053

054

000 001

TABLEMASTER: A Recipe to Advance Table Understanding with Language Models

Anonymous Authors¹

Abstract

Tables serve as a fundamental format for representing structured relational data. While current language models (LMs) excel at many text-based tasks, they still face challenges in table understanding due to the complex characteristics of tabular data, such as their structured nature. In this paper, we aim to enhance LMs for improved table understanding. We identify four key challenges: 1) difficulty in locating target data, 2) deficiency in table semantics, 3) numerical inaccuracies in textual reasoning, and 4) semantic inflexibility in symbolic reasoning. To address these issues, we propose TableMaster, a recipe and comprehensive framework that integrates multiple solutions to overcome these obstacles. TableMaster first extracts relevant table content and verbalizes it with enriched semantic context. Additionally, we introduce adaptive reasoning, a flexible approach that dynamically adjusts between textual and symbolic reasoning, tailoring the reasoning process to each query. Extensive analyses and experiments demonstrate our findings and the effectiveness of TableMaster. On the WikiTO dataset, Table-Master achieves an accuracy of 78.13% using GPT-4o-mini, surpassing existing baselines.

1. Introduction

"Data gains extraordinary power as it transcends the simplicity of one dimension to embrace the richness of higher dimensions."

Tables are widely used in daily life and across various fields, such as healthcare (Ghasemi & Amyot, 2016) and finance (Li et al., 2020), due to their unique ability to efficiently represent two-dimensional relational data. It is crucial to process tabular data with both efficiency and accuracy. Recently, large language models (LLMs) (Gunasekar et al., 2023; OpenAI, 2024; Touvron et al., 2023) have achieved significant progress in the field of natural language processing. They perform well in a wide range of downstream text-based tasks, including language understanding (Minaee et al., 2024; Zhu et al., 2024) and reasoning (Plaat et al., 2024). Naturally, language models (LMs) are increasingly being used to process and understand tabular data (Fang et al., 2024; Zhang et al., 2024b), enabling reasoning for downstream tasks such as table-based question answering (Pasupat & Liang, 2015) and table-based fact verification (Chen et al., 2020).

However, the data structure of tables inherently possess a unique two-dimensional structure that contrasts with the linear text, which dominates the content in language model pretraining corpora. Most advanced LMs are not specifically optimized for processing tabular data. While techniques such as chain-of-thought prompting (Wei et al., 2023) and other reasoning-enhanced methods (Yao et al., 2023) have enabled LMs to perform satisfactorily in reasoning with linear text, significant room for improvement remains in table-based reasoning (Chen, 2023). A notable gap persists in LMs' ability to fully understand tables and effectively reason with tabular data.

Many previous studies have aimed to improve the table understanding capabilities of LMs. One efficient approach is using prompting to adapt LMs for table understanding without requiring fine-tuning, making it applicable to any advanced LM. Recent studies primarily adopt two main strategies to enhance table understanding with LMs. The first strategy involves extracting a sub-table that contains relevant content from the original table to reduce the context size, thereby making it easier for LMs to comprehend. Examples include Dater (Ye et al., 2023) and Chain-of-Table (Wang et al., 2024), among others. The second strategy leverages SQL or Python programs to augment numerical reasoning, locate target data, and enhance table understanding of numerical information, as demonstrated by Binder (Cheng et al., 2023) and LEVER (Ni et al., 2023), etc. However, these studies primarily focus on a single basic aspect

¹Anonymous Institution, Anonymous City, Anonymous Region, Anonymous Country. Correspondence to: Anonymous Author <anon.email@domain.com>.

Preliminary work. Under review by the International Conference on Machine Learning (ICML). Do not distribute.

A Recipe to Advance Table Understanding with Language Models

Challenges





Figure 1. Overview of the challenges and proposed solutions in this work. Tabular data is inherently structured, dense, concise, and numerical. Based on these characteristics, we identify four key challenges. To address them, we propose four targeted solutions. The gray arrows between the characteristics and challenges represent the potential causes of these challenges stemming from specific characteristics. Each proposed solution corresponds to the challenge presented on the left in the same row. TableMaster is a unified recipe developed

to enhance the performance of LMs in table understanding or design complex methods with isolated strategies. There is currently an absence of work that provides a systematic and fundamental analysis of table understanding with language models and proposes comprehensive methods for its

In this paper, we first provide extensive experiments and discussions to identify the challenges in table understanding with language models. To address these challenges, we then introduce TableMaster, a recipe and comprehensive framework that integrates multiple solutions to tackle these issues effectively. In summary, this paper makes the following key

· Challenges of Table Understanding. We observe that tabular data is inherently structured, dense, concise, and numerical. Through empirical analysis, we identify four challenges associated with LMs' table understanding: difficulty in locating target data, deficiency of table semantics, numerical inaccuracies in textual reasoning, and se-100 mantic inflexibility in symbolic reasoning. (Section 3)

102 • A Recipe for Table Understanding. To address these challenges, we propose targeted solutions: table-of-focus, 104 table verbalization, program-aided reasoning, table nor-105 malization, and text-guided symbolic reasoning. Building 106 on these solutions, we introduce a framework as a unified recipe, TableMaster. It also incorporates Adaptive Reasoning (AR), a flexible approach that dynamically ad-109

justs between textual and symbolic reasoning, tailoring the reasoning process to each query. (Section 4)

Proposed Solutions

• Extensive Experiments and Detailed Analyses. We conduct extensive experiments and provide in-depth analyses to support our findings on table understanding with language models. Furthermore, we evaluate and demonstrate the superior performance of TableMaster across three widely used table understanding datasets: WikiTQ, TabFact, and FetaQA. Notably, on the WikiTQ dataset, TableMaster achieves an accuracy of 78.13% based on GPT-40-mini, surpassing existing baselines. (Section 3, Section 5, and Appendix)

2. Related Work

Reasoning with Language Models. It has been observed that language models (LMs) can exhibit reasoning abilities when they are sufficiently large (Wei et al., 2022; Suzgun et al., 2022). LMs are now widely used for various reasoning tasks, such as question answering (Kamalloo et al., 2023), decision making (Yang et al., 2023), and mathematical reasoning (Ahn et al., 2024). At the inference stage, techniques such as chain-of-thought prompting (Wei et al., 2023) are used to trigger step-by-step reasoning processes and improve reasoning performance. Few-shot prompting (Brown et al., 2020), least-to-most prompting (Zhou et al., 2023), and program-of-thought prompting (Chen et al., 2023) have proven effective in specific scenarios. Methods like self-

A Recipe to Advance Table Understanding with Language Models



Figure 2. Experimental analysis of challenges in table understanding with language models. (a) Impact of table size on task difficulty. (b) Effect of verbalized tables with enriched semantic context. (c) Performance comparison of different reasoning methods on calculation-required versus non-calculation questions. (d) Performance differences when processing normalized versus noisy tables.

133 consistency (Wang et al., 2023b) and structuring the reason-134 ing process in forms like trees (Yao et al., 2023) or graphs 135 (Besta et al., 2024; Cao, 2024a) are also useful for more 136 complex reasoning tasks. Recently, many works have fo-137 cused on using reinforcement learning (Lightman et al., 138 2023; Uesato et al., 2022) to improve the reasoning abilities 139 of LMs during training. Our work focuses on inference-time 140 improvements and proposes a general framework applicable 141 to all kinds of LMs for table understanding and reasoning. 142

130

131 132

164

Fine-Tuning LMs for Table Understanding. Several stud-143 144 ies have focused on fine-tuning language models to enhance 145 their understanding of tabular data. For example, based 146 on the masked language modeling approach introduced in 147 BERT (Devlin et al., 2019), models like TaPas (Herzig et al., 2020), Pasta (Gu et al., 2022), and TUTA (Wang et al., 148 149 2021) propose specialized pre-training methods to improve 150 LMs' ability to process tables. Similarly, TAPEX (Liu et al., 2022) pre-trains an encoder-decoder model to function as a 151 152 SQL executor, enabling better table comprehension. Recent 153 advancements, such as TableLlama (Zhang et al., 2024a), 154 TableGPT (Zha et al., 2023), and StructLLM (Zhuang et al., 155 2024), leverage open-sourced decoder-only models like 156 Llama (Touvron et al., 2023) to pre-train larger models optimized for various downstream table-related tasks. 157

Adapting LMs for Table Understanding Without Fine-Tuning. Other studies focus on adapting LMs to tablerelated tasks without requiring fine-tuning. For instance, Binder (Cheng et al., 2023), LEVER (Ni et al., 2023), and PoTable (Mao et al., 2024) generate SQL or Python programs, extending the capabilities of LMs to analyze tabular data. Dater (Ye et al., 2023), TabSQLify (Nahid & Rafiei, 2024a), ReAcTable (Zhang et al., 2023), TAP4LLM (Sui et al., 2024), and Tree-of-Table (Ji et al., 2024) introduce different methods to construct sub-tables, modifying the tabular context for improved understanding. Chain-of-Table (Wang et al., 2024) generalizes various table operations, dynamically generating reasoning chains to create sub-tables. MIX-SC (Liu et al., 2024b) employs table normalization and leverages self-consistency, combining results from Python agents and textual reasoning to enhance performance. SpreadsheetEncoder (Dong et al., 2024) is specifically designed to interpret tabular data within spreadsheet environments. Our work also follows this direction to focus on adapting LMs without fine-tuning. We identify key challenges in table understanding and address them through our proposed method, which can be applied to any advanced LMs.

3. Challenges in Table Understanding

As illustrated in Figure 1, we identify and analyze the challenges in table understanding with language models (LMs) through the experiments shown in Figure 2 and related discussions. Additionally, we propose targeted solutions to address these challenges. The detailed settings of the challenge analysis experiment are provided in Appendix B.

Tabular Data Characteristics. Tabular data differs from regular text, which is typically linear and sequential, due to its **structured** nature. Although tabular data can be 165 represented as sequential text, it is fundamentally a twodimensional array of cells. Each cell primarily contains 167 text, but the cells are interconnected and share relation-168 ships with one another. Typically, cells within the same 169 column represent the same feature or type, while cells in the 170 same row correspond to a single data instance. Tables are 171 highly efficient for data representation, often containing a 172 large amount of information, making them inherently data-173 intensive. Moreover, the text in tables is typically concise, 174 consisting of simple words and phrases rather than contin-175 uous sentences, leading to sparse semantic context. Lastly, 176 tables frequently include substantial amounts of numerical 177 data, such as dates, times, scores, and measurements, which 178 often require specialized processing.

180 **3.1. Difficulty in Locating Target Data**

179

181 When LMs encounter tabular data, they often struggle to 182 locate the target data relevant to a given query, leading to 183 misunderstandings. This challenge arises because tabular 184 data is inherently data-intensive, typically containing large 185 volumes of information. Additionally, the structured nature 186 of tabular data makes it challenging for LMs to interpret in-187 dividual cell contents within the broader context of headers 188 and other structural information. This issue can lead to long-189 context hallucination (Huang et al., 2024). Moreover, LMs 190 are prone to neglecting information in the middle of the con-191 text (Liu et al., 2024a), making it even harder to locate target data and further impairing their overall comprehension of 193 the table. (Figure 1 - C1)

195 As shown in Figure $\underline{2}(a)$, we present the changes in table 196 understanding accuracy across four different table size met-197 rics: row count, column count, area size, and token count, 198 ranging from small to extra-large tables. Row count repre-199 sents the number of data entries, while column count reflects 200 the number of dimensions or attributes per entry. Area size 201 is the product of row count and column count, and token 202 count refers to table sizes from the perspective of LMs. All 203 figures indicate that, regardless of the model used, overall 204 performance tends to decline as table size increases. For weaker LMs, the performance drop is more pronounced. 206

To address this, we propose let LMs focusing on specific parts of the table by explicitly constructing a focused subtable that includes only the relevant information needed for the given context. We define this as the **table-of-focus**. By narrowing the scope, table understanding becomes significantly easier, which aligns with both our previous findings and intuition. (Figure $\underline{1} - S1$)

3.2. Table Semantic Deficiency

214

215

216

217

218

219

Tabular data is typically concise, with most cells containing simple words or phrases. Additionally, for each data entry in a row, some descriptive information may reside outside the row, such as in the top header or other structural elements. Understanding a cell in isolation is challenging and often requires a deeper comprehension of the structural relationships within the table. This leads to the problem of **sparse semantic context**, which is fundamentally different from the rich semantic context found in most data used during LMs' pretraining (Dong et al., 2022). The semantic deficiency in tables makes it difficult for LMs to effectively understand and process tabular data. (Figure 1 - C2)

As shown in Figure 2(b), the *Table* represents the case where the LM is provided only with the table input, while the *Table+Verbal* indicates the table along with an additional description, which we refer to as a **verbalized table**. This description is generated by the LMs themselves, whereas *verbal plus* refers to a description produced by more advanced LMs, which can be considered a ground-truth. We observe that verbalization helps LMs perform better on certain tables, leading to a slight overall performance improvement. This effect is more pronounced in weaker LMs, resulting in a 1.5% increase in accuracy. Additionally, the quality of the description plays a crucial role in improvement.

To address this issue, we propose a solution where tables are first verbalized into sequential, natural text as a description and then provided to LMs alongside the original table before they directly tackle table-related tasks. It is similar to table2text (Parikh et al., 2020). This transformation enriches the semantic context, making the data more aligned with the LMs' pretraining, thereby enhancing their ability to effectively understand and process tabular data. (Figure 1 - S2)

3.3. Numerical Inaccuracy in Textual Reasoning

Tabular data often contains numerical values, such as dates, times, scores, and other recorded numbers, and is typically intensive. However, when LMs are used to process numerical data in textual reasoning, they often face significant limitations. LMs are prone to arithmetic calculation errors, especially when dealing with large numbers. LMs are also inefficient at handling iterative processes, particularly when the number of iteration steps is large (Chen et al., 2023). (Figure $\underline{1}$ - C3)

As shown in Figure 2(c), questions that do not require calculations are relatively easier, allowing textual reasoning to achieve a strong performance of 72.4%. However, when calculations are required, performance drops significantly, falling below that of the enhanced symbolic reasoning introduced later. Specifically, textual reasoning suffers a 20.1% decline, whereas enhanced symbolic reasoning experiences a more moderate drop of 7.6%.

Symbolic methods offer a promising solution to these chal-





Figure 3. The framework of *TableMaster*. It comprises three stages: (1) table structure understanding, where the table's structure is analyzed, and a table-of-focus is constructed through row and column lookup; (2) table content understanding, where the table-of-focus is reconstructed based on the question, and its information is verbalized to enhance the semantic context; and (3) table reasoning for question answering, where an adaptive reasoning strategy determines whether to use textual reasoning or text-guided symbolic reasoning to derive the final answer. The dashed arrows indicate optional workflows, such as the table-of-focus re-construction and incorporating text-guided symbolic reasoning.

lenges and have been explored extensively in prior research (Cheng et al., 2023; Ni et al., 2023; Mao et al., 2024). Using symbolic tools, such as SQL or Python programs in combination with LMs, provides an effective approach to handling numerical data in tabular formats. (Figure <u>1</u> - S3)

3.4. Semantic Inflexibility in Symbolic Reasoning

 Symbolic methods excel at arithmetic calculations. However, when prompting LMs to generate code for program of thought reasoning, the performance is suboptimal. Instead of truly understanding the context and generating problemsolving code, LMs often rely on memorized code from the pretraining stage (Yang et al., 2024). We refer to this limitation as **semantic inflexibility**. In table understanding, this challenge is exacerbated by the table's complex structure and concise text content. In real-world scenarios, noisy tables further hinder LMs' symbolic reasoning capabilities. Consequently, while symbolic reasoning with numerical data is highly accurate, the generated code may be incorrect due to issues in program logic or data handling, leading to errors or unintended results. (Figure <u>1</u> - C4)

As shown in Figure 2(c), basic symbolic reasoning performs worse overall, regardless of whether calculations are required. It indicates that basic symbolic reasoning with current LMs is ineffective. Furthermore, as illustrated in Figure 2(d), when processing the same content in a noisy format, symbolic reasoning suffers a larger performance drop of 31.8%, compared to a 20.5% decline for textual reasoning. This highlights the semantic inflexibility of symbolic reasoning when handling noisy tables.

To address this, we first normalize the table structure and
content, ensuring that each column follows a consistent format. We then propose a solution where LMs first engage
in textual reasoning before generating symbolic reasoning

programs. This preliminary textual reasoning step serves as a guide for subsequent symbolic reasoning, improving alignment with the task context. Our approach can be seen as encouraging LMs to think more thoroughly before reasoning, aligning with techniques like plan-and-solve (Wang et al., 2023a). By incorporating textual reasoning as a foundation, we enhance the accuracy and contextual relevance of symbolic reasoning. As demonstrated in Figure 2(c), this method achieves a higher accuracy of 59.1% for calculationrequired questions. (Figure $\underline{1}$ - S4)

4. *TableMaster*: A Recipe for Table Understanding

Based on findings in Section 3, we introduce a recipe and comprehensive framework, *TableMaster*, as shown in Figure 3. It integrates the propose solution proposed in Section 3 into a unified recipe for table understanding. The framework encompasses three key processes: Table Structure Understanding, Table Content Understanding, and Table Reasoning for QA.

4.1. Task Formulation

In table understanding, the objective is to determine an answer A given a table \mathbb{T} and a question or statement Q related to it. The table T is represented as a two-dimensional array of cells,

$$\mathbb{T}_{m \times n} = \begin{bmatrix} C_{1,1} & C_{1,2} & \dots \\ C_{2,1} & C_{i,j} & \dots \\ \vdots & \vdots & \ddots \end{bmatrix}$$

, where $C_{i,j}$ denotes the cell in the *i*-th row and *j*-th column, with the table consisting of *m* rows and *n* columns. In tablebased question answering tasks, *Q* represents a question, and *A* is the expected answer in natural language. In tablebased fact verification tasks, *Q* is a statement about the table's contents, and $A \in \{\text{True, False}\}\)$ is a Boolean value indicating whether the statement is correct. Therefore, the goal is to develop a system \mathcal{F} that can predict the answer accurately based on the table and the given question or statement, formalized as $\mathcal{F}(\mathbb{T}, Q) = A$.

