain town in Jianyang District, Nanping City. Hu Mou and Liu XX watched, and Song XX used a blade to cut the right pocket of the victim Liu Mou's coat and stole the small yellow envelope containing 1,500 yuan in cash. After getting off the bus, Song XX got 750 yuan of the stolen money, and Hu Mou and Liu XX each got 375 yuan of the stolen money. 2. At about 11:00 on January 20, 2016, the defendants Liu XX, Song XX, and Hu Mou took the Ouyang special bus to a certain section of the road. They saw the victim Zhang Mou 1 sleeping, so Hu Mou and Liu XX watched, and Song XX cut open Zhang Mou 1's coat pocket with a blade, and stole 7,000 yuan in cash tied with cow tendons in the pocket. Afterwards, Song XX received 3,500 yuan, and Hu Mou and Liu XX each received 1,750 yuan.)

To solve this, let's break down the program logic and how it works step by step:

Step-by-step breakdown:

- 1. First Call to LLM (sub\_question1):
- Question: What is the amount of cash stolen in this theft?
- Case: Focus on accurately finding the amount of stolen cash in the description, ensuring no confusion with other theft events.
- Description: The theft took place on January 2, 2016, where Song XX cut open the victim's coat pocket and stole a small yellow envelope containing 1,500 yuan.

Answer: Based on the description, the amount of cash stolen is 1,500 yuan.

answer1 = 1500

- 2. Second Call to LLM (sub\_question2):
- Question: What is the amount of cash stolen in this theft?
- Case: Again, accurately find the amount of stolen cash from the description, ensuring no confusion with other theft events.
- Description: On January 20, 2016, Song XX stole 7,000 yuan from the victim Zhang Mou 1's coat pocket. Answer: Based on the description, the amount of cash stolen is 7,000 yuan.

answer2 = 7000

- 3. Third Call to LLM (sub\_question3):
- Question: What is the total amount involved in all theft crimes?
- Case: Perform the addition operation correctly to avoid errors and omissions, using the amounts obtained from the previous answers.
- Description: The description contains both 'answer1' (1,500 yuan) and 'answer2' (7,000 yuan).

Answer: To calculate the total amount, we add 'answer1' and 'answer2':

1,500 + 7,000 = 8,500 yuan.

answer3 = 8500

4. Final Answer: The total amount involved in all theft crimes is 8,500 yuan.

Table 31: The output of executing the pseudo-program knowledge graph in Table 30.

```
def LLM(case, question, description):
```

LLM\_answer=Given the question, the case is an example about the question. Please study this example and answer the question based on the description.

return LLM\_answer

sub\_question1="Does the contract mention any content regarding usage permissions?"

sub\_ase1="Pay attention to carefully reading the text and find the guiding statements related to content usage permissions."

sub\_description1="See the Permission to use your content section for more about your rights in your content, and how your content is used in our services."

answer1=LLM(sub\_case1,sub\_question1,sub\_description1)

sub\_question2="Does the contract mention the removal of content?"

sub\_case2="Pay attention to the guidance prompts in the text, accurately locate and remove the relevant parts of the content, and do not confuse them with other parts."

sub\_description2="See the Removing your content section to learn why and how we might remove user-generated content from our services."

answer2=LLM(sub\_case2,sub\_question2,sub\_description2)

sub\_question3="How is intellectual property infringement stipulated in the contract?"

sub\_case3="Pay attention to the relevant statements regarding the handling measures after intellectual property infringement, and understand the prescribed infringement handling process and methods."

#Answering sub\_question3 requires relying on the answers of sub\_question1 and sub\_question2, namely answer1 and answer2

sub\_description3=answer1+answer2+"If you think that someone is infringing your intellectual property rights, you can send us notice of the infringement and well take appropriate action. For example, we suspend or close the Google Accounts of repeat copyright infringers as described in our Copyright Help Centre."

answer3=LLM(sub\_case3,sub\_question3,sub\_description3)

sub\_question4="Will Google help me deal with situations where content is used without permission?"
sub\_case4="Pay attention to accurately understanding Google's measures for handling infringement in sub\_question3."

