# TREE-OF-TABLE: UNLEASHING THE POWER OF LLMS FOR ENHANCED LARGE-SCALE TABLE UNDERSTANDING

**Deyi Ji**[1,2]**, Lanyun Zhu**[3]**, Siqi Gao**[2]**, Peng Xu**[2]**, Hongtao Lu**[4]**, Jieping Ye**[2]**, Feng Zhao**[1]
[1]University of Science and Technology of China    [2]Alibaba Group
[3]Singapore University of Technology and Design    [4]Shanghai Jiao Tong University

## ABSTRACT

The ubiquity and value of tables as semi-structured data across various domains necessitate advanced methods for understanding their complexity and vast amounts of information. Despite the impressive capabilities of large language models (LLMs) in advancing the natural language understanding frontier, their application to large-scale tabular data presents significant challenges, specifically regarding table size and complex intricate relationships. Existing works have shown promise with small-scale tables but often flounder when tasked with the complex reasoning required by larger, interconnected tables found in real-world scenarios. To address this gap, we introduce "Tree-of-Table", a novel approach designed to enhance LLMs' reasoning capabilities over large and complex tables. Our method employs Table Condensation and Decomposition to distill and reorganize relevant data into a manageable format, followed by the construction of a hierarchical Table-Tree that facilitates tree-structured reasoning. Through a meticulous Table-Tree Execution process, we systematically unravel the tree-structured reasoning chain to derive the solutions. Experiments across diverse datasets, including WikiTQ, TableFact, FeTaQA, and BIRD, demonstrate that Tree-of-Table sets a new benchmark with superior performance, showcasing remarkable efficiency and generalization capabilities in large-scale table reasoning.

## 1    INTRODUCTION

Tables, as a pivotal form of semi-structured data, ubiquitously underpin numerous aspects of daily life and professional domains, ranging from open data repositories and web pages to critical applications in financial analysis, risk management, health monitoring, and business reporting (Cafarella et al., 2008). The advent of large language models (LLMs) (OpenAI, 2023; Chen, 2023; Jiang et al., 2023; Imani et al., 2023; Anil et al., 2023; Zhu et al., 2024a; Valmeekam et al., 2022; Zhu et al., 2024c) has opened new vistas for understanding and reasoning with tabular data, marking a significant stride in the realm of natural language understanding (Nahid & Rafiei, 2024; Chen et al., 2024; Sui et al., 2024b;a; Ye et al., 2023; Cheng et al., 2022; Jin & Lu, 2023). This intersection is not only instrumental in enhancing the comprehension of tables but also vital for powering a plethora of downstream tasks such as table-based fact verification (Chen et al., 2019) and question answering (Li et al., 2024). Unlike their unstructured text counterparts, tables provide a dense, structured format through the interaction of rows and columns, offering a rich source of information. However, the same structural characteristics pose unique challenges for language models, as they necessitate advanced levels of reasoning over both the textual and numerical data contained within. Given the increasing reliance on tables for data representation and the complexities involved in their interpretation, investigating the integration of LLMs for improved large-scale table understanding has emerged as an essential and compelling research avenue, drawing heightened interest from the global academic and industrial research communities.

Existing methods for table understanding have shown substantial progress in comprehending small-scale tables (Cheng et al., 2022; Ye et al., 2023; Wang et al., 2024). However, these approaches often falter when applied to the more complex and larger tables frequently encountered in real-world scenarios. This gap between academia and practical applications stems from a variety of limitations
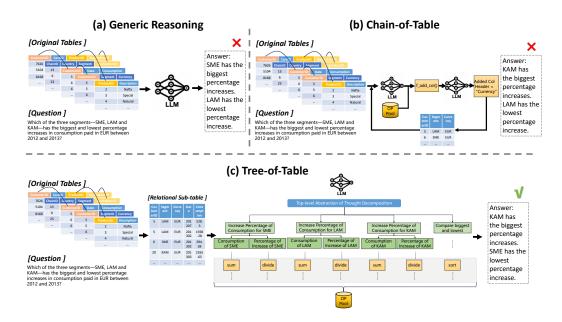
Figure 1: Comparison of (a) Generic Reasoning, (b) Chain-of-Table (Wang et al., 2024), and the proposed (c) Tree-of-Table methods when confronted with large-scale relational tables. Generic Reasoning often struggles with the increased context and complexity, leading to inefficient processing and potential loss of critical information. Chain-of-Table, while more structured with linear thought chain, still faces challenges with the scale and intricacy of data. In contrast, Tree-of-Table showcases a structured and hierarchical reasoning process that adeptly handles large-scale tables, significantly enhancing comprehension and efficiency compared to previous methods, particularly in managing the complexity of expansive tabular data.

inherent in current methodologies. One significant challenge is the limited contextual capacity of today's language models. As tables increase in size, the amount of information that must be processed and understood grows exponentially due to the intricate interactions between rows and columns. This complexity makes it difficult for models to capture and reason about all the necessary information in one go, significantly impeding their understanding capabilities. When faced with complex question-answering logic that spans lengthy chains, pinpointing, extracting, and comprehending key table information becomes an immense challenge.

To address these issues, two main approaches have generally been adopted, as shown in Figure 1. The first involves using only the schema information of tables and employing program-aided methods, such as generating SQL-based answers from questions (Rajkumar et al., 2022b;a; Shi et al., 2020; Pönighaus, 1995; Katsogiannis-Meimarakis & Koutrika, 2023). While this approach avoids directly inputting entire tables, the resulting SQL statements can be lengthy and prone to errors, leading to suboptimal performance. The second strategy involves decomposing tables into multiple sub-tables. Methods like Dater (Ye et al., 2023) attempt to manage larger tables by initially inputting the entire table before breaking it down, which is impractical. The Chain-of-Table (Wang et al., 2024) draws inspiration from the chain-of-thought principle (Wei et al., 2022), performing implicit sub-table extraction. Yet, even this approach is limited to understanding smaller tables. Additionally, traditional table understanding datasets like WikiTQ (Pasupat & Liang, 2015) and TableFact (Chen et al., 2019), which are relatively small, severely restrict the exploration of large-scale table understanding. Fortunately, the introduction of the BIRD (Li et al., 2024) dataset, considered the largest and most complex table understanding dataset to date, highlights the pressing need for improvements. Despite this, due to the reasons mentioned above, existing large language models still exhibit low accuracy on comprehensive, large-scale table datasets like BIRD (Li et al., 2024), signaling a clear necessity for methodological innovations in this area.

Addressing these concerns, we propose "Tree-of-Table", a novel paradigm crafted to optimize LLMs for the task of large-scale table understanding, as shown in Figure 1. By condensing and decomposing tables, our approach distills and systematizes the critical information into a tree-structured model

that resonates with the stepwise reasoning employed by humans. This tree acts as a roadmap, guiding the LLM through the complexities of the table in a logical and organized manner. It provides a structured approach where each node serves a purpose, simplifying the interaction between the LLM and the tabular data. The efficacy of our Tree-of-Table methodology is emphatically validated through rigorous testing across a selection of datasets (including WikiTQ (Pasupat & Liang, 2015), TableFact (Chen et al., 2019), FeTaQA (Nan et al., 2022), and BIRD (Li et al., 2024)) with each presenting its own unique challenges. Consistently achieving top-tier results, Tree-of-Table demonstrates not just its capacity to navigate the intricacies of table reasoning but also its potential to set a new benchmark in the field.

