

# Tabular Data Understanding with LLMs: A Survey of Recent Advances and Challenges

Anonymous ACL submission

## Abstract

Tables have gained significant attention in large language models (LLMs) and multimodal large language models (MLLMs) due to their complex and flexible structure. Unlike linear text inputs, tables are two-dimensional, encompassing formats that range from well-structured database tables to complex, multi-layered spreadsheets, each with different purposes. This diversity in format and purpose has led to the development of specialized methods and tasks, instead of universal approaches—making navigation of table understanding tasks challenging. To address these challenges, this paper introduces key concepts through a taxonomy of tabular input representations and an introduction of table understanding tasks. We highlight several critical gaps in the field that indicate the need for further research: (1) the predominance of retrieval-focused tasks that require minimal reasoning beyond mathematical and logical operations; (2) significant challenges faced by models when processing complex table structures, large-scale tables, length context, or multi-table scenarios; and (3) the limited generalization of models across different tabular representations and formats.

## 1 Introduction

Tables have garnered increasing attention due to advances in large language models (LLMs) and multi-modal large language models (MLLMs), owing to the unique challenges they present. Unlike linear text, tabular data possess an inherently visual, two-dimensional format that requires specialized pipelines to be processed effectively, as shown in Figure 1. Additionally, tables exhibit structural flexibility, serving a wide range of purposes—from well-structured database tables to hierarchical, multi-layered spreadsheets and multimedia-linked info-boxes. These variations in purpose and structure have driven the development of diverse input representations, tasks, and

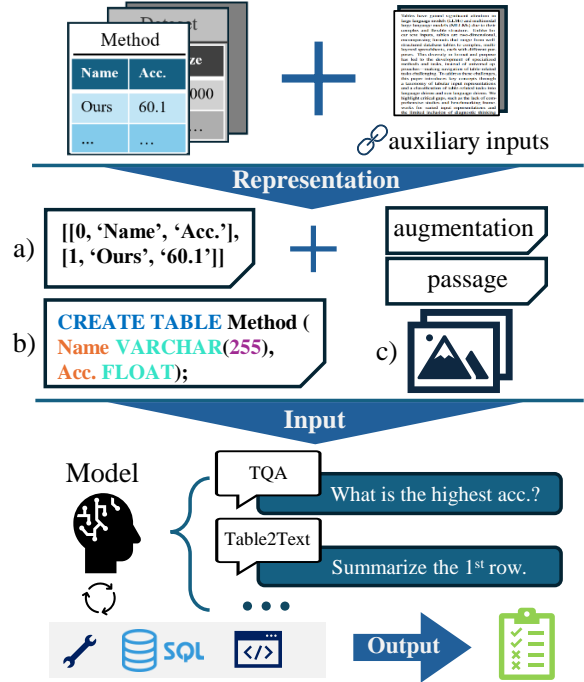


Figure 1: Workflow of table-related tasks in large models. Tables or databases, possibly accompanied by additional input data, are transformed into input representations, which could take the form of (a) serialization, (b) database schema, (c) images, or other format with optional augmentations. These inputs are then processed by models usually leveraging SQL, and other tools to generate task specific outputs.

specialized methods and datasets. However, such specialization often comes at the expense of universality (Zhang et al., 2024a), making it difficult for new researchers to navigate the field effectively. While existing surveys (Fang et al., 2024; Zhang et al., 2024b; Lu et al., 2024; Badaro et al., 2023; Ren et al., 2025) have explored various prompting, training, and transformer-based methods for table processing, there is a need for a comprehensive survey that uncovers new opportunities, focusing on tasks and benchmarks in tabular understanding.

