

Adaptive Prompting: A Dynamic Approach to Temporal Table Reasoning

Anonymous ACL submission

Abstract

Temporal Table Reasoning is a critical challenge for Large Language Models (LLMs), requiring effective prompting techniques to extract relevant insights. Despite existence of multiple prompting methods, their impact on table reasoning remains largely unexplored. Furthermore, the performance of these models varies drastically across different table and context structures, making it difficult to determine an optimal approach. This work investigates multiple prompting technique across diverse table types to determine optimal approaches for different scenarios. We find that performance varies based on *entity type*, *table structure*, *requirement of additional context* and *question complexity*, with "NO" single method consistently outperforming others. To mitigate these challenges, we introduce SEAR, an *adaptive prompting* framework inspired by human reasoning that dynamically adjusts based on context characteristics and integrates a structured reasoning. Our results demonstrate that SEAR achieves superior performance across all table types compared to other baseline prompting techniques. Additionally, we explore the impact of table structure refactoring, finding that a unified representation enhances models reasoning.

1 Introduction

Temporal table reasoning presents a unique challenge, requiring Large Language Models (LLMs) to not only interpret tabular data while understand the embedded temporal relationships. Unlike static tables that provide a fixed snapshot of structured information, temporal tables evolve over time, incorporating event sequences, timestamps and dynamic updates such as new columns and more. Reasoning over such structures essential for financial forecasting, historical trend analysis, medical diagnosis and event based decision making (Gupta et al., 2023; Xiong et al., 2024). However, existing LLMs often struggle with capturing these intricate temporal

| | | | | |
|---------------------|-----------|-------------|--|----------|
| Table 0 | | | Table Pretext - ... operating leases was \$ 100690000 , \$ 92710000 , and \$ 86842000 in fiscal 2006 , 2005 and 2004 ... | |
| Benefit Plan | 2017 | 2016 | | |
| Pension Plan | \$3,856 | \$3,979 | | |
| Health Plan | 11426 | 11530 | | |
| | | | | |
| Table_1 | | | | |
| | 2020 | Thereafter | | |
| Property mortgages | \$703,018 | \$1,656,623 | | amount |
| MRA facilities | — | — | 2007 | 56499000 |
| Sr. unsecured notes | 250000 | 100000 | 2008 | 46899000 |
| Ground leases | 31436 | 703254 | 2009 | 39904000 |
| | | | 2010 | 33329000 |

Question: what was the percentage change in total rental expense under operating leases from July 2 , 2005 to July 1 , 2006?
Dataset: FinQA
Req: Evidence from Text and PoT

| | | |
|--|------------------|------------------|
| Question: What is the sum of Ground leases of 2020, Health Plan of 2016, and Property mortgages and other loans of Thereafter ? Dataset: MultiHiertt Req: Evidence from multiple table and F-COT | | |
| Year | Kit Manufacturer | Shirt Sponsor |
| 1977–1978 | - | National Express |
| 1982–1985 | Umbro | - |
| 1985–1986 | Umbro | Whitbread |
| ... | ... | ... |
| 2008– | Errea | Mira Showers |

Question: what time period had no shirt sponsor?
Dataset: WikiTabQA
Req: Evidence and Direct Answer

| | | |
|------------------------------------|-------------------------|----------------|
| Table Title - Aaron Taylor-Johnson | | |
| Table Subtitle - Film | | |
| Year | Title | Role |
| 2015 | Avengers | Pietro Maimoff |
| | Nocturnal Animals | Ray Marcus |
| 2016 | The Wall | Issac |
| 2017 | Outlaw King | James Douglas |
| 2018 | A Million Little Pieces | - |

Question: What films did Aaron Taylor-Johnson appear in in 2017 and 2018?
Dataset: FeTaQA
Req: Evidence and Decomposition

Figure 1: Examples of Different Table and Contextual structure, taken from different datasets with efficient reasoning method based on specific question, Full Tables are in Appendix C

dependencies, underscoring the need for more effective reasoning frameworks.

Recent research has examined the role of LLMs in table reasoning, demonstrating significant improvements in extracting and analyzing structured data (Zhang et al., 2025). However, studies like Wang and Zhao (2024) highlight persistent challenge in temporal reasoning, as models often struggle to consistently track evolving data and infer event sequences accurately. Moreover, they highlight the limitations of the dominant single-step reasoning technique primarily direct prompting and chain-of-thought reasoning (Wei et al., 2022), which often fail to generalize across different table structures and time sensitive queries.

Despite the growing adoption of LLMs for table reasoning tasks, the impact of different prompting strategies remains largely unexplored. Numerous

prompting techniques have been proposed to enhance LLMs’ reasoning capabilities, yet their effectiveness varies significantly depending on the table structure, question complexity and reasoning requirements. In this study, we evaluate five single-step prompting methods as baselines: Chain-of-Thought (COT), Evidence Extraction, Decomposition, Faithful COT (Radhakrishnan et al., 2023), and Program of Thought (POT) (Chen et al., 2023a). Each of these methods is designed to enhance the logical and numerical reasoning capabilities of LLMs. However, the extent to which these methods influence table reasoning performance, particularly in context of temporal tables, has not been systematically analyzed. Understanding how different prompting strategies affect LLMs’ ability to process structured data is crucial for developing more reliable and adaptable reasoning frameworks.

Tables used in real world applications exhibit diverse structural and contextual characteristics, which significantly impact the effectiveness of different reasoning approaches. Given this variability in table structures, different prompting strategies exhibit varying degrees of effectiveness. Some methods perform well on direct retrieval tasks, while others excel in decomposing complex questions or extraditing relevant evidence as shown in Figure 1.

This study addresses this gap by analyzing the performance of multiple established baselines and a novel adaptive reasoning strategy on a Temporal Tabular Question Answering (TTQA) task. We aim to answer the following research questions:

(RQ1) Given a table and a question, which reasoning strategy should be employed?

(RQ2) Is there a Single reasoning method that can perform well across all types of tabular structure?

(RQ3) Is there a unified representation that can encapsulate all different tabular structures in most effective manner for the TTQA task?

To address these research questions, we conducted experiments on eight distinct tabular structures using multiple state-of-the-art LLMs for the TTQA task. To overcome the limitations of existing baselines we created SEAR (Select-Elaborate-Answer & Reasoning) framework, a novel adaptive prompting strategy designed to systematically address tabular questions with enhanced precision and clarity. SEAR operates in three distinct phases. In the Initial Select phase, it identifies high-level crucial steps, in the subsequent Elaborate phase

it refines these steps by adding detailed instructions, ensuring comprehensive road map. Finally, the Answer & Reasoning phase leverages the structured plan to deliver accurate answers, supported by clean, logical explanations and where necessary, includes integration of Python code for computational tasks.

Furthermore, we combined these three phases to create a single step reasoning strategy, which we call SEAR_Unified. Our results show that SEAR_Unified outperforms all single step baseline reasoning strategies by significant margins, and even standard SEAR exhibits a marginal improvement across all eight distinct tabular structures. This demonstrates the supremacy and efficacy of our proposed reasoning strategy. Additionally, our study also includes detailed analysis of refactoring process, wherein we transform diverse tabular structure into a unified representation ("Refactor"), This approach not only reorganizes the tabular content but also streamlines additional context based on the specific question, enhancing overall clarity and conciseness. Our main contributions are:

- **Benchmarking Prompting Methods:** We evaluate five single-step prompting methods and show that their effectiveness varies based on table structure, entity type, sparsity and question complexity.
- **Adaptive Reasoning Framework:** We introduce SEAR, a multi-step adaptive prompting approach that generalizes well across diverse table structures, also we integrate them into a single unified adaptive prompt SEAR_Unified, outperforming individual methods.
- **Table Structure Refactoring:** We propose refactoring as an enhancement, demonstrating its effectiveness in improving model reasoning by optimized table representation.
- **Comprehensive Evaluation:** We conduct a systematic analysis across various table types, highlighting the impact of different reasoning strategies and structure modifications.

2 Why is Temporal Table Reasoning Challenging?

Tables organize data in a structured format, preserving relationships between values and, in many cases, maintaining temporal order. This structured representation helps in data retrieval and analysis. However, reasoning over tabular data remains a

significant challenge for Large Language Models. When temporal reasoning is involved, the complexity increases further.

One major challenge lies in the diversity of tabular structures. Some tables follow a well-defined row-column format, while others are semi-structured, with implicit relationships between key-value pairs. Some tables are hierarchical, containing multi-level indexing, while others lack a clear hierarchy or even have merged cells. Many tables are interleaved with textual explanations, requiring models to interpret both tabular and textual information jointly. Furthermore, domain-specific tables introduce an additional layer of complexity, as understanding them requires domain knowledge. Tables also vary in representation formats like Markdown, JSON, CSV, and HTML, each demanding different parsing strategies.

Despite the growing interest in table-based reasoning with LLMs, research on temporal reasoning over tabular data remains limited. Datasets such as TempTabQA(Gupta et al., 2023) have been introduced to study this problem, but they often fail to capture the full range of table variations. As a result, models trained on these datasets struggle to generalize. Additionally, datasets frequently contain anomalies that require correction, as highlighted by (Deng et al., 2024).

Symbolic reasoning provides strong logical consistency and performs well for structured tables. Methods like DATER(Ye et al., 2023) and BINDER(Cheng et al., 2023) leverage symbolic logic for table-based reasoning but fail to extend their effectiveness to semi-structured tables. Conversely, C.L.E.A.R(Deng et al., 2024) relies heavily on textual reasoning but does not incorporate logical inference effectively. This limitation highlights the need for a hybrid approach, one that integrates symbolic and textual reasoning dynamically, adapting to the structure of the table and the nature of the question.

To address these challenges, we introduce the Adaptive Prompting Framework. Our method first resolves ambiguities in the tabular context, enhancing the model’s understanding. It then dynamically selects the most effective reasoning strategy based on the question and table structure. This adaptive approach ensures flexibility, allowing the model to leverage both symbolic and textual reasoning depending on the problem at hand.

3 Adaptive Reasoning Framework

Humans intuitively begin by understanding the query objective and analyzing table structures, including cell relationships, headers, and implicit dependencies, while incorporating additional context if available. In temporal tables, this involves identifying both implicit and explicit time-based patterns. Once the problem and context (content and structure wise) are clear, relevant information is retrieved, either directly or by decomposing the task into subproblems based on complexity. Finally, logical and numerical reasoning is applied systematically to arrive at a well-founded conclusion.