## 2 RELATED WORK

### 2.1 TRADITIONAL TABLE UNDERSTANDING.

With the development of deep learning (LeCun et al., 2015; He et al., 2016; Devlin, 2018; Ji et al., 2022; Zhu et al., 2023a;b; Hu et al., 2020; Wang et al., 2021b;a; Xu et al., 2024; Zhu et al., 2024b) and foundational models (Devlin, 2018; Radford et al., 2021; Touvron et al., 2023a; Ji et al., 2024b; Wang et al., 2022) in natural language processing (Mikolov, 2013; Radford et al., 2019; Kaplan et al., 2020; OpenAI, 2023; Bai et al., 2023) and computer vision (He et al., 2016; Alayrac et al., 2022; Ji et al., 2021; Dosovitskiy, 2020; Ji et al., 2023b; Zhu et al., 2021; Ji et al., 2024a; Zhang et al., 2022a; Wang et al., 2023; Ji et al., 2023a), early table understanding generally adopted the approach of generating SQL statements. At the core, traditional methods have focused on generating executable languages like SQL (Rajkumar et al., 2022b; Liu et al., 2021; Eisenschlos et al., 2020; Jiang et al., 2022) to interact with tables. This approach stems from the need to reason over both free-form natural language questions and (semi-)structured tables. While effective in accessing tabular data, these methods often fall short in capturing the nuanced semantics within a table, particularly struggling with web tables that feature free-form text in cells.

### 2.2 PROMPTING LANGUAGE MODELS FOR TABLE UNDERSTANDING.

A novel stride in table understanding has been the application of prompting strategies (Wei et al., 2022; Chen et al., 2022; OpenAI, 2023; Imani et al., 2023; Khot et al., 2022; Zhang et al., 2022b). By generating reasoning steps through in-context learning, models like Chain-of-Thought (Wei et al., 2022) and its evolutions (Yao et al., 2023) break down questions into sub-problems, iteratively solving each to improve comprehension of complex tasks. These methods showcase LLMs' prowess in handling intricate reasoning chains, albeit not being explicitly designed for tabular data. Emerging approaches have sought to extend LLM capabilities beyond text, incorporating external tools to solve reasoning tasks (Cheng et al., 2022; Hsieh et al., 2023; Dhingra et al., 2019; Liu et al., 2023). Generating Python or SQL programs (Cheng et al., 2022; Nahid & Rafiei, 2024; Shi et al., 2020; Pönighaus, 1995) and executing them with interpreters or APIs has shown promise in enhancing arithmetic and table-based reasoning. However, the performance of these program-aided methods sometimes falters in complex table scenarios due to the static nature of tables in the reasoning process. Dater (Ye et al., 2023) dynamically modifies the tabular context to aid in solving table-based tasks, albeit primarily focusing on data pre-processing with limited operations. Subsequently, the Chain-of-Table (Wang et al., 2024) method is inspired by the chain-of-thought (Wei et al., 2022) principle and performs implicit sub-table extraction. Contrarily, our proposed Tree-of-Table approach is inspired by the tree-of-thought principle (Yao et al., 2023), creating adaptive tree-based reasoning chains that exploit the planning capabilities of LLMs for more nuanced and context-specific table reasoning.

### 2.3 TABLE UNDERSTANDING DATASETS.

Datasets like WikiTQ (Pasupat & Liang, 2015) and TableFact (Chen et al., 2019) have been instrumental in developing table understanding methods. These standard benchmarks provide a foundation but are often limited in size and complexity. BIRD (Li et al., 2024) represents a significant leap forward in the field, being one of the largest and most intricate datasets designed for table understanding to date. Spanning across 37 professional domains with a substantial size of 33.4 GB, BIRD offers over 12,000 examples gleaned from real-world databases. Its development involved modifying

open-source relational databases and curating additional ones, all complemented by crowdsourced natural language questions and corresponding SQL queries.

## 3 Tree-of-Table: Unleashing the Power of LLMs

### 3.1 Formulation of Large-Scale Table Understanding

In the domain of table understanding, the core challenge lies in accurately interpreting and extracting information from tabular data in response to a given natural language query or statement. The essence of table understanding can be encapsulated as the task of mapping a natural language question or statement $Q$ to a corresponding output $S$ that accurately reflects the information contained within a table $T$. This table can be characterized by its structure, which includes rows and columns, with each cell representing a specific data point. Formally, a table $T$ can be divided into headers $H$ and data values $D$, where each header in $H$ corresponds to a column in the table and $D$ represents the collective data points contained within these columns.

Furthermore, table understanding involves not just the direct interpretation of tables but also potentially requires external knowledge $K$ and conversion of table data into a format that is amenable to computational models. This is especially relevant for tasks that involve complex reasoning or necessitate an understanding beyond the explicit table content, such as requiring background knowledge or contextual understanding to correctly interpret the question or the data.

In our experiments, we utilize a total of four datasets: WikiTQ (Pasupat & Liang, 2015), TabFact (Chen et al., 2019), FeTaQA (Nan et al., 2022), and BIRD (Li et al., 2024). For WikiTQ, TabFact, and FeTaQA, there is no external knowledge; therefore, the table understanding problem can be defined as finding a function or model $f(\cdot, \theta)$ that satisfies

$$S = f(Q, \langle H, D \rangle | \theta), \tag{1}$$

where $\theta$ represents the model parameters. In contrast, for the BIRD dataset, there is external knowledge used to explain specific terms in the questions, allowing us to define the table understanding problem as

$$S = f(Q, \langle H, D \rangle, K | \theta), \tag{2}$$

where $K$ denotes the external knowledge.

### 3.2 Overview

In this work, we introduce a novel approach named "Tree-of-Table" devised to address the challenge of table reasoning within large-scale table understanding datasets (e.g., BIRD) and real-world applications, as shown in Figure 2 and Figure 3. Our methodology encompasses several steps designed to simplify and enhance the reasoning capabilities of LLMs when confronted with large, interconnected tables. First, we condense the tables based on the specific requirements of the query. This process identifies relevant portions of the tables, thereby reducing the cognitive load on LLMs. We then apply a tree-based decomposition strategy to segment large tables into smaller, manageable units, guided by the relationships among tables, such as foreign keys, and the structure of the query. Next, we construct a "Table-Tree" by reorganizing the condensed information into a hierarchical structure. Each node in this tree represents a logical block of information or a step in the reasoning process, mirroring the cognitive approach of breaking down complex problems into simpler sub-problems. Finally, we perform a sequential traversal of the constructed Table-Tree to derive answers to the queries. This systematic traversal ensures logical progression through each node, allowing for the synthesis of information and insights gained from previous steps. The iterative nature of this reasoning process culminates in a well-informed conclusion.