#Answering sub\_question4 requires relying on the answers of sub\_question3, namely answer3 sub\_description4=answer3

answer4=LLM(sub\_case4,sub\_question4,sub\_description4)

#The final answer final\_answer=answer4 print(final\_answer)

Table 32: The updated pseudo-program knowledge graph on LegalBench. (Question: Will Google help me if I think someone has taken and used content Ive created without my permission? Contract: Some of our services give you the opportunity to make your content publicly available. For example, you might post a product or restaurant review that you wrote, or you might upload a blog post that you created. See the Permission to use your content section for more about your rights in your content, and how your content is used in our services See the Removing your content section to learn why and how we might remove user-generated content from our services If you think that someone is infringing your intellectual property rights, you can send us notice of the infringement and well take appropriate action. For example, we suspend or close the Google Accounts of repeat copyright infringers as described in our Copyright Help Centre.)

Let's go through the program step by step, answering each question based on the cases and descriptions provided:

## Step-by-step Execution:

sub\_question1

Question: Does the contract mention any content regarding usage permissions?

Case: "Pay attention to carefully reading the text and find the guiding statements related to content usage permissions."

Description: "See the Permission to use your content section for more about your rights in your content, and how your content is used in our services."

Analysis for sub\_question1: The description directly refers to a "Permission to use your content" section, which indicates that usage permissions are mentioned in the contract.

Answer1: "Yes, the contract mentions content usage permissions."

## sub\_question2

Question: Does the contract mention the removal of content?

Case: "Pay attention to the guidance prompts in the text, accurately locate and remove the relevant parts of the content, and do not confuse them with other parts."

Description: "See the Removing your content section to learn why and how we might remove user-generated content from our services."

Analysis for sub\_question2: The description references a "Removing your content" section, indicating that the contract does address the removal of content.

Answer2: "Yes, the contract mentions the removal of content."

#### sub\_question3

Question: How is intellectual property infringement stipulated in the contract?

Case: "Pay attention to the relevant statements regarding the handling measures after intellectual property infringement, and understand the prescribed infringement handling process and methods."

Description: This description relies on the answers to sub\_question1 and sub\_question2. Combining those with the following statement: "If you think that someone is infringing your intellectual property rights, you can send us notice of the infringement and we'll take appropriate action. For example, we suspend or close the Google Accounts of repeat copyright infringers as described in our Copyright Help Centre."

Analysis for sub\_question3: The description outlines the process for handling intellectual property infringement, including the option to report infringement and the potential actions (like suspension of accounts). This answers how intellectual property infringement is dealt with in the contract.

Answer3: "The contract specifies that intellectual property infringement is handled by allowing users to report infringement and take actions such as suspending or closing accounts of repeat offenders."

sub\_question4 Question: Will Google help me deal with situations where content is used without permission? Case: "Pay attention to accurately understanding Google's measures for handling infringement in sub\_question3." Description: This relies on answer3: "The contract specifies that intellectual property infringement is handled by allowing users to report infringement and take actions such as suspending or closing accounts of repeat offenders." Analysis for sub\_question4: The description indicates that Google provides help in dealing with content usage without permission through its procedures for reporting and addressing infringement.

Answer4: "Yes, Google will help you deal with situations where content is used without permission by taking actions such as suspending or closing accounts of repeat infringers."

# Final Answer:

The final answer, based on the previous analysis, is:

"Yes, Google will help you deal with situations where content is used without permission by taking actions such as suspending or closing accounts of repeat infringers."

Table 33: The output of executing the pseudo-program knowledge graph in Table 32.

