# Table Question Answering in the Era of Large Language Models:
# A Comprehensive Survey of Tasks, Methods, and Evaluation

**Wei Zhou**[1,3]    **Bolei Ma**[2]    **Annemarie Friedrich**[3]    **Mohsen Mesgar**[1]

[1]Bosch Center for Artificial Intelligence, Renningen, Germany
[2]LMU Munich & Munich Center for Machine Learning, Germany
[3]University of Augsburg, Germany

{wei.zhou3|mohsen.mesgar}@de.bosch.com
bolei.ma@lmu.de    annemarie.friedrich@uni-a.de

## Abstract

Table Question Answering (TQA) aims to answer natural language questions about tabular data, often accompanied by additional contexts such as text passages. The task spans diverse settings, varying in table representation, question/answer complexity, modality involved, and domain. While recent advances in large language models (LLMs) have led to substantial progress in TQA, the field still lacks a systematic organization and understanding of task formulations, core challenges, and methodological trends, particularly in light of emerging research directions such as reinforcement learning. This survey addresses this gap by providing a comprehensive and structured overview of TQA research with a focus on LLM-based methods. We provide a comprehensive categorization of existing benchmarks and task setups. We group current modeling strategies according to the challenges they target, and analyze their strengths and limitations. Furthermore, we highlight underexplored but timely topics that have not been systematically covered in prior research. By unifying disparate research threads and identifying open problems, our survey offers a consolidated foundation for the TQA community, enabling a deeper understanding of the state of the art and guiding future developments in this rapidly evolving area.

## 1 Introduction

Tables are a ubiquitous data format in daily life (Cafarella et al., 2008). Automatically processing and understanding tabular data with (multi-modal) large language models ((M)LLMs) has recently attracted considerable attention from both industry (Katsis et al., 2022; Su et al., 2024) and academia (Pasupat and Liang, 2015; Wolff and Hulsebos, 2025), emerging as a prominent research direction.

Among the various tasks involving tables, including table generation (Gulati and Roysdon, 2023) and table-to-text (Parikh et al., 2020), table question answering (TQA) stands out as one of the most
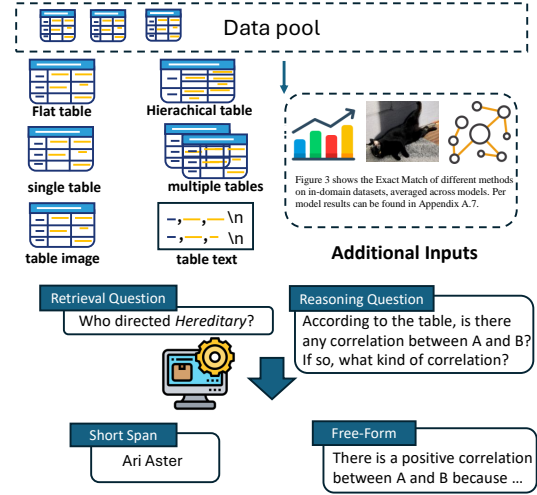


Figure 1: Different table question answering task setups. *Domain:* Either the inputs need to be retrieved from a data pool or directly given. *Table Format*: Tables can exist or be presented in various formats. *Additional Context*: Charts, images, and knowledge graphs can also be involved as inputs. *Question Complexity*: A question can involve retrieving certain cells from a table or require reasoning and analysis to be solved. *Answer Format*: Answers can be in short text spans, consisting only of numbers and entities, or in free-form natural language, with no limitation on types and length.

widely studied (Wu et al., 2025c). The goal of TQA is to answer questions based on tabular data, optionally augmented with additional context such as text passages or images. TQA can be instantiated in diverse settings. As illustrated in Figure 1, the table required to answer a question may be provided directly with the question, or it may be required to first retrieve it from a large corpus. Tables can also vary in format, size, and structural complexity. Moreover, questions may target specific cells or require multi-step/type reasoning over the table content. These variations stem from real-world applications and necessitate different modeling strategies to address the underlying challenges.

With growing interest in TQA, the field has wit-

nessed rapid development, evidenced by the increasing number of works shown in Appendix A.1. This survey not only consolidates key resources and modeling approaches but also distills insights into promising directions for future research.

**Comparing with Existing Surveys.** Most prior surveys on TQA or tabular reasoning focus solely on textual tables (Dong et al., 2022; Jin et al., 2022; Zhang et al., 2024e; Fang et al., 2024b; Ren et al., 2025). Among those covering both image and textual tables, Wu et al. (2025d) concentrate on table representations and related tasks, without discussing modeling approaches, while Tian et al. (2025a) emphasize agentic setups and overlooks fine-tuning methods. Crucially, existing surveys do not provide any overview of TQA task setups, nor do they cover recent advances and emerging themes in the LLM era, such as reinforcement learning, interpretability, and novel evaluation paradigms. Appendix A.2 presents a detailed comparison. To our knowledge, this is the first survey dedicated to TQA in the LLM era, offering timely coverage of contemporary challenges and opportunities.

**Scope.** We include work on both TQA and table fact verification (TFV), as TFV can be reformulated into TQA settings (Lu et al., 2023). We also consider datasets from Text-to-SQL research, since they can serve as TQA benchmarks. As our survey focuses on (M)LLM-based TQA, and the most recent TQA survey (Jin et al., 2022) was published in 2022, we primarily collect modeling papers published since 2022. In total, we review 215 papers, with details of our collection methodology and statistics provided in Appendix A.1.

**Structure.** Section 2 outlines TQA task setups and benchmarks. Section 3 presents modeling approaches grouped by challenges. Section 4 reviews evaluation methodologies and Section 5 discusses emerging topics and future directions.

## 2 Task Setups and Resources

We dissect TQA from five perspectives. Table 2 provides existing TQA datasets categorized by the characteristics of their task setups.

**Table Representation and Format.** Tables exist in both textual and image formats. This representational difference can distinguish different task setups: modeling over textual tables (Zhang et al., 2024g,b; Wang et al., 2024e; Su et al., 2024; He et al., 2025; Zhou et al., 2025f) and table images (Zheng et al., 2024; Zhou et al., 2025a; Jiang et al.,

2025; Yang et al., 2025a; Zhao et al., 2024b). Apart from table representations, table structures (hierarchical vs. flat) and numbers (single vs. multiple) also determine task features. Processing hierarchical tables (Cao et al., 2023; Zhang et al., 2024h; Li et al., 2025c) brings greater challenges of structure understanding than processing flat tables (Liu et al., 2024a; Nahid and Rafiei, 2024b; Zhang et al., 2025a). Similarly, compared to modeling over a single table (Gu et al., 2025; Jin et al., 2025b; Chegini et al., 2025; Yu et al., 2025a), modeling over multiple tables (Zhao et al., 2022; Pal et al., 2023; Zou et al., 2025; Qiu et al., 2024) requires an understanding of inter-table relationships as well as capabilities of processing longer inputs.

**Question Complexity.** A TQA question can require different capabilities to be solved. Zhou et al. (2024a) distinguish retrieval and reasoning questions. The former refers to questions that can be addressed simply by locating relevant cells, while the latter requires additional reasoning to be solved. Based on the level of question complexity, TQA can be categorized into simple and complex setups, with the simple setup involving only retrieval questions (Katsis et al., 2022; Liu et al., 2023a; Wang et al., 2025a) while the complex setup involves numerical (Chen et al., 2021; Zhu et al., 2021; Lu et al., 2022; Tian et al., 2025a), commonsense (Zhang et al., 2023), temporal (Gupta et al., 2023; Shankarampeta et al., 2025) reasoning, optionally with capability of data analysis and plotting (Wu et al., 2024a; He et al., 2024).

**Answer Formats.** Most TQA setups feature short span answers, where an answer is composed of a few tokens or numbers (Pasupat and Liang, 2015; Iyyer et al., 2017; Chen et al., 2019; Cheng et al., 2022a; Wu et al., 2025a). This setup enables easy evaluation, given that one can determine the answer correctness simply by checking if they match the reference answers. Nevertheless, real-world questions might require verbatim answers. For instance, a question asking for data analysis needs several sentences for explanations. This setup with free-form answer format receives more and more attention, given its alignment to real-world queries (Nan et al., 2022; Su et al., 2024; Wu et al., 2024a; Li et al., 2025d; Wu et al., 2025b).

**Modality.** Apart from the standard task setup involving a table and a question as inputs (Zhang et al., 2024b; Liu et al., 2024a; Wang et al., 2024e;

Taxonomy of Task Setups

- Table Format and Representation
  - Representation
    - Text — Ye et al. (2023); Zhang et al. (2024g,b); Wang et al. (2024e); Su et al. (2024); He et al. (2025); Zhou et al. (2025f)
    - Image — Zheng et al. (2024); Zhou et al. (2025a); Jiang et al. (2025); Yang et al. (2025a); Zhao et al. (2024b)
  - Structure
    - Flat — Liu et al. (2024a); Nahid and Rafiei (2024b); Zhang et al. (2025a)
    - Hierarchical — Cao et al. (2023); Zhang et al. (2024h); Li et al. (2025c)
  - Number
    - Multiple — Zhao et al. (2022); Pal et al. (2023); Zou et al. (2025); Qiu et al. (2024)
    - Single — Gu et al. (2025); Jin et al. (2025b); Chegini et al. (2025); Yu et al. (2025a)
- Question Complexity
  - Retrieve — Katsis et al. (2022); Liu et al. (2023a); Wang et al. (2025a)
  - Reasoning — Chen et al. (2021); Zhu et al. (2021); Tian et al. (2025a); Wu et al. (2024a); He et al. (2024)
- Answer Format
  - Short-Span — Pasupat and Liang (2015); Iyyer et al. (2017); Chen et al. (2019); Cheng et al. (2022a)
  - Free-Form — Nan et al. (2022); Su et al. (2024); Wu et al. (2024a); Li et al. (2025d); Wu et al. (2025b)
- Modality
  - Table Only — Zhang et al. (2024b); Liu et al. (2024a); Wang et al. (2024e); Zhou et al. (2025c)
  - Table-Text — Chen et al. (2020b,a); Zhu et al. (2021); Chen et al. (2021); Zhao et al. (2022)
  - Table-Others — Li et al. (2022); Mathur et al. (2024); Titiya et al. (2025); Pramanick et al. (2024); Talmor et al. (2021); Zhao et al. (2024a); Foroutan et al. (2025); Pramanick et al. (2024)
- Domain
  - Close — Pasupat and Liang (2015); Cheng et al. (2022a); Zhou et al. (2024a); Zhang et al. (2025e)
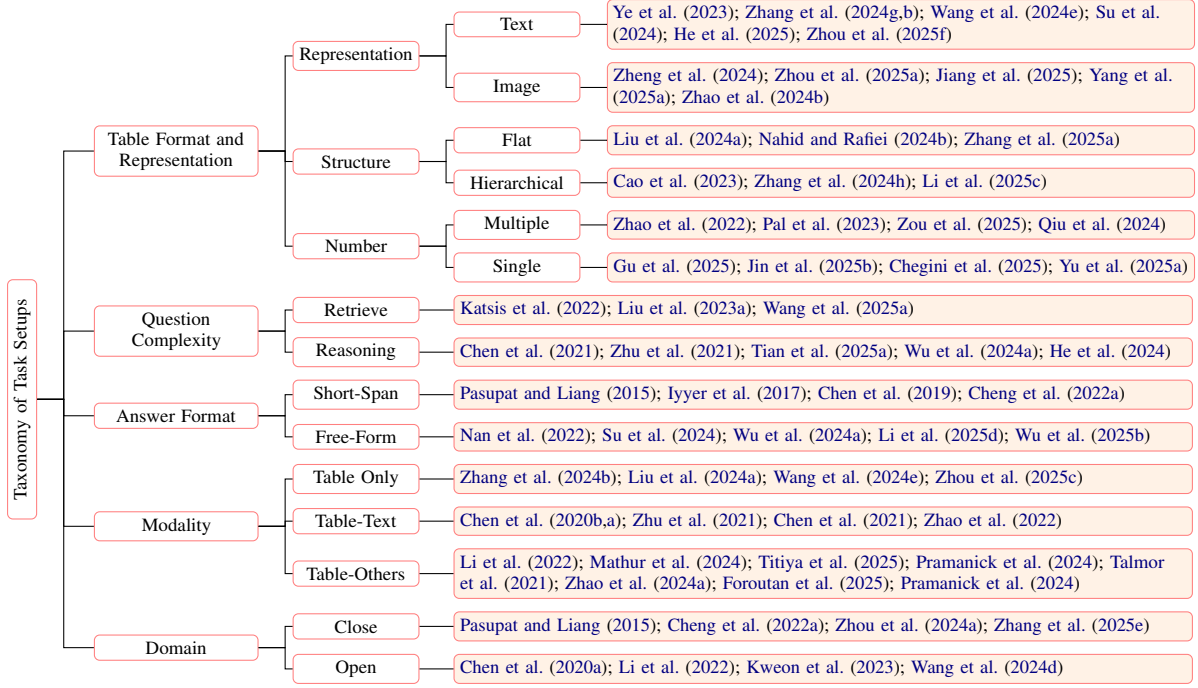  - Open — Chen et al. (2020a); Li et al. (2022); Kweon et al. (2023); Wang et al. (2024d)

Figure 2: A taxonomy of TQA task setups. We list representative papers for each setup.

Zhou et al., 2025c), previous work proposes more complex setups involving additional inputs such as passages (Chen et al., 2020b,a; Zhu et al., 2021; Chen et al., 2021; Zhao et al., 2022), images (Li et al., 2022; Mathur et al., 2024; Titiya et al., 2025; Pramanick et al., 2024; Talmor et al., 2021), charts (Zhao et al., 2024a; Foroutan et al., 2025; Pramanick et al., 2024) and knowledge graphs (Christmann et al., 2023; Hu et al., 2024; Huang et al., 2025a). These complex setups involving inputs from multiple modalities align better with real-world settings, where heterogeneous data often accompanies tables. Among all combinations of input types, tables and passages are the most commonly studied. This task setup is referred to as table-text QA.

**Domains.** TQA task setups can be categorized into open-domain TQA (Chen et al., 2020a; Li et al., 2022; Kweon et al., 2023; Strich et al., 2025b) and closed-domain TQA (Pasupat and Liang, 2015; Cheng et al., 2022a; Zhou et al., 2024a; Zhang et al., 2025e), depending on whether relevant inputs for solving a problem are given or not. In open-domain TQA, an input database (usually a table database) is given. A system needs to first retrieve relevant inputs from a pool of candidates and then carry out reasoning over the target inputs to obtain final answers. Compared to close-domain TQA, where target inputs are given, this setup poses additional challenges in locating relevant inputs.

## 3 Modeling

We categorize modeling methods based on the challenges they try to address. In addition, we also analyze their strengths and limitations.

### 3.1 Table Understanding

**Visual Table Modeling.** Visual table understanding involves comprehending both the tables' contents and their structures. A common approach is to *parse table images into texts* with MLLMs (Nguyen et al., 2023; Hormazábal-Lagos et al., 2025). However, Xia et al. (2024) found that although MLLMs demonstrate promising OCR performance on tabular data, they still exhibit limited spatial and formatting recognition capabilities, especially when dealing with large tables (Zheng et al., 2024). Another line of research focuses on *pre-training and fine-tuning MLLMs* (Zhao et al., 2024b; Zheng et al., 2024; Zhou et al., 2025b; Jiang et al., 2025; Yang et al., 2025a; Zhang et al., 2025b). Training datasets of table images have been constructed using existing tabular datasets (Zhao et al., 2024b), by converting textual tables into image representations (Zheng et al., 2024; Zhou et al., 2025b; Jiang et al., 2025), and from scratch (Yang et al., 2025a; Zhang et al., 2025b).