### 3.3 Table Condensation and Decomposition

Addressing the significant challenges posed by large-scale relational tables to LLMs requires a nuanced understanding of the specific difficulties in question. These challenges primarily stem from two aspects: (1) The intricate foreign key relationships among multiple tables, which are commonly defined using SQL syntax, may not be readily interpretable by LLMs due to their complexity and
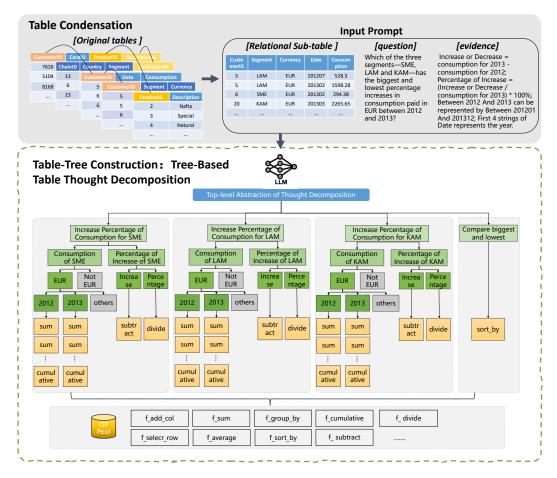
Figure 2: Illustration of the initial phases in the Tree-of-Table methodology, encompassing Table Condensation (the upper part), followed by Table-Tree Construction (the lower part). Starting with a large-scale input table, the process selectively condenses the data, emphasizing task-relevant information. Subsequently, the decomposed elements are methodically reorganized into a Table-Tree, a hierarchical structure designed to streamline and guide the subsequent reasoning process.

the specialized knowledge required to understand relational database schemas. (2) The sheer size of the tables often exceeds the input context limit of LLMs, making it impossible for these models to process the entirety of the data directly. To effectively address these issues, our methodology incorporates two key processes: Table Condensation and Tree-based Decomposition.

### 3.3.1 TABLE CONDENSATION

As shown in the upper part of Figure 2, our initial step involves condensing the tables based on the context of the question $Q$ and any additional evidence provided. Since there are possibly multiple tables, this process employs LLMs to identify one sub-table relevant to $Q$ from them through schema-linking (Lei et al., 2020). Following the identification of relevant schemas/headers, we merge these multiple tables to reduce redundancy and decrease their size,

$$\text{subTable}_Q = f(\text{Schema\_Link}(Q, \{H\}, K)|\theta), \qquad (3)$$

where $\{H\}$ indicates the headers of the tables. This condensation aims to recall one sub-table pertinent to $Q$ and eliminate superfluous table information, thereby enhancing the information density related to $Q$ within the tables and making it more manageable for LLM processing. Through a detailed analysis of the BIRD dataset, we observed that over 70% of questions involved tables whose length exceeded the input limitations of current LLMs. Furthermore, more than 90% of these questions pertained to at least two tables, with 20% involving four or more tables, significantly complicating

the understanding process for LLMs. Post-condensation, we found that the length of tables involved in more than 60% of long questions was reduced below the LLM input limit, and all questions were associated with a singular, condensed table, as shown in the upper part in Figure 2.

### 3.3.2 TREE-BASED DECOMPOSITION

Even with reduced size and complexity post-condensation, the tables might still be too lengthy or intricate for LLMs to handle efficiently, occasionally still surpassing the models' input limits. To mitigate this, we begin by breaking down the question $Q$ into its most general components, delineating the entire problem-solving process into several independent yet sequentially connected steps.

$$S = \mathcal{P}_{\text{decomp}}(\{S_i^1\}, r^1|Q), \quad i < \text{MAXDegree}, \tag{4}$$

where $S$ is the final solution, $\{S_i^1\}$ is the firstly decomposed intermediate sub-solutions towards $S$, $r^1$ is the possible relationship between $\{S_i^1\}$. $\mathcal{P}_{\text{decomp}}$ is the "Thought Decomposition Prompt". "MAXDegree" is the pre-defined maximum degree of the Table-Tree. This decomposition involves mapping out the key stages ($\{S_i^1\}$ and $r^1$) of reasoning required to address $Q$. By doing so, we transform a potentially overwhelming task into a series of manageable sub-tasks, each contributing incrementally to the formulation of the final answer. $\{S_i^1\}$ and $r^1$ also serve as the root node of the first-level subtree we will construct in the Table-Tree, for example, as shown the lower part in Figure 2.

### 3.4 TABLE-TREE CONSTRUCTION

Drawing inspiration from the Tree-of-Thought concept (Yao et al., 2023), our table-tree structure closely resembles how humans naturally approach problem-solving. The illustration of overall construction is showed in the lower part in Figure 2. When faced with a complex problem, people typically employ a "breadth-first" strategy: deconstructing the problem into several general, independent yet interconnected subprocesses and then iteratively refining each subprocess into finer-grained solutions.

### 3.4.1 BREADTH-FIRST THOUGHT GENERATION

Within this framework, we utilize in-context learning to instruct LLMs on dynamically generating thoughts for the question in a breadth-first way. Based on the firstly decomposed $\{S_i^1\}$ and $r^1$ in Eq. 4, the following breadth-first thought generation process can be formulated as,

$$
\begin{aligned}
S &= \mathcal{P}_{\text{decomp}}(\{S_i^1\}, r^1|Q), \quad i < \text{MAXDegree}, \\
S_i^1 &= \mathcal{P}_{\text{decomp}}(\{S_{i,j}^2\}, r_j^2), \quad j < \text{MAXDegree}, \\
&\quad\quad\quad \ldots\ldots \\
S_{i,j,\ldots}^d &= \mathcal{P}_{\text{decomp}}(\{S_{i,j,\ldots,k}^{d+1}\}, r_{i,j,\ldots,k}^{d+1}), \quad k < \text{MAXDegree}, \\
&\quad\quad\quad \ldots\ldots \\
S_{i,j,\ldots,k,\ldots}^{d_{\max}-1} &= \mathcal{P}_{\text{decomp}}(\{S_{i,j,\ldots,k,\ldots,l}^{d_{\max}}\}, r_{i,j,\ldots,k,\ldots,l}^{d_{\max}}), \quad d_{max} <= \text{MAXDepth},
\end{aligned}
\tag{5}
$$

where $S$ is the final solution, $\{S_i^1\}$ is the firstly decomposed intermediate solutions towards $S$, $r^1$ is the possible relationship between $\{S_i^1\}$. $\{S_{i,j,\ldots,k}^{d+1}\}$, $r_{i,j,\ldots,k}^{d+1}$, and so on. Note that $r_{i,j,\ldots,k}^{d+1}$ may be empty in the actual decomposition process if we need not to consider the relationship between $\{S_{i,j,\ldots,k}^{d+1}\}$. $d$ is the depth of thought. "MAXDegree" and "MAXDepth" are the pre-defined maximum degree and depth of the Table-Tree, respectively.