To address the existing gap and assist researchers in navigating table-related tasks, this paper presents a systematic taxonomy of tabular data representations and introduces a broad range of both well-

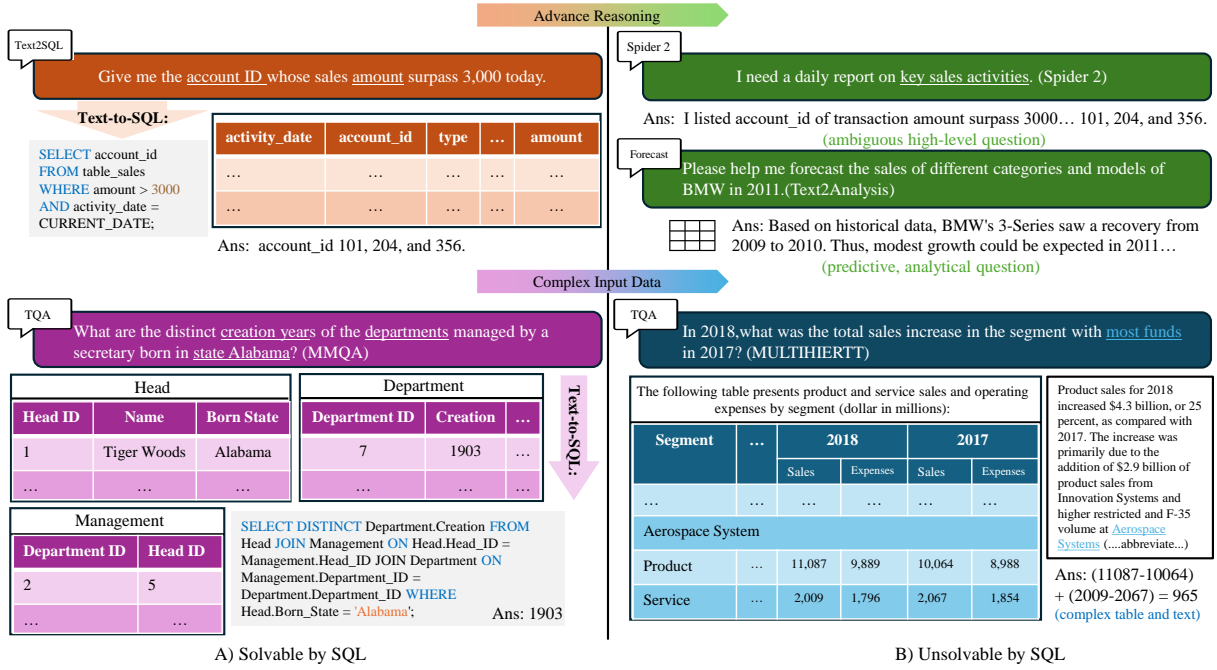


Figure 2: The left side illustrates examples of tasks that can be addressed with SQL-based methods such as typical Text-to-SQL task and a Table QA task from MMQA (Anonymous, 2024). In contrast, the right side presents tasks that demand advanced reasoning or involve complex inputs, such as those found in Spider 2 (Lei et al., 2024), Text2Analysis (He et al., 2024), and MULTIHIERTT (Zhao et al., 2022), which go beyond the capabilities of SQL-based approaches.

established and novel tasks. For instance, we examine *Table QA*, which focuses on answering natural language questions based on table content, and *Table-to-Text*, which involves generating natural language summaries from tabular data. We also highlight innovative tasks such as *leaderboard construction*, which aggregates result tables from scientific papers to provide a comprehensive comparison of methods in one specific field. For well-established tasks, we compile key benchmarks and their associated table formats, categorizing improvements in newer benchmarks relative to earlier ones to highlight emerging research trends.

Furthermore, our survey reveals new opportunities by focusing on tasks and challenges identified in widely used benchmarks. Despite significant progress in prompting and training methods—as highlighted in existing surveys (Lu et al., 2024; Badaro et al., 2023; Ren et al., 2025)—and the robust performance of recent tabular foundational models that integrate tabular data during the pre-training and fine-tuning stages of 72B base models (Su et al., 2024), current table processing benchmarks tend to concentrate on limited reasoning tasks and often rely on simplistic, synthetic tables with inconsistent input representations. While effective for initial evaluations, these benchmarks fall short in assessing the performance of more

advanced methods and models in real-world scenarios that require higher-level reasoning and the processing of complex inputs, ultimately limiting their generalizability and broader applicability.

## 2 Findings and Future Direction

In this section, we outline three key findings that underscore the need for further investigation.