You need to strictly follow the program logic to execute. During execution, the LLM function will be called multiple times, which means you will answer corresponding questions based on different cases and descriptions. Please output the final result of this program text in natural language. def LLM(case, question, description): LLM answer=Given the question, the case is an example about the question. Please study this example and answer the question based on the description. return LLM\_answer sub\_question1="Is the description related to the question 'Can I still buy new energy vehicles'?" sub\_case1="Market inquiry refers to inquiries about the overall situation of the securities market. This includes broad market indices like the Shanghai Composite Index or Shenzhen Component Index, which reflect the performance and conditions of the entire stock market." sub\_description1="Market inquiry" answer1=LLM(sub\_case1,sub\_question1,sub\_description1) sub\_question2="Is the description related to the question 'Can I still buy new energy vehicles'?" sub\_case2="Industry sectors refer to the various parts of the securities market that are classified by industry, such as the financial sector, technology sector, new energy sector, etc. Industry sector inquiries mainly focus on the overall situation of a specific industry, including industry development trends, competition patterns within the industry, the impact of policies on the industry, and the industry's valuation in the market, among many other aspects." sub\_description2="Industry sector inquiry" answer2=LLM(sub\_case2,sub\_question2,sub\_description2) sub\_question3="Is the description related to the question 'Can I still buy new energy vehicles'?" sub\_case3="Individual stock inquiry refers to inquiries related to the stocks of a specific company. This may include questions about the company's financial health, performance, management changes, and stock price trends, as well as technical analysis related to stock price movements." sub description3="Individual stock inquiry answer3=LLM(sub\_case3,sub\_question3,sub\_description3) sub\_question4="Is the description related to the question 'Can I still buy new energy vehicles'?" sub\_case4="Fund inquiry pertains to inquiries about investment funds, which pool money from multiple investors for collective investment. This includes questions about the fund's type (e.g., equity, bond, or money market funds), investment strategy, performance, risk characteristics, management team, and subscription or redemption rules." sub\_description4="Fund inquiry" answer4=LLM(sub\_case4,sub\_question4,sub\_description4) sub\_question5="Is the description related to the question 'Can I still buy new energy vehicles'?" sub\_case5="Customer service inquiry refers to inquiries made by investors to customer service departments of financial institutions regarding issues related to securities trading. These can involve problems with account operations (e.g., account opening, password retrieval), transaction rules (e.g., fees, time), product services, and handling of transaction disputes.' sub\_description5="Customer service inquiry" answer5=LLM(sub\_case5,sub\_question5,sub\_description5) sub\_question6="Which one best fits the intention of the question 'Can I still buy new energy vehicles'?" sub\_case6="Pay attention to comprehensively considering the previous analysis results to ensure accurate and error free judgment, without missing any key information." #Answering sub\_question6 requires relying on the answers of sub\_question1, sub\_question2, sub\_question3, sub question4, and sub question5, namely answer1, answer2, answer3, answer4 and answer5 sub\_description6=answer1+answer2+answer3+answer4+answer5 answer6=LLM(sub\_case6,sub\_question6,sub\_description6) #The final answer final\_answer=answer6 print(final\_answer)

#Given a program text, your role is the LLM function, which has three parameters: case, question, and description. The description is textual, the question requires you to answer based on the description, and the case involves processed scenarios or supplementary information or related considerations to assist you in answering the question.

Table 34: The updated pseudo-program knowledge graph on CFBenchmark. (Question: You are an intention emotion assistant. Please analyze the intention of the following query: [market inquiry, industry sector inquiry, individual stock inquiry, fund inquiry, customer service inquiry]? Query: Can I still buy new energy vehicles?

Let's go through the execution of the program step-by-step, based on the provided logic and descriptions:

#### Step-by-step Execution:

sub question1

Question: Is the description related to the question 'Can I still buy new energy vehicles'?

Case: "Market inquiry refers to inquiries about the overall situation of the securities market. This includes broad market indices like the Shanghai Composite Index or Shenzhen Component Index, which reflect the performance and conditions of the entire stock market. Description: "Market inquiry'

Analysis for sub-question1: The description focuses on market indices and the securities market, which is unrelated to the question about buying new energy vehicles. The question is focused on vehicle availability, not the stock market. Answer1: "No, this description is not related to the question."

# sub\_question2

Question: Is the description related to the question 'Can I still buy new energy vehicles'?

Case: "Industry sectors refer to the various parts of the securities market that are classified by industry, such as the financial sector, technology sector, new energy sector, etc. Industry sector inquiries mainly focus on the overall situation of a specific industry, including industry development trends, competition patterns within the industry, the impact of policies on the industry, and the industry's valuation in the market, among many other aspects.

Description: "Industry sector inquiry"

Analysis for sub\_question2: The description discusses the new energy sector, which is relevant to the question about new energy vehicles. While the focus is on the broader industry trends, it still relates to the availability of vehicles in the new energy sector. Answer2: "Yes, this description is related to the question."