Inspired by this intuitive approach, we propose the SEAR (Select-Elaborate-Answer & Reasoning) a framework designed to dynamically adapt reasoning strategies based on the structure and complexity of the given table. SEAR builds upon existing prompting methods by introducing a structured, multi-phase reasoning process that mirrors human problem-solving. It follows a structured three-step process to improve temporal table reasoning, ensuring systematic problem solving while leveraging In-context learning for adaptability.

Step1: Select Crucial Steps : Identify key reasoning steps without answering directly, creating an efficient problem solving path.

- Problem Understanding: Define the question’s objective and analyze table structure.
- Reasoning Process: Select single or multiple strategies from Extract relevant evidence, decompose complex queries, apply logical steps, and generate Python code if needed.
- Optimization tips: Simplify steps, retrieve direct answers when possible, and use code for numerical operations.

Step 2: Elaborate Crucial Steps : Refine and comprehend selected steps for clarity and effectiveness

- Add contextual details, specify exact table elements, and refine decomposition.
- Ensure a structured and logically coherent flow toward the final answer.

Step 3: Answer & Reasoning : Execute the structured steps to derive an accurate, well-supported answers.

- Follow elaborated steps precisely, referencing extracted evidence.
- Justify answers with logical explanations, when possible directly answer from evidence and integrate Python code for calculations when needed.

By progressively refining reasoning, SEAR ensures adaptability and robustness across diverse table formats and complexities.

Since, standard SEAR is a three step process, it introduces overhead complexity that can impact efficiency. To mitigate this, we consolidated SEAR’s structured reasoning into a unified prompt **SEAR_Unified**, a single step adaptive prompt that dynamically selects, merges and refines reasoning steps based on problem’s complexity and tabular structure. Instead of rigidly following predefined steps, SEAR_Unified reasons on the flow, understanding the query objective, identifying relevant tabular elements and contextual insights, and constructing the most effective problem-solving pathway. It extracts key information, decomposes complex queries when necessary, and applies logical reasoning or computation methods where required. Python code is integrated selectively to handle numerical operations. Finally, SEAR_Unified ensures logical coherence and correctness by validating intermediate steps, summarizing results, and performing error checks, reinforcing accuracy while reducing redundant complexity across diverse table formats. Figure 2 illustrates an example of a prompt used in our method, while Figure 3 depicts the response path followed.

| Categories | fetaqa | finqa | hitab | hybridqa | multi | squall | tatqa | wiki |
|-------------------------------|--------|-------|-------|----------|-------|--------|-------|------|
| Table Structure | 1580 | 961 | 616 | 1528 | 1587 | 774 | 2240 | 1503 |
| Title Clarity | 1582 | 962 | 386 | 1528 | 1587 | 774 | 2244 | 1504 |
| Column/Row Header | 1268 | 919 | 353 | 1229 | 1587 | 774 | 2158 | 1283 |
| Data Formatting | 1329 | 957 | 269 | 1476 | 1585 | 774 | 2124 | 1399 |
| Bolding & Emphasis | 1207 | 934 | 206 | 1460 | 1524 | 347 | 2200 | 478 |
| Other | 328 | 273 | 82 | 468 | 539 | 197 | 696 | 309 |

Table 1: Dataset evaluation for refactoring categories.

To further enhance the efficiency of SEAR and SEAR_Unified, we introduce **table and context refactoring** as an add-on, ensuring that the reasoning process operates on well-structured, unambiguous data. Many real world tables suffer from inconsistencies such as missing titles, unclear headers, and structure misalignment. We hypothesize that incorporating a pre-processing step which refines tables by clarifying headers, aligning data, preserving essential relationship without altering

original content and reducing context to query specific information should enhance the retrieval precision, minimizes reasoning errors and improves adaptability across diverse tabular formats. Table 1 quantifies the various refactoring changes applied on all datasets.

4 Experimental Setup

Datasets. We selected eight diverse tabular datasets spanning structured, semi-structured, hierarchical, and hybrid tables to ensure a comprehensive evaluation. These datasets present challenges such as entity relations, numerical reasoning, and textual integration, making them well-suited for assessing table reasoning in LLMs. To filter temporal reasoning questions, we identified queries containing explicit time-related terms (e.g., day, month, year, decade, season) and domain-specific terms like "fiscal" in financial datasets. This process ensured a focus on questions requiring an understanding of time-based dependencies in tabular data. Below is an overview of the datasets used in this study:

1. **FeTaQA(Nan et al., 2021)** : A Wikipedia-based table QA dataset that requires generating long-form answers by integrating multiple discontinuous facts and reasoning across structured tables. **Temporal Questions: 1,582**
2. **FinQA(Chen et al., 2021)** : A financial QA dataset from reports, requiring expert-verified multi-step numerical reasoning and gold reasoning programs for explainability. **Temporal Questions: 962**
3. **HiTab(Cheng et al., 2022)** : A cross-domain QA and NLG dataset featuring hierarchical tables, analyst-authored questions, and fine-grained annotations for complex numerical reasoning. **Temporal Questions: 897**
4. **HybridQA(Chen et al., 2020b)** : A QA dataset requiring reasoning over Wikipedia tables and linked free-form text, demanding both tabular and textual data for accurate answers. **Temporal Questions: 1,528**
5. **MultiHierTT(Zhao et al., 2022)** : A financial QA benchmark requiring reasoning over multiple hierarchical tables and long unstructured text, with detailed multi-step numerical reasoning annotations. **Temporal Questions: 1,587**

6. **Squall(Shi et al., 2020)** : An extension of WikiTableQuestions with manually created SQL equivalents and fine-grained alignments, supporting structured query reasoning in tabular environments. **Temporal Questions: 774**
7. **TAT-QA(Zhu et al., 2021)** : A financial QA dataset requiring reasoning over both tabular and textual data, involving operations like arithmetic, counting, and sorting for quantitative and qualitative analysis. **Temporal Questions: 2,244**
8. **WikiTableQ(Pasupat and Liang, 2015)** : A Wikipedia-based QA dataset with trivia-style questions requiring factual and numerical reasoning over tables with at least 8 rows and 5 columns. **Temporal Questions: 1,504**

By curating these datasets and extracting temporal reasoning-specific questions, we aim to analyze how different prompting methods perform across diverse table structures and reasoning challenges. The dataset selection ensures a broad coverage of table types, reasoning styles, and domains, making the evaluation framework robust and comprehensive.

Models: We experimented with several state-of-the-art large language models, including GPT-4o-mini, Gemini 1.5 Pro Flash, and LLaMA 3.1 70B. These models represent the latest advancements in both closed-source and open-source LLMs, excelling in natural language understanding and task-oriented generation.

Prompts & Frameworks: Effective prompting improves task comprehension and response quality by providing structured instructions. Some prompts include demonstrations to enhance model performance. We evaluate the following prompting strategies:

- *Chain of Thought (COT)* (Wei et al., 2023): CoT guides models through step-by-step reasoning, promoting structured responses.

- *Evidence Extraction (EE)*: This technique directs LLMs to first identify and extract key supporting information from the input before generating an answer. By explicitly isolating relevant evidence, it enhances factual accuracy and minimizes hallucinations in reasoning.

- *Decomposed Prompting (Decomp)* (Khot et al., 2023): This approach breaks down complex tasks into simpler sub-tasks, each handled by a

specialized prompt. By modularizing reasoning, it enables more precise and interpretable responses, allowing LLMs to tackle intricate problems step by step.

- *Faithful Chain of Thought (F-COT)* (Lyu et al., 2023): F-CoT ensures consistency by maintaining fidelity to the initial prompt throughout response generation.

- *Program of Thought (POT)* (Chen et al., 2023b): PoT offers a predefined sequence of operations for structured, task-specific responses.

To ensure a balanced evaluation, we included both textual and symbolic reasoning prompts. CoT, Evidence Extraction, and Decomposed Prompting guide models through step-by-step interpretation, while PoT and F-CoT generate structured logic for consistent reasoning. This distinction helps assess the impact of different reasoning approaches on temporal table tasks. All methods were evaluated in a few-shot setting.

Evaluation. Evaluating diverse datasets is challenging due to varying answer types, from numerical values to long-form text. A rigid approach may miss semantic correctness, so we combine lexical and contextual metrics for a balanced assessment.

- 1. *Relaxed Exact Match Score(REMS)*: This metric uses an F1-score to measure token overlap between the predicted and gold answer, allowing partial matches for better precision-recall balance. Unlike strict exact match, REMS is more flexible with lexical variations. For numerical answers, it permits a $\pm 5\%$ tolerance after decimal instead of token matching. For example, if the correct answer is 10.64, a prediction of 10.62 is accepted, while 11.64 is not.

Despite its flexibility, REMS does not always reflect true semantic accuracy. High scores indicate strong token alignment, but valid paraphrases can be unfairly penalized. For instance, the correct answer “Barack Obama was the 44th President of the United States” would receive a high score for “Obama was the 44th U.S. President” due to token overlap, but “Obama, a politician, led the U.S.” may score lower despite being factually correct. This limitation makes careful interpretation necessary for low-scoring responses.

- 2. *Contextual Answer Evaluation(CAE)*: CAE is an LLM-based scoring method that assesses responses based on meaning rather than exact token overlap. Using a carefully crafted prompt, it determines whether a response correctly conveys the

intended information. Unlike traditional lexical matching, CAE accounts for paraphrasing and rewording, ensuring a more nuanced assessment of correctness, particularly for complex or free-form answers. The full CAE prompt used for evaluation is provided in Figure 4

3. Hybrid Correctness Score (HCS): To balance both lexical and semantic accuracy, we introduce HCS, which combines REMS and CAE. A response is considered correct if either the REMS score exceeds 80 or CAE validates it as correct. This hybrid approach mitigates the limitations of strict string matching while leveraging LLM-based reasoning for a more comprehensive assessment. By integrating both lexical and contextual evaluation, HCS provides a more reliable measure of answer correctness, ensuring a robust and adaptable evaluation framework for tabular reasoning tasks. In this paper, **all reported scores represent HCS**, ensuring a consistent and comprehensive evaluation standard across all experiments.