### 2.1 Limited Scope Beyond Mathematical Reasoning

Recent work has begun to saturate performance on many widely used benchmarks. For example, question-decomposition pipelines have yielded significant improvements (Gao et al., 2023; Ye et al., 2023; Wang et al., 2024b); the method proposed by Hussain (2025) achieved over 80% accuracy on the Wiki-Table Questions benchmark (Pasupat and Liang, 2015) and more than 93% on TabFact (Chen et al., 2020b), two popular datasets for table QA and fact verification. Moreover, the success of table foundation models—integrating specialized table encoders into large-scale language models pre-trained and fine-tuned on tabular data (Su et al., 2024)—signals a growing trend toward applying tabular methods to larger models. These advances suggest it is time to move beyond data retrieval-based tasks, as most benchmarks rely on detailed

queries that prompt models to extract specific information from tables using logical operations.

Many existing benchmarks are even constructed by first generating SQL queries or sequences of mathematical expressions, which are then translated into natural language query (Pasupat and Liang, 2015; Iyyer et al., 2017; Pal et al., 2023; Anonymous, 2024), or by framing questions whose answers can be fully derived using mathematical functions (Zheng et al., 2023; Zhang et al., 2023d; Zhao et al., 2022; Kweon et al., 2023). While efforts have focused on enhancing task complexity through additional reasoning steps or embedding complex mathematical functions, the core structure of these tasks remains fundamentally unchanged. As shown in Figure 2, such descriptive questions can be solved relatively easily by text-to-SQL methods when tables are well-structured.

Notably, recent work (Majumder et al., 2024) has further pushed the boundaries by emphasizing higher-order reasoning skills. For example, He et al. (2024) introduced tasks that extend beyond basic descriptive analysis, such as insight identification, similar to what is shown in Figure 3, which demands diagnostic thinking; forecasting, which requires predictive thinking; and chart creation from ambiguous queries, a task that requires prescriptive thinking—selecting the appropriate chart type and determining optimal intervals to produce visually appealing figures. In these tasks, models cannot simply rely on finding synonyms or related attributes in the table to perform data retrieval. Instead, they must understand the overall context of the table and the user’s intent to address the query.

A similar direction is explored by Spider 2 (Lei et al., 2024), which introduces questions requiring higher levels of reasoning. Unlike benchmarks such as Spider (Yu et al., 2018) and its extensions, which introduce marginal difficulties by swapping explicit schema names with synonyms or rephrasing utterances (Deng et al., 2021; Gan et al., 2021a), Spider 2 presents high-level, intent-driven queries, as illustrated in Figure 2. For example, instead of asking explicitly (e.g., “Give me the account ID whose sales surpass a threshold today”), Spider 2 poses abstract, goal-oriented queries (e.g., “I need a daily report on key sales activities”). These queries challenge models to infer the user’s intent, requiring a deep understanding of both the database schema and the query’s broader context. Furthermore, Dong et al. (2025) introduce multi-turn conversations that teach models to seek clarification

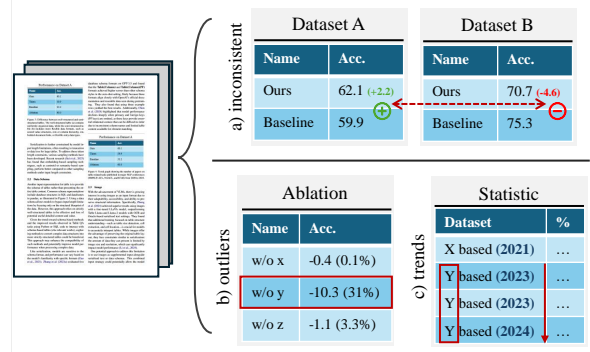


Figure 3: Illustration of the proposed task: Scientific Document Understanding with Tables which require diagnosing implicit knowledge embedded in tabular data, which may not be well addressed in text. Examples include: a) inconsistent results under conditions; b) outliers in values; and c) key trends.

tion whenever a user’s initial query is ambiguous, thereby better mirroring real-world interactions and mitigating the multiple-interpretation issue identified by Pourreza and Rafiei (2023b).