#### sub question3

Question: Is the description related to the question 'Can I still buy new energy vehicles'?

Case: "Individual stock inquiry refers to inquiries related to the stocks of a specific company. This may include questions about the company's financial health, performance, management changes, and stock price trends, as well as technical analysis related to stock price movements. Description: "Individual stock inquiry"

Analysis for sub\_question3: The description focuses on individual company stocks, which is unrelated to the question about buying new energy vehicles. The question pertains to vehicle availability, not individual stock performance.

Answer3: "No, this description is not related to the question."

#### sub question4

Question: Is the description related to the question 'Can I still buy new energy vehicles'?

Case: "Fund inquiry pertains to inquiries about investment funds, which pool money from multiple investors for collective investment. This includes questions about the fund's type (e.g., equity, bond, or money market funds), investment strategy, performance, risk characteristics, management team, and subscription or redemption rules.

Description: "Fund inquiry"

Analysis for sub\_question4: The description concerns investment funds, which is not directly relevant to the question about buying new energy vehicles. The question is focused on vehicles, not investment funds.

Answer4: "No, this description is not related to the question."

# sub\_question5

Question: Is the description related to the question 'Can I still buy new energy vehicles'?

Case: "Customer service inquiry refers to inquiries made by investors to customer service departments of financial institutions regarding issues related to securities trading. These can involve problems with account operations (e.g., account opening, password retrieval), transaction rules (e.g., fees, time), product services, and handling of transaction disputes.'

Description: "Customer service inquiry

Analysis for sub\_question5: The description is about customer service inquiries related to financial services, which is unrelated to the question about buying new energy vehicles. Answer5: "No, this description is not related to the question."

# sub question6

Question: Which one best fits the intention of the question 'Can I still buy new energy vehicles'?

Case: "Pay attention to comprehensively considering the previous analysis results to ensure accurate and error-free judgment, without missing any key information."

Description: The description is a combination of the answers from sub\_question1, sub\_question2, sub\_question3, sub\_question4, and sub\_question5:

- "No, this description is not related to the question."
- "Yes, this description is related to the question."
- "No, this description is not related to the question."
- "No, this description is not related to the question."

"No, this description is not related to the question." Analysis for sub\_question6: The only description that is related to the question is from sub\_question2, which mentions the new energy sector. Therefore, this is the most relevant response

Answer6: "The description in sub\_question2 best fits the intention of the question."

The description "Industry sector inquiry" in sub\_question2, which relates to the new energy sector, best fits the intention of the question "Can I still buy new energy vehicles?'

Table 35: The output of executing the pseudo-program knowledge graph in Table 34.

def LLM(case, question, description):

LLM\_answer=Given the question, the case is an example about the question. Please study this example and answer the question based on the description.

return LLM\_answer

sub\_question1="Is this description consistent with the expression 'Ancient Greek philosophers said that life without reflection is worth less'?"

sub\_case1="It implies that introspection is a prerequisite for a life to have value. The phrase "Only after" emphasizes that without introspection, life has no value. In other words, it presents a rather absolute condition where the value of life is directly tied to the act of introspection and doesn't consider the possibility of value in other contexts that don't involve this specific act."

sub\_description1="Only after introspection can life be of value."
answer1=LLM(sub\_case1,sub\_question1,sub\_description1)

sub\_question2="Is this description consistent with the expression 'Ancient Greek philosophers said that life without reflection is worth less'?"

sub\_case2=" In order for life to possess value, it is necessary to engage in the act of reflection regularly. The use of "To...we must..." indicates a causal or purpose - driven relationship. It emphasizes that reflection is an important activity that one needs to do in order to achieve a valuable life, without completely ruling out that there might be other factors also involved in a valuable life."

sub\_description2="To be valuable in life, we must reflect on it from time to time."
answer2=LLM(sub\_case2,sub\_question2,sub\_description2)

sub\_question3="Is this description consistent with the expression 'Ancient Greek philosophers said that life without reflection is worth less'?"

sub\_case3=" It simply means that a person remains in a state of confusion throughout their entire life. There is no mention of the concepts of reflection or the value of life, as the focus is solely on the state of being confused." sub\_description3="I'm confused for a lifetime."