5 Result and Analysis

In this section, we analyze results using Tables (3, 4, 5) which showcase HCS scores, from multiple perspectives to address our research questions. We have added REMS and CAE scores in Appendix B

Is there a single existing reasoning strategy which works best on all table types? Performance varies depending on table structure, domain, and question complexity. As observed in Gemini 1.5 Flash results (Table 3), COT performs best on HybridQA, Evidence Extraction excels in HiTab, TATQA, FeTaQA and Squall, while Decomposition is most effective for WikiTabQA and FinQA. POT shows the highest performance in MultiHierTT, whereas F-COT does not emerge as the best baseline in any dataset. A similar trend is evident across GPT and LLaMA models as shown in Table 2. Thus, no single prompting method universally outperforms others, as effectiveness is highly dependent on the dataset’s structure and complexity.

Does the Adaptive Reasoning Framework Help?

Table 2 confirms that COT, Evidence Extraction, and Decomposition dominate in most datasets, with POT and F-COT experience improvement in performance for financial and Squall datasets. SEAR dynamically selects its reasoning path, primarily leveraging Evidence Extraction, Decomposition, and Logical Steps (COT) while integrating Python

| | Gemini 1.5 Flash | GPT 4o mini | Llama 3.1 70B |
|--------|------------------------------|--------------------------------|-----------------------------------|
| COT | HybridQA | MultiHierTT TATQA FeTaQA | HiTab HybridQA |
| EE | HiTab TATQA FeTaQA | WikiTabQA HiTab HybridQA | FeTaQA Squall |
| Decomp | Squall WikiTabQA FinQA | FinQA | WikiTabQA MultiHierTT TATQA |
| POT | MultiHierTT | Squall | FinQA |
| F-COT | - | - | - |

Table 2: Dataset for which Baseline reasoning strategy performed best for each model

Program for numerical reasoning. by design, it optimally combines dominant reasoning strategies with computation support. SEAR outperforms baseline in 5 dataset for Gemini, in 2 dataset for GPT, and in 4 datasets for LLaMA. While SEAR consistently improves performance over baseline across multiple models, it does not generalize equally across all datasets.

Does unification of SEAR help? SEAR_Unified optimizes reasoning by merging and refining steps into a single adaptive prompt, reducing overhead while enhancing flexibility. As seen in Table 3, 4, 5, SEAR_Unified outperforms baselines across all datasets for Gemini, while for GPT and LLaMA, it surpasses baselines in 6 datasets, demonstrating its superiority. This highlights SEAR_Unified’s ability to generalize effectively across diverse datasets and models.

| | wiki | multi | hitab | finqa | tatqa | fetaqa | squall | hybridqa |
|---------------------|--------------|--------------|--------------|--------------|--------------|--------------|--------------|--------------|
| COT | 73.60 | 58.79 | 79.04 | 60.08 | 87.30 | 71.30 | 69.90 | 80.76 |
| F-COT | 66.89 | 60.68 | 52.06 | 62.16 | 78.79 | 56.13 | 61.11 | 17.93 |
| Decomp | 78.52 | 61.00 | 75.47 | 62.58 | 91.67 | 67.07 | 67.57 | 74.67 |
| EE | 76.33 | 60.43 | 80.82 | 55.93 | 92.20 | 77.62 | 72.32 | 80.10 |
| PoT | 74.40 | 61.12 | 70.68 | 60.52 | 79.68 | 50.88 | 63.57 | 38.48 |
| Our Approach (SEAR) | | | | | | | | |
| SEAR | 81.45 | 60.18 | 79.71 | 65.90 | 90.02 | 82.87 | 80.23 | 81.15 |
| +Refactor | 82.71 | 58.54 | 81.05 | 65.49 | 89.39 | 84.20 | 78.04 | 65.90 |
| SEAR _U | 82.18 | 61.75 | 82.61 | 68.71 | 92.78 | 79.84 | 81.52 | 82.00 |
| +Refactor | 83.38 | 56.58 | 82.83 | 67.36 | 91.53 | 85.52 | 77.91 | 67.08 |

Table 3: HCS score (in %) for all reasoning strategies across all datasets using Gemini 1.5 Flash, R stands for "Refactoring" and U stands for "Unified".

6 Discussion

The Adaptive Framework consistently generalizes across multiple datasets by dynamically selecting appropriate reasoning paths. Table 6 summarizes the reasoning paths chosen by GPT-4o-mini, showing that Evidence Extraction is always included.

| | wiki | multi | hitab | finqa | tatqa | fetaqa | squall | hybridqa |
|------------|--------------|--------------|--------------|--------------|--------------|--------------|--------------|--------------|
| COT | 78.92 | 57.97 | 77.59 | 64.14 | 92.91 | 84.13 | 67.57 | 78.21 |
| F-COT | 71.61 | 55.32 | 71.35 | 64.97 | 91.04 | 77.81 | 56.46 | 34.62 |
| Decomp | 79.79 | 57.03 | 76.14 | 65.18 | 92.65 | 78.45 | 62.40 | 77.68 |
| EE | 80.12 | 56.77 | 79.38 | 56.03 | 92.81 | 83.88 | 66.67 | 79.58 |
| POT | 79.59 | 57.91 | 76.25 | 56.13 | 90.15 | 72.00 | 72.35 | 61.98 |
| SEAR | 80.19 | 57.40 | 77.37 | 67.26 | 92.42 | 83.38 | 69.64 | 75.33 |
| SEAR_U | 79.92 | 61.00 | 78.93 | 71.10 | 92.91 | 84.89 | 76.74 | 78.27 |
| SEAR + R | 82.91 | 56.65 | 78.82 | 66.94 | 91.84 | 84.77 | 79.33 | 68.72 |
| SEAR_U + R | 84.18 | 59.29 | 80.27 | 69.75 | 91.44 | 84.39 | 79.20 | 70.48 |

Table 4: HCS score (in %) for all reasoning strategies across all datasets using GPT 4o mini, R stands for "Refactoring" and U stands for "Unified".

| | wiki | multi | hitab | finqa | tatqa | fetaqa | squall | hybridqa |
|------------|--------------|--------------|--------------|--------------|--------------|--------------|--------------|--------------|
| COT | 81.05 | 57.59 | 82.95 | 66.22 | 91.00 | 86.03 | 75.45 | 81.66 |
| F-COT | 66.22 | 39.82 | 64.55 | 51.77 | 45.12 | 52.78 | 61.11 | 33.31 |
| Decomp | 82.85 | 59.29 | 81.84 | 65.28 | 93.18 | 84.51 | 73.51 | 80.53 |
| EE | 81.91 | 58.92 | 82.84 | 61.75 | 92.54 | 86.62 | 80.10 | 81.07 |
| POT | 76.53 | 58.98 | 67.56 | 66.42 | 91.40 | 50.44 | 68.22 | 37.76 |
| SEAR | 82.65 | 59.61 | 83.05 | 66.63 | 92.34 | 85.52 | 81.40 | 79.78 |
| SEAR_U | 82.05 | 62.19 | 82.39 | 70.17 | 93.27 | 87.04 | 82.04 | 80.27 |
| SEAR + R | 82.65 | 57.09 | 82.39 | 67.26 | 91.67 | 86.85 | 76.87 | 67.74 |
| SEAR_U + R | 85.11 | 58.16 | 83.05 | 69.67 | 92.89 | 87.23 | 82.49 | 72.16 |

Table 5: HCS score (in %) for all reasoning strategies across all datasets using Llama 3.1 70B, R stands for "Refactoring" and U stands for "Unified".

This step helps the model focus on relevant information, aligning with human intuition (Section 3). For lookup-based questions, Evidence Extraction alone suffices, while more complex tasks require a combination of reasoning methods.

Datasets with long-form answers, such as Fe-TaQA, tend to benefit from textual reasoning methods. As shown in Table 5, for LLaMA 3.1 70B, Fe-TaQA achieves higher accuracy with CoT (84.13%) and Decomposed Prompting (78.45%). This trend is further supported by Table 6, where Evidence Extraction + Decomposed Prompting is the most frequently chosen reasoning path, both being textual techniques. Table 7 reinforces this, showing that 87% of FeTaQA’s reasoning paths rely on textual methods, highlighting their effectiveness for free-form responses.

In contrast, financial datasets like FinQA benefit more from symbolic reasoning due to their reliance on numerical computations. As seen in Table 5, PoT achieves the best performance, with F-CoT also performing well. Table 6 further confirms this, with Evidence Extraction + F-CoT as the most common reasoning path. Similarly, Table 7 shows that 88.25% of FinQA’s reasoning paths involve PoT and F-CoT, emphasizing the strength of symbolic reasoning for computation-heavy datasets.

This pattern extends across datasets, with chosen reasoning paths aligning with their respective

strengths. Table 16 and 17 provide reasoning path analysis for LLaMA 3.1 70B and Gemini-1.5-flash, respectively. By dynamically selecting the most effective reasoning approach based on question type and tabular context, the Adaptive Framework consistently delivers strong performance across diverse table structures and reasoning tasks.

Impact of Table Refactoring. Refactoring tabular data enhances LLM accuracy by improving clarity, structure, and accessibility. Table 1 categorizes key refactoring techniques that aid model interpretation. In ‘Table Structure’, standardizing tables to Markdown format significantly improves performance. For instance, the Squall dataset, originally in JSON, benefits from this transformation. As shown in Table 4, GPT-4o-mini with SEAR + Refactoring (79.33%) outperforms SEAR (69.64%) by 9.69%, and SEAR_U + Refactoring (79.20%) exceeds SEAR_U (76.74%) by 2.46%. Similarly, LLaMA 3.1 70B achieves its highest accuracy (82.49%) with SEAR_U + Refactoring. In ‘Title Clarity’, refining ambiguous or missing table titles improves context. Figure ?? illustrates how adding a player’s name in the title enhances model comprehension. ‘Column/Row Headers’ are refined to eliminate ambiguity and better align entities. ‘Data Formatting’ reduces redundant details, such as excessive decimal places, which can increase hallucinations as context size grows (Liu et al., 2023). Limiting decimals helps models focus and improves accuracy. ‘Bolding and Emphasis’ highlights key details, directing the model’s attention to relevant content. ‘Other’ refinements, such as adding units, removing whitespace, and reformatting text, further enhance readability. The prompt for table refactoring is shown in Figure 6.

| Reasoning Path | Datasets | | | | | | | |
|----------------|-------------|------------|-------|-------------|------------|------------|------------|------------|
| | fetaqa | finqa | hitab | hybridqa | multi | squall | tatqa | wiki |
| EE | 175 | 46 | 476 | 1332 | 194 | 13 | 929 | 703 |
| EE,Decomp | 1365 | 65 | 191 | 28 | 127 | 160 | 249 | 293 |
| EE,F-COT | 23 | 703 | 111 | 5 | 335 | 581 | 547 | 246 |
| EE,POT | 9 | 138 | 107 | 143 | 909 | 14 | 482 | 186 |
| COT,EE | 1 | 1 | 4 | 12 | 5 | - | 5 | 32 |
| COT,EE,Decomp | 8 | 1 | 3 | 2 | - | 1 | 1 | 13 |
| COT,EE,F-COT | 1 | 7 | 1 | - | 5 | 5 | 12 | 17 |
| COT,EE,POT | - | 1 | 4 | 6 | 12 | - | 19 | 14 |
| Total | 1582 | 962 | 897 | 1528 | 1587 | 774 | 2244 | 1504 |

Table 6: Reasoning Path distribution across all datasets for GPT-4o-mini.