In terms of model design, dual vision encoders are employed to capture information at different granularity levels (Zhao et al., 2024b; Zhou et al.,
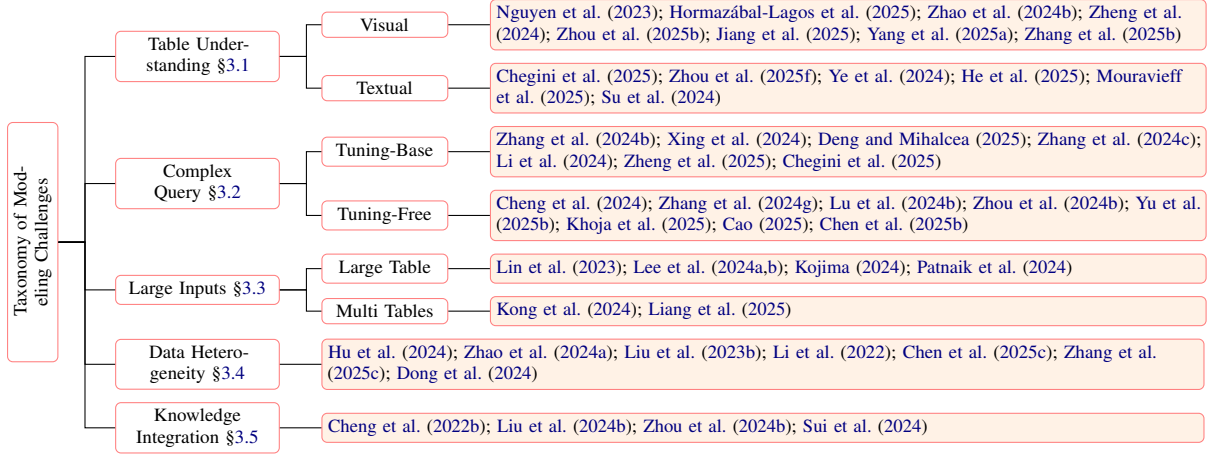
Figure 3: A taxonomy of methods categorized by challenges. We list representative papers for each challenge.

2025b). Zhang et al. (2025b) train a mixture of experts model to capture different types of encoded information, including layout and semantics.

**Textual Table Modeling.** Prior work has explored three main directions for table structure understanding: table representations, architecture modifications, and specialized training tasks. For *representation*, tables have been modeled as relational databases or Pandas DataFrames, where operations are expressed in code (Chegini et al., 2025; Zhou et al., 2025f), or as spreadsheets with formula-based operations (Cao et al., 2025; Wang et al., 2025b). Compared to database and DataFrame formats, which require a pre-defined schema, spreadsheet representations offer greater flexibility. Another approach models tables as hypergraphs (Huang et al., 2025b; Li et al., 2025b; Jin et al., 2025a), where nodes represent cells and edges encode positional relations, enabling explicit structural learning. Tables have also been expressed as natural language tuples (Zhao et al., 2023a; Yang et al., 2025c).

Some methods *encode structure within the model* (He et al., 2025; Mouravieff et al., 2025; Su et al., 2024). For instance, He et al. (2025) serialize tables with special tokens and applies 2D LoRA (Hu et al., 2021) to capture low-rank positional information, while Mouravieff et al. (2025) design a sparse attention mask for tabular data.

Finally, *task-specific objectives* have been proposed to enhance structural reasoning, such as layout transformation inference (Jin et al., 2025b), where models detect changes between original and altered layouts, and cell position generation (Cho et al., 2025), where models predict cell locations.

**Discussion.** Visual table understanding is generally more challenging than textual table understanding (Deng et al., 2024b; Zheng et al., 2024), possibly because it requires additional content interpretation. Notably, OCR-based pipelines do not outperform models directly fine-tuned on table images (Zheng et al., 2024). For details, see Appendix A.3. For textual table understanding, the need to define a fixed table schema may be a drawback for database tables. When using more fine-grained textual representations such as JSON or Markdown, there seems to be no optimal format across datasets and models (Zhang et al., 2024d).

### 3.2 Complex Query

**Tuning-Based.** Methods categorized in this group rely on fine-tuning (M)LLMs to improve reasoning capabilities over tabular data. Diverse training datasets that cover a wide range of reasoning types are crucial for handling complex queries. Such datasets are either collected from existing benchmarks (Zhang et al., 2024b; Xing et al., 2024; Deng and Mihalcea, 2025) or synthesized (Zhang et al., 2024c; Li et al., 2024; Zheng et al., 2025; Chegini et al., 2025). For instance, Zheng et al. (2025) propose selecting training samples based on identified weaknesses and progressively fine-tuning the model. Chegini et al. (2025) collect Python programs generated by large closed-source models and use them to train a smaller open-source LLM. In addition, combining fine-tuned models with tools during inference time has been explored (Wu and Feng, 2024; Mouravieff et al., 2024; Vinayagame et al., 2025). Performance can be further improved through reinforcement learning (Nahid and Rafiei, 2024b; Zhou et al., 2025c; Stoisser et al., 2025).

**Tuning-Free.** In tuning-free approaches, to enable accurate numerical reasoning and reduce hallucinations, agentic workflows that integrate tool usage are often adopted (Cheng et al., 2024; Zhang et al., 2024g; Lu et al., 2024b; Zhou et al., 2024b; Yu et al., 2025b; Khoja et al., 2025; Cao, 2025; Chen et al., 2025b). Typically, models generate and execute Python (Cao et al., 2023; Zhang et al., 2024h; Yu et al., 2025b) or SQL (Abhyankar et al., 2024; Nahid and Rafiei, 2024b; Khoja et al., 2025) code to obtain reasoning results. To improve code generation accuracy, error feedback can be provided to the LLM as a revision signal (Cheng et al., 2024; López Gude et al., 2025; Site et al., 2025). Verification modules have also been introduced to check the correctness of intermediate reasoning (Wang et al., 2024c; Yu et al., 2025a). Most tuning-free methods employ multi-step reasoning to decompose complex questions into simpler sub-problems (Mao et al., 2025; Zhou et al., 2024b; Ji et al., 2024; Zhang et al., 2024a; Deng et al., 2024a; Zhao et al., 2024c; Nguyen et al., 2025). Some also incorporate memory mechanisms to store and reuse past reasoning experiences (Bai et al., 2025; Gu et al., 2025). In terms of prompting strategies, agentic-flow systems often adopt ReAct-style prompting (Zhou et al., 2024b; Yu et al., 2025b; Bai et al., 2025). Zhang et al. (2025g) shows that prompting models to iterate over rows can reduce hallucination. Dixit et al. (2025) find that no single prompting technique consistently outperforms others in temporal table reasoning.

**Discussion.** Tuning-free methods require no training data, but incur longer inference times and higher token costs (Zhou et al., 2025c). In contrast, tuning-based approaches are more efficient at inference but are prone to out-of-domain performance degradation (Deng and Mihalcea, 2025). Both methods demonstrate state-of-the-art performance (Yang et al., 2025d; Abhyankar et al., 2025; Cao, 2025), but involve trade-offs. A promising intermediate strategy might be to fine-tune models for general reasoning capabilities while delegating specific table operations, such as retrieval, to external tools (Wu and Feng, 2024; Zhu et al., 2024).

### 3.3 Large Inputs

The main challenge with large inputs is efficiently identifying relevant information, as processing full tables is often infeasible or ineffective. We review methods for handling large and multiple tables.

**Large Tables.** A common approach is to fine-tune retrievers to identify the most relevant cells for a given question (Lin et al., 2023; Lee et al., 2024a,b; Kojima, 2024; Patnaik et al., 2024). Another strategy leverages LLMs to directly select pertinent table content (Ye et al., 2023; Jiang et al., 2023, 2024; Sui et al., 2024). Some methods define atomic operations for table manipulation (Wang et al., 2024a,e), while others embed both queries and tables for semantic matching (Sui et al., 2024; Yu et al., 2025b). A further line of work employs code generation to execute table-filtering operations (Gemmell and Dalton, 2023; Zhou et al., 2025d; Vyatkin and Oliseenko, 2025).

**Multiple Tables.** LLMs can assist retrieval by enhancing table semantics. For instance, Liang et al. (2025) augment table snippets with LLM-generated questions to produce richer table representations. LLMs can also directly facilitate retrieval. Kong et al. (2024) generate SQL queries to identify relevant tables. Rather than treating each table as an independent document, Zou et al. (2025) represent the table corpus as a hypergraph and select the most relevant subgraph using a multi-stage coarse-to-fine process. Many works adopt dense passage retrieval (Karpukhin et al., 2020) or language model embeddings to encode tables, passages, and questions (Guan et al., 2024; Bardhan et al., 2024).

**Discussion.** Directly using LLMs to retrieve relevant cells or tables can lead to information loss (Zhou et al., 2025d). Fine-tuned sub-table retrievers have shown effectiveness, but their generalizability to diverse table formats remains limited. For both scenarios, retrieval-augmented generation (RAG) offers a viable alternative: tables or cells are embedded into a vector database, and questions are issued as queries (Chen et al., 2024).

### 3.4 Data Heterogeneity

To handle different modalities, existing methods typically follow two directions: (1) employing specialized retrievers and reasoners (Li et al., 2022; Zhao et al., 2024a; Hu et al., 2024; Liu et al., 2023b), or (2) designing unified representations (Dong et al., 2024; Chen et al., 2025c; Zhang et al., 2025c). In the first category, Hu et al. (2024) use a multi-stage knowledge-graph retriever, Zhao et al. (2024a) employ multi-agent retrieval from charts and tables, and Liu et al. (2023b) generate image captions to capture salient visual content; atomic retrieval functions can also target both passages

and tables (Shi et al., 2024; Zhou et al., 2024b). In the second category, unified structures integrate heterogeneous sources: Chen et al. (2025c) use DataFrames to jointly represent tabular and textual data, Zhang et al. (2025c) propose Condition Graphs combining tables and knowledge graphs, and Agarwal et al. (2025) build hybrid graphs from linked entities; tabular content may also be summarized into text for downstream tasks (Bardhan et al., 2024). For reasoning over both tables and text, LLMs can align references across modalities (Luo et al., 2023; Zhang et al., 2025e) or be fine-tuned for domain-specific reasoning, as in Zhu et al. (2024)'s financial QA system, which processes both modalities to produce multi-step reasoning chains combining evidence extraction, logical or equation formulation, and execution. In summary, which method to choose depends on the modalities involved. Constructing graphs is straightforward for tables and text, whereas employing separate retrievers is preferable when modalities are harder to unify, e.g., if information from charts and tables needs to be combined.

## 3.5 Knowledge Integration

External knowledge is often required to answer TQA problems. For example, a question in DataBench (Osés Grijalba et al., 2024) asks: *"What is the total number of rebounds recorded in the dataset where the ball didn't change possession?"* Answering this question requires knowing that *OREB* in the table header denotes *offensive rebounds*, a case where ball possession does not change. To handle such cases, prior work has integrated external resources such as Wikipedia into the reasoning process (Zhou et al., 2024b; Sui et al., 2024). For instance, Zhou et al. (2024b) design an atomic function *Search (arg)*, which returns the first few lines from the Wikipedia page of a specified argument. The retrieved content is then stored in the system's memory for subsequent reference. An alternative strategy is to elicit factual knowledge directly from LLMs (Cheng et al., 2022b; Liu et al., 2024b). However, this approach is susceptible to hallucinations, as LLMs may generate factually incorrect information.

## 4 Evaluation

In this section, we discuss evaluation in terms of task performance, system robustness, and model-generated reasoning as explanations.

## 4.1 Task Performance

Current TQA evaluation primarily focuses on performance, measured by automatic metrics such as Exact Match (EM) and ROUGE. EM suits short, span-based answers, whereas ROUGE, BLEU, or F1 scores are better for free-form responses. While efficient, these metrics often miss subtle mismatches between predicted and gold answers (Wolff and Hulsebos, 2025), e.g., EM may wrongly mark answers as incorrect due to formatting differences (Jan 1 vs. 01-01). To address such issues, outputs are normalized to a canonical form before applying EM (Khoja et al., 2025). This "relaxed EM" improves robustness but can cause inconsistencies, as systems may adopt different normalization rules, leading to misleading cross-system comparisons (Hormazábal-Lagos et al., 2025).

Beyond traditional metrics, some studies employ LLMs as judges (Wu and Feng, 2024; Zhou et al., 2024b; Jiang et al., 2025; Zhang et al., 2025d) or use human evaluation (Zhao et al., 2024c; Ye et al., 2024; Khoja et al., 2025). Dixit et al. (2025) propose the Hybrid Correctness Score, combining $F_1$ with LLM judgments. Wolff and Hulsebos (2025) show that, when calibrated with human annotations, LLM-as-a-judge can offer a reliable evaluation signal for tasks requiring reasoning over tables.

## 4.2 Robustness Evaluation

Robustness to structural or content variations in tables and questions is a key property of TQA systems. Zhao et al. (2023b) present a benchmark for adversarial attacks on table structure/content and question perturbations, showing that state-of-the-art models still struggle. Zhou et al. (2024a) define three robustness dimensions: (1) resilience to table structure changes, (2) resistance to shortcut exploitation, and (3) robustness in numerical reasoning. Their benchmark indicates that pipeline models handle value and positional changes best, whereas LLM-based models are more vulnerable to table shuffling, a trend also observed in other works (Ashury-Tahan et al., 2025; Liu et al., 2024a; Yang et al., 2022). Wolff and Hulsebos (2025) further evaluate LLM robustness on real-world tables with missing or duplicated values, underscoring the need for robustness-oriented evaluation.

## 4.3 Evaluating Explanations and Reasoning

An underexplored aspect of TQA evaluation is assessing explanations and reasoning processes.

Model-generated chains of thought (Wei et al., 2022) are often treated as natural language explanations (Zhao et al., 2024c; Zhou et al., 2025f; Lu et al., 2025), either elicited directly (Zhao et al., 2024c; Zhou et al., 2025f) or derived from executable program outputs (Lu et al., 2025). Nguyen et al. (2024) represent explanations as chains of attribution maps, showing intermediate relevant tables alongside reasoning steps, and propose three evaluation tasks: (1) preference ranking, where judges rank explanation quality; (2) forward simulation, where judges answer using only the explanation; and (3) verification, where judges assess prediction correctness based on the explanation. Both human annotators and LLMs serve as judges. Zhou et al. (2025c) take a different approach, estimating the probability of reaching the correct answer from a reasoning step to quantify each step's contribution to the final outcome.

**Discussion.** Popular datasets like WTQ (Pasupat and Liang, 2015) and TabFact (Chen et al., 2019) may suffer from data contamination (Zhou et al., 2025e), leading to overly optimistic performance estimates. More reliable evaluation requires incorporating additional, uncontaminated datasets. Beyond task performance, dimensions such as robustness and reasoning correctness should be systematically assessed to ensure the development of trustworthy TQA systems.

## 5 Discussion and Future Directions

We discuss emerging topics for future exploration, including table representation, multilinguality, reinforcement learning, multi-modal modeling, interpretability and human-centric setups.

**Table Representation.** Tables can appear in various formats, including structured databases, text, and images. When tables are stored in databases, they can be directly queried using SQL. However, tables in textual or image form are often noisier due to inconsistent formatting (Zhou et al., 2024a), mixed data types (Nahid and Rafiei, 2024a), missing values (Wolff and Hulsebos, 2025), and implicit or incomplete schemas (Zheng et al., 2023). These pose challenges to create models for real-world noisy tables. A growing line of work investigates leveraging both textual and visual representations of tables (Deng et al., 2024b; Zhou et al., 2025e; Liu et al., 2025), capitalizing on the fact that one modality can often be converted to the other (e.g.,

image-to-text via OCR, or text-to-image via HTML rendering). However, most existing approaches adopt ensemble strategies that select the optimal representation based on specific problem features, e.g., by table size (Deng et al., 2024b; Zhou et al., 2025e). This approach lacks flexibility when new features emerge or when interactions among multiple features need to be considered. An alternative is to process textual and visual inputs through dedicated encoders and integrate them within the model, an approach widely used in vision-language tasks (Zhao et al., 2024b; Zhou et al., 2025b). Separate encoders might capture complementary signals, such as layout structure from images and semantic content from text and potentially leading to more robust and generalizable methods.