answer3=LLM(sub\_case3,sub\_question3,sub\_description3)

sub\_question4="Is this description consistent with the expression 'Ancient Greek philosophers said that life without reflection is worth less'?"

sub\_case4="This implies that the purpose of people's lives should be to gain understanding. It doesn't touch on the ideas of reflection and the worth of life as presented in the original maxim. Instead, it centers around the concept of living with the goal of understanding something, without any reference to the importance of reflecting on one's life for its value."

sub\_description4="People should live to understand."
answer4=LLM(sub\_case4,sub\_question4,sub\_description4)

sub\_question5="Which one best fits the expression 'Ancient Greek philosophers said that life without reflection is worth less'?"

sub\_case5="Pay attention to considering the previous analysis results comprehensively and select the best one. Ensure accurate judgment without missing any key information."

#Answering sub\_question6 requires relying on the answers of sub\_question1, sub\_question2, sub\_question3, and sub\_question4, namely answer1, answer2, answer3 and answer4 sub\_description5=answer1+answer2+answer3+answer4

answer5=LLM(sub\_case6,sub\_question6,sub\_description6)

#The final answer final\_answer=answer5 print(final\_answer)

Table 36: The updated pseudo-program knowledge graph on AGIEval. (Question: Ancient Greek philosophers said that life without reflection is worthless. Which of the following options is the closest to the meaning of this maxim? Answer Choices: (A) Only after introspection can life be of value. (B) To be valuable in life, we must

reflect on it from time to time. (C) I'm confused for a lifetime. (D) People should live to understand.)

Let's go through the execution step by step based on the provided program logic. We are essentially evaluating a series of descriptions to see how well they align with the expression "Ancient Greek philosophers said that life without reflection is worth less." Then, we summarize the findings to determine which description fits best with that expression.

Step-by-step Execution:

sub\_question1

Question: Is this description consistent with the expression 'Ancient Greek philosophers said that life without reflection is worth less'?

Case: "It implies that introspection is a prerequisite for a life to have value. The phrase 'Only after' emphasizes that without introspection, life has no value. In other words, it presents a rather absolute condition where the value of life is directly tied to the act of introspection and doesn't consider the possibility of value in other contexts that don't involve this specific act."

Description: "Only after introspection can life be of value."

Analysis for sub\_question1: The description directly correlates with the notion that introspection (or reflection) is essential for life to have value, aligning closely with the expression about life being less valuable without reflection.

Answer1: "Yes, this description is consistent with the expression."

#### sub question2

Question: Is this description consistent with the expression 'Ancient Greek philosophers said that life without reflection is worth less'?

Case: "In order for life to possess value, it is necessary to engage in the act of reflection regularly. The use of 'To...we must...' indicates a causal or purpose-driven relationship. It emphasizes that reflection is an important activity that one needs to do in order to achieve a valuable life, without completely ruling out that there might be other factors also involved in a valuable life."

Description: "To be valuable in life, we must reflect on it from time to time."

Analysis for sub\_question2: The description suggests that reflection is a necessary but not exclusive condition for life to have value. This is slightly different from the original expression, which implies reflection as the sole or primary requirement for value.

Answer2: "Yes, this description is consistent with the expression, but it introduces the possibility of other factors contributing to a valuable life."

# sub\_question3

Question: Is this description consistent with the expression 'Ancient Greek philosophers said that life without reflection is worth less'?

Case: "It simply means that a person remains in a state of confusion throughout their entire life. There is no mention of the concepts of reflection or the value of life, as the focus is solely on the state of being confused."

Description: "I'm confused for a lifetime."

Analysis for sub\_question3: The description focuses on confusion, with no mention of reflection or life's value. It is not aligned with the original expression.

Answer3: "No, this description is not consistent with the expression."

# sub\_question4

Question: Is this description consistent with the expression 'Ancient Greek philosophers said that life without reflection is worth less'?

Case: "This implies that the purpose of people's lives should be to gain understanding. It doesn't touch on the ideas of reflection and the worth of life as presented in the original maxim. Instead, it centers around the concept of living with the goal of understanding something, without any reference to the importance of reflecting on one's life for its value."