283 **4.2. Table Structure Understanding**

282

317

318

319

320

321

322 323

324

325

327

328

329

The goal of table structure understanding is to analyze the table's structure and construct a Table-of-Focus that contains relevant content for the given question. This process reduces context length and simplifies the table as much as possible.

To enhance the efficiency of the framework, we introduce the **table peek** technique. For structure extraction and certain operations, it is often unnecessary to process the entire table; instead, inspecting only the top rows is sufficient. Given a peek size k, the original table $\mathbb{T}_{m \times n}$ is transformed into a peek table $\mathbb{T}_{k \times n}$, where all columns are retained, but the table is truncated to the first k rows.

Given a wild table \mathbb{T}^W , we first normalize it. We begin by 297 determining whether the table is in row-major or column-298 major format. If it is in column-major format, we transpose 299 it using $\mathbb{T} = \text{Transpose}(\mathbb{T}')$. Next, we normalize and clean 300 all columns containing numerical information, ensuring con-301 sistency in formats such as dates and numerical values, mak-302 ing them directly processable in bulk by a program. After 303 this normalization process, we obtain the normalized table 304 \mathbb{T}^{N} . 305

306 We begin by extracting the top headers H and the key col-307 umn. The top headers are used for column lookup, while 308 the key column serves as the subject or unique identifier for 309 each row. Next, we prompt LMs to perform column lookup 310 and row lookup to identify the relevant rows and columns 311 required for the task. Specifically, for column lookup, we 312 first define the set of candidate columns as $\mathbb{C} = \operatorname{Rank}(H)$. 313 LMs will also rank all candidates based on their relevance 314 to the question. We then prompt the LMs to select b relevant 315 columns based on a given question Q: 316

$$C^0 = \text{Column Lookup}(\mathbb{T}^N \mid Q),$$

where $C^0 = \{c_i \mid c_i \in H\}$ and $|C^0| = b$. For row lookup, we instruct the LMs to generate an SQL query to efficiently filter and select *a* relevant rows *R*:

$$R = \text{Row Lookup}(\mathbb{T}^N \mid Q).$$

Using the identified rows and columns, we construct the initial table-of-focus:

$$\mathbb{T}_{a \times b}^{F} = \text{Table Construction}(\mathbb{T}^{N}, C^{0}, R),$$

which contains only the filtered information necessary for the task.

4.3. Table Content Understanding

The goal of table content understanding is to enrich the semantic context of the table.

Studies have shown that LMs can assess whether sufficient information is available to answer a question (Cao, 2024b; Yin et al., 2023). We first prompt the LMs to estimate whether the constructed Table-of-Focus $\mathbb{T}_{a\times b}^F$, containing C^0 , provides enough information to answer the given question Q. If not, additional column attributes from the candidate column set \mathbf{C} are incrementally added from the ranked candidate headers until sufficient information is available or all relevant top headers have been utilized. Subsequently, a total of a' columns from C are selected for further reasoning. We use re-construction to mitigate information loss during the table-of-focus construction process. The detailed re-construction algorithm can be found in Appendix <u>H</u>.

Once the information sufficiency check is passed, we verbalize the table into natural language, adding descriptions to enrich the semantic context and producing a verbalized table:

$$T^{\mathbb{T}} = \text{Verbalization}(\mathbb{T}_{a \times b}^{F})$$

This verbalized table is represented as sequential natural language text T essentially rather than a structured table, preserving rich semantic context while maintaining a concise size. This transformation enhances information density, further facilitating the LMs' reasoning for the given question.

4.4. Table Reasoning for Question Answering

The goal of this stage is to answer table-related questions by understanding the table precisely and calculating accurately.

We employ an **adaptive reasoning** approach. First, we prompt the LMs to determine the most appropriate reasoning strategy S for the given task. In the instruction, for small tables or those without numerical data, the LMs are allowed to perform textual reasoning directly to derive the final result. For larger tables or those containing numerical data, symbolic reasoning with programmatic execution is selected.

$$S =$$
Strategy Assessment $(\mathbb{T}^F, T^{\mathbb{T}}, Q),$

where $S \in \{T, S\}$ represents the chosen reasoning strategy, with T denoting textual reasoning and S denoting symbolic reasoning.

In symbolic reasoning, we first prompt the LMs to perform textual reasoning to generate guidance G without providing the final result. This intermediate reasoning step is then used as input for symbolic reasoning, transitioning to a text-guided symbolic reasoning approach using programmatic methods. This adaptive method dynamically adjusts

330

Table 1. Performance comparison between *TableMaster* and previous work on WikiTQ and TabFact. The values in the table represent accuracy (%). The best result is **bold**, the second-best result is <u>underlined</u>, and the improvement over the previous best result is highlighted in green. '-' indicates that the result values were not reported in the related papers. Our method outperforms all other methods across both datasets and different language models.

Method		WikiTQ			TabFact	
	gpt-3.5-turbo~175B	gpt-40-mini $_{\sim 8B}$	Llama-3.170B	gpt-3.5-turbo $_{\sim 175B}$	gpt-40-mini $_{\sim 8B}$	Llama-3.170B
Text-to-SQL (Rajkumar et al., 2022)	52.90	-	-	64.71	-	-
End-to-End QA (Wang et al., 2024)	51.84	-	-	70.45	-	-
Few-Shot QA (Wang et al., 2024)	52.56	-	-	71.54	-	-
Chain-of-Thought (Wang et al., 2024)	53.48	-	-	65.37	-	-
ReAcTable (Zhang et al., 2023)	52.50	-	-	74.40	-	-
Binder (Cheng et al., 2023)	56.74	58.86	50.51	79.17	84.63	78.16
Dater (Ye et al., 2023)	52.81	58.33	43.53	78.01	80.98	81.57
TabSQLify (Nahid & Rafiei, 2024a)	<u>64.70</u>	57.02	55.78	79.50	78.75	70.70
Chain-of-Table (Wang et al., 2024)	59.94	55.60	62.22	80.20	84.24	85.62
Tree-of-Table (Ji et al., 2024)	61.11	-	-	81.92	-	-
PoTable (Mao et al., 2024)	-	<u>64.73</u>	<u>65.56</u>	-	<u>88.93</u>	<u>87.06</u>
Ours (TableMaster)	68.21 (+3.51)	78.13 (+13.40)	77.95 (+12.39)	83.65 (+1.73)	90.12 (+1.19)	91.16 (+4.10)

based on the table's size, complexity, and the nature of the question, ensuring accurate and reliable results.

$$A = \begin{cases} \text{Chain-of-Thought}(\mathbb{T}^F, T^{\mathbb{T}}, Q), & \text{if } S = \mathcal{T} \\ \mathcal{P}(\text{Program-of-Thought}(\mathbb{T}^F, T^{\mathbb{T}}, Q \mid G)), & \text{if } S = \mathcal{S} \end{cases}$$

where chain-of-thought and program-of-thought are two prompting techniques, \mathcal{P} represents a Python or SQL program executor, A is the final answer for the current table understanding task.

5. Experiments

5.1. Settings

We conduct extensive experiments to evaluate the performance of *TableMaster*. Specifically, we assess its effectiveness across three different table understanding datasets: WikiTQ (Pasupat & Liang, 2015) (table-based question answering), TabFact (Chen et al., 2020) (table-based fact verification), and FetaQA (Nan et al., 2022) (table-based freeform question answering). For WikiTQ and TabFact, following previous work (Wang et al., 2024; Liu et al., 2024b), we use exact match accuracy as the evaluation metric. For FetaQA, we evaluate performance using BLEU (Papineni et al., 2002) and ROUGE (Lin, 2004) scores. Tables are encoded in Markdown format before being input into language models, with or without addresses, depending on the specific case <u>??</u>.

Our experiments utilize OpenAI models. Unless otherwise stated, we set the temperature to 0 to ensure stable output while keeping all other hyperparameters at their default values. The models used in our evaluation include *gpt-4o* (*gpt-4o-0806*), *gpt-4o-mini* (*gpt-4o-mini-0718*), *gpt-3.5-turbo* (*gpt-3.5-turbo-0125*), *o1* (*o1-preview-0912*), and *o1-mini* (*o1-mini-0912*). Additionally, we evaluate our methods on *Llama-3.1-70B* (*Llama-3.1-70B-Instruct*). The exact number of parameters for several LMs (e.g., GPT, o1) has not been publicly disclosed. Most parameter counts are estimates reported to provide context for understanding model performance. For more precise information, please refer to the original or future official documentation (Abacha et al., 2025).

For comparison, we select several strong baselines, including both classic and state-of-the-art methods such as Binder (Cheng et al., 2023), Dater (Ye et al., 2023), and Chain-of-Table (Wang et al., 2024). Performance results for other methods not in this work are cited directly from their original or related papers, with sources indicated alongside the method names in the results table.

Further analysis and additional experiments on *TableMaster* can be found in the Appendix. The prompts used in *TableMaster* can be found in Appendix \underline{N} , while other prompts used in this work are provided in Appendix \underline{O} .

5.2. Main Results

As shown in Table <u>1</u>, our *TableMaster* approach consistently achieves the highest performance across both WikiTQ and TabFact under different backbone models (*gpt-3.5-turbo*, *gpt-4o-mini*, and *Llama-3.1-70B*). On WikiTQ, *TableMaster* outperforms the strongest baselines by +3.51, +13.40, and +12.39 points, respectively. A similar trend is observed on TabFact, with improvements of +1.73, +1.19, and +4.10 points, demonstrating the robustness of our method across diverse large language models. Results on the FetaQA dataset are provided in Appendix <u>C</u>. These results confirm that *TableMaster* not only generalizes well across different base language models but also significantly enhances table understanding and reasoning in complex QA tasks.

Table 2. Ablation results on WikiTQ and TabFact. The values in the table represent accuracy (%), with \bigtriangledown indicating the performance drop. The red text highlights the drop magnitude. Removing any component from *TableMaster* results in a decrease in performance.

Method	WikiTQ	\bigtriangledown	TabFact	∇
TableMaster (gpt-40-mini)	78.13	_	90.12	-
Structure				
w/o Structure Extraction	74.75	(-3.38)	88.98	(-1.14)
w/o Column Lookup	77.00	(-1.13)	90.51	(-0.40)
w/o Row Lookup	76.59	(-1.54)	89.23	(-0.89
w/o Table of Focus	76.40	(-1.73)	89.33	(-0.79
Content				
w/o Re-Construction	75.55	(-2.58)	89.72	(-0.40
w/o Verbalization	75.78	(-2.35)	89.23	(-0.89
Reasoning				
w/o Textual Reasoning	73.85	(-4.28)	88.39	(-1.73
w/o Symbolic Reasoning	76.10	(-2.03)	89.18	(-0.94
w/o Textual Guidance	75.21	(-2.92)	89.67	(-0.44

Notably, methods such as Binder, Dater, TabSQLify, and Chain-of-Table exhibit subpar performance with *gpt-4omini*, in some cases performing worse than with *gpt-3.5turbo*. Our empirical analysis suggests that these methods primarily rely on symbolic approaches to construct subtables, which often fail to leverage the strengths of chainof-thought reasoning in textual contexts. This limitation underscores the necessity of integrating advanced textual reasoning strategies, as effectively demonstrated by our *TableMaster* approach.

416417**5.3. Ablation Study**

385

386

387

388

389

390

395

396

399

400

401

402

403

404

405

406

407

408

409

410

411

412

413

414

415

To analyze the contribution of each component in *TableMas*-*ter*, we conduct an ablation study on WikiTQ and TabFact.
Table 2 presents the results, and the performance drop from
the full model is highlighted in red. The results demonstrate that removing any component leads to a decrease in
accuracy, confirming the importance of each module in the
overall framework.

425 Structure. The structure understanding components play 426 a crucial role in table comprehension. Removing structure 427 extraction results in a notable accuracy drop of 3.38% on 428 WikiTQ and 1.14% on TabFact, indicating that explicitly 429 extracting the table's structure is essential for effective rea-430 soning, as failing to do so can lead to errors in subsequent 431 steps. Among lookup strategies, removing row lookup leads 432 to a 1.54% decrease in WikiTQ accuracy, whereas removing 433 column lookup results in a smaller drop of 1.13%. This sug-434 gests that row-based information retrieval is more critical 435 than column-based lookup, as large tables typically contain 436 a greater number of rows. Additionally, removing the table-437 of-focus reduces performance by 1.73% on WikiTQ and 438 0.79% on TabFact, further emphasizing its important role in 439

structuring relevant table content to extract key information for reasoning.

Content. Table content understanding also significantly influences performance. Eliminating re-construction, which iteratively refines the Table-of-Focus based on the question, results in a 2.58% accuracy drop on WikiTQ and 0.40% on TabFact, highlighting the importance of this process. Similarly, removing table verbalization, which enriches the semantic context of the table by adding descriptive elements, leads to a 2.35% decrease in WikiTQ accuracy. However, its impact on TabFact is minimal (0.23% drop), suggesting that verbalization becomes even more beneficial for complex table understanding tasks.

Reasoning. The reasoning stage exhibits the most significant performance drop when removed. Removing textual reasoning leads to the largest accuracy decline, with a 4.28% drop on WikiTQ and 1.73% on TabFact, underscoring its necessity for complex reasoning tasks. Similarly, removing symbolic reasoning results in a 2.03% and 0.79% drop on WikiTQ and TabFact, respectively, demonstrating that symbolic reasoning enhances numerical and structured table interpretations. Finally, removing textual guidance, which improves the semantic flexibility of symbolic reasoning, reduces accuracy by 2.92% on WikiTQ and 0.44% on Tab-Fact. This highlights that textual guidance is particularly beneficial and important in symbolic reasoning by ensuring alignment with the problem context.

These results reveal that *TableMaster*'s performance relies on a combination of structural understanding, content adaptation, and an effective reasoning strategy. Among these, structural extraction and reasoning components contribute the most to performance, with textual reasoning being the most critical. Additionally, while content understanding components such as re-construction and verbalization provide meaningful improvements, their impact varies across datasets. These findings validate the effectiveness of our framework in handling diverse table understanding tasks.

6. Conclusion

In this paper, we explore table understanding with language models. Given the characteristics of tabular data, we identify key challenges in table understanding. To overcome these challenges, we propose *TableMaster*, a recipe and comprehensive framework that integrates multiple solutions. Extensive analyses and experiments demonstrate our findings and the effectiveness of *TableMaster*. In the future, we plan to extend and refine the framework to improve its performance across diverse practical applications.

Our code can be found at https://anonymous. 4open.science/r/TableMaster-8646.

440 Impact Statement

This paper presents work whose goal is to advance the field of Machine Learning. There are many potential societal consequences of our work, none of which we feel must be specifically highlighted here.

References

441

442

443

444

445

446 447

- Abacha, A. B., wai Yim, W., Fu, Y., Sun, Z., Yetisgen, M.,
 Xia, F., and Lin, T. Medec: A benchmark for medical error detection and correction in clinical notes, 2025.
 URL https://arxiv.org/abs/2412.19260.
- 453 Ahn, J., Verma, R., Lou, R., Liu, D., Zhang, R., and Yin, 454 W. Large language models for mathematical reason-455 ing: Progresses and challenges. In Falk, N., Papi, S., 456 and Zhang, M. (eds.), Proceedings of the 18th Confer-457 ence of the European Chapter of the Association for 458 Computational Linguistics: Student Research Workshop, 459 pp. 225-237, St. Julian's, Malta, March 2024. Asso-460 ciation for Computational Linguistics. URL https: 461 //aclanthology.org/2024.eacl-srw.17.
- 462 Anil, R., Dai, A. M., Firat, O., Johnson, M., Lepikhin, D., 463 Passos, A., Shakeri, S., Taropa, E., Bailey, P., Chen, Z., 464 Chu, E., Clark, J. H., Shafey, L. E., Huang, Y., Meier-465 Hellstern, K., Mishra, G., Moreira, E., Omernick, M., 466 Robinson, K., Ruder, S., Tay, Y., Xiao, K., Xu, Y., Zhang, 467 Y., Abrego, G. H., Ahn, J., Austin, J., Barham, P., Botha, 468 J., Bradbury, J., Brahma, S., Brooks, K., Catasta, M., 469 Cheng, Y., Cherry, C., Choquette-Choo, C. A., Chowd-470 hery, A., Crepy, C., Dave, S., Dehghani, M., Dev, S., 471 Devlin, J., Díaz, M., Du, N., Dyer, E., Feinberg, V., 472 Feng, F., Fienber, V., Freitag, M., Garcia, X., Gehrmann, 473 S., Gonzalez, L., Gur-Ari, G., Hand, S., Hashemi, H., 474 Hou, L., Howland, J., Hu, A., Hui, J., Hurwitz, J., Isard, 475 M., Ittycheriah, A., Jagielski, M., Jia, W., Kenealy, K., 476 Krikun, M., Kudugunta, S., Lan, C., Lee, K., Lee, B., 477 Li, E., Li, M., Li, W., Li, Y., Li, J., Lim, H., Lin, H., 478 Liu, Z., Liu, F., Maggioni, M., Mahendru, A., Maynez, 479 J., Misra, V., Moussalem, M., Nado, Z., Nham, J., Ni, 480 E., Nystrom, A., Parrish, A., Pellat, M., Polacek, M., 481 Polozov, A., Pope, R., Qiao, S., Reif, E., Richter, B., 482 Riley, P., Ros, A. C., Roy, A., Saeta, B., Samuel, R., 483 Shelby, R., Slone, A., Smilkov, D., So, D. R., Sohn, D., 484 Tokumine, S., Valter, D., Vasudevan, V., Vodrahalli, K., 485 Wang, X., Wang, P., Wang, Z., Wang, T., Wieting, J., 486 Wu, Y., Xu, K., Xu, Y., Xue, L., Yin, P., Yu, J., Zhang, 487 Q., Zheng, S., Zheng, C., Zhou, W., Zhou, D., Petrov, 488 S., and Wu, Y. Palm 2 technical report, 2023. URL 489 https://arxiv.org/abs/2305.10403. 490
- Besta, M., Blach, N., Kubicek, A., Gerstenberger, R., Podstawski, M., Gianinazzi, L., Gajda, J., Lehmann, T., Niewiadomski, H., Nyczyk, P., and Hoefler, T. Graph

of thoughts: Solving elaborate problems with large language models. *Proceedings of the AAAI Conference on Artificial Intelligence*, 38(16):17682–17690, March 2024. ISSN 2159-5399. doi: 10.1609/aaai.v38i16. 29720. URL http://dx.doi.org/10.1609/ aaai.v38i16.29720.