**Multilinguality and Low-Resource Settings.** Tables in real-world applications can be text-heavy and may contain content in one or more languages, as in user or product information tables. Nevertheless, most existing TQA datasets and studies focus on English (Pasupat and Liang, 2015; Zhang et al., 2023, 2024b) or other high-resource languages with large speaker populations, such as Chinese (Zheng et al., 2023; Liu et al., 2023a; Zhao et al., 2024a), leaving low-resource and multilingual scenarios underexplored. Recent efforts have introduced datasets for these settings (Minhas et al., 2022; Zhang et al., 2025f; Shu et al., 2025). In terms of modeling, directly applying LLMs yields uneven results. Shu et al. (2025) report the highest performance on Indo-European languages and the lowest on Niger-Congo languages, likely due to disparities in data representation during pretraining. They also find that multilingual fine-tuning does not consistently improve performance for TQA. Translation-based approaches, which convert target language tables to English, also fall short, as their effectiveness depends on translation quality. Core challenges stem from distribution shifts and lexical diversity in multilingual data (Zhang et al., 2025f). Progress will require TQA systems that natively handle low-resource and multilingual inputs rather than relying solely on translation pipelines.

**Reinforcement Learning in TQA.** Reinforcement learning with verifiable rewards (RLVR) (Su et al., 2025) has gained increasing attention due to its success in developing reasoning-oriented models, such as DeepSeek R1 (DeepSeek-AI et al., 2025). Recent studies have also explored RLVR in

TQA (Jin et al., 2025b; Yang et al., 2025d; Jiang et al., 2025; Lei et al., 2025; Cao et al., 2025; Stoisser et al., 2025; Liu et al., 2025). Commonly used reward signals include answer correctness (Yang et al., 2025d; Jiang et al., 2025; Stoisser et al., 2025; Liu et al., 2025), program executability (Jin et al., 2025b; Cao et al., 2025), output formatting (Yang et al., 2025d; Jiang et al., 2025), positional alignment (Lei et al., 2025), and length constraints (Jin et al., 2025b). In contrast to most work that updates models solely based on final outcome rewards, Zhou et al. (2025c) propose a process supervision framework, demonstrating that models trained with intermediate (process) rewards outperform those trained with only final rewards. RLVR has been shown to improve LLMs' reasoning capabilities over tabular data. In particular, RLVR-trained models exhibit better generalizability (Yang et al., 2025d; Cao et al., 2025) and increased robustness to row and column perturbations (Lei et al., 2025) compared to models trained via supervised fine-tuning (SFT). Nevertheless, initializing models with SFT remains crucial for achieving strong performance (Cao et al., 2025).

**Diverse and Multi-Modal Data Modeling.** Many existing TQA studies focus on table-only settings with relatively simple queries (Ye et al., 2023; Ni et al., 2023; Nahid and Rafiei, 2024b; Wang et al., 2024e; Liu et al., 2024a; Zhou et al., 2025d). This setup is suitable for evaluating LLMs' ability to understand table structures. However, it falls short in capturing more complex scenarios that involve multiple modalities and open-domain setups. An increasing body of work has recognized these limitations and proposed more challenging benchmarks (He et al., 2024; Qiu et al., 2024; Wu et al., 2025a; Osés Grijalba et al., 2024; Wu et al., 2024a, 2025b; Zhu et al., 2025). We argue that future research should extend evaluation to these complex datasets.

**Interpretability and Faithfulness.** Instead of directly producing the final answer to a TQA problem (Herzig et al., 2020; Liu et al., 2021; Jiang et al., 2022; Zhang et al., 2025d), an increasing number of approaches also return a reasoning process (Zhao et al., 2024c; Chegini et al., 2025; Nguyen et al., 2024). Such reasoning not only improves system performance, but it also provides human-understandable justifications for how an answer is derived. However, despite appearing plausible, these explanations may not faithfully represent the model's actual decision making process (Turpin et al., 2023; Chen et al., 2025a). Building trustworthy TQA systems is important, especially for high-stakes domains such as medicine (Bardhan et al., 2022). Achieving this requires output reasoning to accurately reflect a model's table understanding capabilities, e.g., faithfully responding with "I don't know" when a table is beyond a model's ability to interpret. We argue that much of the current work on interpretability in TQA focuses on generating post-hoc justifications for answers, rather than genuine explanations that transparently reveal the reasoning process underlying answer derivation.

**Human-Centric and Socially-Aware Setups.** Current research in TQA has largely focused on improving system performance, often overlooking the role of human interaction. However, the ultimate aim of such systems is to empower humans to more effectively interact with tabular data, analogous to the broader goal of developing NLP systems for human and socially aware uses (Hovy and Yang, 2021; Ziems et al., 2024; Yang et al., 2025b). This highlights the importance of human-centric and socially aware design. We identify two complementary research directions: (1) human-centric modeling, which emphasizes representations and benchmarks that capture the unique characteristics of tabular data (e.g., Hu et al., 2024; Ahmad et al., 2025); and (2) socially grounded applications, where TQA systems serve downstream tasks for social good, such as analyzing environmental sustainability reports (Dimmelmeier et al., 2024; Beck et al., 2025), advancing biomedical research (Luo et al., 2022), and supporting decision-making in financial domains (Strich et al., 2025a).

# 6 Conclusion

In this survey, we have reviewed recent advances in TQA with (M)LLMs, covering representative TQA setups, key challenges, and corresponding solutions, along with a discussion of promising future research directions. Looking ahead, we envision a stronger synergy between methods from NLP and related fields, and anticipate that TQA research, both in modeling and evaluation, will increasingly adopt more comprehensive settings. Such settings should account for diverse factors, including human model interaction, system robustness, and real-world applicability.

## Limitations

In this study, we present a survey on TQA with LLMs. Related surveys are discussed in Appendix A.2, and we plan to continuously incorporate additional approaches with more detailed analyses. Despite our best efforts, certain limitations remain.

**References and Methods.** We collected papers published before August 2025, which means that works appearing after this time are not included. We will continue to monitor the literature and update the survey accordingly. The majority of papers were retrieved from venues such as ACL, EMNLP, NAACL, NeurIPS, ICLR, and arXiv, using English-language queries. This approach may have led to the omission of works published in other languages. Furthermore, due to space constraints, we are unable to provide exhaustive technical details for all methods covered in this survey.

**Survey Scope.** This survey focuses exclusively on table question answering. Our objective is to provide a comprehensive and detailed overview of the task's characteristics, challenges, corresponding modeling methods, as well as emerging topics for future research. This focus excludes other table-related tasks, such as table generation and summarization.

## References

Nikhil Abhyankar, Vivek Gupta, Dan Roth, and Chandan K. Reddy. 2024. H-star: Llm-driven hybrid sql-text adaptive reasoning on tables. *ArXiv*, abs/2407.05952.

Nikhil Abhyankar, Vivek Gupta, Dan Roth, and Chandan K. Reddy. 2025. H-STAR: LLM-driven hybrid SQL-text adaptive reasoning on tables. In *Proceedings of the 2025 Conference of the Nations of the Americas Chapter of the Association for Computational Linguistics: Human Language Technologies (Volume 1: Long Papers)*, pages 8841–8863, Albuquerque, New Mexico. Association for Computational Linguistics.

Ankush Agarwal, Chaitanya Devaguptapu, and Ganesh S. 2025. Hybrid graphs for table-and-text based question answering using llms. In *North American Chapter of the Association for Computational Linguistics*.

Chaitanya Agarwal, Vivek Gupta, Anoop Kunchukuttan, and Manish Shrivastava. 2022. Bilingual tabular inference: A case study on Indic languages. In *Proceedings of the 2022 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*, pages

4018–4037, Seattle, United States. Association for Computational Linguistics.

Mohammad S. Ahmad, Zan A. Naeem, Michaël Aupetit, Ahmed Elmagarmid, Mohamed Eltabakh, Xiasong Ma, Mourad Ouzzani, and Chaoyi Ruan. 2025. Hct-qa: A benchmark for question answering on human-centric tables.

Mubashara Akhtar, Chenxi Pang, Andreea Marzoca, Yasemin Altun, and Julian Martin Eisenschlos. 2024. Tanq: An open domain dataset of table answered questions. *Trans. Assoc. Comput. Linguistics*, 13:461–480.

Rana Alshaikh, Israa Alghanmi, and Shelan S. Jeawak. 2025. Aratable: Benchmarking llms' reasoning and understanding of arabic tabular data. *ArXiv*, abs/2507.18442.

Shir Ashury-Tahan, Yifan Mai, C Rajmohan, Ariel Gera, Yotam Perlitz, Asaf Yehudai, Elron Bandel, Leshem Choshen, Eyal Shnarch, Percy Liang, and Michal Shmueli-Scheuer. 2025. The mighty torr: A benchmark for table reasoning and robustness. *ArXiv*, abs/2502.19412.

Ye Bai, Minghan Wang, and Thuy-Trang Vu. 2025. Maple: Multi-agent adaptive planning with long-term memory for table reasoning. *ArXiv*, abs/2506.05813.

Jayetri Bardhan, Anthony Colas, Kirk Roberts, and Daisy Zhe Wang. 2022. DrugEHRQA: A question answering dataset on structured and unstructured electronic health records for medicine related queries. In *Proceedings of the Thirteenth Language Resources and Evaluation Conference*, pages 1083–1097, Marseille, France. European Language Resources Association.

Jayetri Bardhan, Bushi Xiao, and Daisy Zhe Wang. 2024. Ttqa-rs- a break-down prompting approach for multi-hop table-text question answering with reasoning and summarization. *ArXiv*, abs/2406.14732.

Jacob Beck, Anna Steinberg, Andreas Dimmelmeier, Laia Domenech Burin, Emily Kormanyos, Maurice Fehr, and Malte Schierholz. 2025. Addressing data gaps in sustainability reporting: A benchmark dataset for greenhouse gas emission extraction. *Scientific Data*, 12(1):1497.

Michael J. Cafarella, Alon Halevy, Daisy Zhe Wang, Eugene Wu, and Yang Zhang. 2008. Webtables: exploring the power of tables on the web. *Proc. VLDB Endow.*, 1(1):538–549.

Lang Cao. 2025. Tablemaster: A recipe to advance table understanding with language models. *ArXiv*, abs/2501.19378.

Lang Cao, Jingxian Xu, Hanbing Liu, Jinyu Wang, Mengyu Zhou, Haoyu Dong, Shi Han, and Dongmei Zhang. 2025. Fortune: Formula-driven reinforcement learning for symbolic table reasoning in language models. *ArXiv*, abs/2505.23667.

Yihan Cao, Shuyi Chen, Ryan Liu, Zhiruo Wang, and Daniel Fried. 2023. API-assisted code generation for question answering on varied table structures. In *Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing*, pages 14536–14548, Singapore. Association for Computational Linguistics.

Atoosa Chegini, Keivan Rezaei, Hamid Eghbalzadeh, and Soheil Feizi. 2025. RePanda: Pandas-powered tabular verification and reasoning. In *Proceedings of the 63rd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 32200–32212, Vienna, Austria. Association for Computational Linguistics.

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. 2024. Tablerag: Million-token table understanding with language models. *ArXiv*, abs/2410.04739.

Wenhu Chen, Ming-Wei Chang, Eva Schlinger, William Yang Wang, and William W. Cohen. 2020a. Open question answering over tables and text. *ArXiv*, abs/2010.10439.

Wenhu Chen, Hongmin Wang, Jianshu Chen, Yunkai Zhang, Hong Wang, SHIYANG LI, Xiyou Zhou, and William Yang Wang. 2019. Tabfact: A large-scale dataset for table-based fact verification. *ArXiv*, abs/1909.02164.

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

Yanda Chen, Joe Benton, Ansh Radhakrishnan, Jonathan Uesato, Carson E. Denison, John Schulman, Arushi Somani, Peter Hase, Misha Wagner, Fabien Roger, Vladimir Mikulik, Sam Bowman, Jan Leike, Jared Kaplan, and Ethan Perez. 2025a. Reasoning models don't always say what they think. *ArXiv*, abs/2505.05410.

Yibin Chen, Yifu Yuan, Zeyu Zhang, Yan Zheng, Jinyi Liu, Fei Ni, Jianye Hao, Hangyu Mao, and Fuzheng Zhang. 2025b. Sheetagent: Towards a generalist agent for spreadsheet reasoning and manipulation via large language models. In *Proceedings of the ACM on Web Conference 2025*, WWW '25, page 158–177, New York, NY, USA. Association for Computing Machinery.

Yongrui Chen, Junhao He, Linbo Fu, Shenyu Zhang, Rihui Jin, Xinbang Dai, Jiaqi Li, Dehai Min, Nan Hu, Yuxin Zhang, Guilin Qi, Yi Huang, and Tongtong Wu. 2025c. Pandora: A code-driven large language model agent for unified reasoning across diverse structured knowledge. *ArXiv*, abs/2504.12734.

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

Sitao Cheng, Ziyuan Zhuang, Yong Xu, Fangkai Yang, Chaoyun Zhang, Xiaoting Qin, Xiang Huang, Ling Chen, Qingwei Lin, Dongmei Zhang, Saravan Rajmohan, and Qi Zhang. 2024. Call me when necessary: LLMs can efficiently and faithfully reason over structured environments. In *Findings of the Association for Computational Linguistics: ACL 2024*, pages 4275–4295, Bangkok, Thailand. Association for Computational Linguistics.

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

Zhoujun Cheng, Tianbao Xie, Peng Shi, Chengzu Li, Rahul Nadkarni, Yushi Hu, Caiming Xiong, Dragomir R. Radev, Marilyn Ostendorf, Luke S. Zettlemoyer, Noah A. Smith, and Tao Yu. 2022b. Binding language models in symbolic languages. *ArXiv*, abs/2210.02875.

Sanghyun Cho, Hye-Lynn Kim, Jung-Hun Lee, and Hyuk-Chul Kwon. 2025. Swing: Weakly supervised table question answering with self-training via reinforcement learning. In *2025 IEEE International Conference on Big Data and Smart Computing (BigComp)*, pages 273–278.

Philipp Christmann, Rishiraj Saha Roy, and Gerhard Weikum. 2023. Compmix: A benchmark for heterogeneous question answering. *Companion Proceedings of the ACM Web Conference 2024*.