To prevent excessive decomposition of thoughts that could lead to redundant or erroneous reasoning processes, we set a maximum value for the depth $d$, denoted as "MAXDepth", and follow (Yao et al., 2023) to utilize the LM to deliberately reason about end thought states $\{S_{i,j,\ldots,k,\ldots,l}^{d_{\max}}\}, r_{i,j,\ldots,k,\ldots,l}^{d_{\max}}$. Such a deliberate heuristic can be more flexible than programmed rules.
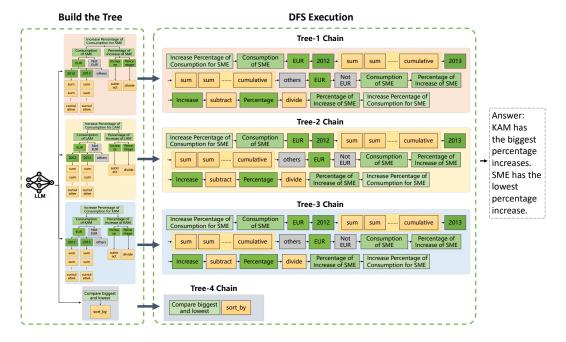
Figure 3: Depiction of the Table-Tree Execution phase within the Tree-of-Table approach. The model traverses the hierarchical Table-Tree, processing each node sequentially from the root to the leaves. At each step, the model integrates the information from the current node with the insights gathered from previous nodes, systematically building upon the reasoning chain to derive the final answer.

### 3.4.2 ITERATIVE CONSTRUCTION

Building upon the breadth-first thought generation, we construct the Table-Tree by iteratively constructing child nodes level by level for $\{S_i^1\}$ and $r^1$, until we reach the leaf nodes at the bottom of the tree. Each sub-thought corresponding to $\{S_{i,j,\dots,k}^{d+1}\}$ and $r_{i,j,\dots,k}^{d+1}$ is regard as intermediate node, which act as "thinking node" representing a subprocess that can be further decomposed. The end thought state $\{S_{i,j,\dots,k,\dots,l}^{d_{\max}}\}, r_{i,j,\dots,k,\dots,l}^{d_{\max}}$ are set to leaf nodes, which functions as "execution nodes". In concrete, leaf nodes are the actionable endpoints of Table-Tree where specific operations such as data retrieval, calculations, or logical evaluations occur based on the parameters defined by their parent nodes. Therefore, we formulate $\{S_{i,j,\dots,k,\dots,l}^{d_{\max}}\}, r_{i,j,\dots,k,\dots,l}^{d_{\max}}$ as

$$S_{i,j,\dots,k,\dots,l}^{d_{\max}}, r_{i,j,\dots,k,\dots,l}^{d_{\max}} = \mathcal{P}_{\text{sample}}(\{\text{OP\_Pool}\}), \quad d_{max} <= \text{MAXDepth}. \qquad (6)$$

where $\mathcal{P}_{\text{sample}}$ is the "Operation Sample Prompt". In selecting the operation pool, we based it on (Wang et al., 2024) and chose the most frequently used table operations from the resource at (Bytescout, 2024).

**Comparison with Chain-of-Table:** Notably, the previous work Chain-of-Table (Wang et al., 2024) processes the entire chain of history for each dynamic planning step and is shown to be effective for relatively small tables, such as WikiTQ. However, as tables grow larger and questions become more complex, maintaining the complete thought chain becomes cumbersome, ultimately decreasing the efficiency of the model. Our Tree-of-Table method addresses this by embracing the tree's inherent "divide and conquer" philosophy (Bentley, 1980) to construct a Table-Tree. Each generation of child nodes relies exclusively on the information from their multi-level parent nodes, without the need for uncle nodes, as illustrated in Figure 2. By this way, we significantly reduce the historical chain's length on which each node's dynamic planning relies, to less than the depth of the tree, thus considerably simplifying the generation process at each level.

Table 1: Comparison of Table Understanding results on WikiTQ, TabFact datasets, with GPT3.5, PaLM2 and LLaMA2.

| Method | WikiTQ | | | TabFact | | |
|---|---|---|---|---|---|---|
| | GPT3.5 | PaLM2 | LLaMA2 | GPT3.5 | PaLM2 | LLaMA2 |
| Text-to-SQL (Rajkumar et al., 2022b) | 52.90 | 52.42 | 36.14 | 64.71 | 68.37 | 64.03 |
| End-to-End QA (Wang et al., 2024) | 51.84 | 60.59 | 23.90 | 70.45 | 77.92 | 44.86 |
| Few-Shot QA (Wang et al., 2024) | 52.56 | 60.33 | 35.52 | 71.54 | 78.06 | 62.01 |
| Binder (Cheng et al., 2022) | 56.74 | 54.88 | 30.92 | 79.17 | 76.98 | 62.76 |
| Chain-of-Thought (Wang et al., 2024) | 53.48 | 60.43 | 36.05 | 65.37 | 79.05 | 60.52 |
| Dater (Ye et al., 2023) | 52.81 | 61.48 | 41.44 | 78.01 | 84.63 | 65.12 |
| Chain-of-Table (Wang et al., 2024) | 59.94 | 67.31 | 42.61 | 80.20 | 86.61 | 67.24 |
| **TREE-OF-TABLE** | **61.11** | **68.77** | **44.01** | **81.92** | **87.88** | **69.33** |

## 3.5 TABLE-TREE EXECUTION

After constructing the Table-Tree, we view it as a proxy task for the entire table understanding procedure. As shown in Figure 3, by traversing and executing operations across this tree, LLMs can implicitly generate tables and save intermediary results, thus enabling a seamless reasoning process. This stage diverges from the construction phase by utilizing a depth-first search approach to execute the thought chain, ensuring a systematic and comprehensive exploration of the tree structure.

**Depth-First Search Execution.** The rationale behind adopting a depth-first search (DFS) (Tarjan, 1972) strategy for Tree Execution lies in its ability to fully explore and resolve each branch of the tree to completion before moving to the next. This method aligns with the logical progression of solving a complex problem by focusing on and completing one aspect of the problem entirely, ensuring that all necessary computations and logical deductions related to a branch are performed before considering alternative or subsequent branches.

**Leveraging Tree Structure for Efficiency.** A key advantage of the tree structure is its inherent ability to enhance reasoning efficiency by logically organizing and compartmentalizing different aspects of the problem-solving process into subtrees. To exploit this benefit to its fullest, we execute the reasoning process subtree by subtree, based on the root node's children. After processing a subtree, we store its result before proceeding. This approach contrasts with linearly merging all subtrees into a single chain for execution. By maintaining the distinction between subtrees and executing them as separate units, we significantly mitigate the risk of intermediary tables becoming excessively large and unwieldy, which in turn, would thwart the reasoning process.