- Brown, T. B., Mann, B., Ryder, N., Subbiah, M., Kaplan, J., Dhariwal, P., Neelakantan, A., Shyam, P., Sastry, G., Askell, A., Agarwal, S., Herbert-Voss, A., Krueger, G., Henighan, T., Child, R., Ramesh, A., Ziegler, D. M., Wu, J., Winter, C., Hesse, C., Chen, M., Sigler, E., Litwin, M., Gray, S., Chess, B., Clark, J., Berner, C., McCandlish, S., Radford, A., Sutskever, I., and Amodei, D. Language models are few-shot learners, 2020. URL https:// arxiv.org/abs/2005.14165.
- Cao, L. GraphReason: Enhancing reasoning capabilities of large language models through a graph-based verification approach. In Dalvi Mishra, B., Durrett, G., Jansen, P., Lipkin, B., Neves Ribeiro, D., Wong, L., Ye, X., and Zhao, W. (eds.), *Proceedings of the 2nd Workshop on Natural Language Reasoning and Structured Explanations* (@ACL 2024), pp. 1–12, Bangkok, Thailand, August 2024a. Association for Computational Linguistics. URL https: //aclanthology.org/2024.nlrse-1.1.
- Cao, L. Learn to refuse: Making large language models more controllable and reliable through knowledge scope limitation and refusal mechanism. In Al-Onaizan, Y., Bansal, M., and Chen, Y.-N. (eds.), *Proceedings of the 2024 Conference on Empirical Methods in Natural Language Processing*, pp. 3628–3646, Miami, Florida, USA, November 2024b. Association for Computational Linguistics. doi: 10.18653/v1/2024.emnlp-main. 212. URL https://aclanthology.org/2024.emnlp-main.212.
- Chen, M., Tworek, J., Jun, H., Yuan, Q., de Oliveira Pinto, H. P., Kaplan, J., Edwards, H., Burda, Y., Joseph, N., Brockman, G., Ray, A., Puri, R., Krueger, G., Petrov, M., Khlaaf, H., Sastry, G., Mishkin, P., Chan, B., Gray, S., Ryder, N., Pavlov, M., Power, A., Kaiser, L., Bavarian, M., Winter, C., Tillet, P., Such, F. P., Cummings, D., Plappert, M., Chantzis, F., Barnes, E., Herbert-Voss, A., Guss, W. H., Nichol, A., Paino, A., Tezak, N., Tang, J., Babuschkin, I., Balaji, S., Jain, S., Saunders, W., Hesse, C., Carr, A. N., Leike, J., Achiam, J., Misra, V., Morikawa, E., Radford, A., Knight, M., Brundage, M., Murati, M., Mayer, K., Welinder, P., Mc-Grew, B., Amodei, D., McCandlish, S., Sutskever, I., and Zaremba, W. Evaluating large language models trained on code, 2021. URL https://arxiv.org/abs/ 2107.03374.

- Chen, W. Large language models are few(1)-shot ta-495 496 ble reasoners. In Vlachos, A. and Augenstein, I. 497 (eds.), Findings of the Association for Computational 498 Linguistics: EACL 2023, pp. 1120-1130, Dubrovnik, 499 Croatia, May 2023. Association for Computational 500 Linguistics. doi: 10.18653/v1/2023.findings-eacl. 501 83. URL https://aclanthology.org/2023. 502 findings-eacl.83.
- 504 Chen, W., Wang, H., Chen, J., Zhang, Y., Wang, H., Li, S.,
 505 Zhou, X., and Wang, W. Y. Tabfact: A large-scale dataset
 506 for table-based fact verification, 2020. URL https:
 507 //arxiv.org/abs/1909.02164.

503

- 508
 509
 509
 510
 510
 511
 511
 512
 512
 512
 514
 515
 516
 517
 517
 518
 519
 519
 510
 510
 510
 511
 512
 512
 512
 512
 512
 514
 515
 515
 516
 517
 517
 518
 518
 519
 519
 510
 510
 510
 511
 512
 512
 512
 512
 512
 512
 512
 512
 512
 512
 512
 512
 512
 512
 512
 512
 512
 512
 512
 512
 512
 512
 512
 512
 512
 512
 512
 512
 512
 512
 512
 512
 512
 512
 512
 512
 512
 512
 512
 512
 512
 512
 512
 512
 512
 512
 512
 512
 512
 512
 512
 512
 512
 512
 512
 512
 512
 512
 512
 512
 512
 512
 512
 512
 512
 512
 512
 512
 512
 512
 512
 512
 512
 512
 512
 512
 512
 512
 512
 512
 512
 512
 512
 512
 512
 512
 512
 512
 512
 512
 512
 512
 512
 512
 512
 512
 512
 512
 512
 512
 512
 512
 512
 512
 512
 512
 512
 512
 512
 512
 512
 512
 512
 512
 512
 512
- 513
 514
 515
 515
 516
 517
 518
 518
 518
 519
 519
 510
 510
 510
 5110
 5111
 512
 512
 512
 513
 514
 514
 515
 515
 516
 517
 518
 518
 518
 518
 518
 519
 510
 510
 510
 511
 512
 512
 513
 514
 515
 515
 516
 517
 518
 517
 518
 518
 518
 518
 518
 518
 518
 518
 518
 518
 518
 518
 518
 518
 518
 518
 518
 518
 518
 518
 518
 518
 518
 518
 518
 518
 518
 518
 518
 518
 518
 518
 518
 518
 518
 518
 518
 518
 518
 518
 518
 518
 518
 518
 518
 518
 518
 518
 518
 518
 518
 518
 518
 518
 518
 518
 518
 518
 518
 518
 518
 518
 518
 518
 518
 518
 518
 518
 518
 518
 518
 518
 518
 518
 518
 518
 518
 518
 518
 518
 518
 518
 518
 518
 518
 518
 518
 518
 518
 518
 518
 518
 518
 518
 518
 518
 518
 518
 518
 518
 518
 518
 518
 518
 518
 518
 518
 518
 518
 518
 518</li
- 519 Devlin, J., Chang, M.-W., Lee, K., and Toutanova, K. Bert:
 520 Pre-training of deep bidirectional transformers for language understanding, 2019. URL https://arxiv.
 522 org/abs/1810.04805.
- Dong, H., Cheng, Z., He, X., Zhou, M., Zhou, A., Zhou,
 F., Liu, A., Han, S., and Zhang, D. Table pre-training:
 A survey on model architectures, pre-training objectives,
 and downstream tasks, 2022. URL https://arxiv.
 org/abs/2201.09745.
- Dong, H., Zhao, J., Tian, Y., Xiong, J., Zhou, M., Lin, Y., 530 Cambronero, J., He, Y., Han, S., and Zhang, D. En-531 coding spreadsheets for large language models. In Al-532 Onaizan, Y., Bansal, M., and Chen, Y.-N. (eds.), Proceed-533 ings of the 2024 Conference on Empirical Methods in 534 Natural Language Processing, pp. 20728–20748, Miami, 535 Florida, USA, November 2024. Association for Compu-536 tational Linguistics. doi: 10.18653/v1/2024.emnlp-main. 537 1154. URL https://aclanthology.org/2024. 538 emnlp-main.1154. 539
- Fang, X., Xu, W., Tan, F. A., Zhang, J., Hu, Z., Qi, Y., Nick-leach, S., Socolinsky, D., Sengamedu, S., and Faloutsos, C. Large language models(llms) on tabular data: Prediction, generation, and understanding a survey, 2024. URL https://arxiv.org/abs/2402.17944.
- Ghasemi, M. and Amyot, D. Process mining in healthcare:
 a systematised literature review. *International Journal of Electronic Healthcare*, 9(1):60–88, 2016.

- Gu, Z., Fan, J., Tang, N., Nakov, P., Zhao, X., and Du, X. Pasta: Table-operations aware fact verification via sentence-table cloze pre-training, 2022. URL https: //arxiv.org/abs/2211.02816.
- Gunasekar, S., Zhang, Y., Aneja, J., Mendes, C. C. T., Giorno, A. D., Gopi, S., Javaheripi, M., Kauffmann, P., de Rosa, G., Saarikivi, O., Salim, A., Shah, S., Behl, H. S., Wang, X., Bubeck, S., Eldan, R., Kalai, A. T., Lee, Y. T., and Li, Y. Textbooks are all you need, 2023. URL https://arxiv.org/abs/2306.11644.
- Herzig, J., Nowak, P. K., Müller, T., rancesco Piccinno, and Eisenschlos, J. M. TAPAS: weakly supervised table parsing via pre-training. *CoRR*, abs/2004.02349, 2020. URL https://arxiv.org/abs/2004.02349.
- Huang, L., Yu, W., Ma, W., Zhong, W., Feng, Z., Wang, H., Chen, Q., Peng, W., Feng, X., Qin, B., and Liu, T. A survey on hallucination in large language models: Principles, taxonomy, challenges, and open questions. ACM Transactions on Information Systems, nov 2024. ISSN 1558-2868. doi: 10.1145/3703155. URL http://dx.doi.org/10.1145/3703155.
- Ji, D., Zhu, L., Gao, S., Xu, P., Lu, H., Ye, J., and Zhao, F. Tree-of-table: Unleashing the power of llms for enhanced large-scale table understanding, 2024. URL https: //arxiv.org/abs/2411.08516.
- Kamalloo, E., Dziri, N., Clarke, C., and Rafiei, D. Evaluating open-domain question answering in the era of large language models. In Rogers, A., Boyd-Graber, J., and Okazaki, N. (eds.), *Proceedings of the 61st Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pp. 5591–5606, Toronto, Canada, July 2023. Association for Computational Linguistics. doi: 10.18653/v1/2023.acl-long.307. URL https: //aclanthology.org/2023.acl-long.307.
- Li, Y., Huang, Z., Yan, J., Zhou, Y., Ye, F., and Liu, X. Gfte: Graph-based financial table extraction, 2020. URL https://arxiv.org/abs/2003.07560.
- Lightman, H., Kosaraju, V., Burda, Y., Edwards, H., Baker, B., Lee, T., Leike, J., Schulman, J., Sutskever, I., and Cobbe, K. Let's verify step by step, 2023. URL https: //arxiv.org/abs/2305.20050.
- Lin, C.-Y. ROUGE: A package for automatic evaluation of summaries. In *Text Summarization Branches Out*, pp. 74–81, Barcelona, Spain, July 2004. Association for Computational Linguistics. URL https: //aclanthology.org/W04-1013/.
- Liu, N. F., Lin, K., Hewitt, J., Paranjape, A., Bevilacqua, M., Petroni, F., and Liang, P. Lost in the middle: How

language models use long contexts. *Transactions of the Association for Computational Linguistics*, 12:157–173,
2024a. doi: 10.1162/tacl_a_00638. URL https://
aclanthology.org/2024.tacl-1.9.

Liu, Q., Chen, B., Guo, J., Ziyadi, M., Lin, Z., Chen, W.,
and Lou, J.-G. Tapex: Table pre-training via learning
a neural sql executor, 2022. URL https://arxiv.
org/abs/2107.07653.

559 Liu, T., Wang, F., and Chen, M. Rethinking tabular data 560 understanding with large language models. In Duh, 561 K., Gomez, H., and Bethard, S. (eds.), Proceedings of 562 the 2024 Conference of the North American Chapter of 563 the Association for Computational Linguistics: Human 564 Language Technologies (Volume 1: Long Papers), pp. 565 450-482, Mexico City, Mexico, June 2024b, Associa-566 tion for Computational Linguistics. doi: 10.18653/v1/ 567 2024.naacl-long.26. URL https://aclanthology. 568 org/2024.naacl-long.26. 569

Mao, Q., Liu, Q., Li, Z., Cheng, M., Zhang, Z., and Li,
R. Potable: Programming standardly on table-based
reasoning like a human analyst, 2024. URL https:
//arxiv.org/abs/2412.04272.

574 Maynez, J., Agrawal, P., and Gehrmann, S. Benchmarking 575 large language model capabilities for conditional gen-576 eration. In Rogers, A., Boyd-Graber, J., and Okazaki, 577 N. (eds.), Proceedings of the 61st Annual Meeting of 578 the Association for Computational Linguistics (Volume 579 1: Long Papers), pp. 9194–9213, Toronto, Canada, 580 July 2023. Association for Computational Linguistics. 581 doi: 10.18653/v1/2023.acl-long.511. URL https: 582 //aclanthology.org/2023.acl-long.511/. 583

- Minaee, S., Mikolov, T., Nikzad, N., Chenaghlu, M., Socher,
 R., Amatriain, X., and Gao, J. Large language models:
 A survey, 2024. URL https://arxiv.org/abs/
 2402.06196.
- Nahid, M. and Rafiei, D. TabSQLify: Enhancing reason-589 ing capabilities of LLMs through table decomposition. 590 In Duh, K., Gomez, H., and Bethard, S. (eds.), Pro-591 ceedings of the 2024 Conference of the North Ameri-592 can Chapter of the Association for Computational Lin-593 guistics: Human Language Technologies (Volume 1: 594 Long Papers), pp. 5725-5737, Mexico City, Mexico, 595 June 2024a. Association for Computational Linguistics. 596 doi: 10.18653/v1/2024.naacl-long.320. URL https:// 597 aclanthology.org/2024.naacl-long.320. 598
- Nahid, M. M. H. and Rafiei, D. NormTab: Improving symbolic reasoning in LLMs through tabular data normalization. In Al-Onaizan, Y., Bansal, M., and Chen, Y.-N. (eds.), *Findings of the Association for Computational Linguistics: EMNLP 2024*, pp. 3569–3585, Miami, Florida,

USA, November 2024b. Association for Computational Linguistics. doi: 10.18653/v1/2024.findings-emnlp. 203. URL https://aclanthology.org/2024. findings-emnlp.203/.

- Nan, L., Hsieh, C., Mao, Z., Lin, X. V., Verma, N., Zhang, R., Kryściński, W., Schoelkopf, H., Kong, R., Tang, X., Mutuma, M., Rosand, B., Trindade, I., Bandaru, R., Cunningham, J., Xiong, C., Radev, D., and Radev, D. Fe-TaQA: Free-form table question answering. *Transactions of the Association for Computational Linguistics*, 10:35–49, 2022. doi: 10.1162/tacl_a_00446. URL https://aclanthology.org/2022.tacl-1.3/.
- Ni, A., Iyer, S., Radev, D., Stoyanov, V., tau Yih, W., Wang, S. I., and Lin, X. V. Lever: Learning to verify languageto-code generation with execution, 2023. URL https: //arxiv.org/abs/2302.08468.
- OpenAI. Gpt-4 technical report, 2024. URL https://arxiv.org/abs/2303.08774.
- Papineni, K., Roukos, S., Ward, T., and Zhu, W.-J. Bleu: a method for automatic evaluation of machine translation. In Isabelle, P., Charniak, E., and Lin, D. (eds.), *Proceedings of the 40th Annual Meeting of the Association for Computational Linguistics*, pp. 311–318, Philadelphia, Pennsylvania, USA, July 2002. Association for Computational Linguistics. doi: 10.3115/1073083.1073135. URL https://aclanthology.org/P02-1040/.
- Parikh, A. P., Wang, X., Gehrmann, S., Faruqui, M., Dhingra, B., Yang, D., and Das, D. Totto: A controlled table-to-text generation dataset, 2020. URL https: //arxiv.org/abs/2004.14373.
- Pasupat, P. and Liang, P. Compositional semantic parsing on semi-structured tables. In Zong, C. and Strube, M. (eds.), Proceedings of the 53rd Annual Meeting of the Association for Computational Linguistics and the 7th International Joint Conference on Natural Language Processing (Volume 1: Long Papers), pp. 1470–1480, Beijing, China, July 2015. Association for Computational Linguistics. doi: 10.3115/v1/P15-1142. URL https://aclanthology.org/P15–1142.
- Plaat, A., Wong, A., Verberne, S., Broekens, J., van Stein, N., and Back, T. Reasoning with large language models, a survey, 2024. URL https://arxiv.org/abs/ 2407.11511.
- Rajkumar, N., Li, R., and Bahdanau, D. Evaluating the text-to-sql capabilities of large language models, 2022. URL https://arxiv.org/abs/2204.00498.
- Sui, Y., Zou, J., Zhou, M., He, X., Du, L., Han, S., and Zhang, D. TAP4LLM: Table provider on

sampling, augmenting, and packing semi-structured 605 606 data for large language model reasoning. In Al-607 Onaizan, Y., Bansal, M., and Chen, Y.-N. (eds.), Find-608 ings of the Association for Computational Linguistics: EMNLP 2024, pp. 10306-10323, Miami, Florida, 609 USA, November 2024. Association for Computational 610 611 Linguistics. doi: 10.18653/v1/2024.findings-emnlp. 612 603. URL https://aclanthology.org/2024. 613 findings-emnlp.603.