7 Related Work

Tabular Reasoning. LLMs have been widely applied to tabular reasoning tasks such as question an-

| Dataset | COT | | EE | | Decomp | | POT | | F-COT | |
|----------|-----|------|------|-----|-------------|--------------|------------|--------------|------------|--------------|
| | # | % | # | % | # | % | # | % | # | % |
| fetaqa | 10 | 0.63 | 1582 | 100 | 1373 | 86.79 | 9 | 0.57 | 24 | 1.52 |
| finqa | 10 | 1.03 | 962 | 100 | 66 | 6.86 | 139 | 14.45 | 710 | 73.8 |
| hitab | 12 | 1.34 | 897 | 100 | 194 | 21.63 | 111 | 12.38 | 112 | 12.49 |
| hybridqa | 20 | 1.31 | 1528 | 100 | 30 | 1.96 | 149 | 9.75 | 5 | 0.33 |
| multi | 22 | 1.39 | 1587 | 100 | 127 | 8.01 | 921 | 58.03 | 132 | 8.32 |
| squall | 6 | 0.78 | 774 | 100 | 161 | 20.8 | 14 | 1.81 | 586 | 75.71 |
| tatqa | 37 | 1.65 | 2244 | 100 | 250 | 11.14 | 501 | 22.33 | 559 | 24.91 |
| wiki | 76 | 5.05 | 1504 | 100 | 306 | 20.35 | 200 | 13.3 | 263 | 17.49 |

Table 7: Distribution of reasoning methods across all the datasets for GPT-4o-mini.

swering, semantic parsing, and table-to-text generation (Chen et al., 2020a; Gupta et al., 2020; Zhang et al., 2020; Zhang and Balog, 2020). Early approaches like TAPAS (Herzig et al., 2020), TaBERT (Yin et al., 2020), and TABBIE (Iida et al., 2021) improve table comprehension by integrating tabular and textual embeddings, allowing models to better process structured information. Other methods, such as Table2Vec (Zhang et al., 2019) and TabGCN (Pramanick and Bhattacharya, 2021), explore alternative tabular representations, enhancing LLMs’ ability to infer relationships between table elements. However, these methods primarily focus on structured tables and do not explicitly address temporal reasoning, which introduces additional complexity when reasoning over tabular data.

Symbolic Reasoning for Tables. Recent work has explored symbolic reasoning for structured tables with predefined schemas, improving logical inference and data consistency (Cheng et al., 2023; Ye et al., 2023; Wang et al., 2024). These methods rely on well-defined structures to extract and process information effectively. However, they struggle with semi-structured and hierarchical tables, where relationships between data points are implicit rather than explicitly defined. Unlike structured tables, these formats require reasoning beyond simple schema-based lookups, often incorporating row-level key-value associations, nested relationships, and missing values. Additionally, temporal reasoning in such tables demands an understanding of time-based dependencies, which current symbolic approaches fail to capture.

Other Reasoning Frameworks. C.L.E.A.R (Deng et al., 2024) demonstrated strong temporal reasoning on domain-specific semi-structured tables by integrating domain knowledge into responses. However, it relies solely on textual reasoning, ignoring robust symbolic approaches, and lacks scalability to other table formats. Similarly, Meta-Reasoning Prompting (MRP)(Gao et al., 2024) selects the optimal reasoning strategy through a two-step process but does not com-

bine reasoning techniques for complex tasks. In contrast, our approach integrates both textual and symbolic reasoning to enhance performance across diverse table types while dynamically selecting the best reasoning path. Moreover, our SEAR-Unified prompt streamlines this into a single-step process, ensuring efficiency and consistency across different table structures.

8 Conclusion and Future Work

This paper introduces SEAR, an adaptive reasoning strategy for LLMs to tackle TTQA tasks, along with its consolidated version, SEAR_Unified. Additionally, we take a step toward a unified table representation by incorporating table refactoring as an enhancement. Our study provides a comprehensive analysis of various reasoning strategies across eight diverse datasets, benchmarking SEAR and SEAR_Unified against multiple baselines.

Results demonstrate that SEAR, SEAR_Unified and with Table Refactoring significantly outperforms popular LLM reasoning methods, with SEAR_Unified surpassing SEAR itself, showcasing its ability to optimize and streamline reasoning with minimal overhead. This highlights capability of modern LLMs to dynamically adjust reasoning within a single prompt, reducing the need for explicit multi-step processes. Our findings reinforce the importance of adaptive reasoning and structured table representation, paving the way for further advancements in LLM-based temporal table reasoning.

While SEAR-based approaches have significantly improved Temporal Table QA, several areas remain open for further exploration. In this work, we have explored Markdown as a unified tabular representation, but future research should investigate alternative formats such as JSON, CSV, or HTML, which could further enhance model adaptability across diverse table structures. Currently, all experiments have been conducted using In-Context Learning (ICL), which limits scalability and efficiency. Future work should explore lightweight adaptive reasoning techniques that could also incorporate self-refinement loops, as the flexibility of SEAR_Unified has shown clear advantages over the rigid reasoning pathways of standard SEAR. Lastly, evaluating SEAR-based approaches on more domains such as medical and scientific evolution dataset which will further validate robustness of Adaptive reasoning strategies

for LLMs.

Limitations

While our study has yielded interesting observations, it’s crucial to acknowledge its limitations. A closer look at the HCS scores in Table 3, 4, 5, reveals that while improvements are observed for datasets with single table contexts, datasets containing multiple tables, such as MultiHierTT and Hybrid tables, show a decline in performance with SEAR-based approaches. This highlights a key limitation of our Table Refactoring method, suggesting that restructuring strategies may need further refinement to handle multi-table contexts effectively. Additionally, scalability remains a concern, as our approach relies on In-Context Learning (ICL), which may not scale effectively for large table datasets. The reliance on ICL-based reasoning can lead to performance bottlenecks.

Ethics Statement

We confirm that our work adheres to the highest ethical standards in research and publication. We will publicly release our code and filtered datasets to enable the research community to validate and build upon our findings. We are committed to the responsible and fair use of computational linguistics methodologies. The claims in our paper accurately reflect the experimental results. While using black-box large language models introduces some stochasticity, we mitigate this by maintaining a fixed temperature. We utilize an AI assistive tools for writing while ensuring absence of bias. We provide comprehensive details on annotations, dataset splits, models used, and prompting methods tried, ensuring the reproducibility of our work.

References

- Wenhu Chen, Xueguang Ma, Xinyi Wang, and William W. Cohen. 2023a. [Program of thoughts prompting: Disentangling computation from reasoning for numerical reasoning tasks](#). *Preprint*, arXiv:2211.12588.
- Wenhu Chen, Xueguang Ma, Xinyi Wang, and William W. Cohen. 2023b. [Program of thoughts prompting: Disentangling computation from reasoning for numerical reasoning tasks](#). *Preprint*, arXiv:2211.12588.
- Wenhu Chen, Hongmin Wang, Jianshu Chen, Yunkai Zhang, Hong Wang, Shiyang Li, Xiyu Zhou, and William Yang Wang. 2020a. [Tabfact: A large-scale](#)

[dataset for table-based fact verification](#). *Preprint*, arXiv:1909.02164.

- Wenhu Chen, Hanwen Zha, Zhiyu Chen, Wenhan Xiong, Hong Wang, and William Yang Wang. 2020b. [HybridQA: A dataset of multi-hop question answering over tabular and textual data](#). In *Findings of the Association for Computational Linguistics: EMNLP 2020*, pages 1026–1036, Online. Association for Computational Linguistics.

- Zhiyu Chen, Wenhu Chen, Charese Smiley, Sameena Shah, Iana Borova, Dylan Langdon, Reema Moussa, Matt Beane, Ting-Hao Huang, Bryan Routledge, and William Yang Wang. 2021. [FinQA: A dataset of numerical reasoning over financial data](#). In *Proceedings of the 2021 Conference on Empirical Methods in Natural Language Processing*, pages 3697–3711, Online and Punta Cana, Dominican Republic. Association for Computational Linguistics.

- Zhoujun Cheng, Haoyu Dong, Zhiruo Wang, Ran Jia, Jiaqi Guo, Yan Gao, Shi Han, Jian-Guang Lou, and Dongmei Zhang. 2022. [HiTab: A hierarchical table dataset for question answering and natural language generation](#). In *Proceedings of the 60th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 1094–1110, Dublin, Ireland. Association for Computational Linguistics.

- Zhoujun Cheng, Tianbao Xie, Peng Shi, Chengzu Li, Rahul Nadkarni, Yushi Hu, Caiming Xiong, Dragomir Radev, Mari Ostendorf, Luke Zettlemoyer, Noah A. Smith, and Tao Yu. 2023. [Binding language models in symbolic languages](#). *Preprint*, arXiv:2210.02875.

- Irwin Deng, Kushagra Dixit, Vivek Gupta, and Dan Roth. 2024. [Enhancing temporal understanding in llms for semi-structured tables](#). *Preprint*, arXiv:2407.16030.

- Peizhong Gao, Ao Xie, Shaoguang Mao, Wenshan Wu, Yan Xia, Haipeng Mi, and Furu Wei. 2024. [Meta reasoning for large language models](#). *Preprint*, arXiv:2406.11698.