## 4 EXPERIMENTS

### 4.1 DATASETS, METRICS, IMPLEMENTATION DETAILS

We evaluate our method on both three small table understanding benchmarks: WikiTQ (Pasupat & Liang, 2015), FeTaQA (Nan et al., 2022), and TabFact (Chen et al., 2019), and one large-scale dataset: BIRD (Li et al., 2024). For WikiTQ and TabFact, we employ the standard denotation accuracy metric. The nature of FeTaQA and BIRD for requiring elaborate responses prompts us to assess performance through a variety of metrics including BLEU, ROUGE-1, ROUGE-2, and ROUGE-L (Lin, 2004) to capture different facets of response quality. For our experiments, following the previous works (Wang et al., 2024), we leverage the computational prowess of advanced language models, namely PaLM 2 (Anil et al., 2023), GPT 3.5 (OpenAI, 2023), and LLaMA 2 (Touvron et al., 2023b). To facilitate in-context learning, we incorporate a few-shot approach using demo samples from the training set within the prompts, ensuring the models can effectively learn from limited examples.

### 4.2 MAIN RESULTS

In our experimental evaluation, we comprehensively compare our proposed approach, Tree-of-Table, with several renowned baselines and state-of-the-art methodologies across both small and large-scale tabular datasets, including WikiTQ, TableFact, FeTaQA, and BIRD, in Table 1 and Table 2. Our analysis is designed to assess the effectiveness of Tree-of-Table in facilitating complex table understanding and reasoning tasks, particularly highlighting its performance in challenging

Table 2: Comparison of Table Understanding results on FetaQA and BIRD datasets.

| Method | FeTaQA | | | | BIRD | | | |
|---|---|---|---|---|---|---|---|---|
| | BLEU | ROUGE-1 | ROUGE-2 | ROUGE-L | BLEU | ROUGE-1 | ROUGE-2 | ROUGE-L |
| FT(T5-large) (Ye et al., 2023) | 30.54 | 0.63 | 0.41 | 0.53 | - | - | - | - |
| End-to-End QA (Wang et al., 2024) | 28.37 | 0.63 | 0.41 | 0.53 | 9.90 | 0.44 | 0.18 | 0.43 |
| Codex (Chen et al., 2021) | 27.96 | 0.62 | 0.40 | 0.52 | - | - | - | - |
| Dater (Ye et al., 2023) | 29.47 | 0.63 | 0.41 | 0.53 | 10.65 | 0.44 | 0.18 | 0.43 |
| Chain-of-Table (Wang et al., 2024) | 32.61 | 0.66 | 0.44 | 0.56 | 12.12 | 0.49 | 0.22 | 0.48 |
| **TREE-OF-TABLE** | **34.73** | **0.68** | **0.46** | **0.58** | **15.70** | **0.53** | **0.26** | **0.52** |



(a) Generalization Ability under Different Table Sizes
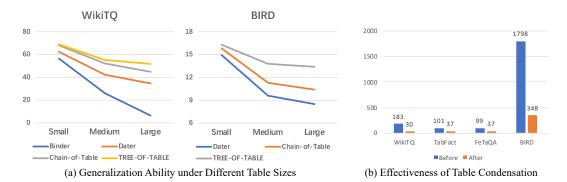
(b) Effectiveness of Table Condensation

Figure 4: Ablation study: (a) Generalization Ability under Different Table Sizes. (b) Effectiveness of Table Condensation.

scenarios involving large-scale tables. The results in Table 1, demonstrate that Tree-of-Table not only significantly outperforms early Generic Reasoning methods (End-to-End QA, Few-Shot QA, Chain-of-Thought) and Program-aid Reasoning methods (Text-to-SQL, Dater, Binder) but also surpasses the current state-of-the-art method, Chain-of-Table. This superiority was consistent across multiple LLMs including GPT3.5, PalM2, and LLAMA2. Our Tree-of-Table approach showcased robust enhancements in reasoning over both small and large tables. Specifically, in datasets with larger tables like BIRD, the benefits of Tree-of-Table were even more pronounced, suggesting that our tree-based method is particularly suited for complex, large-scale reasoning tasks where effective condensation, decomposition, and subtree execution strategies are critical. The results solidify Tree-of-Table as an effective architecture for table understanding, indicating that the tree-shaped mental model is highly adaptive and successful in managing the intricacies of multidimensional reasoning tasks. The framework's hierarchical structure, focused on breadth-first expansion and depth-based optimization, enables a level of complexity reduction and computational efficiency that linear and less structured approaches struggle to match.

## 4.3 ABLATION STUDY

**Generalization Ability under Different Table Sizes.** Here, we evaluate the generalization capacity of our method across tables of varying sizes. Large tables present considerable comprehension challenges to models, as their capacity to grapple with extended contexts and prompts is severely tested. In Figure 4, we provide a detailed comparison of 4 methodologies (Binder, Dater, Chain-of-Table, and Tree-of-Table) across two datasets (WikiTQ and BIRD). The performance metrics in table understanding tasks evidently deteriorate as the size of the tables increases. This degradation in performance reflects the inherent difficulties associated with large table comprehension, confirming that it remains an exceptionally challenging problem area. However, it is notable that the decline in performance with increasing table size is much more gradual for Tree-of-Table as compared to other methods. Especially on the large-scale table dataset BIRD, Tree-of-Table demonstrates superior robustness and generalization ability. Even as table size

Table 3: The Node Number and Height of Tree Chains.

| Dataset | WikiTQ | TabFact | BIRD |
|---|---|---|---|
| Chain-of-Table: Chain Length | 4 | 4 | 11 |
| Tree-of-Table: Tree Height | 3 | 3 | 7 |
| Tree-of-Table: Node Number | 6 | 7 | 18 |

scaled up, Tree-of-Table maintaines a
level of performance that was not only
better than its counterparts but also displayed less variance in its results.

**Node Number and Height of Tree Chains.** The height of the Table-Tree reflects the depth of the model's reasoning chain, indicative of the complexity of the reasoning process. In contrast, the node number corresponds to the length of the thought chain, representing the number of discrete reasoning steps taken by the model. In Table 3, our comparative analysis reveals that across both smaller and larger datasets, the average height of the Table-Trees is generally less than that of the Chain-of-Table. This indicates that the reasoning process in the Tree-of-Table method tends to require fewer levels of hierarchical reasoning to arrive at a solution. Additionally, the average length of the Table-Trees remains within a reasonable range, suggesting a good balance between depth and breadth in our tree structures.

**Comparison of Table Format Encoding.** The encoding format of tables plays a vital role in how effectively a model can interpret and manipulate table data. Early research has indicated that the specific form of table encoding can significantly impact the model's performance in table understanding tasks. Here, we follow the lead of prior work, comparing the effects of four distinct encoding formats on the final performance of table understanding: PIPE, HTML, TSV (Tab Separated Values), and Markdown. As shown in Table 4, the Markdown format leads to the highest performance among the tested encoding styles. The benefits of Markdown may be likely attributed to its readability, clear structure, and straightforward syntax, all of which align well with the parsing capabilities of LLMs.

**Efficiency Analysis.** In this context, efficiency refers to the model's capability to achieve its goals with the least amount of computational resources—specifically, the number of samples it needs to generate to arrive at a correct answer. To substantiate the efficiency of our Tree-of-Table methodology, we scrutinize how it compares with existing methods in terms of the number of required generated samples to solve tasks. For a comprehensive analysis, we compared Tree-of-Table against notable methods such as Dater and Chain-of-Table on BIRD. As depicted in Table 5, our analysis demonstrates that Tree-of-Table consistently requires the fewest generated samples to reach accurate answers across all evaluated datasets. This starkly contrasts with other approaches, which comparatively require more samples to achieve similar levels of performance.