DeepSeek-AI, Daya Guo, Dejian Yang, Haowei Zhang, Jun-Mei Song, Ruoyu Zhang, Runxin Xu, Qihao Zhu, Shirong Ma, Peiyi Wang, Xiaoling Bi, Xiaokang Zhang, Xingkai Yu, Yu Wu, Z. F. Wu, Zhibin Gou, Zhihong Shao, Zhuoshu Li, Ziyi Gao, Aixin Liu, Bing Xue, Bing-Li Wang, Bochao Wu, Bei Feng, Chengda Lu, Chenggang Zhao, Chengqi Deng, Chenyu Zhang, Chong Ruan, Damai Dai, Deli Chen, Dong-Li Ji, Erhang Li, Fangyun Lin, Fucong Dai, Fuli Luo, Guangbo Hao, Guanting Chen, Guowei Li, H. Zhang, Han Bao, Hanwei Xu, Haocheng Wang, Honghui Ding, Huajian Xin, Huazuo Gao, Hui Qu, Hui Li, Jianzhong Guo, Jiashi Li, Jiawei Wang, Jingchang Chen, Jingyang Yuan, Junjie Qiu, Junlong Li, Jiong Cai, Jiaqi Ni, Jian Liang, Jin Chen, Kai Dong, Kai Hu, Kaige Gao, Kang Guan, Kexin Huang, Kuai Yu, Lean Wang, Lecong Zhang, Liang

Zhao, Litong Wang, Liyue Zhang, Lei Xu, Leyi Xia, Mingchuan Zhang, Minghua Zhang, M. Tang, Meng Li, Miaojun Wang, Mingming Li, Ning Tian, Panpan Huang, Peng Zhang, Qiancheng Wang, Qinyu Chen, Qiushi Du, Ruiqi Ge, Ruisong Zhang, Ruizhe Pan, Runji Wang, R. J. Chen, Ruiqi Jin, Ruyi Chen, Shanghao Lu, Shangyan Zhou, Shanhuang Chen, Shengfeng Ye, Shiyu Wang, Shuiping Yu, Shunfeng Zhou, Shuting Pan, S. S. Li, Shuang Zhou, Shao-Kang Wu, Tao Yun, Tian Pei, Tianyu Sun, T. Wang, Wangding Zeng, Wanjia Zhao, Wen Liu, Wenfeng Liang, Wenjun Gao, Wen-Xia Yu, Wentao Zhang, Wangding Xiao, Wei An, Xiaodong Liu, Xiaohan Wang, Xiaokang Chen, Xiaotao Nie, Xin Cheng, Xin Liu, Xin Xie, Xingchao Liu, Xinyu Yang, Xinyuan Li, Xuecheng Su, Xuheng Lin, X. Q. Li, Xiangyu Jin, Xi-Cheng Shen, Xiaosha Chen, Xiaowen Sun, Xiaoxiang Wang, Xinnan Song, Xinyi Zhou, Xianzu Wang, Xinxia Shan, Y. K. Li, Y. Q. Wang, Y. X. Wei, Yang Zhang, Yanhong Xu, Yao Li, Yao Zhao, Yaofeng Sun, Yaohui Wang, Yi Yu, Yichao Zhang, Yifan Shi, Yi Xiong, Ying He, Yishi Piao, Yisong Wang, Yixuan Tan, Yiyang Ma, Yiyuan Liu, Yongqiang Guo, Yuan Ou, Yuduan Wang, Yue Gong, Yu-Jing Zou, Yujia He, Yunfan Xiong, Yu-Wei Luo, Yu mei You, Yuxuan Liu, Yuyang Zhou, Y. X. Zhu, Yanping Huang, Yao Li, Yi Zheng, Yuchen Zhu, Yunxiang Ma, Ying Tang, Yukun Zha, Yuting Yan, Zehui Ren, Zehui Ren, Zhangli Sha, Zhe Fu, Zhean Xu, Zhenda Xie, Zhen guo Zhang, Zhewen Hao, Zhicheng Ma, Zhigang Yan, Zhiyu Wu, Zihui Gu, Zijia Zhu, Zijun Liu, Zi-An Li, Ziwei Xie, Ziyang Song, Zizheng Pan, Zhen Huang, Zhipeng Xu, Zhongyu Zhang, and Zhen Zhang. 2025. Deepseek-r1: Incentivizing reasoning capability in llms via reinforcement learning. *ArXiv*, abs/2501.12948.

Irwin Deng, Kushagra Dixit, Vivek Gupta, and Dan Roth. 2024a. Enhancing temporal understanding in llms for semi-structured tables. *ArXiv*, abs/2407.16030.

Naihao Deng and Rada Mihalcea. 2025. Rethinking table instruction tuning. In *Findings of the Association for Computational Linguistics: ACL 2025*, pages 21757–21780, Vienna, Austria. Association for Computational Linguistics.

Naihao Deng, Zhenjie Sun, Ruiqi He, Aman Sikka, Yulong Chen, Lin Ma, Yue Zhang, and Rada Mihalcea. 2024b. Tables as texts or images: Evaluating the table reasoning ability of LLMs and MLLMs. In *Findings of the Association for Computational Linguistics: ACL 2024*, pages 407–426, Bangkok, Thailand. Association for Computational Linguistics.

Yang Deng, Wenqiang Lei, Wenxuan Zhang, Wai Lam, and Tat-Seng Chua. 2022. PACIFIC: Towards proactive conversational question answering over tabular and textual data in finance. In *Proceedings of the 2022 Conference on Empirical Methods in Natural Language Processing*, pages 6970–6984, Abu Dhabi, United Arab Emirates. Association for Computational Linguistics.

Andreas Dimmelmeier, Hendrik Doll, Malte Schierholz, Emily Kormanyos, Maurice Fehr, Bolei Ma, Jacob Beck, Alexander Fraser, and Frauke Kreuter. 2024. Informing climate risk analysis using textual information - a research agenda. In *Proceedings of the 1st Workshop on Natural Language Processing Meets Climate Change (ClimateNLP 2024)*, pages 12–26, Bangkok, Thailand. Association for Computational Linguistics.

Kushagra Dixit, Abhishek Rajgaria, Harshavardhan Kalalbandi, Dan Roth, and Vivek Gupta. 2025. No universal prompt: Unifying reasoning through adaptive prompting for temporal table reasoning. *ArXiv*, abs/2506.11246.

Haoyu Dong, Zhoujun Cheng, Xinyi He, Mengyu Zhou, Anda Zhou, Fan Zhou, Ao Liu, Shi Han, and Dongmei Zhang. 2022. Table pre-training: A survey on model architectures, pre-training objectives, and downstream tasks. In *Proceedings of the Thirty-First International Joint Conference on Artificial Intelligence, IJCAI-22*, pages 5426–5435. International Joint Conferences on Artificial Intelligence Organization. Survey Track.

Haoyu Dong, Haochen Wang, Anda Zhou, and Yue Hu. 2024. Ttc-quali: A text-table-chart dataset for multimodal quantity alignment. In *Proceedings of the 17th ACM International Conference on Web Search and Data Mining*, WSDM '24, page 181–189, New York, NY, USA. Association for Computing Machinery.

Xi Fang, Weijie Xu, Fiona Anting Tan, Jiani Zhang, Ziqing Hu, Yanjun Qi, Scott Nickleach, Diego Socolinsky, Srinivasan Sengamedu, and Christos Faloutsos. 2024a. Large language models(llms) on tabular data: Prediction, generation, and understanding – a survey.

Xi Fang, Weijie Xu, Fiona Anting Tan, Jiani Zhang, Ziqing Hu, Yanjun Qi, Scott Nickleach, Diego Socolinsky, Srinivasan H. Sengamedu, and Christos Faloutsos. 2024b. Large language models(llms) on tabular data: Prediction, generation, and understanding - a survey. *ArXiv*, abs/2402.17944.

Negar Foroutan, Angelika Romanou, Matin Ansaripour, Julian Martin Eisenschlos, Karl Aberer, and Rémi Lebret. 2025. Wikimixqa: A multimodal benchmark for question answering over tables and charts. *ArXiv*, abs/2506.15594.

Somraj Gautam, Abhishek Bhandari, and Gaurav Harit. 2025. TabComp: A dataset for visual table reading comprehension. In *Findings of the Association for Computational Linguistics: NAACL 2025*, pages 5773–5780, Albuquerque, New Mexico. Association for Computational Linguistics.

Carlos Gemmell and Jeff Dalton. 2023. ToolWriter: Question specific tool synthesis for tabular data. In *Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing*, pages 16137–16148, Singapore. Association for Computational Linguistics.

Akash Ghosh, Venkata Sahith Bathini, Niloy Ganguly, Pawan Goyal, and Mayank Singh. 2024. How robust are the QA models for hybrid scientific tabular data? a study using customized dataset. In *Proceedings of the 2024 Joint International Conference on Computational Linguistics, Language Resources and Evaluation (LREC-COLING 2024)*, pages 8258–8264, Torino, Italia. ELRA and ICCL.

Jiawei Gu, Ziting Xian, Yuanzhen Xie, Ye Liu, Enjie Liu, Ruichao Zhong, Mochi Gao, Yunzhi Tan, Bo Hu, and Zang Li. 2025. Toward structured knowledge reasoning: Contrastive retrieval-augmented generation on experience. In *Findings of the Association for Computational Linguistics: ACL 2025*, pages 23891–23910, Vienna, Austria. Association for Computational Linguistics.

Che Guan, Mengyu Huang, and Peng Zhang. 2024. Mfort-qa: Multi-hop few-shot open rich table question answering. *Proceedings of the 2024 10th International Conference on Computing and Artificial Intelligence*.

Manbir Gulati and Paul Roysdon. 2023. Tabmt: Generating tabular data with masked transformers. In *Advances in Neural Information Processing Systems*, volume 36, pages 46245–46254. Curran Associates, Inc.

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

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

Xinyi He, Yihao Liu, Mengyu Zhou, Yeye He, Haoyu Dong, Shi Han, Zejian Yuan, and Dongmei Zhang. 2025. TableLoRA: Low-rank adaptation on table structure understanding for large language models. In *Proceedings of the 63rd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 22376–22391, Vienna, Austria. Association for Computational Linguistics.

Xinyi He, Mengyu Zhou, Xinrun Xu, Xiaojun Ma, Rui Ding, Lun Du, Yan Gao, Ran Jia, Xu Chen, Shi Han, Zejian Yuan, and Dongmei Zhang. 2023. Text2analysis: A benchmark of table question answering with advanced data analysis and unclear queries. In *AAAI Conference on Artificial Intelligence*.

Xinyi He, Mengyu Zhou, Xinrun Xu, Xiaojun Ma, Rui Ding, Lun Du, Yan Gao, Ran Jia, Xu Chen, Shi Han, Zejian Yuan, and Dongmei Zhang. 2024. Text2analysis: a benchmark of table question answering with advanced data analysis and unclear queries. In *Proceedings of the Thirty-Eighth AAAI Conference on Artificial Intelligence and Thirty-Sixth Conference on Innovative Applications of Artificial Intelligence and Fourteenth Symposium on Educational Advances in Artificial Intelligence*, AAAI'24/IAAI'24/EAAI'24. AAAI Press.

Jonathan Herzig, Pawel Krzysztof Nowak, Thomas Müller, Francesco Piccinno, and Julian Eisenschlos. 2020. TaPas: Weakly supervised table parsing via pre-training. In *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*, pages 4320–4333, Online. Association for Computational Linguistics.

Maximiliano Hormazábal-Lagos, Álvaro Bueno Saez, Pedro Alonso Doval, Jorge Alcalde Vesteiro, and Héctor Cerezo-Costas. 2025. Explicit-qa: Explainable code-based image table question answering. *ArXiv*, abs/2507.11694.

Dirk Hovy and Diyi Yang. 2021. The importance of modeling social factors of language: Theory and practice. In *Proceedings of the 2021 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*, pages 588–602, Online. Association for Computational Linguistics.

J. Edward Hu, Yelong Shen, Phillip Wallis, Zeyuan Allen-Zhu, Yuanzhi Li, Shean Wang, and Weizhu Chen. 2021. Lora: Low-rank adaptation of large language models. *ArXiv*, abs/2106.09685.

Mengkang Hu, Haoyu Dong, Ping Luo, Shi Han, and Dongmei Zhang. 2024. KET-QA: A dataset for knowledge enhanced table question answering. In *Proceedings of the 2024 Joint International Conference on Computational Linguistics, Language Resources and Evaluation (LREC-COLING 2024)*, pages 9705–9719, Torino, Italia. ELRA and ICCL.

Sirui Huang, Yanggan Gu, Zhonghao Li, Xuming Hu, Li Qing, and Guandong Xu. 2025a. StructFact: Reasoning factual knowledge from structured data with large language models. In *Findings of the Association for Computational Linguistics: ACL 2025*, pages 7521–7552, Vienna, Austria. Association for Computational Linguistics.

Sirui Huang, Hanqian Li, Yanggan Gu, Xuming Hu, Qing Li, and Guandong Xu. 2025b. Hyperg: Hypergraph-enhanced llms for structured knowledge. *ArXiv*, abs/2502.18125.

Mohit Iyyer, Wen-tau Yih, and Ming-Wei Chang. 2017. Search-based neural structured learning for sequential question answering. In *Proceedings of the 55th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 1821–1831, Vancouver, Canada. Association for Computational Linguistics.

Sujay Kumar Jauhar, Peter Turney, and Eduard Hovy. 2016. Tables as semi-structured knowledge for question answering. In *Proceedings of the 54th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 474–483, Berlin, Germany. Association for Computational Linguistics.

Deyi Ji, Lanyun Zhu, Siqi Gao, Peng Xu, Hongtao Lu, Jieping Ye, and Feng Zhao. 2024. Tree-of-table: Unleashing the power of llms for enhanced large-scale table understanding. *ArXiv*, abs/2411.08516.

Jinhao Jiang, Kun Zhou, Zican Dong, Keming Ye, Xin Zhao, and Ji-Rong Wen. 2023. 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*, pages 9237–9251, Singapore. Association for Computational Linguistics.

Jun-Peng Jiang, Yu Xia, Hai-Long Sun, Shiyin Lu, Qing-Guo Chen, Weihua Luo, Kaifu Zhang, De-Chuan Zhan, and Han-Jia Ye. 2025. Multimodal tabular reasoning with privileged structured information. *ArXiv*, abs/2506.04088.

Ruya Jiang, Chun Wang, and Weihong Deng. 2024. Seek and solve reasoning for table question answering. *ArXiv*, abs/2409.05286.

Zhengbao Jiang, Yi Mao, Pengcheng He, Graham Neubig, and Weizhu Chen. 2022. 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*, pages 932–942, Seattle, United States. Association for Computational Linguistics.

Nengzheng Jin, Joanna Siebert, Dongfang Li, and Qingcai Chen. 2022. A survey on table question answering: Recent advances.

Rihui Jin, Yu Li, Guilin Qi, Nan Hu, Yuan-Fang Li, Jiaoyan Chen, Jianan Wang, Yongrui Chen, Dehai Min, and Sheng Bi. 2025a. Hegta: Leveraging heterogeneous graph-enhanced large language models for few-shot complex table understanding. *Proceedings of the AAAI Conference on Artificial Intelligence*, 39(23):24294–24302.

Rihui Jin, Zheyu Xin, Xing Xie, Zuoyi Li, Guilin Qi, Yongrui Chen, Xinbang Dai, Tongtong Wu, and Gholamreza Haffari. 2025b. Table-r1: Self-supervised and reinforcement learning for program-based table reasoning in small language models. *ArXiv*, abs/2506.06137.

Changwook Jun, Jooyoung Choi, Myoseop Sim, Hyun Kim, Hansol Jang, and Kyungkoo Min. 2022. Korean-specific dataset for table question answering. In *Proceedings of the Thirteenth Language Resources and Evaluation Conference*, pages 6114–6120, Marseille, France. European Language Resources Association.

Vladimir Karpukhin, Barlas Oguz, Sewon Min, Patrick Lewis, Ledell Wu, Sergey Edunov, Danqi Chen, and Wen-tau Yih. 2020. Dense passage retrieval for open-domain question answering. In *Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing (EMNLP)*, pages 6769–6781, Online. Association for Computational Linguistics.

Yannis Katsis, Saneem Chemmengath, Vishwajeet Kumar, Samarth Bharadwaj, Mustafa Canim, Michael Glass, Alfio Gliozzo, Feifei Pan, Jaydeep Sen, Karthik Sankaranarayanan, and Soumen Chakrabarti. 2022. AIT-QA: Question answering dataset over complex tables in the airline industry. In *Proceedings of the 2022 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies: Industry Track*, pages 305–314, Hybrid: Seattle, Washington + Online. Association for Computational Linguistics.