- Vivek Gupta, Pranshu Kandoi, Mahek Vora, Shuo Zhang, Yujie He, Ridho Reinanda, and Vivek Srikrumar. 2023. [TempTabQA: Temporal question answering for semi-structured tables](#). In *Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing*, pages 2431–2453, Singapore. Association for Computational Linguistics.

- Vivek Gupta, Maitrey Mehta, Pegah Nokhiz, and Vivek Srikrumar. 2020. [INFOTABS: Inference on tables as semi-structured data](#). In *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*, pages 2309–2324, Online. Association for Computational Linguistics.

- Jonathan Herzig, Pawel Krzysztof Nowak, Thomas Müller, Francesco Piccinno, and Julian Eisenschlos. 2020. [Tapas: Weakly supervised table parsing via](#)

| | | |
|-----|---|-----|
| 788 | pre-training . In <i>Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics</i> . Association for Computational Linguistics. | 845 |
| 789 | | 846 |
| 790 | | 847 |
| 791 | Hiroshi Iida, Dung Thai, Varun Manjunatha, and Mohit Iyyer. 2021. Tabbie: Pretrained representations of tabular data . <i>Preprint</i> , arXiv:2105.02584. | 848 |
| 792 | | 849 |
| 793 | | 850 |
| 794 | 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 . <i>Preprint</i> , arXiv:2210.02406. | 851 |
| 795 | | 852 |
| 796 | | 853 |
| 797 | | 854 |
| 798 | | 855 |
| 799 | Nelson F. Liu, Kevin Lin, John Hewitt, Ashwin Paranjape, Michele Bevilacqua, Fabio Petroni, and Percy Liang. 2023. Lost in the middle: How language models use long contexts . <i>Preprint</i> , arXiv:2307.03172. | 856 |
| 800 | | 857 |
| 801 | | 858 |
| 802 | | 859 |
| 803 | Qing Lyu, Shreya Havaldar, Adam Stein, Li Zhang, Delip Rao, Eric Wong, Marianna Apidianaki, and Chris Callison-Burch. 2023. Faithful chain-of-thought reasoning . In <i>Proceedings of the 13th International 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)</i> , pages 305–329, Nusa Dua, Bali. Association for Computational Linguistics. | 860 |
| 804 | | 861 |
| 805 | | 862 |
| 806 | | 863 |
| 807 | | 864 |
| 808 | | 865 |
| 809 | | 866 |
| 810 | | 867 |
| 811 | | 868 |
| 812 | | 869 |
| 813 | Linyong Nan, Chiachun Hsieh, Ziming Mao, Xi Victoria Lin, Neha Verma, Rui Zhang, Wojciech Kryściński, Nick Schoelkopf, Riley Kong, Xiangru Tang, Murori Mutuma, Ben Rosand, Isabel Trindade, Renusree Bandaru, Jacob Cunningham, Caiming Xiong, and Dragomir Radev. 2021. Fetaqa: Free-form table question answering . <i>Preprint</i> , arXiv:2104.00369. | 870 |
| 814 | | 871 |
| 815 | | 872 |
| 816 | | 873 |
| 817 | | 874 |
| 818 | | 875 |
| 819 | | 876 |
| 820 | Panupong Pasupat and Percy Liang. 2015. Compositional semantic parsing on semi-structured tables . <i>Preprint</i> , arXiv:1508.00305. | 877 |
| 821 | | 878 |
| 822 | | 879 |
| 823 | Aniket Pramanick and Indrajit Bhattacharya. 2021. Joint learning of representations for web-tables, entities and types using graph convolutional network . In <i>Proceedings of the 16th Conference of the European Chapter of the Association for Computational Linguistics: Main Volume</i> , pages 1197–1206, Online. Association for Computational Linguistics. | 880 |
| 824 | | 881 |
| 825 | | 882 |
| 826 | | 883 |
| 827 | | 884 |
| 828 | | 885 |
| 829 | | 886 |
| 830 | Ansh Radhakrishnan, Karina Nguyen, Anna Chen, Carol Chen, Carson Denison, Danny Hernandez, Esin Durmus, Evan Hubinger, Jackson Kernion, Kamilė Lukošiuūtė, Newton Cheng, Nicholas Joseph, Nicholas Schiefer, Oliver Rausch, Sam McCandlish, Sheer El Showk, Tamera Lanham, Tim Maxwell, Venkatesa Chandrasekaran, Zac Hatfield-Dodds, Jared Kaplan, Jan Brauner, Samuel R. Bowman, and Ethan Perez. 2023. Question decomposition improves the faithfulness of model-generated reasoning . <i>Preprint</i> , arXiv:2307.11768. | 887 |
| 831 | | 888 |
| 832 | | 889 |
| 833 | | 890 |
| 834 | | 891 |
| 835 | | 892 |
| 836 | | 893 |
| 837 | | 894 |
| 838 | | 895 |
| 839 | | 896 |
| 840 | | 897 |
| 841 | Tianze Shi, Chen Zhao, Jordan Boyd-Graber, Hal Daumé III, and Lillian Lee. 2020. On the potential of lexico-logical alignments for semantic parsing to SQL queries . In <i>Findings of the Association</i> | 898 |
| 842 | | 899 |
| 843 | | 900 |
| 844 | | |
| | <i>for Computational Linguistics: EMNLP 2020</i> , pages 1849–1864, Online. Association for Computational Linguistics. | |
| | Yuqing Wang and Yun Zhao. 2024. TRAM: Benchmarking temporal reasoning for large language models . In <i>Findings of the Association for Computational Linguistics: ACL 2024</i> , pages 6389–6415, Bangkok, Thailand. Association for Computational Linguistics. | |
| | Zilong Wang, Hao Zhang, Chun-Liang Li, Julian Martin Eisenschlos, Vincent Perot, Zifeng Wang, Lesly Miculicich, Yasuhisa Fujii, Jingbo Shang, Chen-Yu Lee, and Tomas Pfister. 2024. Chain-of-table: Evolving tables in the reasoning chain for table understanding . <i>Preprint</i> , arXiv:2401.04398. | |
| | Jason Wei, Xuezhi Wang, Dale Schuurmans, Maarten Bosma, Brian Ichter, Fei Xia, Ed Chi, Quoc Le, and Denny Zhou. 2023. Chain-of-thought prompting elicits reasoning in large language models . <i>Preprint</i> , arXiv:2201.11903. | |
| | 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 <i>Proceedings of the 36th International Conference on Neural Information Processing Systems, NIPS '22</i> , Red Hook, NY, USA. Curran Associates Inc. | |
| | Siheng Xiong, Ali Payani, Ramana Kompella, and Faramarz Fekri. 2024. Large language models can learn temporal reasoning . In <i>Proceedings of the 62nd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)</i> , pages 10452–10470, Bangkok, Thailand. Association for Computational Linguistics. | |
| | Yunhu Ye, Binyuan Hui, Min Yang, Binhua Li, Fei Huang, and Yongbin Li. 2023. Large language models are versatile decomposers: Decompose evidence and questions for table-based reasoning . <i>Preprint</i> , arXiv:2301.13808. | |
| | Pengcheng Yin, Graham Neubig, Wen tau Yih, and Sebastian Riedel. 2020. Tabert: Pretraining for joint understanding of textual and tabular data . <i>Preprint</i> , arXiv:2005.08314. | |
| | Li Zhang, Shuo Zhang, and Krisztian Balog. 2019. Table2vec: Neural word and entity embeddings for table population and retrieval . In <i>Proceedings of the 42nd International ACM SIGIR Conference on Research and Development in Information Retrieval, SIGIR '19</i> . ACM. | |
| | Shuo Zhang and Krisztian Balog. 2020. Web table extraction, retrieval, and augmentation: A survey . <i>ACM Trans. Intell. Syst. Technol.</i> , 11(2):13:1–13:35. | |
| | Shuo Zhang, Zhuyun Dai, Krisztian Balog, and Jamie Callan. 2020. Summarizing and exploring tabular data in conversational search . In <i>Proceedings of the 43rd International ACM SIGIR Conference on Research and Development in Information Retrieval</i> , | |

SIGIR '20, pages 1537–1540, New York, NY, USA. Association for Computing Machinery.

Xuanliang Zhang, Dingzirui Wang, Longxu Dou, Qingfu Zhu, and Wanxiang Che. 2025. [A survey of table reasoning with large language models](#). *Front. Comput. Sci.*, 19(9).

Yilun Zhao, Yunxiang Li, Chenying Li, and Rui Zhang. 2022. [MultihierTT: Numerical reasoning over multi hierarchical tabular and textual data](#). *Preprint*, arXiv:2206.01347.

Fengbin Zhu, Wenqiang Lei, Youcheng Huang, Chao Wang, Shuo Zhang, Jiancheng Lv, Fuli Feng, and Tat-Seng Chua. 2021. [Tat-qa: A question answering benchmark on a hybrid of tabular and textual content in finance](#). *Preprint*, arXiv:2105.07624.

A Example Appendix

B REMS & CAE Results

C Figure 1: Expanded Tables

| Year | Amount (\$) |
|-------------|-------------|
| 2007 | 56499000 |
| 2008 | 46899000 |
| 2009 | 39904000 |
| 2010 | 33329000 |
| 2011 | 25666000 |
| Later Years | 128981000 |

Table 15: Aggregate Minimum Lease Payments, Lease Payments FinQA

Leases: for FinQA example

Although Sysco normally purchases assets, it has obligations under capital and operating leases for certain distribution facilities, vehicles, and computers. Total rental expense under operating leases was \$100,690,000, \$92,710,000, and \$86,842,000 in fiscal 2006, 2005, and 2004, respectively. Contingent rentals, subleases, and assets and obligations under capital leases are not significant. Aggregate minimum lease payments by fiscal year under existing non-capitalized long-term leases are as follows:

Total Debt Overview: for FinQA example

Total debt at July 1, 2006 was \$1,762,692,000, of which approximately 75 was at fixed rates averaging 6.0 with an average life of 19 years, and the remainder was at floating rates averaging 5.2. Certain loan agreements contain typical debt covenants to protect noteholders, including provisions to maintain the company’s long-term debt to total capital ratio below a specified level. Sysco was in compliance with all debt covenants at July 1, 2006.