## 2.2 Lack of Robustness on Input Complexity

Another area of opportunity in current table-related research is enhancing model robustness when processing complex input scenarios, including intricate table structures, long tables, lengthy texts, and multi-table contexts—challenges that have minimal impact on human performance (Anonymous, 2024; Pal et al., 2023). Benchmarks such as HiTab (Cheng et al., 2022) and MULTIHIERTT (Zhao et al., 2022) have been instrumental in highlighting these challenges. HiTab features hierarchical multidimensional tables, while MULTIHIERTT further incorporates lengthy texts where answers may be embedded, as well as multi-table scenarios. Both benchmarks report model performances below 50%, compared to a human accuracy of around 83% on MULTIHIERTT. Similarly, benchmarks like MultiTableQA (Pal et al., 2023) and MMQA (Anonymous, 2024), which focus on multi-table question answering from well-structured databases such as those in the Spider benchmark, provide valuable insights into current model limitations. For instance, in MMQA the strongest model evaluated, o1-preview (OpenAI, 2024), achieves an exact match score slightly above 50%, while human performance reaches approximately 89%.

## Scientific Document Understanding with Tables.

Scientific documents provide a rich test bed for information extraction and table extraction (Bai et al., 2024; Yang et al., 2022; Kardas et al., 2020). These papers typically contain complex ablation,

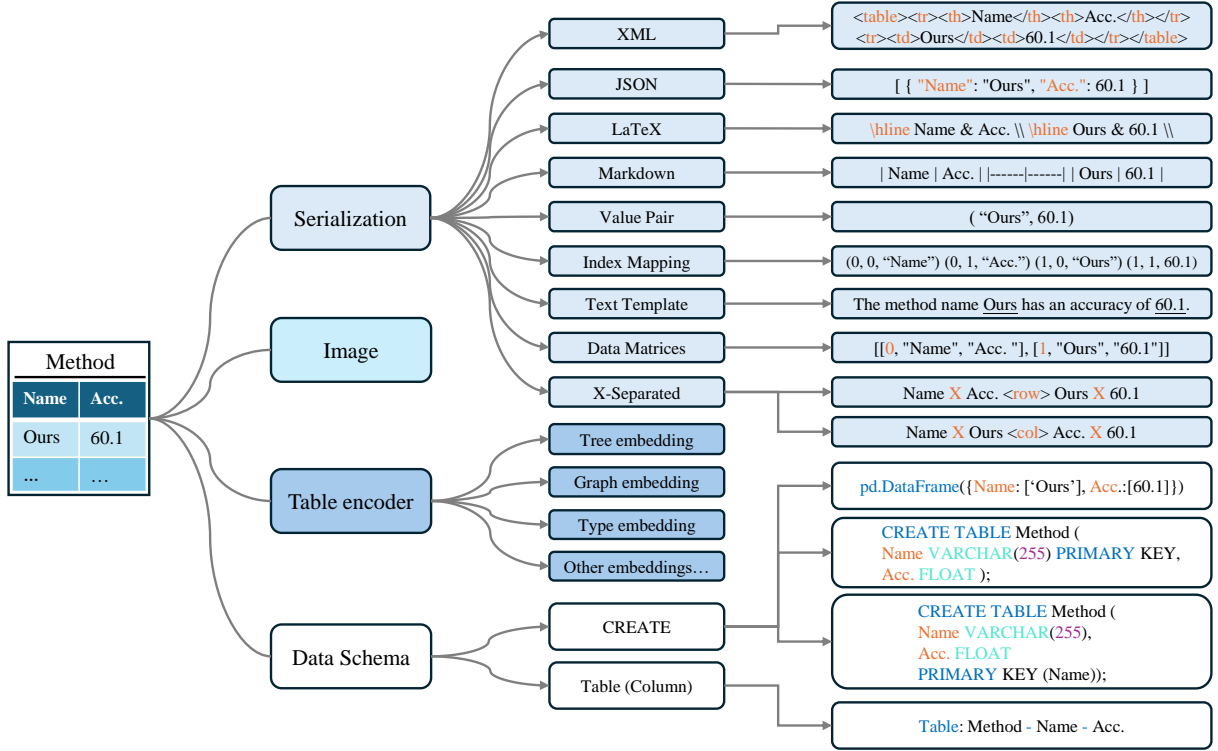


Figure 4: Taxonomy of table input representation methods, encompassing serialization, image, specialized table encoders, and data schema. Examples illustrating each representation type are shown on the right.

analysis, and method-comparison tables alongside extensive textual discussion, all of which demand sophisticated reasoning for accurate interpretation (Zhang et al., 2023c; Asai et al., 2024). Building on this foundation, future work can harness scientific-document data to develop higher-level table-reasoning systems that demand a broad repertoire of skills—such as trend detection, diagnostic assessment, and forecasting (see Figure 3).