Rohit Khoja, Devanshu Gupta, Yanjie Fu, Dan Roth, and Vivek Gupta. 2025. Weaver: Interweaving sql and llm for table reasoning. *ArXiv*, abs/2505.18961.

Yoonsik Kim, Moonbin Yim, and Ka Yeon Song. 2024. Tablevqa-bench: A visual question answering benchmark on multiple table domains. *ArXiv*, abs/2404.19205.

Atsushi Kojima. 2024. Sub-table rescorer for table question answering. In *Proceedings of the 2024 Joint International Conference on Computational Linguistics, Language Resources and Evaluation (LREC-COLING 2024)*, pages 15422–15427, Torino, Italia. ELRA and ICCL.

Kezhi Kong, Jiani Zhang, Zhengyuan Shen, Balasubramaniam Srinivasan, Chuan Lei, Christos Faloutsos, Huzefa Rangwala, and George Karypis. 2024. Opentab: Advancing large language models as open-domain table reasoners. *ArXiv*, abs/2402.14361.

Sunjun Kweon, Yeonsu Kwon, Seonhee Cho, Yohan Jo, and Edward Choi. 2023. Open-WikiTable : Dataset for open domain question answering with complex reasoning over table. In *Findings of the Association for Computational Linguistics: ACL 2023*, pages 8285–8297, Toronto, Canada. Association for Computational Linguistics.

Chia-Hsuan Lee, Oleksandr Polozov, and Matthew Richardson. 2021. KaggleDBQA: Realistic evaluation of text-to-SQL parsers. In *Proceedings of the 59th Annual Meeting of the Association for Computational Linguistics and the 11th International Joint Conference on Natural Language Processing (Volume 1: Long Papers)*, pages 2261–2273, Online. Association for Computational Linguistics.

Wonjin Lee, Kyumin Kim, Sungjae Lee, Jihun Lee, and Kwang In Kim. 2024a. Piece of table: A divide-and-conquer approach for selecting sub-tables in table question answering. *ArXiv*, abs/2412.07629.

Younghun Lee, Sungchul Kim, Ryan A. Rossi, Tong Yu, and Xiang Chen. 2024b. Learning to reduce: Towards improving performance of large language models on structured data. *ArXiv*, abs/2407.02750.

Fangyu Lei, Jinxiang Meng, Yiming Huang, Tinghong Chen, Yun Zhang, Shizhu He, Jun Zhao, and Kang Liu. 2025. Reasoning-table: Exploring reinforcement learning for table reasoning. *ArXiv*, abs/2506.01710.

Ce Li, Xiaofan Liu, Zhiyan Song, Ce Chi, Chen Zhao, Jingjing Yang, Zhendong Wang, Kexin Yang, Boshen Shi, Xing Wang, Chao Deng, and Junlan Feng. 2025a. Treb: A comprehensive benchmark for evaluating table reasoning capabilities of large language models. *ArXiv*, abs/2506.18421.

Jinyang Li, Binyuan Hui, Ge Qu, Jiaxi Yang, Binhua Li, Bowen Li, Bailin Wang, Bowen Qin, Ruiying Geng, Nan Huo, Xuanhe Zhou, Chenhao Ma, Guoliang Li, Kevin C.C. Chang, Fei Huang, Reynold Cheng, and Yongbin Li. 2023. Can llm already serve as a database interface? a big bench for large-scale database grounded text-to-sqls. In *Proceedings of the 37th International Conference on Neural Information Processing Systems*, NIPS '23, Red Hook, NY, USA. Curran Associates Inc.

Liyao Li, Chao Ye, Wentao Ye, Yifei Sun, Zhe Jiang, Haobo Wang, Gang Chen, and Junbo Zhao. 2025b. Injecting learnable table features into LLMs.

Peng Li, Yeye He, Dror Yashar, Weiwei Cui, Song Ge, Haidong Zhang, Danielle Rifinski Fainman, Dongmei Zhang, and Surajit Chaudhuri. 2024. Table-gpt: Table fine-tuned gpt for diverse table tasks. *Proc. ACM Manag. Data*, 2(3).

Qianlong Li, Chen Huang, Shuai Li, Yuanxin Xiang, Deng Xiong, and Wenqiang Lei. 2025c. GraphOTTER: Evolving LLM-based graph reasoning for complex table question answering. In *Proceedings of the 31st International Conference on Computational Linguistics*, pages 5486–5506, Abu Dhabi, UAE. Association for Computational Linguistics.

Xiao Li, Yawei Sun, and Gong Cheng. 2021. Tsqa: Tabular scenario based question answering. *ArXiv*, abs/2101.11429.

Yongqi Li, Wenjie Li, and Liqiang Nie. 2022. MM-CoQA: Conversational question answering over text, tables, and images. In *Proceedings of the 60th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 4220–4231, Dublin, Ireland. Association for Computational Linguistics.

Zheng Li, Yang Du, Mao Zheng, and Mingyang Song. 2025d. MiMoTable: A multi-scale spreadsheet benchmark with meta operations for table reasoning. In *Proceedings of the 31st International Conference on Computational Linguistics*, pages 2548–2560, Abu Dhabi, UAE. Association for Computational Linguistics.

Hsing-Ping Liang, Che-Wei Chang, and Yao-Chung Fan. 2025. Improving table retrieval with question generation from partial tables. In *Proceedings of the 4th Table Representation Learning Workshop*, pages 217–228, Vienna, Austria. Association for Computational Linguistics.

Weizhe Lin, Rexhina Blloshmi, Bill Byrne, Adria de Gispert, and Gonzalo Iglesias. 2023. An inner table retriever for robust table question answering. In *Proceedings of the 61st Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 9909–9926, Toronto, Canada. Association for Computational Linguistics.

Chuang Liu, Junzhuo Li, and Deyi Xiong. 2023a. TabCQA: A tabular conversational question answering dataset on financial reports. In *Proceedings of the 61st Annual Meeting of the Association for Computational Linguistics (Volume 5: Industry Track)*, pages 196–207, Toronto, Canada. Association for Computational Linguistics.

Qian Liu, Bei Chen, Jiaqi Guo, Zeqi Lin, and Jian-Guang Lou. 2021. Tapex: Table pretraining via learning a neural sql executor. *ArXiv*, abs/2107.07653.

Tianyang Liu, Fei Wang, and Muhao Chen. 2024a. Rethinking tabular data understanding with large language models. In *Proceedings of the 2024 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies (Volume 1: Long Papers)*, pages 450–482, Mexico City, Mexico. Association for Computational Linguistics.

Weihao Liu, Fangyu Lei, Tongxu Luo, Jiahe Lei, Shizhu He, Jun Zhao, and Kang Liu. 2023b. Mmhqa-icl: Multimodal in-context learning for hybrid question answering over text, tables and images. *ArXiv*, abs/2309.04790.

Yujian Liu, Jiabao Ji, Tong Yu, Ryan Rossi, Sungchul Kim, Handong Zhao, Ritwik Sinha, Yang Zhang, and Shiyu Chang. 2024b. Augment before you try: Knowledge-enhanced table question answering via table expansion. *ArXiv*, abs/2401.15555.

Zhenghao Liu, Haolan Wang, Xinze Li, Qiushi Xiong, Xiaocui Yang, Yu Gu, Yukun Yan, Qi Shi, Fangfang Li, Ge Yu, and Maosong Sun. 2025. Hippo: Enhancing the table understanding capability of large language models through hybrid-modal preference optimization. *ArXiv*, abs/2502.17315.

Adrián López Gude, Roi Santos Ríos, Francisco Prado Valiño, Ana Ezquerro, and Jesús Vilares. 2025. LyS at SemEval 2025 task 8: Zero-shot code generation for tabular QA. In *Proceedings of the 19th International Workshop on Semantic Evaluation (SemEval-2025)*, pages 1282–1288, Vienna, Austria. Association for Computational Linguistics.

Pan Lu, Liang Qiu, Kai-Wei Chang, Ying Nian Wu, Song-Chun Zhu, Tanmay Rajpurohit, Peter Clark, and A. Kalyan. 2022. Dynamic prompt learning via policy gradient for semi-structured mathematical reasoning. *ArXiv*, abs/2209.14610.

Weizheng Lu, Jiaming Zhang, Jing Zhang, and Yueguo Chen. 2024a. Large language model for table processing: A survey. *Frontiers Comput. Sci.*, 19:192350.

Xinyuan Lu, Liangming Pan, Qian Liu, Preslav Nakov, and Min-Yen Kan. 2023. SCITAB: A challenging benchmark for compositional reasoning and claim verification on scientific tables. In *Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing*, pages 7787–7813, Singapore. Association for Computational Linguistics.

Xinyuan Lu, Liangming Pan, Yubo Ma, Preslav Nakov, and Min-Yen Kan. 2024b. Tart: An open-source tool-augmented framework for explainable table-based reasoning. In *North American Chapter of the Association for Computational Linguistics*.

Xinyuan Lu, Liangming Pan, Yubo Ma, Preslav Nakov, and Min-Yen Kan. 2025. TART: An open-source tool-augmented framework for explainable table-based reasoning. In *Findings of the Association for Computational Linguistics: NAACL 2025*, pages 4323–4339, Albuquerque, New Mexico. Association for Computational Linguistics.

Man Luo, Sharad Saxena, Swaroop Mishra, Mihir Parmar, and Chitta Baral. 2022. Biotabqa: Instruction learning for biomedical table question answering.

Tongxu Luo, Fangyu Lei, Jiahe Lei, Weihao Liu, Shihu He, Jun Zhao, and Kang Liu. 2023. Hrot: Hybrid prompt strategy and retrieval of thought for table-text hybrid question answering. *ArXiv*, abs/2309.12669.

Qingyang Mao, Qi Liu, Zhi Li, Mingyue Cheng, Zheng Zhang, and Rui Li. 2025. Potable: Towards systematic thinking via stage-oriented plan-then-execute reasoning on tables.

Suyash Vardhan Mathur, Jainit Sushil Bafna, Kunal Kartik, Harshita Khandelwal, Manish Shrivastava, Vivek Gupta, Mohit Bansal, and Dan Roth. 2024. Knowledge-aware reasoning over multimodal semi-structured tables. In *Findings of the Association for Computational Linguistics: EMNLP 2024*, pages 14054–14073, Miami, Florida, USA. Association for Computational Linguistics.

Bhavnick Minhas, Anant Shankhdhar, Vivek Gupta, Divyanshu Aggarwal, and Shuo Zhang. 2022. XInfoTabS: Evaluating multilingual tabular natural language inference. In *Proceedings of the Fifth Fact Extraction and VERification Workshop (FEVER)*, pages 59–77, Dublin, Ireland. Association for Computational Linguistics.

Raphaël Mouravieff, Benjamin Piwowarski, and Sylvain Lamprier. 2024. Learning relational decomposition of queries for question answering from tables. In *Proceedings of the 62nd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 10471–10485, Bangkok, Thailand. Association for Computational Linguistics.

Raphaël Mouravieff, Benjamin Piwowarski, and Sylvain Lamprier. 2025. Structural deep encoding for table question answering. In *Findings of the Association for Computational Linguistics: ACL 2025*, pages 2389–2402, Vienna, Austria. Association for Computational Linguistics.

Md Mahadi Hasan Nahid and Davood Rafiei. 2024a. Normtab: Improving symbolic reasoning in llms through tabular data normalization. *ArXiv*, abs/2406.17961.

Md Mahadi Hasan Nahid and Davood Rafiei. 2024b. Tabsqlify: Enhancing reasoning capabilities of llms through table decomposition. *ArXiv*, abs/2404.10150.

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. 2022. FeTaQA: Free-form table question answering. *Transactions of the Association for Computational Linguistics*, 10:35–49.

Giang (Dexter) Nguyen, Ivan Brugere, Shubham Sharma, Sanjay Kariyappa, Anh Totti Nguyen, and Freddy Lécué. 2024. Interpretable llm-based table question answering. *ArXiv*, abs/2412.12386.

Phuc Nguyen, Nam Tuan Ly, Hideaki Takeda, and Atsuhiro Takasu. 2023. Tabiqa: Table questions answering on business document images. *ArXiv*, abs/2303.14935.

Thi-Nhung Nguyen, Hoang Ngo, D.Q. Phung, Thuy-Trang Vu, and Dat Quoc Nguyen. 2025. Planning for success: Exploring llm long-term planning capabilities in table understanding. *Proceedings of the 29th Conference on Computational Natural Language Learning*.

Ansong Ni, Srini Iyer, Dragomir Radev, Ves Stoyanov, Wen-tau Yih, Sida I. Wang, and Xi Victoria Lin. 2023. Lever: learning to verify language-to-code generation with execution. In *Proceedings of the 40th International Conference on Machine Learning*, ICML'23. JMLR.org.

Jorge Osés Grijalba, L. Alfonso Ureña-López, Eugenio Martínez Cámara, and Jose Camacho-Collados. 2024. Question answering over tabular data with DataBench: A large-scale empirical evaluation of

LLMs. In *Proceedings of the 2024 Joint International Conference on Computational Linguistics, Language Resources and Evaluation (LREC-COLING 2024)*, pages 13471–13488, Torino, Italia. ELRA and ICCL.

Vaishali Pal, Evangelos Kanoulas, Andrew Yates, and Maarten de Rijke. 2024. Table question answering for low-resourced Indic languages. In *Proceedings of the 2024 Conference on Empirical Methods in Natural Language Processing*, pages 75–92, Miami, Florida, USA. Association for Computational Linguistics.

Vaishali Pal, Andrew Yates, Evangelos Kanoulas, and Maarten de Rijke. 2023. MultiTabQA: Generating tabular answers for multi-table question answering. In *Proceedings of the 61st Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 6322–6334, Toronto, Canada. Association for Computational Linguistics.

Chaoxu Pang, Yixuan Cao, Chunhao Yang, and Ping Luo. 2024. Uncovering limitations of large language models in information seeking from tables. In *Findings of the Association for Computational Linguistics: ACL 2024*, pages 1388–1409, Bangkok, Thailand. Association for Computational Linguistics.

Ankur Parikh, Xuezhi Wang, Sebastian Gehrmann, Manaal Faruqui, Bhuwan Dhingra, Diyi Yang, and Dipanjan Das. 2020. ToTTo: A controlled table-to-text generation dataset. In *Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing (EMNLP)*, pages 1173–1186, Online. Association for Computational Linguistics.

Panupong Pasupat and Percy Liang. 2015. 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)*, pages 1470–1480, Beijing, China. Association for Computational Linguistics.

Sohan Patnaik, Heril Changwal, Milan Aggarwal, Sumita Bhatia, Yaman Kumar, and Balaji Krishnamurthy. 2024. Cabinet: Content relevance based noise reduction for table question answering. *ArXiv*, abs/2402.01155.

Shraman Pramanick, Rama Chellappa, and Subhashini Venugopalan. 2024. Spiqa: A dataset for multimodal question answering on scientific papers. *ArXiv*, abs/2407.09413.

Zipeng Qiu, You Peng, Guangxin He, Binhang Yuan, and Chen Wang. 2024. Tqa-bench: Evaluating llms for multi-table question answering with scalable context and symbolic extension. *ArXiv*, abs/2411.19504.

Weijieying Ren, Tianxiang Zhao, Yuqing Huang, and Vasant Honavar. 2025. Deep learning within tabular data: Foundations, challenges, advances and future directions. *ArXiv*, abs/2501.03540.

Abhilash Shankarampeta, Harsh Mahajan, Tushar Kataria, Dan Roth, and Vivek Gupta. 2025. TRANSIENTTABLES: Evaluating LLMs' reasoning on temporally evolving semi-structured tables. In *Proceedings of the 2025 Conference of the Nations of the Americas Chapter of the Association for Computational Linguistics: Human Language Technologies (Volume 1: Long Papers)*, pages 6526–6544, Albuquerque, New Mexico. Association for Computational Linguistics.