- Suzgun, M., Scales, N., Schärli, N., Gehrmann, S., Tay,
 Y., Chung, H. W., Chowdhery, A., Le, Q. V., Chi, E. H.,
 Zhou, D., and Wei, J. Challenging big-bench tasks and
 whether chain-of-thought can solve them, 2022. URL
 https://arxiv.org/abs/2210.09261.
- Touvron, H., Lavril, T., Izacard, G., Martinet, X., Lachaux,
 M.-A., Lacroix, T., Rozière, B., Goyal, N., Hambro,
 E., Azhar, F., Rodriguez, A., Joulin, A., Grave, E., and
 Lample, G. Llama: Open and efficient foundation language models, 2023. URL https://arxiv.org/
 abs/2302.13971.
- 627 Uesato, J., Kushman, N., Kumar, R., Song, F., Siegel, N.,
 628 Wang, L., Creswell, A., Irving, G., and Higgins, I. Solv629 ing math word problems with process- and outcome630 based feedback, 2022. URL https://arxiv.org/
 631 abs/2211.14275.
- 632 Wang, L., Xu, W., Lan, Y., Hu, Z., Lan, Y., Lee, R. K.-633 W., and Lim, E.-P. Plan-and-solve prompting: Improv-634 ing zero-shot chain-of-thought reasoning by large lan-635 guage models. In Rogers, A., Boyd-Graber, J., and 636 Okazaki, N. (eds.), Proceedings of the 61st Annual Meet-637 ing of the Association for Computational Linguistics (Vol-638 ume 1: Long Papers), pp. 2609–2634, Toronto, Canada, 639 July 2023a. Association for Computational Linguistics. 640 doi: 10.18653/v1/2023.acl-long.147. URL https: 641 //aclanthology.org/2023.acl-long.147. 642
- Wang, X., Wei, J., Schuurmans, D., Le, Q., Chi, E., Narang,
 S., Chowdhery, A., and Zhou, D. Self-consistency improves chain of thought reasoning in language models, 2023b. URL https://arxiv.org/abs/2203.
 11171.
- Wang, Z., Dong, H., Jia, R., Li, J., Fu, Z., Han, S., and
 Zhang, D. Tuta: Tree-based transformers for generally
 structured table pre-training. In *Proceedings of the 27th ACM SIGKDD Conference on Knowledge Discovery & Data Mining*, pp. 1780–1790, 2021.
- Wang, Z., Zhang, H., Li, C.-L., Eisenschlos, J. M., Perot,
 V., Wang, Z., Miculicich, L., Fujii, Y., Shang, J., Lee,
 C.-Y., and Pfister, T. Chain-of-table: Evolving tables in
 the reasoning chain for table understanding, 2024. URL
 https://arxiv.org/abs/2401.04398.

- Wei, J., Tay, Y., Bommasani, R., Raffel, C., Zoph, B., Borgeaud, S., Yogatama, D., Bosma, M., Zhou, D., Metzler, D., Chi, E. H., Hashimoto, T., Vinyals, O., Liang, P., Dean, J., and Fedus, W. Emergent abilities of large language models, 2022. URL https: //arxiv.org/abs/2206.07682.
- Wei, J., Wang, X., Schuurmans, D., Bosma, M., Ichter, B., Xia, F., Chi, E., Le, Q., and Zhou, D. Chain-ofthought prompting elicits reasoning in large language models, 2023. URL https://arxiv.org/abs/ 2201.11903.
- Yang, S., Nachum, O., Du, Y., Wei, J., Abbeel, P., and Schuurmans, D. Foundation models for decision making: Problems, methods, and opportunities, 2023. URL https://arxiv.org/abs/2303.04129.
- Yang, Y., Xiong, S., Payani, A., Shareghi, E., and Fekri, F. Can LLMs reason in the wild with programs? In Al-Onaizan, Y., Bansal, M., and Chen, Y.-N. (eds.), *Findings of the Association for Computational Linguistics: EMNLP 2024*, pp. 9806–9829, Miami, Florida, USA, November 2024. Association for Computational Linguistics. doi: 10.18653/v1/2024.findings-emnlp. 573. URL https://aclanthology.org/2024. findings-emnlp.573/.
- Yao, S., Yu, D., Zhao, J., Shafran, I., Griffiths, T. L., Cao, Y., and Narasimhan, K. Tree of thoughts: Deliberate problem solving with large language models, 2023. URL https://arxiv.org/abs/2305.10601.
- Ye, Y., Hui, B., Yang, M., Li, B., Huang, F., and Li, Y. Large language models are versatile decomposers: Decompose evidence and questions for table-based reasoning, 2023. URL https://arxiv.org/abs/2301.13808.
- Yin, Z., Sun, Q., Guo, Q., Wu, J., Qiu, X., and Huang, X. Do large language models know what they don't know? In Rogers, A., Boyd-Graber, J., and Okazaki, N. (eds.), *Findings of the Association for Computational Linguistics: ACL 2023*, pp. 8653–8665, Toronto, Canada, July 2023. Association for Computational Linguistics. doi: 10.18653/v1/2023.findings-acl. 551. URL https://aclanthology.org/2023. findings-acl.551.
- Zha, L., Zhou, J., Li, L., Wang, R., Huang, Q., Yang, S., Yuan, J., Su, C., Li, X., Su, A., Zhang, T., Zhou, C., Shou, K., Wang, M., Zhu, W., Lu, G., Ye, C., Ye, Y., Ye, W., Zhang, Y., Deng, X., Xu, J., Wang, H., Chen, G., and Zhao, J. Tablegpt: Towards unifying tables, nature language and commands into one gpt, 2023. URL https://arxiv.org/abs/2307.08674.

- Kang, T., Yue, X., Li, Y., and Sun, H. Tablellama: Towards
 open large generalist models for tables, 2024a. URL
 https://arxiv.org/abs/2311.09206.
- Kang, X., Wang, D., Dou, L., Zhu, Q., and Che, W. A survey of table reasoning with large language models, 2024b.
 URL https://arxiv.org/abs/2402.08259.
- ⁶⁶⁷
 ⁶⁶⁸
 ⁶⁶⁹
 ⁶⁶⁹
 ⁶⁷⁰
 ⁶⁷⁰
 ⁶⁷¹
 ⁶⁷¹
 ⁶⁷¹
 ⁶⁷²
 ⁶⁷³
 ⁶⁷⁴
 ⁶⁷⁵
 ⁶⁷⁵
 ⁶⁷⁶
 ⁶⁷⁶
 ⁶⁷⁷
 ⁶⁷⁷
 ⁶⁷⁸
 ⁶⁷⁸
 ⁶⁷⁹
 ⁶⁷⁹
 ⁶⁷⁹
 ⁶⁷⁰
 ⁶⁷¹
 ⁶⁷¹
 ⁶⁷¹
 ⁶⁷²
 ⁶⁷³
 ⁶⁷⁴
 ⁶⁷⁵
 ⁶⁷⁵
 ⁶⁷⁶
 ⁶⁷⁶
 ⁶⁷⁶
 ⁶⁷⁷
 ⁶⁷⁷
 ⁶⁷⁸
 ⁶⁷⁸
 ⁶⁷⁹
 ⁶⁷⁹
 ⁶⁷⁹
 ⁶⁷⁹
 ⁶⁷⁹
 ⁶⁷⁰
 ⁶⁷¹
 ⁶⁷¹
 ⁶⁷¹
 ⁶⁷²
 ⁶⁷³
 ⁶⁷⁴
 ⁶⁷⁵
 ⁶⁷⁵
 ⁶⁷⁵
 ⁶⁷⁶
 ⁶⁷⁶
 ⁶⁷⁶
 ⁶⁷⁷
 ⁶⁷⁶
 ⁶⁷⁶
 ⁶⁷⁶
 ⁶⁷⁷
 ⁶⁷⁷
 ⁶⁷⁸
 ⁶⁷⁸
 ⁶⁷⁹
 ⁶⁷⁹
 ⁶⁷⁹
 ⁶⁷⁹
 ⁶⁷⁹
 ⁶⁷⁹
 ⁶⁷⁰
 ⁶⁷⁰
 ⁶⁷¹
 ⁶⁷¹
 ⁶⁷¹
 ⁶⁷²
 ⁶⁷⁵
 ⁶⁷⁵
 ⁶⁷⁵
 ⁶⁷⁵
 ⁶⁷⁶
 ⁶⁷⁶
 ⁶⁷⁶
 ⁶⁷⁷
 ⁶⁷⁶
 ⁶⁷⁶
 ⁶⁷⁶
 ⁶⁷⁶
 ⁶⁷⁶
 ⁶⁷⁷
 ⁶⁷⁶
 <
- Zhou, D., Schärli, N., Hou, L., Wei, J., Scales, N., Wang, X.,
 Schuurmans, D., Cui, C., Bousquet, O., Le, Q., and Chi,
 E. Least-to-most prompting enables complex reasoning in
 large language models, 2023. URL https://arxiv.
 org/abs/2205.10625.
- Zhu, Y., Moniz, J. R. A., Bhargava, S., Lu, J., Piraviperu-mal, D., Li, S., Zhang, Y., Yu, H., and Tseng, B.-H. Can large language models understand context? In Graham, Y. and Purver, M. (eds.), Findings of the Association for Computational Linguistics: EACL 2024, pp. 2004–2018, St. Julian's, Malta, March 2024. Association for Compu-tational Linguistics. URL https://aclanthology. org/2024.findings-eacl.135.
- Zhuang, A., Zhang, G., Zheng, T., Du, X., Wang, J.,
 Ren, W., Huang, S. W., Fu, J., Yue, X., and Chen,
 W. Structlm: Towards building generalist models for
 structured knowledge grounding, 2024. URL https:
 //arxiv.org/abs/2402.16671.

A Recipe to Advance Table Understanding with Language Models

٨		
A	Limitations, Extendability, and Future Works	15
	A.1 Technical Refinement	15
	A.2 Downstream Applications	15
B	Detailed Settings of Challenge Analysis Experiments	16
С	Experimental Results on the FetaQA Dataset	16
D	Table Understanding Baselines	17
E	Performance Analysis Under Different Table Sizes	18
F	Performance Analysis Under Different Table Peek Sizes	20
G	Efficiency Analysis of TableMaster	21
	G.1 Theoretical Analysis	21
	G.2 Empirical Analysis	21
H	Detailed Algorithm of Table-of-Focus Re-Construction	22
I	Analysis of Table-of-Focus Re-Construction	23
J	Analysis of Adaptive Reasoning	23
K	Information Missing and Table Reasoning with Full Table	25
L	Case Study of Table Verbalization	25
M	Case Study of TableMaster	27
N	Prompt Design in TableMaster	28
0	Prompts Used in Analysis Experiments	37
Р	Notion Table	42

A. Limitations, Extendability, and Future Works

Although we conduct extensive experiments and in-depth analysis, this work still has certain limitations. However, we believe that *TableMaster* possesses extensibility, allowing for future refinements. These improvements may include technical refinement as well as optimizing its application in downstream applications.

A.1. Technical Refinement

Wild Table. In our experiments, the tables in the three datasets we use are already cleaned; therefore, we do not explicitly implement table normalization in our evaluation experiments. However, we conduct analysis experiments to highlight the importance of table normalization for handling wild tables. In practical scenarios, various tools are available for table normalization. Regular expression matching can be employed for formatting, and small language models can also be leveraged to efficiently process and normalize tables (Nahid & Rafiei, 2024b).

Hierarchy Table. In our work, we assume all tables are flat, allowing for straightforward utilization and extraction of structural information. However, many real-world tables are hierarchical, where data is organized in a tree structure, making table structure understanding more challenging. We envision two possible solutions: converting hierarchical tables into flat tables or designing a tree-based structure extraction method to effectively locate target data.

Table-of-Focus Construction. In designing the Table-of-Focus, we employ two efficient methods: LM prompting for column lookup and SQL generation for row lookup. The Table-of-Focus is then constructed based on the results of these two lookups. Many previous works (Ji et al., 2024; Wang et al., 2024) have introduced complex approaches for extracting relevant sub-tables. In contrast, our method remains intentionally simple, prioritizing efficiency and adaptability. We believe that in the future, more advanced techniques may emerge to further enhance the extraction of key information.

Table Verbalization. To facilitate the implementation of *TableMaster*, we utilize language models themselves to verbalize the table. However, the quality of the generated text is not optimal due to the challenges of open-ended text generation. Several existing studies, such as Table-to-Text (Parikh et al., 2020), have explored this sub-task. In the future, we can enhance performance and efficiency by replacing this step with specifically trained small language models, which could further improve the semantic density of the verbalized table.

Adaptive Reasoning. Adaptive reasoning can be unstable, as language models may not always select the optimal strategy. We further explore this issue in Appendix J. In the future, training a dedicated machine learning model to guide LMs in selecting the most effective reasoning strategy could improve stability and performance.

Information Missing. The construction of the Table-of-Focus involves a trade-off between precision and recall. If recall is insufficient, essential information may be missing for final reasoning, while low precision can render the extracted content less useful. Although we use re-construction to mitigate information loss during the Table-of-Focus construction process, our analysis reveals that some information missing persist in row lookup. We further investigate this issue in Appendix <u>K</u>.

Efficiency. Efficiency is crucial in table processing and table understanding. To enhance efficiency, we incorporate the table peek technique, which reduces the context that language models need to process at certain steps. We further explore this technique in Appendix <u>F</u> and analyze the overall efficiency in Appendix <u>G</u>. In real-world applications, for optimal efficiency, we consider replacing certain steps with specialized small language models, balancing the trade-off between efficiency and performance .

A.2. Downstream Applications

Web Tables. The web contains a vast number of structured tables, including Wikipedia tables, government reports, and other online tabular data. Extracting and reasoning over these tables is crucial for applications such as fact verification, web search, and knowledge graph construction. *TableMaster* enhances the ability to interpret, query, and reason over these tables, enabling more accurate and context-aware information retrieval.

Spreadsheets. Spreadsheets are widely used in business, finance, and scientific research for data management and analysis. Traditional spreadsheet tools require manual formula creation and human intervention to derive insights. In contrast, *TableMaster* can automate tasks such as data summarization, trend analysis, anomaly detection, and reasoning-based computations. By integrating with tools like Microsoft Excel and Google Sheets, *TableMaster* enables intelligent spreadsheet interactions, allowing users to query data using natural language and receive precise, structured responses.

Databases. Structured databases store vast amounts of relational data, typically accessed through SQL queries or predefined interfaces. However, many users lack SQL proficiency, posing barriers to efficient data retrieval. *TableMaster*, with its Tableof-Focus mechanism, facilitates the quick understanding of large databases, enabling seamless querying of relational data without the need for manual SQL query writing. Additionally, it enhances database reasoning tasks, including knowledge extraction, making structured data more accessible to non-technical users.

In real-world applications, different scenarios have varying requirements, and it may not be necessary to incorporate all
 aspects of *TableMaster*. Instead, certain components can be adapted or selectively applied based on specific needs.

Finally, as discussed above, there is still much work to be done in the future to further enhance language model-based table
understanding. We hope this work serves as a recipe of comprehensive references on current state-of-the-art methods and
provides guidance for future advancements in this field.

B. Detailed Settings of Challenge Analysis Experiments

We conduct extensive experiments to analyze the challenges of table understanding with language models (LMs). Specifically, we perform challenge analysis experiments on the WikiTQ dataset (Pasupat & Liang, 2015), which consists of 4,344 data instances. Following previous work (Wang et al., 2024; Liu et al., 2024b), we use the exact match of the final answer as the evaluation metric to measure accuracy. Our experiments utilize OpenAI models. Unless otherwise stated, we set the temperature to 0 to ensure stable output while keeping all other hyperparameters at their default values. For each model, we use the following versions: *gpt-4o* (*gpt-4o-0806*), *gpt-4o-mini* (*gpt-4o-mini-0718*), *gpt-3.5-turbo* (*gpt-3.5-turbo-0125*), and *o1* (*o1-preview-0912*).

Effect of Table Size (Figure 2(a)). We evaluate how table size impacts task difficulty using a direct prompting approach (Prompt 21) with *gpt-4o*, *gpt-4o-mini* and *gpt-3.5-turbo* to generate answers. We categorize table size based on four metrics: row count, column count, area size (computed as the product of row and column counts), and token count (measured using the *cl100k_base* encoding). The tables are divided into four size categories—small, medium, large, and extra-large—strictly partitioned into quartiles from the smallest to the largest. We then analyze results by splitting performance based on table size.

Effect of Verbalization (Figure 2(b)). We investigate the impact of enriching semantic context through verbalized tables by comparing three approaches. In the *Table* setting, the LM processes the raw table directly using direct prompting (Prompt 21). In *Table + Verbal*, the table is first verbalized using the LM itself (Prompt 24), and both the original and verbalized tables are then provided as input. Lastly, in *Table + Verbal Plus*, the verbalized table is generated using *gpt-4o*, further enhancing the semantic richness of the input.

Comparison of Reasoning Methods (Figure 2(c)). We compare different reasoning approaches—textual reasoning (Prompt 22), symbolic reasoning (Prompt 23), and text-guided symbolic reasoning (Prompt 25)—on calculation-required versus non-calculation questions using *gpt-3.5-turbo*. To classify WikiTQ questions into calculation-required or not, we use *o1* (Prompt 28), identifying 2,692 calculation-required questions and 1,652 non-calculation questions. The results are then analyzed based on this classification.

Impact of Noisy Tables (Figure 2(d)). We investigate how performance varies between normalized and noisy tables. To generate noisy tables, we use *o1* (Prompt 29), instructing it to introduce noise into table contents while preserving actual values and diversifying entries within columns. Additionally, each table has a 50% chance of being randomly transposed from the default row-major format to the column-major format. We then filter the generated tables through a combination of human verification and *o1* checks to ensure that answers remain derivable from the noisy tables. After filtering, 2,565 noisy tables remain. We evaluate textual reasoning (Prompt 22) and symbolic reasoning (Prompt 23) on both the noisy and original normalized tables using *gpt-4o-mini*.