The fair value of Sysco’s total long-term debt is estimated based on the quoted market prices for the same or similar issues or on the current rates offered to the company for debt of the same remaining maturities. The fair value of total long-term debt approximated \$1,669,999,000 at July 1, 2006 and \$1,442,721,000 at July 2, 2005, respectively. As of July 1, 2006 and July 2, 2005, letters of credit outstanding were \$60,000,000 and \$76,817,000, respectively.

| | wiki | | multi | | hitab | | finqa | | tatqa | | fetaqa | | squall | | hybridqa | |
|------------|--------------|--------------|--------------|--------------|--------------|--------------|--------------|--------------|--------------|--------------|--------------|--------------|--------------|--------------|--------------|--------------|
| | REMS | CAE | REMS | CAE | REMS | CAE | REMS | CAE | REMS | CAE | REMS | CAE | REMS | CAE | REMS | CAE |
| COT | 77.31 | 76.66 | 57.49 | 49.72 | 75.26 | 74.58 | 58.94 | 57.07 | 81.68 | 87.48 | 28.38 | 84.13 | 66.27 | 65.25 | 74.07 | 76.51 |
| F-COT | 67.85 | 67.82 | 49.39 | 51.35 | 41.44 | 69.79 | 60.78 | 61.12 | 67.36 | 86.76 | 40.46 | 77.69 | 52.97 | 53.36 | 29.78 | 32.79 |
| Decomp | 77.69 | 76.60 | 56.12 | 49.02 | 73.19 | 73.36 | 60.40 | 58.21 | 86.13 | 87.25 | 28.71 | 78.45 | 61.07 | 59.30 | 74.71 | 74.87 |
| EE | 78.57 | 77.86 | 56.32 | 48.27 | 76.16 | 76.92 | 50.94 | 46.88 | 90.22 | 88.06 | 28.42 | 83.82 | 65.55 | 64.60 | 75.85 | 76.96 |
| POT | 76.28 | 75.93 | 53.41 | 53.12 | 41.92 | 73.47 | 51.88 | 52.49 | 66.88 | 86.10 | 29.71 | 72.00 | 65.90 | 69.12 | 58.66 | 60.27 |
| SEAR | 78.32 | 76.60 | 54.70 | 50.98 | 67.36 | 74.58 | 62.52 | 60.91 | 81.94 | 85.83 | 29.53 | 83.38 | 67.56 | 60.72 | 72.07 | 73.63 |
| SEAR_U | 77.50 | 77.53 | 56.39 | 56.84 | 71.78 | 76.70 | 62.87 | 67.57 | 88.31 | 89.75 | 31.06 | 84.89 | 72.26 | 73.77 | 74.96 | 75.85 |
| SEAR + R | 80.51 | 79.39 | 54.04 | 51.10 | 68.40 | 75.92 | 61.88 | 60.08 | 81.63 | 85.87 | 29.71 | 84.39 | 76.85 | 74.03 | 65.89 | 66.03 |
| SEAR_U + R | 81.14 | 81.25 | 55.54 | 55.51 | 72.13 | 77.59 | 62.43 | 66.53 | 86.56 | 88.23 | 30.47 | 84.70 | 76.21 | 76.87 | 66.96 | 67.74 |

Table 8: REMS & CAE score (in %) for all reasoning strategies across all datasets using GPT 4o mini, R stands for "Refactoring" and U stands for "Unified"

| | wiki | | multi | | hitab | | finqa | | tatqa | | fetaqa | | squall | | hybridqa | |
|------------|--------------|--------------|--------------|--------------|--------------|--------------|--------------|--------------|--------------|--------------|--------------|--------------|--------------|--------------|--------------|--------------|
| | REMS | CAE | REMS | CAE | REMS | CAE | REMS | CAE | REMS | CAE | REMS | CAE | REMS | CAE | REMS | CAE |
| COT | 71.86 | 71.28 | 57.29 | 39.26 | 73.97 | 74.25 | 58.00 | 39.29 | 80.81 | 85.34 | 28.25 | 71.24 | 69.44 | 69.66 | 77.29 | 76.57 |
| F-COT | 64.76 | 57.51 | 58.36 | 47.83 | 35.68 | 49.34 | 60.60 | 34.20 | 64.86 | 74.88 | 37.05 | 55.69 | 60.40 | 60.73 | 17.86 | 15.97 |
| Decomp | 76.26 | 75.00 | 58.90 | 41.84 | 71.70 | 72.44 | 60.72 | 32.22 | 84.23 | 85.12 | 29.91 | 67.07 | 66.01 | 65.98 | 72.94 | 69.31 |
| EE | 74.24 | 72.81 | 59.02 | 42.41 | 74.61 | 76.43 | 54.54 | 30.46 | 86.14 | 86.27 | 28.63 | 77.62 | 71.89 | 72.03 | 74.12 | 68.72 |
| POT | 72.65 | 66.69 | 60.00 | 47.01 | 41.37 | 67.54 | 54.90 | 58.10 | 66.74 | 75.61 | 26.73 | 50.88 | 62.66 | 62.99 | 38.18 | 33.84 |
| SEAR | 79.08 | 78.19 | 57.15 | 54.69 | 74.93 | 76.81 | 59.90 | 61.02 | 75.07 | 83.87 | 28.75 | 82.87 | 76.14 | 68.60 | 77.61 | 78.08 |
| SEAR_U | 79.32 | 80.32 | 59.27 | 57.34 | 78.53 | 79.38 | 63.16 | 65.59 | 82.70 | 86.68 | 31.57 | 79.77 | 77.29 | 79.59 | 77.11 | 79.84 |
| SEAR + R | 80.27 | 78.46 | 55.32 | 52.30 | 75.08 | 77.37 | 59.88 | 60.50 | 73.57 | 84.54 | 28.97 | 84.20 | 76.13 | 72.09 | 62.43 | 62.24 |
| SEAR_U + R | 80.78 | 81.32 | 53.09 | 53.62 | 78.94 | 79.60 | 61.98 | 63.83 | 82.20 | 85.65 | 32.89 | 85.52 | 75.16 | 75.97 | 62.96 | 64.86 |

Table 9: REMS & CAE score (in %) for all reasoning strategies across all datasets using Gemini 1.5 Flash, R stands for "Refactoring" and U stands for "Unified"

| | wiki | | multi | | hitab | | finqa | | tatqa | | fetaqa | | squall | | hybridqa | |
|------------|--------------|--------------|--------------|--------------|--------------|--------------|--------------|--------------|--------------|--------------|--------------|--------------|--------------|--------------|--------------|--------------|
| | REMS | CAE | REMS | CAE | REMS | CAE | REMS | CAE | REMS | CAE | REMS | CAE | REMS | CAE | REMS | CAE |
| COT | 79.20 | 78.86 | 56.91 | 48.71 | 80.77 | 81.38 | 60.91 | 60.81 | 83.69 | 86.10 | 28.07 | 86.03 | 73.21 | 73.39 | 79.10 | 79.78 |
| F-COT | 63.02 | 62.43 | 37.21 | 37.30 | 37.35 | 61.76 | 48.14 | 48.44 | 59.67 | 61.72 | 25.34 | 52.72 | 56.53 | 58.01 | 30.30 | 31.28 |
| Decomp | 80.71 | 80.78 | 58.39 | 52.24 | 78.71 | 80.72 | 60.50 | 59.77 | 86.62 | 86.41 | 29.36 | 84.51 | 71.00 | 71.58 | 79.98 | 77.75 |
| EE | 80.30 | 79.79 | 57.70 | 48.27 | 81.42 | 80.05 | 57.03 | 53.53 | 89.09 | 87.70 | 28.63 | 86.62 | 78.12 | 77.78 | 78.33 | 78.73 |
| POT | 74.74 | 73.34 | 56.47 | 55.14 | 37.05 | 65.44 | 62.44 | 61.75 | 65.02 | 87.17 | 20.25 | 50.44 | 63.43 | 64.73 | 35.63 | 35.60 |
| SEAR | 80.69 | 78.79 | 57.76 | 50.79 | 75.45 | 78.60 | 61.40 | 60.40 | 84.67 | 88.41 | 29.47 | 85.52 | 78.74 | 72.22 | 76.43 | 77.29 |
| SEAR_U | 78.91 | 79.26 | 60.02 | 58.03 | 75.12 | 79.38 | 63.30 | 66.01 | 89.20 | 86.36 | 34.15 | 87.04 | 78.74 | 80.62 | 77.11 | 78.24 |
| SEAR + R | 80.17 | 78.46 | 54.97 | 48.02 | 75.77 | 78.37 | 62.00 | 61.43 | 81.71 | 86.99 | 29.53 | 86.85 | 73.95 | 70.67 | 67.35 | 70.75 |
| SEAR_U + R | 82.53 | 82.05 | 56.15 | 52.68 | 76.19 | 77.70 | 61.66 | 66.03 | 86.58 | 86.47 | 34.83 | 87.17 | 79.01 | 80.68 | 67.11 | 67.80 |

Table 10: REMS & CAE score (in %) for all reasoning strategies across all datasets using Llama 3.1 70B, R stands for "Refactoring" and U stands for "Unified"

SEAR_UNIFIED PROMPT

Instruction
You are a adaptive-reasoner with the capabilities to select or merge steps to create the most appropriate reasoning pathway based on the tabular question provided by the user. You can even develop new reasoning steps by combining the new steps or learning from illustrations to create new pathways depending on the provided problem.

Steps for Adaptive Reasoning:
Each section has multiple approaches, you do not have to use all the approaches. Understand their use-cases and then pick minimal relevant steps to create your own optimal approach to answer the question.

Problem Understanding:
- Determine the objective: Identify the goal or desired outcome of the reasoning process.
- Understand the problem: Comprehend the nature and scope of the problem.

Reasoning Process:
- Step-by-step reasoning: Approach the problem logically, ensuring clarity at each step or stage.
- Extract relevant information: Gather all necessary data and details pertinent to the problem, by extracting relevant rows, columns and textual information.
- Decomposition of problem into sub-problems: Break down the main question into smaller and more manageable sub questions.
- Individually answer each sub-problem with reasoning: Apply logical steps to solve each sub question separately.
- Write a single Python program for solving the problem: Create a detailed unified Python script with comments describing the steps and stages.
- Individually write a Python program for each sub-problem: Develop separate Python scripts for each sub-problem, ensuring modularity and clarity.

Conclusion:
- Summarize findings: Combine the results from each step or sub question to give the final answer as Final Answer: {{[Answer]}}.
- Combine Python code: If necessary, integrate the individual Python scripts into a cohesive program at the end. Print the final answer as Final Answer: {{[Answer]}} , end your code with a comment "#Done".