Qi Shi, Han Cui, Haofeng Wang, Qingfu Zhu, Wanxiang Che, and Ting Liu. 2024. Exploring hybrid question answering via program-based prompting. In *Proceedings of the 62nd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 11035–11046, Bangkok, Thailand. Association for Computational Linguistics.

Daixin Shu, Jian Yang, Zhenhe Wu, Xianjie Wu, Xianfu Cheng, Xiangyuan Guan, Yanghai Wang, Pengfei Wu, Tingyang Yang, Hualei Zhu, Wei Zhang, Ge Zhang, Jiaheng Liu, and Zhoujun Li. 2025. M3tqa: Massively multilingual multitask table question answering.

Anshul Singh, Christian Biemann, and Jan Strich. 2025. Mtabvqa: Evaluating multi-tabular reasoning of language models in visual space. *ArXiv*, abs/2506.11684.

Atakan Site, Emre Hakan Erdemir, and Gülşen Eryiğit. 2025. Itunlp at semeval-2025 task 8: Question-answering over tabular data: A zero-shot approach using llm-driven code generation.

Josefa Lia Stoisser, Marc Boubnovski Martell, and Julien Fauqueur. 2025. Sparks of tabular reasoning via text2sql reinforcement learning. *ArXiv*, abs/2505.00016.

Jan Strich, Enes Kutay Isgorur, Maximilian Trescher, Chris Biemann, and Martin Semmann. 2025a. $T^2$-ragbench: Text-and-table benchmark for evaluating retrieval-augmented generation.

Jan Strich, Enes Kutay Isgorur, Maximilian Trescher, Christian Biemann, and Martin Semmann. 2025b. T2-ragbench: Text-and-table benchmark for evaluating retrieval-augmented generation. *ArXiv*, abs/2506.12071.

Aofeng Su, Aowen Wang, Chaonan Ye, Chengcheng Zhou, Ga Zhang, Gang Chen, Guangcheng Zhu, Haobo Wang, Haokai Xu, Hao Chen, Haoze Li, Haoxuan Lan, Jiaming Tian, Jing Yuan, Junbo Zhao, Junlin Zhou, Kaizhe Shou, Liangyu Zha, Lin Long, Liyao Li, Peng Wu, Qi Zhang, Qingyi Huang, Sa Yang, Tao Zhang, Wen-Yuan Ye, Wufang Zhu, Xiaomeng Hu, Xijun Gu, Xinjie Sun, Xiang Li, Yuhang Yang, and Zhiqing Xiao. 2024. Tablegpt2: A large multimodal model with tabular data integration. *ArXiv*, abs/2411.02059.

Yi Su, Dian Yu, Linfeng Song, Juntao Li, Haitao Mi, Zhaopeng Tu, Min Zhang, and Dong Yu. 2025. Crossing the reward bridge: Expanding rl with verifiable rewards across diverse domains. *ArXiv*, abs/2503.23829.

Yuan Sui, Jiaru Zou, Mengyu Zhou, Xinyi He, Lun Du, Shi Han, and Dongmei Zhang. 2024. 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*, pages 10306–10323, Miami, Florida, USA. Association for Computational Linguistics.

Alon Talmor, Ori Yoran, Amnon Catav, Dan Lahav, Yizhong Wang, Akari Asai, Gabriel Ilharco, Hannaneh Hajishirzi, and Jonathan Berant. 2021. Multimodalqa: Complex question answering over text, tables and images. *ArXiv*, abs/2104.06039.

Jiaming Tian, Liyao Li, Wentao Ye, Haobo Wang, Lingxin Wang, Lihua Yu, Zujie Ren, Gang Chen, and Junbo Zhao. 2025a. Toward real-world table agents: Capabilities, workflows, and design principles for llm-based table intelligence.

Shi-Yu Tian, Zhi Zhou, Wei Dong, Ming Yang, Kun-Yang Yu, Zi-Jian Cheng, Lan-Zhe Guo, and Yu-Feng Li. 2025b. Automated text-to-table for reasoning-intensive table qa: Pipeline design and benchmarking insights. *ArXiv*, abs/2505.19563.

Prasham Titiya, Jainil Trivedi, Chitta Baral, and Vivek Gupta. 2025. Mmtbench: A unified benchmark for complex multimodal table reasoning. *ArXiv*, abs/2505.21771.

Miles Turpin, Julian Michael, Ethan Perez, and Sam Bowman. 2023. Language models don't always say what they think: Unfaithful explanations in chain-of-thought prompting. *ArXiv*, abs/2305.04388.

Vishnou Vinayagame, Gregory Senay, and Luis Martí. 2025. Matata: Weakly supervised end-to-end mathematical tool-augmented reasoning for tabular applications.

A. A. Vyatkin and V. D. Oliseenko. 2025. Generating pandas code for big table question answering using large language models. In *2025 XXVIII International Conference on Soft Computing and Measurements (SCM)*, pages 164–166.

Hanjun Wang, Wenda Liu, Qun Wang, Jiaming Xu, Shuai Nie, Yuchen Liu, and Runyu Shi. 2024a. Light-table-chain: Simplifying and enhancing chain-based table reasoning. In *2024 10th International Conference on Computer and Communications (ICCC)*, pages 269–275.

Lanrui Wang, Mingyu Zheng, Hongyin Tang, Zheng Lin, Yanan Cao, Jingang Wang, Xunliang Cai, and Weiping Wang. 2025a. Needleinatable: Exploring long-context capability of large language models towards long-structured tables. *ArXiv*, abs/2504.06560.

Yuqi Wang, Lyuhao Chen, Songcheng Cai, Zhijian Xu, and Yilun Zhao. 2024b. Revisiting automated evaluation for long-form table question answering. In *Proceedings of the 2024 Conference on Empirical Methods in Natural Language Processing*, pages 14696–14706, Miami, Florida, USA. Association for Computational Linguistics.

Yuxiang Wang, Jianzhong Qi, and Junhao Gan. 2024c. Accurate and regret-aware numerical problem solver for tabular question answering. In *AAAI Conference on Artificial Intelligence*.

Zhensheng Wang, Wenmian Yang, Kun Zhou, Yiquan Zhang, and Weijia Jia. 2024d. Retqa: A large-scale open-domain tabular question answering dataset for real estate sector. In *AAAI Conference on Artificial Intelligence*.

Zhongyuan Wang, Richong Zhang, and Zhijie Nie. 2025b. General table question answering via answer-formula joint generation. *ArXiv*, abs/2503.12345.

Zilong Wang, Hao Zhang, Chun-Liang Li, Julian Martin Eisenschlos, Vincent Perot, Zifeng Wang, Lesly Miculicich, Yasuhisa Fujii, Jingbo Shang, Chen-Yu Lee, and Tomas Pfister. 2024e. Chain-of-table: Evolving tables in the reasoning chain for table understanding. *ArXiv*, abs/2401.04398.

Jason Wei, Xuezhi Wang, Dale Schuurmans, Maarten Bosma, Ed H. Chi, F. Xia, Quoc Le, and Denny Zhou. 2022. Chain of thought prompting elicits reasoning in large language models. *ArXiv*, abs/2201.11903.

Cornelius Wolff and Madelon Hulsebos. 2025. How well do llms reason over tabular data, really? *ArXiv*, abs/2505.07453.

Jian Wu, Linyi Yang, Dongyuan Li, Yuliang Ji, Manabu Okumura, and Yue Zhang. 2025a. MMQA: Evaluating LLMs with multi-table multi-hop complex questions. In *The Thirteenth International Conference on Learning Representations*.

Pengzuo Wu, Yuhang Yang, Guangcheng Zhu, Chao Ye, Hong Gu, Xu Lu, Ruixuan Xiao, Bowen Bao, Yijing He, Liangyu Zha, Wentao Ye, Junbo Zhao, and Haobo Wang. 2025b. RealHiTBench: A comprehensive realistic hierarchical table benchmark for evaluating LLM-based table analysis. In *Findings of the Association for Computational Linguistics: ACL 2025*, pages 7105–7137, Vienna, Austria. Association for Computational Linguistics.

Xianjie Wu, Jian Yang, Linzheng Chai, Ge Zhang, Jiaheng Liu, Xinrun Du, Di Liang, Daixin Shu, Xianfu Cheng, Tianzhen Sun, Guanglin Niu, Tongliang Li, and Zhoujun Li. 2024a. Tablebench: A comprehensive and complex benchmark for table question answering. *ArXiv*, abs/2408.09174.

Xiaofeng Wu, Alan Ritter, and Wei Xu. 2025c. Tabular data understanding with llms: A survey of recent advances and challenges. *ArXiv*, abs/2508.00217.

Xiaofeng Wu, Alan Ritter, and Wei Xu. 2025d. Tabular data understanding with llms: A survey of recent advances and challenges.

Xueqing Wu, Rui Zheng, Jingzhen Sha, Te-Lin Wu, Hanyu Zhou, Mohan Tang, Kai-Wei Chang, Nanyun Peng, and Haoran Huang. 2024b. Daco: Towards application-driven and comprehensive data analysis via code generation. *ArXiv*, abs/2403.02528.

Zirui Wu and Yansong Feng. 2024. Protrix: Building models for planning and reasoning over tables with sentence context. In *Conference on Empirical Methods in Natural Language Processing*.

Shiyu Xia, Junyu Xiong, Haoyu Dong, Jianbo Zhao, Yuzhang Tian, Mengyu Zhou, Yeye He, Shi Han, and Dongmei Zhang. 2024. Vision language models for spreadsheet understanding: Challenges and opportunities. In *Proceedings of the 3rd Workshop on Advances in Language and Vision Research (ALVR)*, pages 116–128, Bangkok, Thailand. Association for Computational Linguistics.

Junjie Xing, Yeye He, Mengyu Zhou, Haoyu Dong, Shi Han, Dongmei Zhang, and Surajit Chaudhuri. 2024. Table-llm-specialist: Language model specialists for tables using iterative generator-validator fine-tuning. *ArXiv*, abs/2410.12164.

Bohao Yang, Yingji Zhang, Dong Liu, Andr'e Freitas, and Chenghua Lin. 2025a. Does table source matter? benchmarking and improving multimodal scientific table understanding and reasoning. *ArXiv*, abs/2501.13042.

Diyi Yang, Dirk Hovy, David Jurgens, and Barbara Plank. 2025b. Socially aware language technologies: Perspectives and practices. *Computational Linguistics*, 51:689–703.

Jingfeng Yang, Aditya Gupta, Shyam Upadhyay, Luheng He, Rahul Goel, and Shachi Paul. 2022. TableFormer: Robust transformer modeling for table-text encoding. In *Proceedings of the 60th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 528–537, Dublin, Ireland. Association for Computational Linguistics.

Zhen Yang, Ziwei Du, Minghan Zhang, Wei Du, Jie Chen, Zhen Duan, and Shu Zhao. 2025c. Triples as the key: Structuring makes decomposition and verification easier in LLM-based tableQA. In *The Thirteenth International Conference on Learning Representations*.

Zheyuan Yang, Lyuhao Chen, Arman Cohan, and Yilun Zhao. 2025d. Table-r1: Inference-time scaling for table reasoning. *ArXiv*, abs/2505.23621.

Junyi Ye, Mengnan Du, and Guiling Wang. 2024. Dataframe qa: A universal llm framework on dataframe question answering without data exposure. *ArXiv*, abs/2401.15463.

Yunhu Ye, Binyuan Hui, Min Yang, Binhua Li, Fei Huang, and Yongbin Li. 2023. Large language models are versatile decomposers: Decomposing evidence and questions for table-based reasoning. In *Proceedings of the 46th International ACM SIGIR Conference on Research and Development in Information Retrieval*, SIGIR '23, page 174–184, New York, NY, USA. Association for Computing Machinery.

Peiying Yu, Guoxin Chen, and Jingjing Wang. 2025a. Table-critic: A multi-agent framework for collaborative criticism and refinement in table reasoning. In *Proceedings of the 63rd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 17432–17451, Vienna, Austria. Association for Computational Linguistics.

Tao Yu, Rui Zhang, Kai Yang, Michihiro Yasunaga, Dongxu Wang, Zifan Li, James Ma, Irene Li, Qingning Yao, Shanelle Roman, Zilin Zhang, and Dragomir Radev. 2018. Spider: A large-scale human-labeled dataset for complex and cross-domain semantic parsing and text-to-SQL task. In *Proceedings of the 2018 Conference on Empirical Methods in Natural Language Processing*, pages 3911–3921, Brussels, Belgium. Association for Computational Linguistics.

Xiaohan Yu, Pu Jian, and Chong Chen. 2025b. Tablerag: A retrieval augmented generation framework for heterogeneous document reasoning. *ArXiv*, abs/2506.10380.

Han Zhang, Yuheng Ma, and Hanfang Yang. 2024a. Alter: Augmentation for large-table-based reasoning. In *North American Chapter of the Association for Computational Linguistics*.

Han Zhang, Yuheng Ma, and Hanfang Yang. 2025a. ALTER: Augmentation for large-table-based reasoning. In *Proceedings of the 2025 Conference of the Nations of the Americas Chapter of the Association for Computational Linguistics: Human Language Technologies (Volume 1: Long Papers)*, pages 179–198, Albuquerque, New Mexico. Association for Computational Linguistics.

Junwen Zhang, Pu Chen, and Yin Zhang. 2025b. Table-moe: Neuro-symbolic routing for structured expert reasoning in multimodal table understanding. *ArXiv*, abs/2506.21393.

Tianshu Zhang, Xiang Yue, Yifei Li, and Huan Sun. 2024b. TableLlama: Towards open large generalist models for tables. In *Proceedings of the 2024 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies (Volume 1: Long Papers)*, pages 6024–6044, Mexico City, Mexico. Association for Computational Linguistics.

Wen Zhang, Long Jin, Yushan Zhu, Jiaoyan Chen, Zhiwei Huang, Junjie Wang, Yin Hua, Lei Liang,

and Huajun Chen. 2025c. Trustuqa: a trustful framework for unified structured data question answering. In *Proceedings of the Thirty-Ninth AAAI Conference on Artificial Intelligence and Thirty-Seventh Conference on Innovative Applications of Artificial Intelligence and Fifteenth Symposium on Educational Advances in Artificial Intelligence*, AAAI'25/IAAI'25/EAAI'25. AAAI Press.

Xiaokang Zhang, Sijia Luo, Bohan Zhang, Zeyao Ma, Jing Zhang, Yang Li, Guanlin Li, Zijun Yao, Kangli Xu, Jinchang Zhou, Daniel Zhang-Li, Jifan Yu, Shu Zhao, Juanzi Li, and Jie Tang. 2025d. TableLLM: Enabling tabular data manipulation by LLMs in real office usage scenarios. In *Findings of the Association for Computational Linguistics: ACL 2025*, pages 10315–10344, Vienna, Austria. Association for Computational Linguistics.

Xiaokang Zhang, Jing Zhang, Zeyao Ma, Yang Li, Bohan Zhang, Guanlin Li, Zijun Yao, Kangli Xu, Jinchang Zhou, Daniel Zhang-Li, Jifan Yu, Shu Zhao, Juan-Zi Li, and Jie Tang. 2024c. Tablellm: Enabling tabular data manipulation by llms in real office usage scenarios. *ArXiv*, abs/2403.19318.

Xuanliang Zhang, Dingzirui Wang, Longxu Dou, Baoxin Wang, Dayong Wu, Qingfu Zhu, and Wanxiang Che. 2024d. Flextaf: Enhancing table reasoning with flexible tabular formats. *ArXiv*, abs/2408.08841.