C. Experimental Results on the FetaQA Dataset

PaLM 2 has been deprecated (Anil et al., 2023) and is no longer accessible. Therefore, we use a comparable language
model, *gpt-4o*, to conduct experiments on FetaQA and compare the results with previous methods. Additionally, we use 20
exemplars for few-shot in-context learning to align with the dataset's format.

Table <u>3</u> shows that *TableMaster* improves free-form question answering performance on FetaQA compared to the base

878 879

877

871 872

Table 3. Performance comparison on FetaQA. The values are multiplied by 100, and the percentage improvement represents the performance gain compared to the end-to-end QA of the base model. The results demonstrate that *TableMaster* achieves strong performance in long-form question answering.

Methods	BLEU	ROUGE-1	ROUGE-2	ROUGE-L
Fine-Tuning (T5-large) (Ye et al., 2023)	30.54	63	41	53
End-to-End QA (Codex) (Chen et al., 2021)	27.96	62	40	52
End-to-End QA (PaLM 2) (Wang et al., 2024)	28.37	63	41	53
Dater (PaLM 2) (Ye et al., 2023)	29.47	63	41	53
Chain-of-Table (PaLM 2) (Wang et al., 2024)	32.61 (+14.9%)	66 (+4.8%)	44 (+7.3%)	56 (+5.7%)
End-to-End QA (gpt-4o) Ours (Tablemaster - gpt4o)	24.91 28.94 (+16.2%)	62.05 66.06 (+6.5%)	41.29 45.29 (+9.7%)	50.36 54.56 (+ 8.3%)

Question: How did Napolitano perform compared to the other candidates?



Figure 4. An example (fetaqa-164) from the FetaQA dataset where the result is accurate, but the evaluation metric assigns a low score.

End-to-End QA model, achieving improvements of 16.2% in BLEU and 6.5% in ROUGE-1. These improvements surpass those of Chain-of-Table when compared to its respective End-to-End QA baseline.

916 However, the improvement of *TableMaster* over baseline methods remains marginal, with some values even falling below 917 those of previous approaches in absolute terms. We believe this does not fully reflect the model's actual performance in 918 free-form QA. We attribute this to the n-gram text similarity metrics used in ROUGE-1/2/L (Lin, 2004), which are known to 919 be insensitive to improvements gained from in-context learning (Maynez et al., 2023). These metrics struggle to capture 920 stylistic and structural enhancements in free-form text generation. Since models rely on instructions and a limited number of 921 examples, they may not fully adapt to the expected output format, leading to an underestimation of performance gains.

To further investigate this, we analyze a specific case, FetaQA-164, as shown in Figure <u>4</u>. In this instance, the BLEU and ROUGE metrics assign low scores, as only two words match in the entire sentence. However, manual review confirms that the generated answer is indeed correct—these two words are the most important, and the overall meaning of the response is both accurate and superior to the ground truth. This highlights the limitations of ROUGE in evaluating free-form QA and suggests that qualitative analysis is essential for a more comprehensive assessment of model improvements. Nonetheless, based on quantitative analysis, *TableMaster* is overall effective.

D. Table Understanding Baselines

To better facilitate future research, we evaluate different reasoning methods across various base models. Table <u>4</u> presents the accuracy results of our reproduced baselines on WikiTQ and TabFact, comparing different base LLMs and reasoning methods. The table includes evaluations on *o1-preview* (\sim 300B), *o1-mini* (\sim 100B), *gpt-4o* (\sim 200B), *gpt-4o-mini* (\sim 8B),

A Recipe to Advance Table Understanding with Language Models

Base LLM	Method	WikiTQ	TabFact
o1-preview~300B	Direct	84.60	92.05
o1-mini~100B	Direct	83.49	91.35
	Direct	73.07	84.73
ant la	Chain of Thought	83.98	91.90
gpt-40~200B	Program of Thought	74.63	90.02
	TableMaster (gpt-40)	84.55	94.52
	Direct	59.53	71.25
ant la mini	Chain of Thought	72.97	87.40
gpt-40-mm _{∼8B}	Program of Thought	61.83	85.18
	TableMaster (gpt-40-mini)	78.13	90.12
	Direct	56.58	70.90
254 1	Chain of Thought	59.92	69.52
gpt-3.3-turbo $\sim 175B$	Program of Thought	50.32	68.82
	TableMaster (gpt-3.5-turbo)	68.21	83.65

935

949 950 951

977

Table 4. Results of our reproduced baselines on WikiTQ and TabFact. The values in the table represent accuracy (%).

and *gpt-3.5-turbo* (\sim 175B). Each model is tested with various reasoning strategies, including Direct, chain of thought, and Program of Thought, alongside our proposed *TableMaster*.

Across all base models, TableMaster consistently achieves the highest accuracy. For gpt-40, TableMaster reaches 84.55% on
WikiTQ and 94.52% on TabFact, outperforming both chain of thought (83.98%, 91.90%) and Program of Thought (74.63%,
90.02%). Similarly, for gpt-40-mini, TableMaster achieves 78.13% on WikiTQ and 90.12% on TabFact, significantly
improving over the Direct method (59.53%, 71.25%) and surpassing chain of thought (72.97%, 87.40%).

The performance gap is even more pronounced for gpt-3.5-turbo, where TableMaster reaches 68.21% on WikiTQ and 83.65% on TabFact, significantly outperforming both chain of thought (59.92%, 69.52%) and Program of Thought (50.32%, 68.82%). Interestingly, we observe that while *TableMaster* 's improvement is limited on gpt-40, the weaker the base model, the greater the performance improvement. While o1-preview and o1-mini achieve high accuracy with the Direct method (84.60%, 92.05% for o1-preview and 83.49%, 91.35% for o1-mini), the results of *TableMaster* on gpt-40 demonstrate that our method is capable of achieving state-of-the-art performance across different LL architectures.

Additionally, we find that chain of thought reasoning is highly effective, achieving strong accuracy across models. Even a simple chain of thought approach outperforms previous methods that rely solely on symbolic reasoning (Mao et al., 2024), indicating that chain of thought should be retained as a key component in the reasoning framework.

These results confirm that *TableMaster* enhances table reasoning performance across various LLMs, effectively outperforming both direct prompting and traditional reasoning strategies, particularly in cases where table complexity and reasoning demands are higher.

978 E. Performance Analysis Under Different Table Sizes

Table <u>5</u> presents a performance comparison across different table sizes, categorized into small (<2k tokens), medium (2k \sim 4k tokens), and large (>4k tokens). The results compare several methods, including Binder (Cheng et al., 2023), Dater (Ye et al., 2023), and Chain-of-Table (Wang et al., 2024), against *TableMaster*. All methods are evaluated using *gpt-3.5-turbo*, with additional results of *TableMaster* provided for *gpt-4o-mini*.

Across all table sizes, *TableMaster* consistently outperforms baseline methods. Specifically, for *gpt-3.5-turbo*, *TableMaster* achieves the highest performance in all table size categories, scoring 69.01% on small tables, 58.00% on medium tables, and 56.73% on large tables. This demonstrates its ability to maintain robust performance even as table size increases, significantly outperforming Binder, Dater, and Chain-of-Table, especially on medium and large tables, where the performance gap becomes more pronounced.

A Recipe to Advance	Table	Understanding w	with Language	Models
---------------------	-------	-----------------	---------------	--------

Method	Table Size (Token)				
	Small (<2k)	Medium (2k \sim 4k)	Large (>4k)		
Binder (Cheng et al., 2023)	56.54	26.13	6.41		
Dater (Ye et al., 2023)	62.50	42.34	34.62		
Chain-of-Table (Wang et al., 2024)	68.13	52.25	44.87		
TableMaster (gpt-3.5-turbo)	69.01	58.00	56.73		
TableMaster (gpt-4o-mini)	78.71	70.50	70.19		



Figure 5. Performance Comparison Across Table Sizes (Row Count, Column Count, Area Size, Token Count).

Furthermore, *TableMaster* with *gpt-4o-mini* achieves even stronger performance, with accuracy scores of 78.71% (small tables), 70.50% (medium tables), and 70.19% (large tables). These results highlight that leveraging stronger base models further enhances *TableMaster* 's effectiveness, making it particularly well-suited for large-scale table reasoning tasks. Notably, when transitioning from medium to large tables, *TableMaster* (*gpt-4o-mini*) experiences only a 0.31% performance drop (from 70.50% to 70.19%), demonstrating its strong capability in handling increasing table complexity. This minimal decline contrasts sharply with other methods, which show significantly larger drops, further reinforcing the scalability and robustness of *TableMaster* in processing large-scale tabular data.

Figure 5 illustrates the accuracy trends of different models across various table sizes, categorized based on row count, column count, area size, and token count. The models evaluated in this study include *gpt-3.5-turbo* (*gpt35*), *gpt-4o-mini* (*gpt4m*), *TableMaster* (*gpt35*), and *TableMaster* (*gpt4m*). The results provide insights into how these models handle increasing table complexity and size, revealing the comparative strengths and limitations of each approach. The size split in this study is strictly partitioned into quartiles, ranging from the smallest to the largest tables.

Row Count. The top-left plot analyzes accuracy trends as row count increases. *TableMaster (gpt4m)* consistently outperforms other models, maintaining high accuracy levels even with an increasing number of rows. In contrast, *gpt-3.5-turbo (gpt35)* starts with the highest accuracy, peaking in the 11–15 row range before experiencing a decline as row count further increases. Smaller models such as *gpt35* and *gpt4m* exhibit a sharper decline, highlighting the challenge of 1045 processing larger tables with more rows.

Column Count. The top-right plot examines model performance as column count increases. *TableMaster (gpt4m)* again achieves strong performance, peaking at around five columns before showing a slight decline. This result highlights the effectiveness of **table-of-focus re-construction**, demonstrating that column re-selection can effectively adapt to scenarios with many columns. While *gpt35* initially maintains the highest accuracy, other models experience a steeper drop as the number of columns increases. These trends suggest that column-heavy tables pose greater challenges for reasoning compared to row-heavy tables, likely due to the increased dimensional complexity and interdependencies between attributes.

Area Size. The bottom-left plot evaluates the relationship between accuracy and table area size, calculated as the product of row and column counts. *TableMaster (gpt4m)* reaches peak performance in the mid-range (96–167 area size) before slightly declining for larger tables. *gpt35* initially performs well but deteriorates as table area size increases, while *gpt4m* and *gpt35* show a noticeable decline overall, reinforcing that larger tables significantly impact accuracy across models.

Token Count. The bottom-right plot assesses accuracy as a function of table token count, which reflects the amount of textual information models need to process. *TableMaster (gpt4m)* consistently achieves the highest accuracy, followed by *TableMaster (gpt35)*. A general downward trend is observed across all models as token count increases, indicating that larger input lengths negatively affect performance. Notably, *gpt35* experiences the sharpest drop, suggesting its lower capacity for handling long-context table data compared to *gpt4m*.

1063 Overall, these findings confirm that *TableMaster* is highly scalable and generalizable across different table sizes, consistently 1064 outperforming previous methods, particularly in handling larger and more complex tables. Its robust performance and 1065 gradual decline in accuracy as table size increases make it a reliable and efficient solution for table-based reasoning tasks. 1066

1069 Plot of Row Sizes in WikiTQ Datas Accuracy vs Peek Siz 0.05 78 0.04 1073 Accuracy (%) 24 Density Density 0.02 72 0.01 72.3 1079 0.00 70 24 20 60 80 Number of Rows per Table 100 120 10 25 50 Peek Size (a) (b) 1082

1067
1068F. Performance Analysis Under Different Table Peek Sizes

1083 1084 *Figure 6.* The row count distribution in the WikiTQ dataset and the analysis of accuracy variation with different peek sizes.

We propose the concept of table peek, which enhances the efficiency of *TableMaster* for table understanding tasks by reducing the context that language models need to process at certain steps.

To analyze the effectiveness of this approach, we first examine the row count distribution in the WikiTQ dataset, as shown in Figure 6(a). To improve visualization, we remove 72 extreme outliers with exceptionally large row counts. The resulting density plot illustrates that the majority of tables contain fewer than 20 rows, with a pronounced peak around 10 rows. As the number of rows increases, the density gradually declines, indicating that large tables are relatively uncommon. Although a small number of tables exceed 100 rows, their frequency is minimal.

The line graph in Figure <u>6(b)</u> illustrates the variation in accuracy with different peek sizes, where the peek size determines the number of rows considered during processing. Initially, accuracy is relatively low when only a small number of rows (e.g., 2–4) are used, reaching its minimum at a peek size of 4. We hypothesize that this occurs because, at a peek size of the table includes only the top headers and a single example row, which may provide a clear structure for the language model to follow. However, at a peek size of 4, the table includes three example rows, potentially causing the language model to overfit the first few rows and misinterpret the overall table structure. This misalignment may lead to ineffective SQL 1102 1103 1104 full table while still maintaining strong performance. 1105 1106 G. Efficiency Analysis of TableMaster 1107 1108 **G.1.** Theoretical Analysis 1109 1110 1111 1112 1113 1114 define the computational cost in terms of the total area size of the table that the language model processes. 1115 Below are the main components of our efficiency analysis: 1116 1117 • Structure extraction: $k \times n$ 1118 1119 • Row lookup: $k \times n$ 1120 1121 • Column lookup: n 1122 • Table-of-Focus Re-Construction $a \times b \times e$ 1123 1124 • Table Verbalization: $a \times b$, 1125 1126 • Reasoning Strategy Assessment: $a \times b$, 1127 1128 1129 weighted as 2) 1130 1131 Here, k represents the size of table peek, and e represents the number of table-of-focus re-constructions after information 1132 1133 cost is given by: 1134 1135 1136 1137

A Recipe to Advance Table Understanding with Language Models

generation for row lookup, resulting in a temporary drop in accuracy. 1100

1101 As the peek size increases, accuracy improves significantly, showing a sharp rise up to 25 rows. Beyond this point, the accuracy continues to improve but at a slower rate, eventually reaching its peak when the entire table is utilized ('All'). This trend suggests that a moderate peek size can effectively balance efficiency and accuracy, eliminating the need to process the

Efficiency is a critical factor in table-understanding methods. We analyze the efficiency of *TableMaster* theoretically, following the notations introduced in Section 4. Our analysis considers the length of the table input as the primary computational cost, excluding any additional prompts or external information, and does not account for output length. This is because, in most cases, the output is relatively short compared to the large volume of data in the table. Specifically, we

• **Reasoning:** $1.5 a \times b$ (where the factor 1.5 accounts for textual processes weighted as 1 and symbolic processes

estimation. a and b denote the dimensions of the table-of-focus $\mathbb{T}_{a \times b}$. Combining these components, the total computational

$$\text{Total Cost} = (2k+1) \times n + (e+2.5) \times (a \times b). \tag{1}$$

G.2. Empirical Analysis 1138

1139 The bar chart in Figure 7 illustrates the change in table area size before and after table condensation for the WikiTQ and 1140 TabFact datasets. The y-axis represents the table size, while the x-axis categorizes the datasets. Each dataset has two bars: 1141 the blue bar represents the original table size, and the orange bar represents the condensed table size after table-of-focus 1142 construction. WikiTQ exhibits a significant reduction in table size, approximately 1:3, with the condensed table being much 1143 smaller than the original. In contrast, TabFact also undergoes condensation but to a lesser extent, around 1:2. This suggests 1144 that WikiTQ tables require more substantial structural modifications to focus on relevant content, while TabFact tables need 1145 comparatively less condensation. 1146

As shown in Equation 1, the theoretical cost is independent of the number of rows m, while $a \times b$ reflects the size of the 1147 small sub-table, which is influenced by the estimated table condensation ratio 2.5. As stated in Table 6, the reconstruction 1148 occurs 1.5 averagely, so e is typically 1.5. In an ideal scenario, if the peek size is negligible, the cost is approximately 1149 $1.6 \times (m \times n)$. In the worst-case scenario, where the entire table must be examined and all content is required, the cost 1150 reaches $6 \times (m \times n)$ approximately. The estimation range for each table is 1.6 to 6 times the original table size. 1151

1152 Recent advancements in table understanding, such as Chain-of-Table (Wang et al., 2024) and Tree-of-Table (Ji et al., 2024), 1153 involve a step-by-step evolution of tables through a long chain of transformations. In each new step, both the original table 1154

A Recipe to Advance Table Understanding with Language Models



Figure 7. Changes in Table Condensation After Table-of-Focus Construction in Table Structure Understanding.

and the newly generated sub-table must be processed by language models. Additionally, their iterative process is complex, unstable, and difficult to analyze theoretically. In contrast, our approach is general and comprehensive, avoiding the trivial overhead of sub-table extraction. Instead, it focuses on holistic reasoning while maintaining ideal efficiency.

1177 H. Detailed Algorithm of Table-of-Focus Re-Construction

1155

1156

1161

1163 1164

1167

1169

1171

1176

Here, we provide a detailed description of the Table-of-Focus Re-Construction algorithm, as shown in Algorithm <u>1</u>.