### 2.3 Limited Generalization Across Tabular Representations

Despite recent advances, current models still struggle to generalize across diverse tabular representations. Their performance on commonly used benchmarks can vary by up to 5% depending on how closely input formats align with the data encountered during pretraining (Sui et al., 2024), as similarly observed by Gao et al. (2023) in the Text-to-SQL domain. Benchmarks highlight this issue by relying on a variety of input representations chosen based on convenience and accessibility. As demonstrated in our collection of major benchmarks (see Tables 1, 2, and 3), tabular representations for the same type of task lack universality. Even when categorized under the same format, such as JSON, the internal structures can vary greatly (Aly et al., 2021; Chen et al., 2020c), further complicating

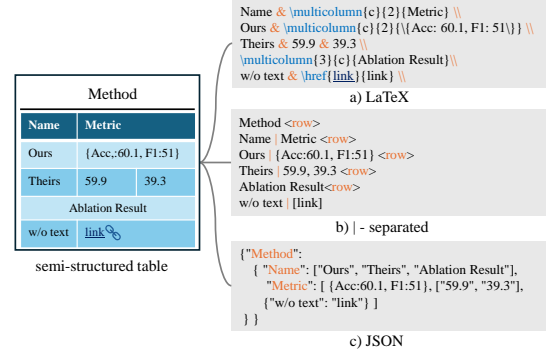


Figure 5: Comparison of serialization methods for semi-structured tables: a) LaTeX, b) X-separated, and c) JSON. Each method has its strengths and weaknesses in handling aspects such as nested value structures, row or column hierarchies, embedded document links, and flexible data types.

performance evaluations and introducing bias.

Efforts to address these inconsistencies are emerging. For example, Lei et al. (2023) provides standardized serialization options such as Markdown and flattened text, though additional formats remain underexplored. Another line of research (Zheng et al., 2024) focuses on visual representations of complex tables—such as Table Cell Locating and Merged Cell Detection—to generate serialized versions from images. Integrating these tasks into fine-tuning pipelines has proven beneficial.

Future research could explore serialization-to-serialization tasks, where models transform one



format (e.g., JSON) into another (e.g., LaTeX or Markdown). Integrating such task could enhance models’ robustness to varied input styles and create opportunities for fine-tuning across multiple representations. Additionally, limited investigation has been conducted into the effectiveness of different representations for complex tables. For instance, LaTeX’s `\multicolumn` command effectively captures hierarchical structures, whereas other formats may ignore this type of relationship during serialization process, as Figure 5 shown.

### 3 Modalities of Table Representation

In this section, we introduce key tabular representations that are essential for enabling large models to process table data effectively. Since these models require one-dimensional input formats, structured, two-dimensional tables must be converted accordingly. This transformation, however, often results in the loss of valuable structural information. To address these challenges, various methods have been developed, including serialization, database schema representations, image-based formats, and specialized table encoders, as illustrated in Figure 4. Recent studies (Sui et al., 2024; Zhang et al., 2023a) demonstrate that model performance is sensitive to the chosen input representation, underscoring the data-dependent nature of current approaches to processing tabular data. Unfortunately, many existing benchmarks rely on representations selected primarily for convenience (Sundararajan et al., 2024), lacking of robust, unbiased comparisons.

#### 3.1 Serialization

Serialization has long been a common method for representing tabular data, transforming tables into serialized text. Its primary advantages lie in compatibility with standard models and ease of access to existing formats, such as HTML or Markdown tables on the web, LaTeX tables in PDF documents, and JSON or key-value pairs in code environments (see Figure 4). Most current benchmarks rely on serialization, as illustrated in Tables 1, 2, and 3. Below, we highlight several noteworthy papers:

**Sensitivity of Input Design.** Models are not only sensitive to different serialization formats, but variations in input design can also cause significant fluctuations in performance across table interaction tasks such as table partitioning, cell lookup, and reverse lookup (Sui et al., 2024). For example, omitting marked partitions or altering the input or-

der has resulted in performance drops of up to 20%, while removing example shots has led to accuracy deteriorations of as much as 50%.