Description: "People should live to understand."

Analysis for sub\_question4: The description speaks about understanding as the goal of life, but it doesn't emphasize reflection or introspection as the path to valuing life, making it inconsistent with the original expression.

Answer4: "No, this description is not consistent with the expression."

# sub\_question5

Question: Which one best fits the expression 'Ancient Greek philosophers said that life without reflection is worth less'?

Case: "Pay attention to considering the previous analysis results comprehensively and select the best one. Ensure accurate judgment without missing any key information."

Description: The description is a combination of the answers from sub\_question1, sub\_question2, sub\_question3, and sub\_question4: "Yes, this description is consistent with the expression. Yes, this description is consistent with the expression, but it introduces the possibility of other factors contributing to a valuable life. No, this description is not consistent with the expression. No, this description is not consistent with the expression." Analysis for sub\_question5: The best match comes from sub\_question1 and sub\_question2, which both emphasize the importance of reflection in a valuable life. Sub\_question1 suggests a direct, strong connection between reflection and life's value, while sub\_question2 allows for other factors but still ties value to reflection.

Answer5: "The descriptions from sub\_question1 and sub\_question2 best fit the expression, but sub\_question1 aligns more directly with the original maxim."

# Final Answer

The description that best fits the expression "Ancient Greek philosophers said that life without reflection is worth less" is the one in sub\_question1, as it strongly links introspection to life's value without introducing any exceptions or other factors.

Table 37: The output of executing the pseudo-program knowledge graph in Table 36.

```
#include <iostream>
#include <string>
using namespace std;
string LLM(string case, string question, string description) {
  string LLM_answer = Given the question, the case is an example about the question. Please study this example
and answer the question based on the description;
  return LLM_answer;
int main() {
  string sub_question1 = "How many lunch trays can the person carry at once?";
  string sub_case1 = "Description:Tom is helping the restaurant staff collect dishes, but he can only carry 6 plates
at a time.
  Answer:Tom can carry 6 plates at a time.";
  string sub_description1 = "Roger can only carry 4 trays at a time.";
  string answer1 = LLM(sub_case1, sub_question1, sub_description1);
  string sub_question2 = "How many trays does the person need to pick up from the first table?";
  string sub_case2 = "Description:Tom needs to pick up 15 plates from one table and 9 plates from another.
  Answer:Tom needs to pick up 15 plates from the first table.";
  string sub_description2 = "Roger had to pick up 10 trays from one table and 2 trays from another.";
  string answer2 = LLM(sub_case2, sub_question2, sub_description2);
  string sub_question3 = "How many trays does the person need to pick up from the second table?";
  string sub_case3 = "Description:Tom needs to pick up 15 plates from one table and 9 plates from another.
  Answer:Tom needs to pick up 9 plates from the second table.";
  string sub_description3 = "Roger had to pick up 10 trays from one table and 2 trays from another.";
  string answer3 = LLM(sub_case3, sub_question3, sub_description3);
  string sub_question4 = "How many trays does the person need to carry in total?";
  string sub_case4 = "Description:Tom needs to pick up 15 plates from the first table and 9 plates from the second
  Answer: 15+9=24. Therefore Tom needs to carry 24 plates in total.";
// Answering sub_question4 requires relying on the answers of sub_question2 and sub_question3
  string sub_description4 = answer2 + answer3;
  string answer4 = LLM(sub_case4, sub_question4, sub_description4);
  string sub_question5 = "How many trips will the person need to make to carry all the trays?";
  string sub_case5 = "Description:Tom can carry 6 plates at a time and Tom needs to carry 24 plates in total.
  Answer: 24/6=4. Therefore Tom needs to carry 24 plates in total.";
// Answering sub_question5 requires relying on the answers of sub_question1 and sub_question4
  string sub_description5 = answer1 + answer4;
  string answer5 = LLM(sub_case5, sub_question5, sub_description5);
// The final answer
  string final_answer = answer5;
  cout « final_answer « endl;
  return 0;
```

Table 38: The updated pseudo-program knowledge graph in the form of C++ on MultiArith. (Question: Roger was helping the cafeteria workers pick up lunch trays, but he could only carry 4 trays at a time. If he had to pick up 10 trays from one table and 2 trays from another, how many trips will he make?)