1180			
1181	Alg	orithm 1 Algorithm of Table-of-Focus Re-Co	nstruction
1182	Ree	quire: \mathbb{T} : The original table	
1183	Ree	quire: Q : The question	
1184	Ree	quire: R: Selected rows	
1185	Ree	quire: C^0 : Initially selected columns	
1186	Ree	quire: \mathbb{C} : Ranked candidate column indices	
1187	Ens	sure: \mathbb{T}^{F} : Final table-of-focus	
1188	Ens	sure: C: Updated selected columns	
1189	1:	Initialize $C^{candidate} \leftarrow \{c \in \mathbb{C} \mid c \notin C^0\}$	
1190	2:	Initialize $C \leftarrow \operatorname{Copy}(C^0)$	
1191	3:	while true do	
1192	4:	$\mathbb{T}^{F} \leftarrow \text{extractTable}(\mathbb{T}, R, C)$	
1193	5:	$E \leftarrow estimateInformation(\mathbb{T}^{F}, Q)$	
1194	6:	if E or $len(C^{candidate}) = \emptyset$ then	
1195	7:	break	
1196	8:	else	
1197	9:	$c \leftarrow \text{popFront}(C^{candidate})$	{Select the next candidate column}
1198	10:	$C \leftarrow C \cup \{c\}$	
1199	11:	end if	
1200	12:	end while	
1201	13:	return \mathbb{T}^{F}, C	
1202			

The Table-of-Focus Re-Construction Algorithm iteratively refines a table by selecting relevant columns to form the final table-of-focus \mathbb{T}^F . It starts by initializing the set of candidate columns $C^{candidate}$ that are not part of the initially selected columns C^0 , and copies C^0 to initialize C. In each iteration, it extracts a sub-table \mathbb{T}^F using the current selected columns and estimates whether the extracted sub-table contains sufficient information to answer the given question Q. If the information is sufficient E = True or no more candidate columns remain, the process terminates. Otherwise, the next ranked candidate column is selected and added to C, repeating the process. The algorithm ultimately returns the refined table \mathbb{T}^F and

the updated set of selected columns, ensuring an efficient and structured approach to dynamically refining a table while 1211 balancing relevance and minimal table size.

1212 1213

1215 1216

1224

1225

1246

1259

I. Analysis of Table-of-Focus Re-Construction 1214

Dataset	Initial Columns Final Columns		Columns	Addec	l Columns	
	Number (#)	Percentage (%)	Number (#)	Percentage (%)	Number (#)	Percentage (%)
TabFact	2.44	40.74	3.34	54.64	0.90	13.91
WikiTQ	2.87	47.67	4.72	75.91	1.85	28.23

Table 6 presents Column Selection Statistics before and after Table-of-Focus Re-Construction for two datasets: TabFact and WikiTQ. The table measures how many columns were initially selected, how many remained after refinement, and how many were newly added during the reconstruction process. 1226

The table is structured into three main sections: Initial Columns, Final Columns, and Added Columns. Each section includes 1228 two metrics: the number of columns and the percentage of total columns in the dataset. The Initial Columns represent 1229 the starting number of columns before any refinement. The Final Columns show the number of columns retained after 1230 the reconstruction process. The Added Columns indicate the number of additional columns incorporated to enhance table 1231 comprehension.

1232 For the TabFact dataset, the number of Initial Columns is 2.44, covering 40.74% of the table's total columns. After the 1233 reconstruction process, the Final Columns increase to 3.34, covering 54.64%. This means that 0.90 additional columns were 1234 introduced averagely, which accounts for 13.91% of the total columns. For the WikiTQ dataset, the pattern is similar but with higher values. The Initial Columns start at 2.87, representing 47.67% of the total table. After reconstruction, the Final 1236 Columns expand to 4.72, covering 75.91% of the table's total structure. This increase results from 1.85 additional columns, 1237 which make up 28.23% of the total columns. 1238

1239 Overall, this mechanism has been proven to be effective while remaining lightweight. The table demonstrates that Table-of-1240 Focus Re-Construction slightly increases the number of selected columns, with a more pronounced effect in the WikiTO 1241 dataset compared to TabFact. This suggests that WikiTQ tables require a greater degree of expansion to ensure adequate 1242 information coverage, whereas TabFact tables undergo a more moderate refinement process. 1243

1244 J. Analysis of Adaptive Reasoning 1245

	Table 7. Performance of	of Different Reasoning Method	s Across Base LLMs	
Base LLM	Method	Calculation Required #2692	No Calculation Required #1652	Overall #4344
	Textual Reasoning	81.17	88.56	83.98
gpt-40	Symbolic Reasoning	74.59	74.70	74.63
	Text-Guided Symbolic Reasoning	76.49	77.36	76.82
	Textual Reasoning	67.50	81.90	72.97
gpt-40-mini	Symbolic Reasoning	61.55	62.29	61.83
	Text-Guided Symbolic Reasoning	67.24	71.43	68.83
	Textual Reasoning	52.27	72.40	59.92
gpt-3.5-turbo	Symbolic Reasoning	43.28	61.80	50.32
	Text-Guided Symbolic Reasoning	59.10	66.65	61.97

Table 7 compares different reasoning methods-textual reasoning, symbolic reasoning, and text-guided symbolic reason-1260 ing-across various LLMs under calculation-required and no-calculation-required scenarios. This experiment is conducted 1261 using gpt-4o-mini on the WikiTQ dataset. 1262

For gpt-40, textual reasoning achieves the highest accuracy (83.98% overall), excelling in both calculation-required (81.17%)

and no-calculation-required (88.56%) cases. Symbolic reasoning performs worse (74.63% overall), while text-guided symbolic reasoning offers slight improvements (76.82%). For *gpt-4o-mini*, a similar trend is observed, with textual reasoning maintaining the highest accuracy (72.97% overall), followed by text-guided symbolic reasoning (68.83%), and symbolic reasoning performing the worst (61.83%). For *gpt-3.5-turbo*, performance drops significantly, with textual reasoning at 59.92%, symbolic reasoning struggling at 50.32%, and text-guided symbolic reasoning achieving the best results (61.97%),

1270 indicating that symbolic guidance benefits weaker models.

Symbolic reasoning is consistently outperformed by textual reasoning, while text-guided symbolic reasoning surpasses
textual reasoning only in *gpt-3.5-turbo* under calculation-required scenarios. One reason for this is that not all calculation-required questions necessarily benefit from symbolic reasoning; for simple calculations, textual reasoning is more effective.
However, for complex calculation-required questions, text-guided symbolic reasoning is the preferred approach. This
provides a key insight for prompt design of reasoning strategy assessment.

1277 Overall, textual reasoning consistently outperforms symbolic reasoning across all models, while text-guided symbolic 1278 reasoning helps mitigate weaker numerical capabilities in smaller models. These results suggest that adaptive reasoning 1279 should prioritize textual approaches, incorporating symbolic methods selectively for numerical calculations in weaker 1280 models.

1281 1282

1283	Table 8. Performance and Inference Times for Different Methods					
1284	Method	Accuracy (%)	Inference Times (#)			
1285	Chain of Thought	72.97	1			
1286	Program of Thought	61.83	1			
1287	Text-Guided Program of Thought	68.83	1			
1289	Self-Consistency (5 CoT)	74.98	3			
1290	Self-Consistency (5 PoT)	63.97	3			
1291	Mix Self-Consistency (3+3)	76.70	6			
1292	Mix Self-Consistency (5+5)	77.46	10			
1293	Self-Eval	70.58	2			
1294	Adaptive Reasoning (POT)	71.18	1			
1295	Adaptive Reasoning (TPOT)	74.08	1			
1296	Adaptive Reasoning (POT - Upper Bound)	82.99	1			
1297	Adaptive Reasoning (TPOT - Upper Bound)	85.06	1			
1000						

1298 1299

Table <u>8</u> compares the performance (accuracy %) and inference times of various reasoning methods, including chain of thought (CoT), program of thought (PoT), text-guided program of thought (TPoT), self-consistency, and adaptive reasoning. This experiment is conducted using *gpt-4o-mini* on the WikiTQ dataset.

Among single-pass methods (1 inference), chain of thought achieves 72.97% accuracy, outperforming program of thought (61.83%) and text-guided program of thought (68.83%). This suggests that CoT is more effective than pure symbolic reasoning when only one inference is allowed.

Self-consistency methods, which perform multiple inferences to improve reliability, achieve better results. Five-shot CoT self-consistency reaches 74.98%, while five-shot PoT self-consistency lags behind at 63.97%. As introduced in (Liu et al., 2024b), mixed self-consistency (3 CoT + 3 PoT) and (5+5) further improve accuracy to 76.70% and 77.46%, respectively, at the cost of increased inference time (6 and 10 passes). Self-evaluation (self-eval) first performs CoT and PoT inferences (Prompt <u>26</u>), then selects the better result, achieving 70.58% with 2 inferences.

Adaptive reasoning achieves competitive performance while maintaining single-pass efficiency. PoT-based adaptive reasoning reaches 71.18%, while TPOT-based adaptive reasoning, which combines textual and text-guided symbolic

1315 methods, improves to 74.08%. The upper-bound performance of these adaptive strategies (assuming perfect strategy

1316 selection) reaches 82.99% (PoT) and 85.06% (TPOT), significantly outperforming all other methods, highlighting the

1317 importance of textual guidance and strategy selection.

For the selection distribution between CoT and PoT (TPoT):

Self-eval: 1,962 PoT and 2,382 CoT

1322 1323

1324 1325

1326

1327 1328

1338 1339

1366

- Adaptive reasoning (PoT): 1,590 PoT and 2,754 CoT
- Adaptive reasoning (TPoT): 1,590 PoT and 2,754 CoT
- Adaptive reasoning (PoT upper bound): 435 PoT and 3,909 CoT
- Adaptive reasoning (TPoT upper bound): 525 PoT and 3,819 CoT

These results suggest that language models should prioritize textual reasoning and reserve symbolic reasoning for more complex numerical calculations where it provides a clear advantage.

Overall, self-consistency enhances accuracy but requires multiple inferences, whereas adaptive reasoning effectively balances
 accuracy and efficiency. To further improve strategy assessment, we will explore ways to approach this upper bound in
 future work. This demonstrates that well-designed adaptive reasoning strategies can rival more computationally expensive
 self-consistency methods while maintaining efficiency.

1336 1337 **K. Information Missing and Table Reasoning with Full Table**

Table 9. Performance comparison of reasoning with and without the full table on WikiTQ and TabFact.

Method	WikiTQ	TabFact
PoTable (Previous SOTA) (Mao et al., 2024)	64.73	88.93
TableMaster w/ Full Table in Reasoning	78.13	90.12
TableMaster w/o Full Table in Reasoning	77.23 (-0.90)	89.58 (-0.54)

As discussed in our limitations, the table-of-focus process may sometimes lead to the loss or omission of key relevant information. This issue is inevitable when attempting to locate specific data. If no relevant data exists within the selected portion, the reasoning result will naturally be incorrect.

In our experiments, we found that when using the table-of-focus and its verbalized representation for reasoning, 265 out of 4,344 questions in the WikiTQ dataset had no available answers. This led to a performance drop, as the language model responded with an inability to provide an answer. To address this, we replaced the table-of-focus with the original full table, combined with verbalized table-of-focus as input in those questions. The performance under this adjustment is shown in Table 9, reaching 77.23% in WikiTQ. When we directly replaced the table-of-focus with the full table for all questions in WikiTQ, the performance increased to 78.13%, resulting in a slight improvement of 0.9%. Two results are similar.

We believe this approach does not contradict previous steps such as structure extraction and table-of-focus selection. These steps remain valuable, as the extracted target data is retained in the verbalized table, where the information density is higher and semantic context is richer. The language model prioritizes this high-density information, and if it is insufficient, it can then reference the global information from the full table. This demonstrates the complementary nature of the full table and the verbalized table-of-focus. From an efficiency perspective, it is preferable to use the sub-table for reasoning initially and only switch to the full table when necessary.

To highlight the performance of *TableMaster*, we report the best scores of 78.13 and 90.12 in the main results table.
Regardless, our method consistently outperforms the previous state-of-the-art, PoTable (Mao et al., 2024), on both WikiTQ
and TabFact.

1367 L. Case Study of Table Verbalization

Table verbalization brings a slight overall improvement in table understanding and is particularly effective in cases where deeper comprehension of the table's context is required to answer questions accurately.

Figure <u>8</u> presents a case study on the impact of table verbalization in helping language models reason about structured data.
The setup includes a table listing U.S. congressional representatives, their terms in office, political affiliations, and notable events. The question posed is: *How many people stayed at least 3 years in office?*, with a ground-truth answer of 4.

A Recipe to Advanc	e Table U	Inderstanding	with I	anguage	Models
--------------------	-----------	---------------	--------	---------	--------

375	Table	Question	: how many people stay	yed at least 3 years in of	ffice? Ground-truth: 4				
1376	Party	Name	Took office	Left office	Party	Notes/Events			
1377	11	William McCreery	March 4, 1803	March 3, 1809	Democratic Republican		LM Reasoning		
1378	12	Alexander McKim	March 4, 1809	March 3, 1815	Democratic Republican		J		
1379	13	William Pinkney	March 4, 1815	April 18, 1816	Democratic Republican	Resigned to	Pr	ediction: 5	
1380	14	Peter Little	September 2, 1816	March 3, 1823	Democratic Republican				
1381	14	Peter Little	March 4, 1823	March 3, 1825	Jacksonian DR				
1382	14	Peter Little	March 4, 1825	March 3, 1829	Adams				
1383	15	Benjamin C. Howard	March 4, 1829	March 3, 1833	Jacksonian				
1384									
1385	Verhal	ized Table		Table Verba	lization				
1386				•			1		
1387	William	McCreery, a Democra er McKim, also a Dem	itic-Republican, serveo locratic-Republican, w	d in office from March ho held the position fr	4, 1803, to March 3, 1809, rom March 4, 1809, to Marc	followed by ch 3, 1815.	LM Reasoning		
1388	William	Pinkney then took off	ice on March 4, 1815, t	out resigned on April 1	8, 1816, to accept a diplon	natic role as			
1389	Minister as a Der	Plenipotentiary to Ru nocratic-Republican.	ssia. Peter Little succe then continued in offic	eeded him, serving from ce under different polit	m September 2, 1816, to M tical affiliations: as a Jacks	larch 3, 1823, Ionian	Pr	rediction: 4	
1390	Democr	atic-Republican from	March 4, 1823, to Marc	ch 3, 1825, and later a	s an Adams supporter fron	n March 4,			
1391	1825, to 1829, to	March 3, 1829. His te March 3, 1833.	nure was followed by I	Benjamin C. Howard, a	a Jacksonian, who served f	rom March 4,			
1392							1		
1393									
1394		Eim	ma 9 Casa study	on the immediate of t	able verbalization T	The date is fre	m the WiltimO d	atasat	
1395		Figu	Te 8. Case study (on the impact of t	able verbalization.	ne data is in	m the wiking da	ataset.	
1396									
1397	When t	he table is inpu	t directly, the mo	odel incorrectly	predicts 5, as it mi	stakenly cou	unts rows rather	than identify	ing unique
1398	individ	uals. This sugg	ests that the mo	del relies on si	mple row counting	g instead of	truly understand	ding the data	. However
1399	with th	e verbalized tab	ole, the model ac	curately interpr	ets the description	s, grasps the	e actual meaning	g, and correct	ly answers
1400	with 4.		,	5 1	1	<i>, C</i> 1	c c		5
1401									
1402									
1403									
1404									
1405									
1406									
1407									
1408									
1409									
1410									
1411									
1412									
1413									
1414									
1415									
1416									
1417									
1418									
1419									
1420									
1421									
1422									
1423									
1424									
1425									
1426									
1427									
1428									
1429									

1430 M. Case Study of *TableMaster*

1478

1484

As shown in Figure 9, we present a case study of *TableMaster* to provide a detailed illustration of its workflow. The process begins with SQL generation, followed by structure extraction, where key columns such as Nominated work are identified. Next, row and column lookup selects relevant data, leading to table-of-focus construction, which refines the table to retain only essential information. Table verbalization then converts structured data into a human-readable summary, and a textual overview highlights Leona Lewis's major awards and nominations. Additionally, a step-by-step guide explains the counting process, providing a structured approach to symbolic reasoning. The reasoning and execution phase includes code execution, where a Python snippet correctly counts the occurrences of Won, yielding the correct prediction of 20.

This case study demonstrates *TableMaster* 's effectiveness in extracting, processing, and reasoning over structured data,
 ultimately enabling accurate table-based question answering.



Figure 9. Case study of TableMaster. The data is from the WikiTQ dataset.

N. Prompt Design in TableMaster Prompt for TableMaster – Structure Extraction ## Obiective You are provided with a text representation of a table in string format, detailing the content of each cell. Your task is to identify and extract the Top Header and Key Column of the table. ## Table Definition The table is represented by cell-value pairs, where each pair consists of a cell address and a value of the content in that cell, separated by a comma (e.g., 'A1, Year'). Multiple cells are separated by '|' (e.g., 'A1, Year | A2, Profit'). Cells may contain empty values, represented as 'A1, |A2, Profit'. ## Table {table} ## Instructions 1. Top Header: The section at the top of the table, often spanning multiple columns horizontally, that describes the primary information presented in the table. 2. Key Column: A column where the values best represent the subject or key identifier for each row in the table, typically containing row labels or keys (e.g., year, date, number, name, etc.). 3. You should extract the top headers with address and value, like ['A1,Year', 'A2,Profit', ...]. 4. key_column_index should be like 'A' or 'B' ... 5. The key column should contain meaningful values instead of id. ## Response Format The response should be in JSON format: ```json {{ "topheaders": ["address1, header1", "address2, header2", ...], "key_column_index": "column1", }}

A Recipe to Advance Table Understanding with Language Models

Figure 10. Prompt for structure extraction in *TableMaster*. Blue text indicates placeholders for variables within the prompt. The prompt guides the language model in extracting the table's structure.