Table 4: Comparison of Table Format Encoding.

| Table Formatting | WikiTQ |
|---|---|
| HTML | 68.01 |
| TSV | 68.12 |
| PIPE | 69.34 |
| MarkDown | 69.77 |

Table 5: Efficiency Analysis.

| Method | Generate Samples |
|---|---|
| Dater | 300 |
| Chain-of-Table | 120 |
| TREE-OF-TBALE | 90 |

**Effectiveness of Table Condensation.** Finaly, we validate the efficacy of the proposed Table Condensation component in reducing table sizes, making them more amenable to LLMs for reasoning. By condensing tables, we aim to filter out irrelevant information, thereby boosting the signal-to-noise ratio and allowing the model to focus on the most pertinent data. We conducte a comparative analysis of the number of table cells before and after applying Table Condensation across four datasets: WikiTQ, TabFact, FeTaQA, and BIRD. Figure 4 (b) highlights the stark contrast in table sizes before and after the application of Table Condensation. Across all examined datasets, there is a significant reduction in the number of table cells post-condensation. This reduction demonstrates the effectiveness of our method in shrinking table dimensions, ensuring that tables remain within a tractable size range and contain information that is highly relevant to the task at hand.

## 5   CONCLUSION

In this paper, we address the profound challenge of advancing table understanding with LLMs, specifically in the domain of large and complex tabular datasets. Our innovative approach, Tree-of-Table, integrates table condensation and decomposition with a hierarchical reasoning construct that aligns with human cognitive processes to tackle intricate problem-solving tasks. Our extensive experiments conduct across various datasets, including WikiTQ, TableFact, FeTaQA, and BIRD,

demonstrate that Tree-of-Table not only achieves state-of-the-art performance but also presents remarkable improvements in efficiency and generalizability.

**Limitations.** While Tree-of-Table has shown exceptional results in enhancing the efficiency and effectiveness of large language models in processing extensive tabular data, the tree-based reasoning may need to require careful calibration to balance depth and breadth effectively—a task that necessitates fine-tuning and may impose certain limitations on adaptability.

**Broader Impact.** The broader impact of Tree-of-Table is multi-faceted, extending across academic, industrial, and societal domains. Academically, our work contributes a significant leap forward in the intersection of table understanding and natural language processing, providing a reference point for future research and development in this area. In industry, the application of Tree-of-Table can revolutionize the way organizations interact with large datasets. By simplifying the complexity and enhancing the reasoning capabilities of models with tabular data, Tree-of-Table can facilitate more informed decision-making, enhance prediction systems, and optimize data-driven strategies across various sectors such as finance, healthcare, and logistics. From a societal perspective, improving the accessibility and comprehension of large-scale data has the potential to democratize information. By enabling a more nuanced understanding of data presented in tabular form, Tree-of-Table can contribute to greater transparency and empower individuals to make better data-informed decisions.

## REFERENCES

Zilong Wang, Hao Zhang, Chun-Liang Li, Julian Martin Eisenschlos, Vincent Perot, Zifeng Wang, Lesly Miculicich, Yasuhisa Fujii, Jingbo Shang, Chen-Yu Lee, et al. Chain-of-table: Evolving tables in the reasoning chain for table understanding. *arXiv preprint arXiv:2401.04398*, 2024.

Michael J. Cafarella, Alon Halevy, Daisy Zhe Wang, Eugene Wu, and Yang Zhang. Webtables: Exploring the power of tables on the web. *Proc. VLDB Endow.*, 1(1):538–549, aug 2008. ISSN 2150-8097. doi: 10.14778/1453856.1453916.

OpenAI. Gpt-4 technical report. *ArXiv*, abs/2303.08774, 2023.

Wenhu Chen. Large language models are few (1)-shot table reasoners. In *Findings of the Association for Computational Linguistics: EACL 2023*, pp. 1120–1130, 2023.

Jinhao Jiang, Kun Zhou, Zican Dong, Keming Ye, Wayne Xin Zhao, and Ji-Rong Wen. Structgpt: A general framework for large language model to reason over structured data. In *Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing*, pp. 9237–9251, 2023.

Haoran Wang, Licheng Jiao, Fang Liu, Lingling Li, Xu Liu, Deyi Ji, and Weihao Gan. Learning social spatio-temporal relation graph in the wild and a video benchmark. *IEEE Transactions on Neural Networks and Learning Systems*, 34(6):2951–2964, 2021b.

Shima Imani, Liang Du, and Harsh Shrivastava. MathPrompter: Mathematical reasoning using large language models. In *Proceedings of the 61st Annual Meeting of the Association for Computational Linguistics (Volume 5: Industry Track)*, pp. 37–42, Toronto, Canada, July 2023. Association for Computational Linguistics. doi: 10.18653/v1/2023.acl-industry.4.

Rohan Anil, Andrew M Dai, Orhan Firat, Melvin Johnson, Dmitry Lepikhin, Alexandre Passos, Siamak Shakeri, Emanuel Taropa, Paige Bailey, Zhifeng Chen, et al. Palm 2 technical report. *arXiv preprint arXiv:2305.10403*, 2023.

Lanyun Zhu, Tianrun Chen, Deyi Ji, Jieping Ye, and Jun Liu. Llafs: When large language models meet few-shot segmentation. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp. 3065–3075, 2024a.

Karthik Valmeekam, Alberto Olmo, Sarath Sreedharan, and Subbarao Kambhampati. Large language models still can't plan (a benchmark for llms on planning and reasoning about change). In *NeurIPS 2022 Foundation Models for Decision Making Workshop*, 2022.

Md Mahadi Hasan Nahid and Davood Rafiei. Tabsqlify: Enhancing reasoning capabilities of llms through table decomposition. *arXiv preprint arXiv:2404.10150*, 2024.

Si-An Chen, Lesly Miculicich, Julian Martin Eisenschlos, Zifeng Wang, Zilong Wang, Yanfei Chen, Yasuhisa Fujii, Hsuan-Tien Lin, Chen-Yu Lee, and Tomas Pfister. TableRAG: Million-token table understanding with language models. In *The Thirty-eighth Annual Conference on Neural Information Processing Systems*, 2024.

Yuan Sui, Jiaru Zou, Mengyu Zhou, Xinyi He, Lun Du, Shi Han, and Dongmei Zhang. TAP4LLM: Table provider on sampling, augmenting, and packing semi-structured data for large language model reasoning. In *Findings of the Association for Computational Linguistics: EMNLP 2024*, 2024b.

Yuan Sui, Mengyu Zhou, Mingjie Zhou, Shi Han, and Dongmei Zhang. Table meets llm: Can large language models understand structured table data? a benchmark and empirical study. In *Proceedings of the 17th ACM International Conference on Web Search and Data Mining*, pp. 645–654, 2024a.

Lanyun Zhu, Deyi Ji, Tianrun Chen, Peng Xu, Jieping Ye, and Jun Liu. Ibd: Alleviating hallucinations in large vision-language models via image-biased decoding. *arXiv preprint arXiv:2402.18476*, 2024c.