Error Detection:
- Review each step or sub-problem: Ensure each step or sub-problem has been addressed thoroughly and correctly.
- Ensure logical flow: Verify that the reasoning process flows logically from one step to the next.
- Check Python program for syntax and errors: Confirm that the final Python program is syntactically correct and free of errors.

"Helpful Tips for Creating Appropriate and Optimal Approach":
- Understand what is asked in the question, mention all the steps required to answer the question and why each step is necessary.
- If the question can be broken into smaller and more manageable sub questions, always decompose the question into relevant sub questions.
- If there are "calculations involved you must use python code" for performing calculations and reaching the final answer.
- If the question is directly answerable by direct look up from the tabular data or from the extracted evidence then provide a direct answer.

Table:
Context:

Race Results Overview

This table showcases the results of various athletes who participated in different heats, including their times and nationalities.

| Rank | Heat | Name | Nationality | Time | Notes |
|------|------|--------------------------|----------------|-------|-------|
| 1 | 1 | Salem Al-Yami | Saudi Arabia | 10.55 | Q |
| 2 | 1 | Hiroyasu Tsuchie | Japan | 10.64 | Q |
| 3 | 1 | Khaled Yousef Al-Obaidli | Qatar | 10.68 | Q |
| 4 | 1 | Chintake De Zoysa | Sri Lanka | 10.78 | q |
| 5 | 1 | Suminda Mendis | Sri Lanka | 10.82 | q, PB |
| 6 | 1 | Vissanu Sophanich | Thailand | 10.87 | |
| 1 | 2 | Gennadiy Chernovol | Kazakhstan | 10.59 | Q |
| 2 | 2 | Yuta Kanno | Japan | 10.64 | Q |
| 3 | 2 | Shen Yunbao | China | 10.72 | Q |
| 4 | 2 | Tsai Meng-Lin | Chinese Taipei | 10.74 | q |
| 5 | 2 | Tan Kok Lim | Malaysia | 10.83 | q |
| 6 | 2 | Ahmad Hudeib Al-Mamari | Oman | 10.97 | |
| 1 | 3 | Jamal Al-Saffar | Saudi Arabia | 10.57 | Q |
| 2 | 3 | Reanchai Srihawong | Thailand | 10.72 | Q |
| 3 | 3 | Shin Jung-Ki | South Korea | 10.79 | [2] |
| 4 | 3 | Chen Tien-Wen | Chinese Taipei | 10.74 | q |
| 5 | 3 | To Wai Lok | Hong Kong | 10.92 | |
| 6 | 3 | Poh Seng Song | Singapore | 11.10 | SB |
| 1 | 4 | Chen Haijian | China | 10.65 | Q |
| 2 | 4 | Saad Faraj Al-Shahwani | Qatar | 10.67 | Q |
| 3 | 4 | Azmi Ibrahim | Malaysia | 10.78 | Q |
| 4 | 4 | Chiang Wai Hung | Malaysia | 10.89 | |
| 5 | 4 | Nguyen Thanh Hai | Vietnam | 11.16 | PB |
| 6 | 4 | Phipop Rasme Prum Keo | Cambodia | 11.70 | PB |

Question : which country finished earlier, thailand or japan?

Optimal Approach:

Figure 2: SEAR_UNIFIED Prompt Example

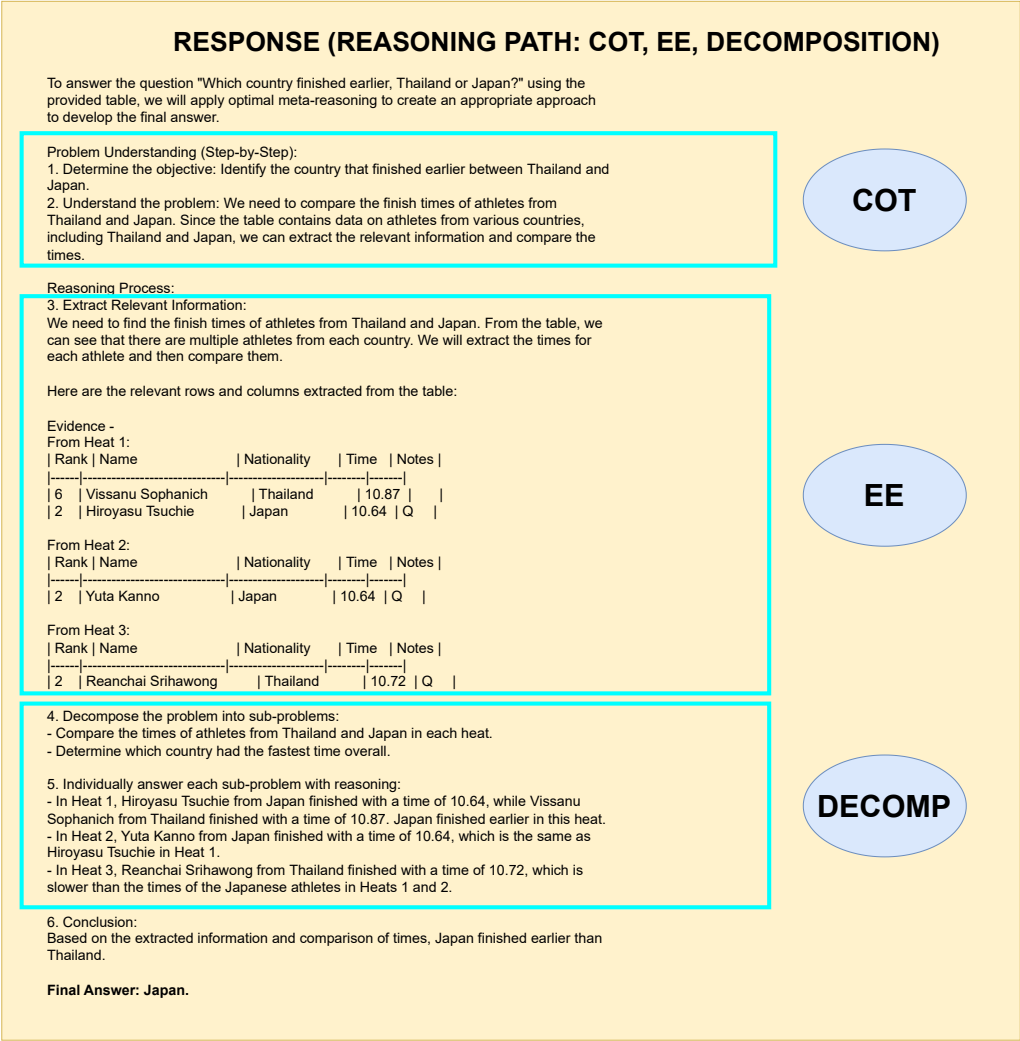


Figure 3: The figure illustrates the response path followed by SEAR_Unified Prompting. The reference prompt is provided in Figure 2

| Year | Kit Manufacturer | Shirt Sponsor | Back of Shirt Sponsor | Short Sponsor |
|-----------|------------------|--------------------------|-------------------------|----------------------|
| 1977–1978 | - | National Express | - | - |
| 1982–1985 | Umbro | - | - | - |
| 1985–1986 | Umbro | Whitbread | - | - |
| 1986–1988 | Henson | Duraflex | - | - |
| 1988–1989 | - | Gulf Oil | - | - |
| 1991–1993 | Technik | Gulf Oil | - | - |
| 1993–1994 | Club Sport | Gulf Oil | - | - |
| 1994–1995 | Klüb Sport | Empress | - | - |
| 1995–1996 | Matchwinner | Empress | - | - |
| 1996–1997 | UK | Endsleigh Insurance | - | - |
| 1997–1999 | Errea | Endsleigh Insurance | - | - |
| 1999–2004 | Errea | Towergate Insurance | - | - |
| 2004–2008 | Errea | Bence Building Merchants | - | - |
| 2008– | Errea | Mira Showers | - | - |
| 2009–2011 | Errea | Mira Showers | PSU Technology Group | - |
| 2011–2013 | Errea | Mira Showers | Barr Stadia | Gloucestershire Echo |
| 2013– | Errea | Mira Showers | Gloucestershire College | Gloucestershire Echo |

Table 11: Historical Sponsorship and Kit Manufacturer Data, WikiTabQA example

Input :

You are an expert LLM evaluator tasked with assessing the accuracy of model responses against gold standard answers. Your role is to determine if the core content and intent of the model's response align with the gold answer, considering various answer formats and implicit information.

Key Guidelines

- **Understand the question's essence**, including specific operations or units mentioned.
- **Compare model responses** to gold answers, focusing on key information.
- **Allow a small margin of error** ($\pm 0.1\%$) for numerical answers.
- **Recognize correct answers in different formats**, such as percentages and decimals.
- **Consider implicit information and context** in responses.
- **For list-type answers:**
 - Evaluate based on content rather than order.
 - If more than **two elements are missing** (context-dependent), evaluate as incorrect.
- **Assess mathematical answers** based on value range unless a specific value is required.
- **Check for appropriate units** in mathematical answers.

Final Judgment

Provide a "Yes" or "No" judgment without explanation unless explicitly requested.