Let's walk through the program to understand how the 'LLM' function is used to answer the questions and derive the final answer.

Breakdown of Execution:

- 1. First Call to LLM (sub\_question1):
- Question: How many lunch trays can the person carry at once?
- Case: Tom can carry 6 plates at a time.
- Description: Roger can carry 4 trays at a time.
- Answer: Since Roger can carry 4 trays at a time, the answer is simply 4.

answer1 = "4"

- 2. Second Call to LLM (sub\_question2):
- Question: How many trays does the person need to pick up from the first table?
- Case: Tom needs to pick up 15 plates from one table.
- Description: Roger needs to pick up 10 trays from one table.
- Answer: Roger needs to pick up 10 trays from the first table.

answer2 = "10"

- 3. Third Call to LLM (sub\_question3):
- Question: How many trays does the person need to pick up from the second table?
- Case: Tom needs to pick up 9 plates from another table.
- Description: Roger needs to pick up 2 trays from another table.
- Answer: Roger needs to pick up 2 trays from the second table.

answer3 = "2"

- 4. Fourth Call to LLM (sub\_question4):
- Question: How many trays does the person need to carry in total?
- Case: Tom needs to pick up 15 plates from one table and 9 plates from the second table.
- Description: The total number of trays that Roger needs to carry is the sum of the trays picked up from both tables. From the previous answers:
- 10 trays from the first table
- 2 trays from the second table
- Total = 10 + 2 = 12 trays

answer4 = "12"

- 5. Fifth Call to LLM (sub\_question5):
- Question: How many trips will the person need to make to carry all the trays?
- Case: Tom can carry 6 plates at a time and needs to carry 24 plates in total.
- Description: Roger can carry 4 trays at a time and needs to carry 12 trays in total. To calculate the number of trips, we divide the total number of trays by the number of trays he can carry at a time:
- -12/4 = 3 trips

answer5 = "3"

Final Answer: Roger needs to make 3 trips to carry all the trays.

Output: 3

Table 39: The output of executing the pseudo-program knowledge graph in Table 38.

```
public class Main {
  public static String LLM(String case, String question, String description) {
    String LLM_answer = Given the question, the case is an example about the question. Please study this example
and answer the question based on the description;
    return LLM_answer;
  }
  public static void main(String[] args) {
  String sub_question1 = "How many lunch trays can the person carry at once?";
  String sub_case1 = "Description:Tom is helping the restaurant staff collect dishes, but he can only carry 6 plates
  Answer:Tom can carry 6 plates at a time.";
  String sub_description1 = "Roger can only carry 4 trays at a time.";
  String answer1 = LLM(sub_case1, sub_question1, sub_description1);
  String sub_question2 = "How many trays does the person need to pick up from the first table?";
  String sub_case2 = "Description:Tom needs to pick up 15 plates from one table and 9 plates from another.
  Answer: Tom needs to pick up 15 plates from the first table.";
  String sub_description2 = "Roger had to pick up 10 trays from one table and 2 trays from another.";
  String answer2 = LLM(sub_case2, sub_question2, sub_description2);
  String sub_question3 = "How many trays does the person need to pick up from the second table?";
  String sub_case3 = "Description:Tom needs to pick up 15 plates from one table and 9 plates from another.
  Answer: Tom needs to pick up 9 plates from the second table.";
  String sub_description3 = "Roger had to pick up 10 trays from one table and 2 trays from another.";
  String answer3 = LLM(sub_case3, sub_question3, sub_description3);
  String sub_question4 = "How many travs does the person need to carry in total?":
  String sub_case4 = "Description:Tom needs to pick up 15 plates from the first table and 9 plates from the second
table.
  Answer: 15+9=24. Therefore Tom needs to carry 24 plates in total.";
  // Answering sub_question4 requires relying on the answers of sub_question2 and sub_question3
  String sub_description4 = answer2 + answer3;
  String answer4 = LLM(sub_case4, sub_question4, sub_description4);
  String sub_question5 = "How many trips will the person need to make to carry all the trays?";
  String sub_case5 = "Description:Tom can carry 6 plates at a time and Tom needs to carry 24 plates in total.
  Answer: 24/6=4. Therefore Tom needs to carry 24 plates in total.";
// Answering sub_question5 requires relying on the answers of sub_question1 and sub_question4
  String sub_description5 = answer1 + answer4;
  String answer5 = LLM(sub_case5, sub_question5, sub_description5);
// The final answer
  String string final_answer = answer5;
  System.out.println(final_answer);
```