Xuanliang Zhang, Dingzirui Wang, Longxu Dou, Qingfu Zhu, and Wanxiang Che. 2024e. A survey of table reasoning with large language models. *Frontiers Comput. Sci.*, 19:199348.

Xuanliang Zhang, Dingzirui Wang, Longxu Dou, Qingfu Zhu, and Wanxiang Che. 2024f. A survey of table reasoning with large language models.

Xuanliang Zhang, Dingzirui Wang, Baoxin Wang, Longxu Dou, Xinyuan Lu, Keyan Xu, Dayong Wu, and Qingfu Zhu. 2025e. SCITAT: A question answering benchmark for scientific tables and text covering diverse reasoning types. In *Findings of the Association for Computational Linguistics: ACL 2025*, pages 3859–3881, Vienna, Austria. Association for Computational Linguistics.

Xuanliang Zhang, Dingzirui Wang, Keyan Xu, Qingfu Zhu, and Wanxiang Che. 2025f. Multitat: Benchmarking multilingual table-and-text question answering. *ArXiv*, abs/2502.17253.

Xuanliang Zhang, Dingzirui Wang, Keyan Xu, Qingfu Zhu, and Wanxiang Che. 2025g. Rot: Enhancing table reasoning with iterative row-wise traversals. *ArXiv*, abs/2505.15110.

Yuji Zhang, Qingyun Wang, Cheng Qian, Jiateng Liu, Chenkai Sun, Denghui Zhang, Tarek F. Abdelzaher, ChengXiang Zhai, Preslav Nakov, and Heng Ji. 2025h. Atomic reasoning for scientific table claim verification. *ArXiv*, abs/2506.06972.

Yunjia Zhang, Jordan Henkel, Avrilia Floratou, Joyce Cahoon, Shaleen Deep, and Jignesh M. Patel. 2024g. Reactable: Enhancing react for table question answering. *Proc. VLDB Endow.*, 17(8):1981–1994.

Zhehao Zhang, Yan Gao, and Jian-Guang Lou. 2024h. $e^5$: Zero-shot hierarchical table analysis using augmented LLMs via explain, extract, execute, exhibit and extrapolate. In *Proceedings of the 2024 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies (Volume 1: Long Papers)*, pages 1244–1258, Mexico City, Mexico. Association for Computational Linguistics.

Zhehao Zhang, Xitao Li, Yan Gao, and Jian-Guang Lou. 2023. CRT-QA: A dataset of complex reasoning question answering over tabular data. In *Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing*, pages 2131–2153, Singapore. Association for Computational Linguistics.

Bowen Zhao, Tianhao Cheng, Yuejie Zhang, Ying Cheng, Rui Feng, and Xiaobo Zhang. 2024a. Ct2c-qa: Multimodal question answering over chinese text, table and chart. In *Proceedings of the 32nd ACM International Conference on Multimedia*, MM '24, page 3897–3906, New York, NY, USA. Association for Computing Machinery.

Bowen Zhao, Changkai Ji, Yuejie Zhang, Wen He, Yingwen Wang, Qing Wang, Rui Feng, and Xiaobo Zhang. 2023a. Large language models are complex table parsers. In *Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing*, pages 14786–14802, Singapore. Association for Computational Linguistics.

Weichao Zhao, Hao Feng, Qi Liu, Jingqun Tang, Shubo Wei, Binghong Wu, Lei Liao, Yongjie Ye, Hao Liu, Houqiang Li, and Can Huang. 2024b. Tabpedia: Towards comprehensive visual table understanding with concept synergy. *ArXiv*, abs/2406.01326.

Yilun Zhao, Lyuhao Chen, Arman Cohan, and Chen Zhao. 2024c. TaPERA: Enhancing faithfulness and interpretability in long-form table QA by content planning and execution-based reasoning. In *Proceedings of the 62nd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 12824–12840, Bangkok, Thailand. Association for Computational Linguistics.

Yilun Zhao, Yunxiang Li, Chenying Li, and Rui Zhang. 2022. MultiHiertt: Numerical reasoning over multi hierarchical tabular and textual data. In *Proceedings of the 60th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 6588–6600, Dublin, Ireland. Association for Computational Linguistics.

Yilun Zhao, Chen Zhao, Linyong Nan, Zhenting Qi, Wenlin Zhang, Xiangru Tang, Boyu Mi, and Dragomir Radev. 2023b. RobuT: A systematic study

of table QA robustness against human-annotated adversarial perturbations. In *Proceedings of the 61st Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 6064–6081, Toronto, Canada. Association for Computational Linguistics.

Mingyu Zheng, Xinwei Feng, Qingyi Si, Qiaoqiao She, Zheng Lin, Wenbin Jiang, and Weiping Wang. 2024. Multimodal table understanding. In *Proceedings of the 62nd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 9102–9124, Bangkok, Thailand. Association for Computational Linguistics.

Mingyu Zheng, Zhifan Feng, Jia Wang, Lanrui Wang, Zheng Lin, Hao Yang, and Weiping Wang. 2025. TableDreamer: Progressive and weakness-guided data synthesis from scratch for table instruction tuning. In *Findings of the Association for Computational Linguistics: ACL 2025*, pages 7290–7315, Vienna, Austria. Association for Computational Linguistics.

Mingyu Zheng, Yang Hao, Wenbin Jiang, Zheng Lin, Yajuan Lyu, QiaoQiao She, and Weiping Wang. 2023. IM-TQA: A Chinese table question answering dataset with implicit and multi-type table structures. In *Proceedings of the 61st Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 5074–5094, Toronto, Canada. Association for Computational Linguistics.

Victor Zhong, Caiming Xiong, and Richard Socher. 2017. Seq2sql: Generating structured queries from natural language using reinforcement learning. *ArXiv*, abs/1709.00103.

Bangbang Zhou, Zuan Gao, Zixiao Wang, Boqiang Zhang, Yuxin Wang, Zhineng Chen, and Hongtao Xie. 2025a. SynTab-LLaVA: Enhancing Multimodal Table Understanding with Decoupled Synthesis . In *2025 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 24796–24806, Los Alamitos, CA, USA. IEEE Computer Society.

Bangbang Zhou, Zuan Gao, Zixiao Wang, Boqiang Zhang, Yuxin Wang, Zhineng Chen, and Hongtao Xie. 2025b. Syntab-llava: Enhancing multimodal table understanding with decoupled synthesis. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 24796–24806.

Wei Zhou, Mohsen Mesgar, Heike Adel, and Annemarie Friedrich. 2024a. FREB-TQA: A fine-grained robustness evaluation benchmark for table question answering. In *Proceedings of the 2024 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies (Volume 1: Long Papers)*, pages 2479–2497, Mexico City, Mexico. Association for Computational Linguistics.

Wei Zhou, Mohsen Mesgar, Heike Adel, and Annemarie Friedrich. 2025c. Ppt: A process-based preference learning framework for self improving table question answering models. *ArXiv*, abs/2505.17565.

Wei Zhou, Mohsen Mesgar, Heike Adel, and Annemarie Friedrich. 2025d. Ritt: A retrieval-assisted framework with image and text table representations for table question answering. *Proceedings of the 4th Table Representation Learning Workshop*.

Wei Zhou, Mohsen Mesgar, Heike Adel, and Annemarie Friedrich. 2025e. Texts or images? a fine-grained analysis on the effectiveness of input representations and models for table question answering. In *Findings of the Association for Computational Linguistics: ACL 2025*, pages 2307–2318, Vienna, Austria. Association for Computational Linguistics.

Wei Zhou, Mohsen Mesgar, Annemarie Friedrich, and Heike Adel. 2024b. Efficient multi-agent collaboration with tool use for online planning in complex table question answering. In *North American Chapter of the Association for Computational Linguistics*.

Wei Zhou, Mohsen Mesgar, Annemarie Friedrich, and Heike Adel. 2025f. G-MACT at SemEval-2025 task 8: Exploring planning and tool use in question answering over tabular data. In *Proceedings of the 19th International Workshop on Semantic Evaluation (SemEval-2025)*, pages 726–742, Vienna, Austria. Association for Computational Linguistics.

Fengbin Zhu, Wenqiang Lei, Youcheng Huang, Chao Wang, Shuo Zhang, Jiancheng Lv, Fuli Feng, and Tat-Seng Chua. 2021. TAT-QA: A question answering benchmark on a hybrid of tabular and textual content in finance. In *Proceedings of the 59th Annual Meeting of the Association for Computational Linguistics and the 11th International Joint Conference on Natural Language Processing (Volume 1: Long Papers)*, pages 3277–3287, Online. Association for Computational Linguistics.

Fengbin Zhu, Ziyang Liu, Fuli Feng, Chao Wang, Moxin Li, and Tat Seng Chua. 2024. Tat-llm: A specialized language model for discrete reasoning over financial tabular and textual data. In *Proceedings of the 5th ACM International Conference on AI in Finance*, ICAIF '24, page 310–318, New York, NY, USA. Association for Computing Machinery.

Junnan Zhu, Jingyi Wang, Bohan Yu, Xiaoyu Wu, Junbo Li, Lei Wang, and Nan Xu. 2025. Tableeval: A real-world benchmark for complex, multilingual, and multi-structured table question answering. *ArXiv*, abs/2506.03949.

Caleb Ziems, William Held, Omar Shaikh, Jiaao Chen, Zhehao Zhang, and Diyi Yang. 2024. Can large language models transform computational social science? *Computational Linguistics*, 50(1):237–291.

Jiaru Zou, Dongqi Fu, Sirui Chen, Xinrui He, Li, Yada Zhu, Jiawei Han, and Jingrui He. 2025. Gtr: Graph-table-rag for cross-table question answering. *ArXiv*, abs/2504.01346.

# A Appendix

## A.1 Survey Scope

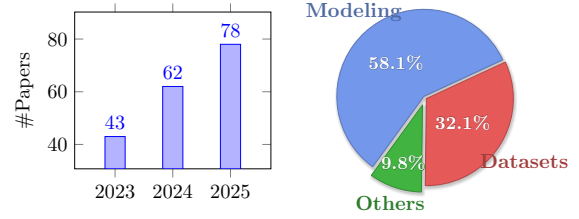We clarify the scope of this survey and describe the process used to collect the papers reviewed.

**Tasks.** This survey primarily focuses on TQA. We also include table fact verification, as this task can be readily reformulated into TQA. For example, by appending a question such as "Is the statement true or false?" to the statement being validated. Another related task is text-to-SQL, for which we mainly discuss relevant datasets, since they can also serve as benchmarks for TQA. We further consider SQL generation as an approach to solving TQA, but we do not provide a detailed review of methods specifically targeting SQL generation. Tasks that are not explicitly covered in this survey include table prediction, table generation, and table summarization, as these differ from TQA in terms of problem formulation and objectives.

**Models.** Given our focus on TQA in the era of LLMs, we primarily include work that leverages these models. For modeling papers, we restrict our collection to works published after 2022, as earlier studies are comprehensively reviewed in the survey by Jin et al. (2022).

**Paper Collection.** We search for papers using the keywords table/tabular reasoning, table/tabular question answering, and table/tabular understanding on both arXiv and the ACL Anthology. We include works published up to August 1st, and expand the collection by identifying relevant papers cited in the related work sections of the retrieved publications. In total, we compile a corpus of 215 papers. Figure 4 presents the distribution of papers by year and theme, illustrating a clear upward trend in research on table question answering.

## A.2 Comparing with Other Surveys

We compare our survey with recent work on table question answering (TQA) in Table 1. *TQA Setups* indicates whether a given survey proposes a TQA-specific task taxonomy or discusses multiple task setups within TQA. *Input Modality* specifies the types of input modalities considered in TQA. *Lan* denotes the languages of the benchmarks covered in the survey, and *Eval* indicates whether evaluation methodologies are discussed. *RLVR* and *ITPT* refer to reinforcement learning with verifiable rewards and interpretability, respectively. As shown



(a) #collected paper per year.    (b) Paper types distribution.

Figure 4: Statistics of the collected paper. We show the number of collected paper by year as well as the distribution of different types of paper.



Figure 5: Performance of (M)LLMs in textual and image-based table understanding. FT-Model denotes fine-tuned models, specifically TableLlaVA-7B (Zheng et al., 2024) and TableLlaMA-7B (Zhang et al., 2024b). OCR refers to configurations in which image tables are first converted to text via optical character recognition (OCR) and then processed using TableLlaMA-7B.

in Table 1, our survey differs from prior work in the following ways: (1) We provide a fine-grained discussion of diverse TQA task setups and the modalities involved. (2) We include benchmarks covering languages beyond English. (3) We present an up-to-date review of modeling approaches with (M)LLMs and evaluation paradigms. (4) We discuss recent advances and emerging themes in the era of LLMs.

## A.3 Method Comparison

Figure 5 compares the performance of (multi-modal) large language models ((M)LLMs)) in textual and image-based table understanding. The results are drawn from Deng et al. (2024b) and Zheng et al. (2024).

## A.4 TQA Datasets

Table 2 shows existing TQA datasets, categorized by features introduced in Section 2. Licenses are

| Survey | Publication Year | TQA Setups | Input Modality | Lan | Modeling with LLMs | Eval | RLVR | ITPT | Summary |
|---|---|---|---|---|---|---|---|---|---|
| Dong et al. (2022) | 2022 | ✗ | text | en | ✗ | ✗ | ✗ | ✗ | table-pretraining |
| Jin et al. (2022) | 2022 | ✓ | text | en | ✗ | ✗ | ✗ | ✗ | benchmarks & pre-LLM modeling |
| Zhang et al. (2024f) | 2024 | ✗ | text | en | ✓ | ✗ | ✗ | ✗ | LLM-based modeling |
| Fang et al. (2024a) | 2024 | ✗ | text | en | ✓ | p | ✗ | ✓ | table prediction, generation and understanding |
| Lu et al. (2024a) | 2024 | ✗ | text, image | en | ✓ | p,r | ✗ | ✗ | tasks discussed via a data lifecycle aspect |
| Wu et al. (2025d) | 2025 | ✗ | text, image | en | ✗ | p,r | ✗ | ✗ | table representations and tasks |
| Tian et al. (2025a) | 2025 | ✗ | text, image | en | ✓ | p,r | ✗ | ✗ | LLM agents |
| Ren et al. (2025) | 2025 | ✗ | text | en | ✓ | ✗ | ✗ | ✗ | general table modeling with deep learning |
| Ours | 2025 | ✓ | text, image, KB | various | ✓ | p,r,e | ✓ | ✓ | fine-grained TQA task taxonomy, up-to-date modeling and discussion |

Table 1: Comparing this work with recent surveys pertinent to table question answering from various perspectives. *TQA Setups*: if a fine-grained TQA task taxonomy is given. *Lan*: language of sourced benchmarks. *Eval*: Evaluation discussions. p,r,e stands for performance, robustness and explanation, respectively. *RLVR*: reinforcement learning with verifiable reward. *ITPT*: interpretability.

subject to each individual dataset. Please refer to
the original datasets for more information.