Yunhu Ye, Binyuan Hui, Min Yang, Binhua Li, Fei Huang, and Yongbin Li. Large language models are versatile decomposers: Decompose evidence and questions for table-based reasoning. *arXiv preprint arXiv:2301.13808*, 2023.

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

Ziqi Jin and Wei Lu. Tab-cot: Zero-shot tabular chain of thought. *arXiv preprint arXiv:2305.17812*, 2023.

Wenhu Chen, Hongmin Wang, Jianshu Chen, Yunkai Zhang, Hong Wang, Shiyang Li, Xiyou Zhou, and William Yang Wang. Tabfact: A large-scale dataset for table-based fact verification. In *International Conference on Learning Representations*, 2019.

Jinyang Li, Binyuan Hui, Ge Qu, Jiaxi Yang, Binhua Li, Bowen Li, Bailin Wang, Bowen Qin, Ruiying Geng, Nan Huo, et al. Can llm already serve as a database interface? a big bench for large-scale database grounded text-to-sqls. *Advances in Neural Information Processing Systems*, 36, 2024.

Nitarshan Rajkumar, Raymond Li, and Dzmitry Bahdanau. Evaluating the text-to-sql capabilities of large language models. *arXiv preprint arXiv:2204.00498*, 2022b.

Deyi Ji, Wenwei Jin, Hongtao Lu, and Feng Zhao. Pptformer: Pseudo multi-perspective transformer for uav segmentation. *International Joint Conference on Artificial Intelligence*, pp. 893–901, 2024a.

Nitarshan Rajkumar, Raymond Li, and Dzmitry Bahdanau. Evaluating the text-to-sql capabilities of large language models. *arXiv preprint arXiv:2204.00498*, 2022a.

Tianze Shi, Chen Zhao, Jordan Boyd-Graber, Hal Daumé III, and Lillian Lee. On the potential of lexico-logical alignments for semantic parsing to sql queries. *Findings of the Association for Computational Linguistics: EMNLP 2020*, 2020.

Richard Pönighaus. 'favourite'sql-statements—an empirical analysis of sql-usage in commercial applications. In *International Conference on Information Systems and Management of Data*, pp. 75–91. Springer, 1995.

George Katsogiannis-Meimarakis and Georgia Koutrika. A survey on deep learning approaches for text-to-sql. *The VLDB Journal*, pp. 1–32, 2023.

Jason Wei, Xuezhi Wang, Dale Schuurmans, Maarten Bosma, Fei Xia, Ed Chi, Quoc V Le, Denny Zhou, et al. Chain-of-thought prompting elicits reasoning in large language models. *Advances in Neural Information Processing Systems*, 35:24824–24837, 2022.

Panupong Pasupat and Percy Liang. Compositional semantic parsing on semi-structured tables. In *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.

Deyi Ji, Feng Zhao, Hongtao Lu, Mingyuan Tao, and Jieping Ye. Ultra-high resolution segmentation with ultra-rich context: A novel benchmark. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp. 23621–23630, 2023b.

Linyong Nan, Chiachun Hsieh, Ziming Mao, Xi Victoria Lin, Neha Verma, Rui Zhang, Wojciech Kryściński, Hailey Schoelkopf, Riley Kong, Xiangru Tang, Mutethia Mutuma, Ben Rosand, Isabel Trindade, Renusree Bandaru, Jacob Cunningham, Caiming Xiong, Dragomir Radev, and Dragomir Radev. FeTaQA: Free-form table question answering. *Transactions of the Association for Computational Linguistics*, 10:35–49, 2022. doi: 10.1162/tacl_a_00446.

Yann LeCun, Yoshua Bengio, and Geoffrey Hinton. Deep learning. *nature*, 521(7553):436–444, 2015.

Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image recognition. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pp. 770–778, 2016.

Jacob Devlin. Bert: Pre-training of deep bidirectional transformers for language understanding. *arXiv preprint arXiv:1810.04805*, 2018.

Qianxiong Xu, Xuanyi Liu, Lanyun Zhu, Guosheng Lin, Cheng Long, Ziyue Li, and Rui Zhao. Hybrid mamba for few-shot segmentation. *arXiv preprint arXiv:2409.19613*, 2024.

Alec Radford, Jong Wook Kim, Chris Hallacy, Aditya Ramesh, Gabriel Goh, Sandhini Agarwal, Girish Sastry, Amanda Askell, Pamela Mishkin, Jack Clark, et al. Learning transferable visual models from natural language supervision. In *International conference on machine learning*, pp. 8748–8763. PMLR, 2021.

Hugo Touvron, Thibaut Lavril, Gautier Izacard, Xavier Martinet, Marie-Anne Lachaux, Timothée Lacroix, Baptiste Rozière, Naman Goyal, Eric Hambro, Faisal Azhar, et al. Llama: Open and efficient foundation language models. *arXiv preprint arXiv:2302.13971*, 2023a.

Deyi Ji, Feng Zhao, Lanyun Zhu, Wenwei Jin, Hongtao Lu, and Jieping Ye. Discrete latent perspective learning for segmentation and detection. In *Forty-first International Conference on Machine Learning*, pp. 21719–21730, 2024b.

Peng Wang, An Yang, Rui Men, Junyang Lin, Shuai Bai, Zhikang Li, Jianxin Ma, Chang Zhou, Jingren Zhou, and Hongxia Yang. Ofa: Unifying architectures, tasks, and modalities through a simple sequence-to-sequence learning framework. In *International conference on machine learning*, pp. 23318–23340. PMLR, 2022.

Tomas Mikolov. Efficient estimation of word representations in vector space. *arXiv preprint arXiv:1301.3781*, 3781, 2013.

Alec Radford, Jeffrey Wu, Rewon Child, David Luan, Dario Amodei, Ilya Sutskever, et al. Language models are unsupervised multitask learners. *OpenAI blog*, 1(8):9, 2019.

Jared Kaplan, Sam McCandlish, Tom Henighan, Tom B Brown, Benjamin Chess, Rewon Child, Scott Gray, Alec Radford, Jeffrey Wu, and Dario Amodei. Scaling laws for neural language models. *arXiv preprint arXiv:2001.08361*, 2020.

Jinze Bai, Shuai Bai, Yunfei Chu, Zeyu Cui, Kai Dang, Xiaodong Deng, Yang Fan, Wenbin Ge, Yu Han, Fei Huang, et al. Qwen technical report. *arXiv preprint arXiv:2309.16609*, 2023.

Haoran Wang, Licheng Jiao, Fang Liu, Lingling Li, Xu Liu, Deyi Ji, and Weihao Gan. Ipgn: Interactiveness proposal graph network for human-object interaction detection. *IEEE Transactions on Image Processing*, 30:6583–6593, 2021a.

Jean-Baptiste Alayrac, Jeff Donahue, Pauline Luc, Antoine Miech, Iain Barr, Yana Hasson, Karel Lenc, Arthur Mensch, Katherine Millican, Malcolm Reynolds, et al. Flamingo: a visual language model for few-shot learning. *Advances in neural information processing systems*, 35:23716–23736, 2022.