Prompt for TableMaster – Column Ranking ## Objective You are provided with information of a table and a question related to the table. Your task is to rank the column indices based on the relevance to the question. ## Table Definition The table is represented by cell-value pairs, where each pair consists of a cell address and a value of the content in that cell, separate by a comma (e.g., 'A1, Year'). Multiple cells are separated by ' ' (e.g., 'A1, Year A2, Profit'). Cells may contain empty values, represented as 'A1, A2, Profit'. ## Table Information Table: (table) Top Headers: (topheaders) ## Question (question) ## nstructions 1. The column indices must only contain letters, like ['A', 'B', 'C',]. 2. You should first rank all the column indices. ## Response Format The response should be in JSON format: "join (''ranked_column_indices": [''column indexA", ''column indexB",]		
## Objective You are provided with information of a table and a question related to the table. Your task is to rank the column indices based on the relevance to the question. ## Table Definition The table is represented by cell-value pairs, where each pair consists of a cell address and a value of the content in that cell, separate by a commo (e.g., YA1,Year). Multiple cells are separated by '['(e.g., YA1,Year]A2,Profit'). Cells may contain empty values, represented as 'A1, [A2,Profit'. ## Table Information Table: (table) Top Headers: (topheaders) ## Question (question) ## Instructions 1. The column indices must only contain letters, like ['A', 'B', 'C',]. 2. You should first rank all the column indices. ## Response Format The response should be in JSON format: "json ('''ranked_column_indices": [''column indexA", ''column indexB",]		
Prompt for TableMaster – Column Ranking ## Objective You are provided with information of a table and a question related to the table. Your task is to rank the column indices based on the relevance to the question. ## Table Definition The table is represented by cell-value pairs, where each pair consists of a cell address and a value of the content in that cell, separat by a comma (e.g., 'A1,Year'). Multiple cells are separated by 'I' (e.g., 'A1,Year A2,Profit'). Cells may contain empty values, represented as 'A1, A2,Profit'). Cells may contain empty values, represented as 'A1, A2,Profit'. ## Table Information Table: (table) Top Headers: {topheaders} ## Question (question) ## Instructions 1. The column indices must only contain letters, like ['A', 'B', 'C',]. 2. You should first rank all the column indices based on the relevance to the question. 3. Your output should contain all the column indices. ## Response Format The response should be in JSON format: "json (["ranked_column_indices": ["column indexA", "column indexB",]		
Prompt for TableMaster – Column Ranking ## Objective You are provided with information of a table and a question related to the table. Your task is to rank the column indices based on the relevance to the question. ## Table Definition The table is represented by cell-value pairs, where each pair consists of a cell address and a value of the content in that cell, separat by a comma (e.g., 'A1,Year'). Multiple cells are separated by 'I' (e.g., 'A1,Year/A2,Profit'). Cells may contain empty values, represented as 'A1, [A2,Profit']. ## Table Information Table: {table} Top Headers: (topheaders) ## Question {question] ## Instructions 1. The column indices must only contain letters, like ['A', 'B', 'C',]. 2. You sould first rank all the column indices. ## Response Format The response should be in JSON format: "joon {{ ""ranked_column_indices": ["column indexA", "column indexB",] }		
## Objective You are provided with information of a table and a question related to the table. Your task is to rank the column indices based on the relevance to the question. ## Table Definition The table is represented by cell-value pairs, where each pair consists of a cell address and a value of the content in that cell, separate by a comma (e.g., 'A1, Year'). Multiple cells are separated by '[' (e.g., 'A1, Year]A2, Profit'). Cells may contain empty values, represented as 'A1, [A2, Profit'). Cells may contain on mapty values, represented as 'A1, [A2, Profit'. ## Table Information Table: (rable) Top Headers: {topheaders] ## Question (question) ## Instructions 1. The column indices must only contain letters, like ['A', 'B', 'C',]. 2. You should first rank all the column indices. ## Response Format The response should be in JSON format: "join "ranked_column_indices": ["column indexA", "column indexB",]		
## Objective You are provided with information of a table and a question related to the table. Your task is to rank the column indices based on the relevance to the question. ## Table Definition The table Definition The table is represented by cell-value pairs, where each pair consists of a cell address and a value of the content in that cell, separate by a comma (e.g., 'A1, Year'). Multiple cells are separated by ' ' (e.g., 'A1, Year A2, Profit'). Cells may contain empty values, represented as 'A1, A2, Profit'). Cells may contain empty values, represented as 'A1, A2, Profit'). Table Information Table: (table) Top Headers: {topheaders} ## Question {question} ## Instructions 1. The column indices must only contain letters, like ['A', 'B', 'C',]. 2. You should first rank all the column indices based on the relevance to the question. 3. Your output should contain all the column indices. ## Response Format The response should be in JSON format: "joan ['' ranked_column_indices": ["column indexA", "column indexB",]		
<pre>Prompt for TableMaster – Column Ranking ## Objective You are provided with information of a table and a question related to the table. Your task is to rank the column indices based on the relevance to the question. ## Table Definition The table is represented by cell-value pairs, where each pair consists of a cell address and a value of the content in that cell, separate by a comma (e.g., 'A1,Year'). Multiple cells are separated by 'I' (e.g., 'A1,Year A2,Profit'). Cells may contain empty values, represented as 'A1, A2,Profit'. ## Table Information Table: (table) Top Headers: (topheaders) ## (Question (question) ## Instructions 1. The column indices must only contain letters, like ['A', 'B', 'C',]. 2. You should first rank all the column indices based on the relevance to the question. 3. Your output should contain all the column indices. ## Response Format The response should be in JSON format: ''joan { "ranked_column_indices": ["column indexA", "column indexB",] }</pre>		
Prompt for TableMaster – Column Ranking ## Objective You are provided with information of a table and a question related to the table. Your task is to rank the column indices based on the relevance to the question. ## Table Definition The table is represented by cell-value pairs, where each pair consists of a cell address and a value of the content in that cell, separate by a comma (e.g., 'A1, Year'). Multiple cells are separated by ' ' (e.g., 'A1, Year A2, Profit'). Cells may contain empty values, represented as 'A1, A2, Profit'. ## Table Information Table: ta		
<pre>## Objective You are provided with information of a table and a question related to the table. You rask is to rank the column indices based on the relevance to the question. ## Table Definition The table is represented by cell-value pairs, where each pair consists of a cell address and a value of the content in that cell, separat by a comma (e.g., 'A1, Year'). Multiple cells are separated by ' ' (e.g., 'A1,Year A2,Profit'). Cells may contain empty values, represented as 'A1, A2,Profit'. ## Table Information Table: (table) Top Headers: {topheaders} ## Question (question) ## Instructions 1. The column indices must only contain letters, like ['A', 'B', 'C',]. 2. You should first rank all the column indices. ## Response Format The response should be in JSON format: ``json {{ "ranked_column_indices": ["column indexA", "column indexB",] } </pre>	Ī	Prompt for TableMaster – Column Ranking
You are provided with information of a table and a question related to the table. You r task is to rank the column indices based on the relevance to the question. ## Table Definition The table is represented by cell-value pairs, where each pair consists of a cell address and a value of the content in that cell, separat by a comma (e.g., 'A1,Year'). Multiple cells are separated by ' (e.g., 'A1,Year A2,Profit'). Cells may contain empty values, represented as 'A1, A2,Profit'. ## Table Information Table: (table) Top Headers: {topheaders} ## Question {question} ## Instructions 1. The column indices must only contain letters, like ['A', 'B', 'C',]. 2. You should first rank all the column indices. ## Response Format The response should be in JSON format: ``json {{ "ranked_column_indices": ["column indexA", "column indexB",] }	Γ	## Objective
<pre>Your task is to rank the column indices based on the relevance to the question. ## Table Definition The table is represented by cell-value pairs, where each pair consists of a cell address and a value of the content in that cell, separat by a comma (e.g., 'A1,Year'). Multiple cells are separated by ' ' (e.g., 'A1,Year A2,Profit'). Cells may contain empty values, represented as 'A1, A2,Profit'. ## Table Information Table: {table is represented by cell-value pairs, values, represented as 'A1, A2,Profit'. ## Table Information Table: {table is represented to y ' ' (e.g., 'A1,Year', 'A2,Profit'). Cells may contain empty values, represented as 'A1, A2,Profit'. ## Table Information Table: {table is represented to y ' ' (e.g., 'A1,Year', 'A2,Profit'). Top Headers: {topheaders} ## Question {question {uestion Subset Y# Question Subset Y# Instructions Y* Contain letters, like ['A', 'B', 'C',]. You should first rank all the column indices based on the relevance to the question. Your output should contain all the column indices. ## Response Format The response should be in JSON format: ```json {{ "ranked_column_indices": ["column indexA", "column indexB",] } </pre>	1	<i>(ou are provided with information of a table and a question related to the table.</i>
<pre>## Table Definition The table is represented by cell-value pairs, where each pair consists of a cell address and a value of the content in that cell, separat by a comma (e.g., 'A1, Year'). Multiple cells are separated by ' ' (e.g., 'A1, Year A2, Profit'). Cells may contain empty values, represented as 'A1, A2, Profit'. ## Table Information Table: {table} Top Headers: {topheaders} ## Question {question} 1. The column indices must only contain letters, like ['A', 'B', 'C',]. 2. You should first rank all the column indices based on the relevance to the question. 3. Your output should contain all the column indices. ## Response Format The response should be in JSON format: ```;son {{ "ranked_column_indices": ["column indexA", "column indexB",] }</pre>	1	our task is to rank the column indices based on the relevance to the question.
<pre>/// Note Column indices ": ["column indexA", "column indexB",] }/// a column_indices ": ["column indexA", "column indexB",] }// a column_indices ": ["column indexA", "column indexB",] </pre>		## Table Definition
by a comma (e.g., 'A1,Year'). Multiple cells are separated by ' ' (e.g., 'A1,Year A2,Profit'). Cells may contain empty values, represented as 'A1, A2,Profit'. ## Table Information Table: {table} Top Headers: {topheaders} ## Question {question} ## Instructions 1. The column indices must only contain letters, like ['A', 'B', 'C',]. 2. You should first rank all the column indices based on the relevance to the question. 3. Your output should contain all the column indices. ## Response Format The response should be in JSON format: "json {{ "ranked_column_indices": ["column indexA", "column indexB",] }		Fhe table is represented by cell-value pairs, where each pair consists of a cell address and a value of the content in that cell, separat
Multiple cells are separated by ' '(e.g., 'A1, Year/A2, Profit'). Cells may contain empty values, represented as 'A1, /A2, Profit'. ## Table Information Table: {table} Top Headers: {topheaders} ## Question {question} ## Instructions 1. The column indices must only contain letters, like ['A', 'B', 'C',]. 2. You should first rank all the column indices based on the relevance to the question. 3. Your output should contain all the column indices. ## Response Format The response should be in JSON format: "json {''ranked_column_indices": ["column indexA", "column indexB",] }	1	yy a comma (e.g., 'A1,Year').
<pre>## Table Information Table: {table} Top Headers: {topheaders} ## Question {question} ## Instructions 1. The column indices must only contain letters, like ['A', 'B', 'C',]. 2. You should first rank all the column indices based on the relevance to the question. 3. Your output should contain all the column indices. ## Response Format The response should be in JSON format: ``json {{ "ranked_column_indices": ["column indexA", "column indexB",] }}</pre>		Aultiple cells are separated by ' ' (e.g., 'A1,Year A2,Profit').
<pre>## Table Information Table: {table} Top Headers: {topheaders} ## Question {question} ## Instructions 1. The column indices must only contain letters, like ['A', 'B', 'C',]. 2. You should first rank all the column indices based on the relevance to the question. 3. Your output should contain all the column indices. ## Response Format The response should be in JSON format: ``json {{ "ranked_column_indices": ["column indexA", "column indexB",] }}</pre>	(elis may contain empty values, represented as "A1, A2,Profit".
Table: {table} Top Headers: {topheaders} ## Question {question} ## Instructions 1. The column indices must only contain letters, like ['A', 'B', 'C',]. 2. You should first rank all the column indices based on the relevance to the question. 3. Your output should contain all the column indices. ## Response Format The response should be in JSON format: ``json {{ "ranked_column_indices": ["column indexA", "column indexB",] }	ŧ	## Table Information
<pre>{table} Top Headers: {topheaders} ## Question {question {question} ## Instructions 1. The column indices must only contain letters, like ['A', 'B', 'C',]. 2. You should first rank all the column indices based on the relevance to the question. 3. Your output should contain all the column indices. ## Response Format The response should be in JSON format: ```json {{ "ranked_column_indices": ["column indexA", "column indexB",] }</pre>		Table:
Top Headers: {topheaders} ## Question {question} ## Instructions 1. The column indices must only contain letters, like ['A', 'B', 'C',]. 2. You should first rank all the column indices based on the relevance to the question. 3. Your output should contain all the column indices. ## Response Format The response should be in JSON format: ```json {{ "ranked_column_indices": ["column indexA", "column indexB",] }	1	table}
<pre>## Question {question} ## Instructions 1. The column indices must only contain letters, like ['A', 'B', 'C',]. 2. You should first rank all the column indices based on the relevance to the question. 3. Your output should contain all the column indices. ## Response Format The response Format The response should be in JSON format: ```json {{ "ranked_column_indices": ["column indexA", "column indexB",] }</pre>		Fop Headers: {topheaders}
<pre>## Question {question {question {question } ## Instructions 1. The column indices must only contain letters, like ['A', 'B', 'C',]. 2. You should first rank all the column indices based on the relevance to the question. 3. Your output should contain all the column indices. ## Response Format The response should be in JSON format: ```json {{ "ranked_column_indices": ["column indexA", "column indexB",] }}</pre>		
<pre>## Instructions 1. The column indices must only contain letters, like ['A', 'B', 'C',]. 2. You should first rank all the column indices based on the relevance to the question. 3. Your output should contain all the column indices. ## Response Format The response should be in JSON format: ```json {{ "ranked_column_indices": ["column indexA", "column indexB",] }}</pre>	1	/# Question
<pre>## Instructions 1. The column indices must only contain letters, like ['A', 'B', 'C',]. 2. You should first rank all the column indices based on the relevance to the question. 3. Your output should contain all the column indices. ## Response Format The response should be in JSON format: ```json {{ "ranked_column_indices": ["column indexA", "column indexB",] }}</pre>		
 The column indices must only contain letters, like ['A', 'B', 'C',]. You should first rank all the column indices based on the relevance to the question. Your output should contain all the column indices. ## Response Format The response should be in JSON format: ```json {{ "ranked_column_indices": ["column indexA", "column indexB",] }}	1	## Instructions
2. You should jist fank all the column indices based on the relevance to the question. 3. Your output should contain all the column indices. ## Response Format The response should be in JSON format: ```json {{ "ranked_column_indices": ["column indexA", "column indexB",] }}		1. The column indices must only contain letters, like ['A', 'B', 'C',].
## Response Format The response should be in JSON format: ```json {{ "ranked_column_indices": ["column indexA", "column indexB",] }}		2. Fou should just rank un the column malces based on the relevance to the question. 3. Your output should contain all the column indices.
<pre>## Response Format The response should be in JSON format: ```json {{ "ranked_column_indices": ["column indexA", "column indexB",] }}</pre>		
The response should be in JSON format: ""json {{ "ranked_column_indices": ["column indexA", "column indexB",] }}	\$	## Response Format
<pre>[son {{ "ranked_column_indices": ["column indexA", "column indexB",] }}</pre>		The response should be in JSON format:
"ranked_column_indices": ["column indexA", "column indexB",] }}		json v
}}	ľ	"ranked_column_indices": ["column indexA", "column indexB",]
	j	}

Figure 11. Prompt for column ranking in *TableMaster*. Blue text indicates placeholders for variables within the promp guides the language model to rank the priority of all columns based on the given table, top headers, and related question.

Ρ	rompt for TableMaster – Column Lookup
##	
Yc	u are provided with information of a table and a question related to the table.
Yc	ur task is to lookup the column indices that are needed to answer the question based on the table.
ш	
## Th	: Table Definition Table is represented by cell-value pairs, where each pair consists of a cell address and a value of the content in that cell, ser
by	a comma (e.g., 'A1,Year').
M	ultiple cells are separated by ' ' (e.g., 'A1,Year A2,Profit').
Ce	lls may contain empty values, represented as 'A1, A2,Profit'.
##	Table Information
Та	ble:
{to	nble}
та	n Headers: (tenheaders)
10	p neulers. {lopineulers}
##	Question
{ q	uestion}
##	Instructions
1.	The column indices must only contain letters, like ['A', 'B', 'C',].
2.	Your output of the column indices should not any contain number, like ['A1', 'B2', 'C1',].
3.	Your output of the column indices should not contain the column name.
4.	You should select the column that are relevant and necessary to answer the question.
##	t Response Format
Th	e response should be in JSON format:
•••	ison
{{	
11	selectea_column_indices": ["column indexA", "column indexB",]
<i>ss</i>	

Figure 12. Prompt for column lookup in *TableMaster*. Blue text indicates placeholders for variables within the prompt. The prompt guidesthe language model to select relevant columns based on the given table, top headers, and related question.

# Objective 'ou are provided with informatic 'our task is to generate a SQL qu	on of a table and a question related to the table. Very that can be used to find the rows that answer the question.
# Table Information Part of Table: table}	
# Question question}	
# Instructions . The SQL query must be in the where Table is the table name, X . If the information is not enoug . Do not give complex sql query . Use this SQL query only to sele	format of `SELECT XXX, FROM Table WHERE XXX`, XX is the column name, and WHERE is the criteria. h to answer the question, you should return a sql to select all rows. just simple query to select rows. ct relevant rows, not for getting the final answer.
# Response Format Provide the response in the follo ``json { "sal": "SELECT XXX FROM T	ving JSON format:

Figure 13. Prompt for SQL generation for row lookup in TableMaster. Blue text indicates placeholders for variables within the prompt.

The prompt guides the language model to generate SQL for selecting relevant rows based on the given table and related question.

Prompt for TableMaster – Table Verbalization

Objective

You are provided with a table in string format. Your task is to convert the table into a detailed text description.