Table 40: The updated pseudo-program knowledge graph in the form of **Java** on MultiArith. (Question: Roger was helping the cafeteria workers pick up lunch trays, but he could only carry 4 trays at a time. If he had to pick up 10 trays from one table and 2 trays from another, how many trips will he make?)

Let's break down the program step-by-step and understand how the 'LLM'function is used to answer the questions and derive the final result.

Breakdown of Execution:

- 1. First Call to LLM (sub\_question1):
- Question: How many lunch trays can the person carry at once?
- Case: Tom can carry 6 plates at a time.
- Description: Roger can carry 4 trays at a time.
- Answer: Based on the description, Roger can carry 4 trays at a time.

So, answer1 = "4 trays"

- 2. Second Call to LLM (sub\_question2):
- Question: How many trays does the person need to pick up from the first table?
- Case: Tom needs to pick up 15 plates from one table.
- Description: Roger needs to pick up 10 trays from one table and 2 trays from another.
- Answer: Roger needs to pick up 10 trays from the first table.

So, answer2 = "10 trays"

- 3. Third Call to LLM (sub\_question3):
- Question: How many trays does the person need to pick up from the second table?
- Case: Tom needs to pick up 15 plates from one table and 9 plates from another.
- Description: Roger needs to pick up 10 trays from one table and 2 trays from another.
- Answer: Roger needs to pick up 2 trays from the second table.

So, answer3 = "2 trays"

- 4. Fourth Call to LLM (sub\_question4):
- Question: How many trays does the person need to carry in total?
- Case: Tom needs to pick up 15 plates from the first table and 9 plates from the second table.
- Description: By adding up the trays from the first and second tables, we get the total number of trays Roger needs to carry:
- 10 trays from the first table
- 2 trays from the second table
- Total = 10 + 2 = 12 trays

So, answer4 = "12 trays"

- 5. Fifth Call to LLM (sub\_question5):
- Question: How many trips will the person need to make to carry all the trays?
- Case: Tom can carry 6 plates at a time and needs to carry 24 plates in total.
- Description: Roger can carry 4 trays at a time and needs to carry 12 trays in total. To calculate the number of trips, we divide the total number of trays by the number of trays he can carry at a time:
- -12/4 = 3 trips

So, answer5 = "3 trips"

Final Answer: Roger needs to make 3 trips to carry all the trays.

Output: 3 trips

Table 41: The output of executing the pseudo-program knowledge graph in Table 40.

You are tasked with generating new questions that share the same underlying reasoning logic as a given original question but differ in content. Follow these steps carefully:

- 1. **Input**: You will be provided with an original question from an open-source dataset.
- 2. **Task**: For each original question, generate **4 new questions** that meet the following criteria:
  - Same Reasoning Logic: The new questions must require the same reasoning steps or problem-solving approach as the original.
  - **Different Content**: The surface content (e.g., entities, scenarios, or context) should be altered to create distinct variations.
  - High Quality: Ensure the new questions are clear, grammatically correct, and logically valid.
- 3. **Output Format**: Return the generated questions in this structured format:

Original Question: [Insert original question here]

New Questions:

- a. [Generated question 1]
- b. [Generated question 2]
- c. [Generated question 3]
- d. [Generated question 4]

# 4. Example:

Original Question: "If a train travels 300 miles in 5 hours, what is its average speed?"

New Questions:

- a. "If a car travels 240 kilometers in 4 hours, what is its average speed?"
- b. "A cyclist covers 45 miles in 1.5 hours. What is their average speed?" c. "A plane flies 1,800 miles in 6 hours. Calculate its average speed."
- d. "A runner completes a 10-kilometer race in 50 minutes. What is their average speed in km/h?"

Table 42: The prompt for using LLMs to assist in generating datasets.