Alexey Dosovitskiy. An image is worth 16x16 words: Transformers for image recognition at scale. *arXiv preprint arXiv:2010.11929*, 2020.

Lanyun Zhu, Deyi Ji, Shiping Zhu, Weihao Gan, Wei Wu, and Junjie Yan. Learning statistical texture for semantic segmentation. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp. 12537–12546, June 2021.

Haotian Zhang, Pengchuan Zhang, Xiaowei Hu, Yen-Chun Chen, Liunian Li, Xiyang Dai, Lijuan Wang, Lu Yuan, Jenq-Neng Hwang, and Jianfeng Gao. Glipv2: Unifying localization and vision-language understanding. *Advances in Neural Information Processing Systems*, 35:36067–36080, 2022a.

Xinlong Wang, Xiaosong Zhang, Yue Cao, Wen Wang, Chunhua Shen, and Tiejun Huang. Seggpt: Segmenting everything in context. *arXiv preprint arXiv:2304.03284*, 2023.

Deyi Ji, Feng Zhao, and Hongtao Lu. Guided patch-grouping wavelet transformer with spatial congruence for ultra-high resolution segmentation. *International Joint Conference on Artificial Intelligence*, pp. 920–928, 2023a.

Qian Liu, Bei Chen, Jiaqi Guo, Morteza Ziyadi, Zeqi Lin, Weizhu Chen, and Jian-Guang Lou. TAPEX: Table pre-training via learning a neural sql executor. In *International Conference on Learning Representations*, 2021.

Julian Eisenschlos, Syrine Krichene, and Thomas Müller. Understanding tables with intermediate pre-training. In *Findings of the Association for Computational Linguistics: EMNLP 2020*, pp. 281–296, Online, November 2020. Association for Computational Linguistics. doi: 10.18653/v1/2020.findings-emnlp.27.

Zhengbao Jiang, Yi Mao, Pengcheng He, Graham Neubig, and Weizhu Chen. OmniTab: Pretraining with natural and synthetic data for few-shot table-based question answering. In *Proceedings of the 2022 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*, pp. 932–942, Seattle, United States, July 2022. Association for Computational Linguistics. doi: 10.18653/v1/2022.naacl-main.68.

Wenhu Chen, Xueguang Ma, Xinyi Wang, and William W Cohen. Program of thoughts prompting: Disentangling computation from reasoning for numerical reasoning tasks. *arXiv preprint arXiv:2211.12588*, 2022.

Deyi Ji, Haoran Wang, Hanzhe Hu, Weihao Gan, Wei Wu, and Junjie Yan. Context-aware graph convolution network for target re-identification. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 35, pp. 1646–1654, 2021.

Tushar Khot, Harsh Trivedi, Matthew Finlayson, Yao Fu, Kyle Richardson, Peter Clark, and Ashish Sabharwal. Decomposed prompting: A modular approach for solving complex tasks. In *International Conference on Learning Representations*, 2022.

Zhuosheng Zhang, Aston Zhang, Mu Li, and Alex Smola. Automatic chain of thought prompting in large language models. In *International Conference on Learning Representations*, 2022b.

Lanyun Zhu, Tianrun Chen, Jianxiong Yin, Simon See, and Jun Liu. Learning gabor texture features for fine-grained recognition. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pp. 1621–1631, 2023b.

Hanzhe Hu, Deyi Ji, Weihao Gan, Shuai Bai, Wei Wu, and Junjie Yan. Class-wise dynamic graph convolution for semantic segmentation. In *European Conference on Computer Vision*, pp. 1–17. Springer, 2020.

Shunyu Yao, Dian Yu, Jeffrey Zhao, Izhak Shafran, Thomas L Griffiths, Yuan Cao, and Karthik Narasimhan. Tree of thoughts: Deliberate problem solving with large language models. *arXiv preprint arXiv:2305.10601*, 2023.

Cheng-Yu Hsieh, Chun-Liang Li, Chih-kuan Yeh, Hootan Nakhost, Yasuhisa Fujii, Alex Ratner, Ranjay Krishna, Chen-Yu Lee, and Tomas Pfister. Distilling step-by-step! outperforming larger language models with less training data and smaller model sizes. In *Findings of the Association for Computational Linguistics: ACL 2023*. Association for Computational Linguistics, 2023.

Bhuwan Dhingra, Manaal Faruqui, Ankur Parikh, Ming-Wei Chang, Dipanjan Das, and William Cohen. Handling divergent reference texts when evaluating table-to-text generation. In *Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics*, pp. 4884–4895, 2019.

Nelson F Liu, Kevin Lin, John Hewitt, Ashwin Paranjape, Michele Bevilacqua, Fabio Petroni, and Percy Liang. Lost in the middle: How language models use long contexts. *arXiv preprint arXiv:2307.03172*, 2023.

Lanyun Zhu, Tianrun Chen, Jianxiong Yin, Simon See, and Jun Liu. Addressing background context bias in few-shot segmentation through iterative modulation. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp. 3370–3379, 2024b.

Wenqiang Lei, Weixin Wang, Zhixin Ma, Tian Gan, Wei Lu, Min-Yen Kan, and Tat-Seng Chua. Re-examining the role of schema linking in text-to-sql. In *Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing (EMNLP)*, pp. 6943–6954, 2020.

Bytescout. https://bytescout.com/blog/20-important-sql-queries.html. In *Bytescout*, 2024.

Jon Louis Bentley. Multidimensional divide-and-conquer. *Communications of the ACM*, 23(4): 214–229, 1980.

Deyi Ji, Haoran Wang, Mingyuan Tao, Jianqiang Huang, Xian-Sheng Hua, and Hongtao Lu. Structural and statistical texture knowledge distillation for semantic segmentation. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp. 16876–16885, 2022.

Robert Tarjan. Depth-first search and linear graph algorithms. *SIAM journal on computing*, 1(2): 146–160, 1972.

Lanyun Zhu, Tianrun Chen, Jianxiong Yin, Simon See, and Jun Liu. Continual semantic segmentation with automatic memory sample selection. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp. 3082–3092, 2023a.

Chin-Yew Lin. ROUGE: A package for automatic evaluation of summaries. In *Text Summarization Branches Out*, pp. 74–81, Barcelona, Spain, July 2004. Association for Computational Linguistics.

Hugo Touvron, Louis Martin, Kevin Stone, Peter Albert, Amjad Almahairi, Yasmine Babaei, Nikolay Bashlykov, Soumya Batra, Prajjwal Bhargava, Shruti Bhosale, et al. Llama 2: Open foundation and fine-tuned chat models. *arXiv preprint arXiv:2307.09288*, 2023b.

Mark Chen, Jerry Tworek, Heewoo Jun, Qiming Yuan, Henrique Ponde de Oliveira Pinto, Jared Kaplan, Harri Edwards, Yuri Burda, Nicholas Joseph, Greg Brockman, et al. Evaluating large language models trained on code. *arXiv preprint arXiv:2107.03374*, 2021.