Table 2: Existing TQA Datasets. *MC* denotes multiple-choice. *ret* and *rea* refer to retrieval and reasoning, respectively. *LAN* indicates the language of the dataset. *Q* and *T* represent question and table, respectively. The value *both* in the table suggests that the dataset contains both single and multiple tables, or includes both flat and hierarchical tables. *DB* stands for database.

| Datasets | Source/Domain | Open domain | Answer format | Reasoning | LAN | Size #Q | Size #T | table features Single | table features Flat | table features Format | additional inputs text | additional inputs image | additional inputs chart | additional inputs KB |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| WTQ (Pasupat and Liang, 2015) | Wikipedia | ✗ | spans | ret, rea | EN | 22k | 2.1k | ✓ | ✓ | text | ✗ | ✗ | ✗ | ✗ |
| SQA (Iyyer et al., 2017) | WTQ | ✗ | spans | ret, rea | EN | 17k | | ✓ | ✓ | text | ✗ | ✗ | ✗ | ✗ |
| TabMCQ (Jauhar et al., 2016) | science exam | ✗ | MC | | EN | 9k | 68 | ✓ | ✓ | text | ✗ | ✗ | ✗ | ✗ |
| TabFact (Chen et al., 2019) | Wikipedia | ✗ | MC | ret, rea | EN | 118k | 16k | ✓ | ✓ | text | ✗ | ✗ | ✗ | ✗ |
| WikiSQL (Zhong et al., 2017) | Wikipedia | ✗ | SQL | ret, rea | EN | 80k | 26k | ✓ | ✓ | text | ✗ | ✗ | ✗ | ✗ |
| Spider (Yu et al., 2018) | WikiSQL, Internet | ✗ | SQL | ret, rea | EN | 10k | 200 | both | ✓ | DB | ✗ | ✗ | ✗ | ✗ |
| INFOTABS (Gupta et al., 2020) | Wikipedia | ✗ | MC | ret, rea | EN | 23k | 2.5k | ✓ | ✓ | text | ✗ | ✗ | ✗ | ✗ |
| KaggleDBQA (Lee et al., 2021) | Kaggle | ✗ | SQL | ret,rea | EN | 272 | 8 | both | ✓ | DB | ✗ | ✗ | ✗ | ✗ |
| HiTab (Cheng et al., 2022a) | statistical report | ✗ | spans | ret, rea | EN | 10k | 3.6k | ✓ | ✗ | text | ✗ | ✗ | ✗ | ✗ |
| AIT-QA (Katsis et al., 2022) | airline | ✗ | spans | ret | EN | 515 | 116 | ✓ | both | text | ✗ | ✗ | ✗ | ✗ |
| FeTaQA (Nan et al., 2022) | ToTTo | ✗ | free-form | ret, rea | EN | 10k | 10k | ✓ | ✓ | text | ✗ | ✗ | ✗ | ✗ |
| TabMWP (Lu et al., 2022) | websites | ✗ | MC, spans | numerical | EN | 38k | 37k | ✓ | ✓ | text, image | ✗ | ✗ | ✗ | ✗ |
| XINFOTABS (Minhas et al., 2022) | INFOTABS | ✗ | MC | ret, rea | MLT | 23k | 2.5k | ✓ | ✓ | text | ✗ | ✗ | ✗ | ✗ |
| EI-INFOTABS (Agarwal et al., 2022) | INFOTABS | ✗ | MC | ret, rea | HI | 23k | 2.5k | ✓ | ✓ | text | ✗ | ✗ | ✗ | ✗ |
| KorWikiTQ (Jun et al., 2022) | Wikipedia | ✗ | spans | ret, rea | KO | 70k | | ✓ | ✓ | text | ✗ | ✗ | ✗ | ✗ |
| TempTabTQA (Gupta et al., 2023) | Wikipedia | ✗ | spans | temporal | EN | 11k | 1.2k | ✓ | ✓ | text | ✗ | ✗ | ✗ | ✗ |
| RobuT (Zhao et al., 2023b) | WTQ, SQA, WikiSQL | ✗ | spans | robustness | EN | 143k | | ✓ | ✓ | text | ✗ | ✗ | ✗ | ✗ |
| IM-TQA (Zheng et al., 2023) | reports | ✗ | spans | ret,rea | ZH | 5k | 1.2k | ✓ | ✗ | text | ✗ | ✗ | ✗ | ✗ |
| Text2Analysis (He et al., 2023) | | ✗ | code | ret,rea | EN | 2.2k | 347 | ✓ | ✓ | text | ✗ | ✗ | ✗ | ✗ |
| SciTab (Lu et al., 2023) | SciGen | ✗ | MC | ret,rea | EN | 1.2k | | ✓ | ✓ | text | ✗ | ✗ | ✗ | ✗ |
| CRT (Zhang et al., 2023) | TabFact | ✗ | spans | ret,rea | EN | 728 | 423 | ✓ | ✓ | text | ✗ | ✗ | ✗ | ✗ |
| Tab-CQA (Liu et al., 2023a) | reports | ✗ | spans | ret | ZH | 10k | 7k | ✓ | | text | ✗ | ✗ | ✗ | ✗ |
| BIRD (Li et al., 2023) | internet | ✗ | SQL | ret,rea | EN | 12.7k | 95 | both | ✓ | DB | ✗ | ✗ | ✗ | ✗ |
| LF-TQA (Wang et al., 2024b) | FeTAQA, QTSUMM | ✗ | free-form | ret,rea | EN | 2.9k | | ✓ | ✓ | text | ✗ | ✗ | ✗ | ✗ |
| Indic-TQA (Pal et al., 2024) | Wikipedia | ✗ | spans | ret,rea | BN,Hi | 2m | 21k | ✓ | ✓ | text | ✗ | ✗ | ✗ | ✗ |
| TableBench (Wu et al., 2024a) | existing datasets | ✗ | spans, free-form | ret,rea | EN | 20k | 3.6k | ✓ | ✓ | text | ✗ | ✗ | ✗ | ✗ |
| Databench (Osés Grijalba et al., 2024) | internet | ✗ | spans | ret,rea | EN | 1.8k | 65 | ✓ | ✓ | DB | ✗ | ✗ | ✗ | ✗ |
| DACO (Wu et al., 2024b) | Spider, Kaggle | ✗ | free-form | ret,rea | EN | 1.9k | 440 | ✗ | ✓ | DB | ✗ | ✗ | ✗ | ✗ |
| TabIS (Pang et al., 2024) | ToTTo HiTab | ✗ | MC | ret,rea | EN | 61k | | ✓ | ✓ | DB | ✗ | ✗ | ✗ | ✗ |
| FREB-TQA (Zhou et al., 2024a) | existing datasets | ✗ | spans | robustness | EN | 8.5k | | ✓ | ✓ | DB | ✗ | ✗ | ✗ | ✗ |
| RealTabBench (Su et al., 2024) | existing datasets | ✗ | free-form | robustness | EN,ZH | 6k | 360 | ✓ | both | csv | ✗ | ✗ | ✗ | ✗ |

| Datasets | Source/Domain | Open domain | Answer format | Reasoning | LAN | Size | | table features | | | additional inputs | | | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| | | | | | | #Q | #T | Single | Flat | Format | text | image | chart | KB |
| TabularGSM (Tian et al., 2025b) | GSM8K | ✗ | number | numerical | EN | 3.5k | | ✓ | ✓ | text | ✗ | ✗ | ✗ | ✗ |
| NIAT (Wang et al., 2025a) | WTQ, HiTab AIT-QA | ✗ | spans | ret | EN | 287k | 750 | ✓ | both | text | ✗ | ✗ | ✗ | ✗ |
| HCT-QA (Ahmad et al., 2025) | internet | ✗ | spans | ret,rea | EN,AR | 77k | 6.7k | ✓ | ✗ | text image | ✗ | ✗ | ✗ | ✗ |
| TableEval (Zhu et al., 2025) | reports | ✗ | spans | ret,rea | EN,ZH | 2.3k | 617 | ✓ | both | csv | ✗ | ✗ | ✗ | ✗ |
| SciAtomicBench (Zhang et al., 2025h) | PubTables MatSciTable | ✗ | MC | ret,rea | EN | 2.5k | | ✓ | both | text | ✗ | ✗ | ✗ | ✗ |
| AraTable (Alshaikh et al., 2025) | Internet | ✗ | spans | ret,rea | AR | 615 | 41 | ✓ | ✓ | text | ✗ | ✗ | ✗ | ✗ |
| RealHiTBench (Wu et al., 2025b) | online platforms | ✗ | free-form | ret,rea | EN | 3.7k | 708 | both | ✗ | text image | ✗ | ✗ | ✗ | ✗ |
| TReB (Li et al., 2025a) | existing datasets | ✗ | free-form | ret,rea | EN,ZH | 7.8k | | ✓ | both | text | ✗ | ✗ | ✗ | ✗ |
| M3TQA (Shu et al., 2025) | reports | ✗ | spans | ret,rea | MLT | 46k | 50 | ✓ | both | text | ✗ | ✗ | ✗ | ✗ |
| TableDreamer (Zheng et al., 2025) | synthesized | ✗ | free-form | ret,rea | EN | 27k | | ✓ | ✓ | text | ✗ | ✗ | ✗ | ✗ |
| MMQA (Wu et al., 2025a) | Spider | ✗ | spans | ret,rea | EN | 3.3k | 3.3k | ✗ | ✓ | DB | ✗ | ✗ | ✗ | ✗ |
| MultiTableQA (Zou et al., 2025) | existing datasets | ✗ | spans | ret,rea | EN | 23k | 57k | ✗ | ✓ | text | ✓ | ✗ | ✗ | ✗ |
| TRANSIENTTABLES (Shankarampeta et al., 2025) | Wikipedia | ✗ | spans | temporal | EN | 3.9k | 14k | ✗ | ✓ | text | ✗ | ✗ | ✗ | ✗ |
| MiMoTable (Li et al., 2025d) | internet | ✗ | free-form | ret,rea | EN,ZH | 1.7k | 428 | both | both | csv | ✗ | ✗ | ✗ | ✗ |
| TQA-Bench (Qiu et al., 2024) | world-Bank BIRD | ✗ | MC | ret,rea | EN | 56k | 10 | ✗ | ✓ | DB | ✗ | ✗ | ✗ | ✗ |
| ComTQA (Zhao et al., 2024b) | FinTabNet PubTab1M D | ✗ | spans | ret,rea | EN | 9k | 1.5k | ✓ | both | image | ✗ | ✗ | ✗ | ✗ |
| MMTab (Zheng et al., 2024) | existing datasets | ✗ | spans | ret,rea | EN | 49k | 23k | ✓ | both | image | ✗ | ✗ | ✗ | ✗ |
| MTabVQA (Singh et al., 2025) | existing datasets | ✗ | spans | ret,rea | EN | 3.7k | 8.5k | ✗ | ✓ | image | ✗ | ✗ | ✗ | ✗ |
| MMSci (Yang et al., 2025a) | SciGen | ✗ | spans | ret,rea | EN | 15k | 52k | ✓ | ✓ | image | ✗ | ✗ | ✗ | ✗ |
| TabComp (Gautam et al., 2025) | DocVQA | ✗ | free-form | ret,rea | EN | 30k | 10k | ✓ | both | image | ✗ | ✗ | ✗ | ✗ |
| TableVQA-Bench (Kim et al., 2024) | WTQ, TabFact FinTabNet | ✗ | spans | ret,rea | EN | 1.5k | 894k | ✓ | ✓ | image | ✗ | ✗ | ✗ | ✗ |
| OTT-QA (Chen et al., 2020a) | Wikipedia | ✓ | spans | ret,rea | EN | 45k | | | ✓ | text | ✓ | ✗ | ✗ | ✗ |
| TANQ (Akhtar et al., 2024) | QAMPARI Wikipedia | ✓ | table | ret,rea | EN | 43k | | | ✓ | text | ✓ | ✗ | ✗ | ✗ |
| RETQA (Wang et al., 2024d) | QAMPARI Wikipedia | ✓ | free-form | ret,rea | ZH | 20k | 4.9k | | ✓ | DB | ✓ | ✗ | ✗ | ✗ |
| Open-WikiTable (Kweon et al., 2023) | WTQ WikiSQL | ✓ | spans | ret,rea | EN | 67k | 24k | | ✓ | text | ✗ | ✗ | ✗ | ✗ |
| CompMix (Christmann et al., 2023) | CONVMIX | ✓ | spans | ret,rea | EN | 9.4k | | ✓ | ✓ | text | ✓ | ✗ | ✗ | ✓ |
| $T^2$-RAGBench (Strich et al., 2025b) | existing datasets | ✓ | spans | ret,rea | EN | 32k | | | | text | ✓ | ✗ | ✗ | ✗ |
| MMCoQA (Li et al., 2022) | MMQA | ✓ | spans | ret,rea | EN | 5.7k | 10k | | | text | ✓ | ✓ | ✗ | ✗ |
| MMTBENCH (Titiya et al., 2025) | Internet | ✗ | spans | ret,rea | EN | 4k | 500 | ✓ | both | csv image | ✗ | ✓ | ✓ | ✗ |
| KET-QA (Hu et al., 2024) | HybridTQA | ✗ | spans | ret,rea | EN | 9.4k | 5.7k | ✓ | ✓ | text | ✗ | ✗ | ✗ | ✓ |
| CT2C-QA (Zhao et al., 2024a) | reports | ✗ | spans | ret,rea | ZH | 9.9k | 369 | | | html | ✓ | ✗ | ✓ | ✗ |
| mmtabqa (Mathur et al., 2024) | existing datasets | ✗ | spans | ret,rea | EN | 69k | 259 | ✓ | ✓ | text | ✗ | ✓ | ✗ | ✗ |
| StructFact (Huang et al., 2025a) | existing datasets | ✗ | MC | ret,rea | EN | 13k | | ✓ | ✓ | text | ✓ | ✗ | ✗ | ✓ |

| Datasets | Source/Domain | Open domain | Answer format | Reasoning | LAN | Size | | table features | | | additional inputs | | | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| | | | | | | #Q | #T | Single | Flat | Format | text | image | chart | KB |
| WikiMixQA (Foroutan et al., 2025) | WTabHTML | ✗ | MC | ret,rea | EN | 1k | | ✓ | ✓ | text | ✗ | ✗ | ✓ | ✗ |
| SPIQA (Pramanick et al., 2024) | arXiv | ✗ | free-form | ret,rea | EN | 270k | 117k | both | both | image | ✓ | ✓ | ✓ | ✗ |
| MMQA (Talmor et al., 2021) | Wikipedia | ✗ | spans | ret,rea | EN | 30k | | ✓ | ✓ | text | ✓ | ✓ | ✗ | ✗ |
| SciTabQA (Ghosh et al., 2024) | SciGen | ✗ | spans | ret,rea | EN,AR | 822 | 198 | ✓ | ✓ | text | ✓ | ✗ | ✗ | ✗ |
| MULTITAT (Zhang et al., 2025f) | existing datasets | ✗ | spans | ret,rea | MLT | 250 | | ✓ | ✓ | text | ✓ | ✗ | ✗ | ✗ |
| SciTAT (Zhang et al., 2025e) | arXiv | ✗ | spans free-form | ret,rea | EN | 13k | | ✓ | ✓ | text | ✓ | ✗ | ✗ | ✗ |
| PACIFIC (Deng et al., 2022) | TATQA | ✗ | spans | ret,rea | EN | 19k | | ✓ | ✓ | text | ✓ | ✗ | ✗ | ✗ |
| TAT-QA (Zhu et al., 2021) | reports | ✗ | spans | ret,rea | EN | 16k | | ✓ | ✓ | text | ✓ | ✗ | ✗ | ✗ |
| GeoTSQA (Li et al., 2021) | exams | ✗ | MC | ret,rea | ZH | 1k | | ✗ | | text | ✓ | ✗ | ✗ | ✗ |
| HybridQA (Chen et al., 2020b) | Wikipedia | ✗ | spans | ret | EN | 70k | 13k | ✓ | ✓ | text | ✓ | ✗ | ✗ | ✗ |
| FinQA (Chen et al., 2021) | reports | ✗ | spans | ret,rea | EN | 8k | | ✓ | ✓ | text | ✓ | ✗ | ✗ | ✗ |
| MultiHiertt (Zhao et al., 2022) | FinTabNet | ✗ | spans | ret,rea | EN | 10k | | ✗ | ✗ | text | ✓ | ✗ | ✗ | ✗ |