Figure 4: Prompt for Contextual Answer Evaluation(CAV)

| Year | Title | Role | Director | Notes |
|------|---------------------------------|-----------------------|-----------------------|-----------------------------|
| 2000 | The Apocalypse | Johanan | Raffaele Mertes | - |
| 2002 | Tom & Thomas | Tom Sheppard / Thomas | Esmé Lammers | - |
| 2003 | Behind Closed Doors | Sam Goodwin | Louis Caulfield | - |
| 2003 | Shanghai Knights | Charlie Chaplin | David Dobkin | - |
| 2004 | Dead Cool | George | David Cohen | - |
| 2006 | The Thief Lord | Prosper | Richard Claus | - |
| 2006 | The Illusionist | Young Eisenheim | Neil Burger | - |
| 2006 | Fast Learners | Neil | Christoph Röhl | Short film |
| 2006 | The Best Man | Michael (Aged 15) | Stefan Schwartz | - |
| 2007 | The Magic Door | Flip | Paul Matthews | - |
| 2008 | Dummy | Danny | Matthew Thompson | Nominated — ALFS Award |
| 2008 | Angus, Thongs | Robbie Jennings | Gurinder Chadha | - |
| 2009 | The Greatest | Bennett Brewer | Shana Feste | - |
| 2009 | Nowhere Boy | John Lennon | Sam Taylor-Johnson | Empire Award for Best... |
| 2010 | Kick-Ass | David "Dave" Lizewski | Matthew Vaughn | Nominated — Empire Award... |
| 2010 | Chatroom | William Collins | Hideo Nakata | - |
| 2011 | Albert Nobbs | Joe Mackins | Rodrigo García | - |
| 2012 | Savages | Ben | Oliver Stone | - |
| 2012 | Anna Karenina | Count Vronsky | Joe Wright | Final time credited as... |
| 2013 | Kick-Ass 2 | David "Dave" Lizewski | Jeff Wadlow | First time credited as... |
| 2014 | Captain America: Winter Soldier | Pietro Maximoff | Anthony and Joe Russo | Uncredited cameo |
| 2014 | Godzilla | Lt. Ford Brody | Gareth Edwards | - |
| 2015 | Avengers: Age of Ultron | Pietro Maximoff | Joss Whedon | - |
| 2016 | Nocturnal Animals | Ray Marcus | Tom Ford | Golden Globe Award for... |
| 2017 | The Wall | Isaac | Doug Liman | - |
| 2018 | Outlaw King | James Douglas | David Mackenzie | - |
| 2018 | A Million Little Pieces | James Frey | Sam Taylor-Johnson | - |
| 2020 | Kingsman: The Great Game | - | Matthew Vaughn | Filming |

Table 12: Aaron Taylor-Johnson Filmography, example FeTaQA

Input :

Instruction

You are given the following **Question** and **Context**. The **Context** includes a table that may be incomplete, ambiguous, or poorly structured. Your task is to produce a **cleaned version of the table** that improves its clarity and structure so that it can be correctly used to answer the **Question**.

Guidelines

1. **Do not add, remove, or alter any data.** Only restructure and clarify what is already present.
2. You may improve the **table title** if it is missing or ambiguous:
 - If a title is missing, infer an appropriate one based on the **question** and table content.
 - If the existing title is unclear or misleading, revise it for clarity while keeping its original meaning.
3. You may improve the **table headers** if needed:
 - Rename ambiguous column/row headers for clarity.
 - Ensure column and row labels accurately describe their content.
4. You may fix **structural inconsistencies**:
 - Align misaligned data properly under the correct headers.
 - Ensure row and column structures are uniform.
 - Remove redundant headers or merge split headers where necessary.
5. The data should be kept in the same order whenever possible. However, if **minor reordering of rows or columns** helps fix structural issues, you may do so—**only if it does not change or omit any data**.

Output Format

- Provide only the **cleaned table** as your output in a structured format appropriate for the data in **Markdown format**.
- **Do not add any explanations, reasoning, or commentary.**

Question: {question}

Context: {context}

Now produce just the cleaned table.

Figure 5: Prompt for Refactoring Tables.

Table Refactoring Example

Question: how many passing yards did J.J. Raterink get in 2012?

Initial Table

Title : aft statistics ← Lack of Context About Table

| year | team | passing | cmp | att | pct | yds | td | int | rtg | rushing | att | yds | td | | |
|--------|-------------|---------|-------|------|--------|-----|----|--------|-----|---------|-----|-----|----|--|--|
| 2010 | chicago | 65 | 102 | 63.7 | 767 | 14 | 2 | 112.66 | 8 | 9 | 2 | | | | |
| 2011 | chicago | 64 | 105 | 61.0 | 888 | 16 | 2 | 118.27 | 4 | 8 | 2 | | | | |
| 2011 | kansas city | 311 | 500 | 62.2 | 3,723 | 65 | 17 | 103.28 | 48 | 138 | 5 | | | | |
| 2012 | iowa | 413 | 618 | 66.8 | 4,870 | 93 | 10 | 121.49 | 37 | 110 | 8 | | | | |
| 2013 | iowa | 346 | 575 | 60.2 | 4,015 | 78 | 18 | 102.19 | 32 | 10 | 8 | | | | |
| 2014 | los angeles | 211 | 383 | 55.1 | 2,335 | 38 | 19 | 77.53 | 6 | 5 | 1 | | | | |
| 2014 | iowa | 101 | 163 | 62.0 | 1,320 | 22 | 1 | 118.65 | 37 | 111 | 9 | | | | |
| 2015 | las vegas | 178 | 325 | 54.8 | 1,986 | 35 | 9 | 88.57 | 32 | 19 | 6 | | | | |
| career | | 1,689 | 2,771 | 61.0 | 19,904 | 361 | 78 | 103.65 | 204 | 410 | 41 | | | | |

← Bad Column Headers

Refactored Table

Title: Player Statistics for J.J. Raterink ← Improved Title for better Context

| Year | Team | Passing Completions | Passing Attempts | Completion Percentage | Passing Yards | Touchdowns | Interceptions | Rating | Rushing Attempts |
|------------|-------------|---------------------|------------------|-----------------------|---------------|------------|---------------|--------|------------------|
| 2010 | Chicago | 65 | 102 | 63.7% | 767 | 14 | 2 | 112.66 | 8 |
| 2011 | Chicago | 64 | 105 | 61.0% | 888 | 16 | 2 | 118.27 | 4 |
| 2011 | Kansas City | 311 | 500 | 62.2% | 3,723 | 65 | 17 | 103.28 | 48 |
| 2012 | Iowa | 413 | 618 | 66.8% | 4,870 | 93 | 10 | 121.49 | 37 |
| 2013 | Iowa | 346 | 575 | 60.2% | 4,015 | 78 | 18 | 102.19 | 32 |
| 2014 | Los Angeles | 211 | 383 | 55.1% | 2,335 | 38 | 19 | 77.53 | 6 |
| 2014 | Iowa | 101 | 163 | 62.0% | 1,320 | 22 | 1 | 118.65 | 37 |
| 2015 | Las Vegas | 178 | 325 | 54.8% | 1,986 | 35 | 9 | 88.57 | 32 |
| **Career** | | 1,689 | 2,771 | 61.0% | 19,904 | 361 | 78 | 103.65 | 204 |

← Improved Column Headers

Figure 6: Prompt Example

| Benefit Plan | 2017 | 2016 | 2015 |
|--------------------------|-------------|-------------|-------------|
| Pension Plan | 3856 | 3979 | 2732 |
| Health Plan | 11426 | 11530 | 8736 |
| Other plans | 1463 | 1583 | 5716 |
| Total plan contributions | 16745 | 17092 | 17184 |

Table 13: Benefit Plan Contributions, Benefits, Multi-Hierrt example Table 0

| | 2018 | 2019 | 2020 | 2021 | 2022 | Thereafter | Total |
|------------------------------------|--------|---------|---------|--------|---------|------------|----------|
| Property mortgages and other loans | 153593 | 42289 | 703018 | 11656 | 208003 | 1656623 | 2775182 |
| MRA facilities | 90809 | 0 | 0 | 0 | 0 | 0 | 90809 |
| Revolving credit facility | 0 | 0 | 0 | 0 | 0 | 40000 | 40000 |
| Unsecured term loans | 0 | 0 | 0 | 0 | 0 | 1500000 | 1500000 |
| Senior unsecured notes | 250000 | 0 | 250000 | 0 | 800000 | 100000 | 1400000 |
| Trust preferred securities | 0 | 0 | 0 | 0 | 0 | 100000 | 100000 |
| Capital lease | 2387 | 2411 | 2620 | 2794 | 2794 | 819894 | 832900 |
| Ground leases | 31049 | 31066 | 31436 | 31628 | 29472 | 703254 | 857905 |
| Estimated interest expense | 226815 | 218019 | 184376 | 163648 | 155398 | 281694 | 1229950 |
| Joint venture debt | 200250 | 717682 | 473809 | 449740 | 223330 | 2119481 | 4184292 |
| Total | 954903 | 1011467 | 1645259 | 659466 | 1418997 | 7320946 | 13011038 |

Table 14: Loans and Liabilities, Loans, MultiHiertt example Table 1

| Reasoning Path | fetaqa | finqa | hitab | hybridqa | multi | squall | tatqa | wiki |
|----------------|--------|-------|-------|----------|-------|--------|-------|------|
| EE | 221 | 39 | 561 | 1072 | 356 | 9 | 1040 | 987 |
| EE,Decomp | 553 | 21 | 19 | 8 | 33 | 28 | 59 | 81 |
| EE,F-COT | 571 | 853 | 123 | 35 | 262 | 709 | 391 | 236 |
| EE,POT | 234 | 45 | 194 | 405 | 919 | 25 | 753 | 187 |
| COT,EE | - | - | - | 6 | 5 | - | 1 | 7 |
| COT,EE,Decomp | 3 | - | - | 2 | 10 | 1 | - | 2 |
| COT,EE,F-COT | - | 3 | - | - | - | 2 | - | 4 |
| POT | - | 1 | - | - | 2 | - | - | - |
| Total | 1582 | 962 | 897 | 1528 | 1587 | 774 | 2244 | 1504 |

Table 16: Reasoning Path distribution across all datasets for Llama 3.1 70B.

| Reasoning Path | fetaqa | finqa | hitab | hybridqa | multi | squall | tatqa | wiki |
|----------------|--------|-------|-------|----------|-------|--------|-------|------|
| EE | 982 | 106 | 675 | 1492 | 155 | 112 | 1160 | 875 |
| EE,DecompE | 197 | 16 | 6 | 2 | 87 | 17 | 9 | 186 |
| EE,F-COT | 175 | 796 | 29 | - | 333 | 516 | 49 | 173 |
| EE,POT | 191 | 42 | 186 | 33 | 1010 | 119 | 1025 | 268 |
| COT,EE | 25 | - | - | 1 | - | 1 | 1 | 2 |
| COT,EE,Decomp | 3 | - | - | - | - | 1 | - | - |
| COT,EE,F-COT | 2 | 1 | - | - | - | 6 | - | - |
| COT,EE,POT | 7 | - | 1 | - | - | 1 | - | - |
| Decomp | - | 1 | - | - | - | - | - | - |
| POT | - | - | - | - | 2 | 1 | - | - |
| Total | 1582 | 962 | 897 | 1528 | 1587 | 774 | 2244 | 1504 |

Table 17: Reasoning Path distribution across all datasets for Gemini-1.5-Flash.