Table {table}

Instructions

1. Provide a detailed description of the table, covering all rows and columns.

2. Include every detail and numerical value without omitting or summarizing.

3. Use external knowledge only to enhance clarity, while staying faithful to the table's content.

4. If the table only contains headers and no rows, it should be described as an empty table.

Now, please provide the verbalized description of the table:

Figure 14. Prompt for table verbalization in TableMaster. Blue text indicates placeholders for variables within the prompt. The prompt guides the language model to verbalize the given table by adding detailed descriptions and additional knowledge about the table.

Prompt for TableMaster – Information Estimation

Objective
You are provided with information from a table and a question related to the table.
Your task is to estimate whether the current information of the table can answer the question.
Table Information
Top Headers: {topheader_info}
Table Content:
{table}
Question
{question
{question}
Response Format
The response should be in JSON format:
```json
{{
 "results": True of False
}}

*Figure 15.* Prompt for information estimation in *TableMaster*. Blue text indicates placeholders for variables within the prompt. The
 prompt guides the language model to evaluate the given table's content and determine whether it contains sufficient information to answer
 the provided question

### Prompt for TableMaster – Reasoning Strategy Assessment

| Your task   | is to assess whether answering this question needs mathematical calculation.                          |
|-------------|-------------------------------------------------------------------------------------------------------|
|             |                                                                                                       |
| ## Table    |                                                                                                       |
| {table}     |                                                                                                       |
| ## Questi   | on                                                                                                    |
| {question   |                                                                                                       |
|             |                                                                                                       |
| ## Instruc  |                                                                                                       |
| 1. If the q | iestion can be directly answered using the information in the table, you should respond with "False". |
| 2. If the q | iestion involves counting something, you should respond with `True`.                                  |
| 3. If the q | estion requires calculations based on the data in the table, you should respond with `True`.          |
| ## Respoi   | ise Format                                                                                            |
| The respo   | nse should be in JSON format:                                                                         |
| ```ison     |                                                                                                       |
| {{          |                                                                                                       |
| "results    | "· True of False                                                                                      |
| 11          |                                                                                                       |
| <i>}}</i>   |                                                                                                       |

Figure 16. Prompt for reasoning strategy assessment in *TableMaster*. Blue text indicates placeholders for variables within the prompt.
 The prompt guides the language model to evaluate whether answering the given question requires direct information retrieval, counting, or mathematical calculations based on the table's content. The response determines the subsequent reasoning strategy.

### Prompt for TableMaster – Textual Reasoning

## Objective

 You are provided with a table, a verbalization of the table, and a question related to the table. Your task is to reason step by step to answer the question based on the table.

## Table Definition

The table is represented by cell-value pairs, where each pair consists of a cell address and a value of the content in that cell, separated by a comma (e.g., 'A1,Year'). Multiple cells are separated by '|' (e.g., 'A1,Year|A2,Profit').

Cells may contain empty values, represented as 'A1, |A2, Profit'.

## Table {table}

## Verbalized Table {verbalized\_table}

**Question Answering** 

The answer should be short and simple. It can be a number, a word, or a phrase in the table, but not a full sentence. Your response should end with `Answer: xxx` (answer to the question).

Your response should end with `Answer: true` or `Answer: false` (answer to the question). If the table only contain headers and no rows, it indicates there is no information available for this question, therefore the answer should be "false."

*Now, give me the answer step by step: Question: {question}*  Fact Verification

*Figure 17.* Prompt for textual reasoning in *TableMaster*. Blue text represents placeholders for variables within the prompt, while the grey
 region indicates optional sections to adapt the prompt for question-answering or fact-verification tasks. The prompt guides the language
 model to answer the question step by step.

|   | Prompt for Table Master Taxtual Guidance Concretion                                                                                                              |
|---|------------------------------------------------------------------------------------------------------------------------------------------------------------------|
| _ | Frompt for TableMaster – Textual Guidance Generation                                                                                                             |
|   | ## Objective                                                                                                                                                     |
|   | You are provided with a table, a verbalized table, and a question related to the table.                                                                          |
|   | Your task is to give a step-by-step guidance to answer the question based on the table.                                                                          |
| I |                                                                                                                                                                  |
| I | ## Table Definition The table is represented by cell value pairs, where each pair consists of a cell address and a value of the content in that cell, constrated |
| l | hy a comma (e.a. '41 Year')                                                                                                                                      |
| l | Multiple cells are separated by ' ' (e.a., 'A1.Year   A2.Profit').                                                                                               |
| l | Cells may contain empty values, represented as 'A1, /A2, Profit'.                                                                                                |
|   |                                                                                                                                                                  |
|   | ## Table                                                                                                                                                         |
|   | {table}                                                                                                                                                          |
|   | ## Varbalized Tabla                                                                                                                                              |
|   | {verbalized table}                                                                                                                                               |
|   |                                                                                                                                                                  |
|   |                                                                                                                                                                  |
|   | You do not need to give the answer. You need to give a reasoning process as a guidance that will be used later.                                                  |
|   | ## Response Format                                                                                                                                               |
|   | The response should be a list of steps:                                                                                                                          |
|   | 1. xxx                                                                                                                                                           |
|   | 2. xxx                                                                                                                                                           |
|   |                                                                                                                                                                  |
|   |                                                                                                                                                                  |
|   | Now, give me the guidance to answer the question step by step:                                                                                                   |
|   |                                                                                                                                                                  |
|   |                                                                                                                                                                  |
| u | re 18. Prompt for textual guidance generation in TableMaster. Blue text indicates placeholders for variables within the prom                                     |
| n | npt guides the language model to generate textual guidance that can be utilized for subsequent symbolic reasoning.                                               |
|   |                                                                                                                                                                  |
|   |                                                                                                                                                                  |
|   |                                                                                                                                                                  |
|   |                                                                                                                                                                  |
|   |                                                                                                                                                                  |
|   |                                                                                                                                                                  |
|   |                                                                                                                                                                  |
|   |                                                                                                                                                                  |
|   |                                                                                                                                                                  |
|   |                                                                                                                                                                  |
|   |                                                                                                                                                                  |
|   |                                                                                                                                                                  |
|   |                                                                                                                                                                  |
|   |                                                                                                                                                                  |
|   |                                                                                                                                                                  |

### Prompt for TableMaster – Symbolic Reasoning (Programming of Thought)

| ## Objective |
|--------------|
| Vou are pro  |

You are provided with a table, a verbalized table, a guidance, and a question related to the table. Your task is to generate Python code that answers the question using the table and the quidance as a quide. ## Table Definition The table is represented by cell-value pairs, where each pair consists of a cell address and a value of the content in that cell, separated by a comma (e.g., 'A1, Year'). Multiple cells are separated by '|' (e.g., 'A1, Year | A2, Profit'). Cells may contain empty values, represented as 'A1, |A2, Profit'. ## Table {table} ## Verbalized Table {verbalized\_table} ## Guidance {textual\_guidance} ## Question {question} ## Instructions 1. The actual data of the table is stored in the variable `table` as a list of lists. 2. The result should be store in the variable `answer` as a string and do not need to print it. 3. You need to generate Python code within ```python``` code block. Now, give me the executable python code to answer the question: ```python table = {table\_array}

Figure 19. Prompt for symbolic reasoning in *TableMaster*. Blue text indicates placeholders for variables within the prompt. The prompt guides the language model to generate Python code to answer the question.

### Prompt for TableMaster – Answer Formatting

## Objective

You are provided with a process of text-guided reasoning with programming and a question related to the table. Your task is to answer the question using the reasoning process. ## Table Definition The table is represented by cell-value pairs, where each pair consists of a cell address and a value of the content in that cell, separated by a comma (e.g., 'A1, Year'). Multiple cells are separated by '|' (e.g., 'A1, Year |A2, Profit'). Cells may contain empty values, represented as 'A1, |A2, Profit'. ## Table {table} ## Textual Reasoning Process {textual\_reasoning\_process} ## Programmed Reasoning Process {symbolic\_reasoning\_process} The answer should be short and simple. It can be a number, a word, or a phrase in the table, but not a full sentence. Your response should be in the format of 'Answer: xxx' (answer to the question).

Question: {question} Answer:

*Figure 20.* Prompt for answer formatting in *TableMaster*. Blue text indicates placeholders for variables within the prompt. The prompt guides the language model to format the final answer based on the given table, question, and reasoning process.

#### **O. Prompts Used in Analysis Experiments**

### **Prompt for Direct Baseline**

## Objective You are provided with a table and a question related to the table. Your task is to answer the question directly based on the table. ## Table {table} ## Question {question} The answer should be short and simple. It can be a number, a word, or a phrase in the table, but not a full sentence. Now, answer the question directly: Answer:

Figure 21. Direct prompt for table understanding in analysis experiment. Blue text indicates placeholders for variables within the prompt. The prompt guides the language model to directly give the final answer based on the given table and question.

### **Prompt for Chain of Thought Baseline**

| ## Obj  | jective                                                                                                          |
|---------|------------------------------------------------------------------------------------------------------------------|
| You ar  | re provided with a table and a question related to the table.                                                    |
| Your t  | ask is to answer the question step by step based on the table.                                                   |
| ## Tal  | hle                                                                                                              |
| (table  |                                                                                                                  |
| {luble  | 1                                                                                                                |
| ## Qu   | estion                                                                                                           |
| {quest  | tion}                                                                                                            |
| {quest  | ion}                                                                                                             |
| The ar  | nswer should be short and simple. It can be a number, a word, or a phrase in the table, but not a full sentence. |
| Your r  | esponse should end with `Answer: xxx` (answer to the question).                                                  |
| Now     | answer the question step hyston.                                                                                 |
| 110W, 1 | unswer the question step by step.                                                                                |

Figure 22. Chain of thought prompt for table understanding in analysis experiment. Blue text indicates placeholders for variables within the prompt. The prompt guides the language model to give the answer step by step based on the given table and question. 

| 2015 | the prompt The prompt galaxy are language model to give the answer step of step case |
|------|--------------------------------------------------------------------------------------|
| 2016 |                                                                                      |
| 2017 |                                                                                      |
| 2018 |                                                                                      |
| 2019 |                                                                                      |
| 2020 |                                                                                      |
| 2021 |                                                                                      |
| 2022 |                                                                                      |
| 2023 |                                                                                      |
| 2024 |                                                                                      |
| 2025 |                                                                                      |
| 2026 |                                                                                      |
| 2027 |                                                                                      |
| 2028 |                                                                                      |
| 2029 |                                                                                      |
| 2030 |                                                                                      |
| 2031 |                                                                                      |
| 2032 |                                                                                      |
| 2033 |                                                                                      |
| 2034 |                                                                                      |
|      |                                                                                      |

### **Prompt for Program of Thought Baseline**

You are provided with a table and a question related to the table.

#### ## Objective

Your task is to answer the question based on the table by writing python code as a solution. ## Table

2035 2036

2038 2039

2040

2041

2042

2043 2044

2045 2046

2047

2048

2049

2050

2051

2052

2054

2055

2056 2057

2069

2073

2074

2076

2078

2079

2080 2081

2082

### {table}

## Reasoning Instructions

1. You must use executable python code to solve the question.

- 2. The final answer should be variable named "answer" in the code.
- 3. Do not execute the code in the response.
- 4. The python code should be in the following format:

```python # your code here

> Now, answer the question by writing python code as a solution: Question: {question} ``python

Figure 23. Program of thought prompt for table understanding in analysis experiment. Blue text indicates placeholders for variables within 2058 the prompt. The prompt guides the language model to generate code to derive the answer based on the given table and question. 2059

Prompt for Verbalization Baseline

Objective

You are provided with a table in string format. Your task is to convert the table into a detailed text description.

Table

{table}

Instructions

1. Provide a comprehensive description of the table.

- 2. Include all details and numerical values from the table in your response.
- 3. Do not omit or summarize any information from the table.
- 4. You may use external knowledge to enhance your understanding of the table, but the response must remain faithful to the table's content.

Now, please provide the verbalized description of the table:

Figure 24. Prompt for table verbalization in analysis experiment. Blue text indicates placeholders for variables within the prompt. The 2084 prompt guides the language model to verbalize a table to add detailed description. 2085

- 2086
- 2087

| ## Objective
You are provided with a table, and a question related to the table.
Your task is to give a step-by-step guidance to answer the question based on the table. |
|---|
| ## Table Definition |
| The table is represented by cell-value pairs, where each pair consists of a cell address and a value of the content in that cell, separated by a comma (e.g., 'A1,Year'). |
| Multiple cells are separated by ' ' (e.g., 'A1, Year A2, Profit').
Cells may contain empty values, represented as 'A1, A2, Profit'. |
| ## Table {table} |
| You do not need to give the answer. You need to give a reasoning process as a guidance that will be used later.
Keep the reasoning process concise and clear.
Control the number of steps in the reasoning process in the range of 1-5. |
| ## Response Format |
| The response should be a list of steps:
1. xxx |
| 2. xxx |
| Now, give me the guidance to answer the question step by step: |
| Question: {question} |

Prompt for Reasoning Strategy Evaluation

Question: {question}

Table: {table}

| Method 1 Solution: {cot_prediction} |
|-------------------------------------|
| Method 1 Reasoning: {cot_reasoning} |

Method 2 Solution: {pot_prediction} Method 2 Reasoning: {pot_reasoning}

Please evaluate which method is better. Respond in the following JSON format: {{ "better_method": 1 or 2

}}

Figure 26. Prompt for reasoning strategy evaluation in analysis experiment. Blue text indicates placeholders for variables within the prompt. The prompt guides the language model to select the better reasoning process after table reasoning.

| You are provide
Your task is to | ed with a table and a question related to the table.
assess whether answering this question needs mathematical calculation. |
|------------------------------------|--|
| ## Table | |
| {table} | |
| ## Question | |
| {question} | |
| ## Instructions | |
| 1. If the question | on can be easily answered using the information in the table, respond with False. |
| 2. If the question | on involves comparison, respond with False. |
| 2. When the qu | estion involves counting a substantial number (more than 5) of items or rows, respond with True. |
| 3. If the question of the True | on demands complex calculations or multi-step mathematical operations based on the table's data, the response shou |
| 4. For simple a | ithmetic or small-scale counting that requires minimal computational effort, respond with False. |
| ## Response Fo | ormat |
| The response s | hould be in ISON format: |
| ```json | |
| {{ | |
| "need_calcul | ation": true/false |
| }} | |
| | |

Figure 27. Prompt for reasoning strategy evaluation in analysis experiment. Blue text indicates placeholders for variables within the prompt. The prompt guides the language model to select the better reasoning strategy before table reasoning.

| Promp | t for C | Juestion | Τνι | pe Classificati | on (| Calculation | Requ | uired) |
|-------|---------|-----------------|-----|-----------------|------|-------------|------|--------|
| | | | | | | | | |

| Table:
{table} |
|---|
| Question: |
| {question} |
| Determine whether a calculation is required to answer the question, or if the question can be directly answered using the information in the table. |
| Provide your response in the following JSON format: |
| {{ |
| "need_calculation": true/false |
| |

Figure 28. Prompt for classifying a question type based on whether calculation is required in the analysis experiment. Blue text indicates placeholders for variables within the prompt.

Prompt for Noised Table Generation

Given a table, you need to generate a new table by disrupting the content in the table.

Table:

{table}

Rules:

- Your goal is to make the content in each row within the same column follows a different format to increase diversity as much as possible.

- You cannot change the structure of the table.

- You cannot add or remove any rows or columns.

- You cannot modify the column names in the first row.

- You can only alter the format of the content in each cell, not the actual values.

- You should not make the content in each row within the same column in the same format as much as possible.

Format Change Examples:

- Change a number format from 123456 to 123,456.

- Change a date format from 2024-01-01 to 2024/01/01.

- Simplify or abbreviate text content.

Provide your new table in the following JSON format:
""json
{{
"table": [[...], [...]],

}}

Figure 29. Prompt for generating noised tables in the analysis experiment. Blue text represents placeholders for variables within the prompt. The prompt instructs the language model to add noise by altering the cell content format based on a given table.

P. Notion Table

Table <u>10</u> provides a comprehensive list of the notations used throughout this paper, along with their corresponding descriptions. This table serves as a quick reference to help readers better understand the concepts presented in our work.

| 2259 | | |
|------|------------------|---|
| 2260 | | Table 10. Notation used throughout the paper |
| 2261 | Notation | Description |
| 2262 | Notation | Description |
| 2263 | General | |
| 2264 | Q | Given question or query |
| 2265 | A | Generated answer |
| 2266 | \mathbb{T} | Input table |
| 2267 | \mathbb{T}^W | Wild table before normalization |
| 2268 | \mathbb{T}^N | Normalized table |
| 2269 | \mathbb{T}^F | Table-of-Focus |
| 2270 | $C_{i,j}$ | Cell in the <i>i</i> -th row and <i>j</i> -th column |
| 271 | m, n | Number of rows and columns in the table |
| 272 | Table Stru | cture Understanding |
| 273 | H | Set of ton headers |
| 274 | II
K | Key column serving as row identifier |
| 275 | n
C | Candidate column set |
| 276 | C^0 | Selected relevant columns |
| 277 | P | Selected relevant cours |
| 278 | n k | Deak size for table processing |
| 279 | <u> </u> | reek size for table processing |
| 280 | Table Cont | tent Understanding |
| 281 | $T^{\mathbb{T}}$ | Verbalized table (natural language text) |
| 282 | a',b' | Number of refined columns and rows after reconstruction |
| 283 | Table Reas | soning |
| 284 | S | Selected reasoning strategy |
| 285 | \mathcal{T} | Textual reasoning strategy |
| 286 | S | Symbolic reasoning strategy |
| 287 | G | Textual reasoning guidance |
| 288 | \mathcal{P} | Program executor (Python/SOL) |
| 289 | | |
| 290 | | |
| 291 | | |
| 292 | | |
| 293 | | |
| 294 | | |
| 295 | | |
| 296 | | |
| 297 | | |
| 298 | | |
| 299 | | |
| 300 | | |
| 2201 | | |