**Sampling and Augmentation.** Long or multi-table inputs pose challenges for serialization due to model input length limitations, often resulting in truncation or data loss. To address these constraints, researchers have developed methods for sampling rows or columns that capture the key information in a table. Recent research (Sui et al., 2023) demonstrates that embedding-based sampling techniques, such as centroid and semantic-based sampling, outperform other approaches. Furthermore, they show a balanced combination of augmentation data (e.g., table sizes and keyword explanations) and sampled table text has proven effective in achieving better overall performance within token limits.

#### 3.2 Data Schema

Another input representation for table is to provide the schema of tables rather than presenting the entire table content. Common schema representations include database structures in SQL and dataframes in pandas, as illustrated in Figure 4. Using a data schema allows models to bypass input length limitations by focusing only on the structural blueprint of the data. However, this approach relies on strictly well-structured tables to be effective and loss of potential useful detailed content and value.

**Sensitivity of Input Design.** Like serialization, models are not only sensitive to the schema format, but also its designs: Zhang et al. (2023a) evaluated schema input designs on GPT-3.5 and found that using three example rows yielded the best results. Additionally, they highlighted that model performance declines sharply when primary and foreign keys (PF keys) in the data schema are omitted, which Chen et al. (2024) also mentioned.

**Normalized structure.** Given the trend toward schema-based methods and the improved results observed in Table QA tasks using Python or SQL code to interact with schema-based tables (Wang et al., 2024b; Pourreza and Rafiei, 2023a; Ye et al., 2023), exploring methods to convert complex data structures into more structured tables could be beneficial to enhance the compatibility of such methods.

#### 3.3 Image

With the advancement of MLLMs, there is growing interest in using images as an input format

Benchmark	Sources / Domain	# Q	# T	Passage	Table Format	Output	Directions
WTQ (2015)	Wikipedia	22,033	2,108		HTML	cells	-
SQA (2017)	Wikipedia	17,553	6,066		HTML	cells	Input Complexity
HybridQA (2020c)	Wikipedia	69,611	13,000	✓	JSON	text-span	Input Complexity
FetaQA (2021a)	Wikipedia	-	10,330		Data Matrices	free-form	Answer Format
TAT-QA (2021)	Financial Reports	16,552	7,431	✓	Data Matrices	number	Domain, Input
OTT-QA (2021)	Wikipedia	-	45,841	✓	JSON	text-span	Input, Reasoning
AIT-QA (2022)	Airline Industry	515	113		Data Matrices	cells	Domain, Input
FinQA (2022)	Financial Report	8,281	2,789	✓	Data Matrices	number	Domain Knowledge
MMCoQA (2022)	MMQA (2018)	1,715	10,042	✓	JSON	text-span	Input Complexity
HiTab (2022)	Wikipedia, Statistic	10,672	3,597		Row-Separated	text-span	Input Complexity
MULTIHIERTT (2022)	Financial Report	10,440	2,513	✓	HTML	number	Input, Reasoning
Open-WikiTable (2023)	Wikipedia	67,023	24,680		Row-Separated	text-span, SQL	Answer Format
QTSUMM (2023)	Wikipedia	7,111	2,934		Data Matrices	free-form	Answer Format
TEMPTABQA (2023)	Wikipedia	11,454	1,208		JSON, HTML	text-span	Reasoning Difficulty
CRT-QA (2023d)	TabFact (2020b)	1,000	423		Row-Separated	text-span	Reasoning Difficulty
IM-TQA (2023)	Baidu Encyclopedia	5,000	1,200		Index Mapping	text-span	Input Complexity
TabCQA (2023a)	Financial Report	109,089	7,041		Text Template, Value Pair	text-span	Input Complexity
MultiTabQA (2023)	Spider (2018), Synthetic, TAPEX (2022) Corpus	136,461	-		Row-Separated	sub-table	Answer, Input
TABMWP (2023a)	Online Learning Web	38,431	37,544		Row-Separated, Spreadsheet, Image	free-form	Reasoning Difficulty
FREB-TQA (2024)	WTQ, WikiSQL (2017), SQA, TAT-QA	75,205	8,590		Data Matrices	text-span	Input, Reasoning
Text2Analysis (2024)	Data Analysis Libraries	2,249	347		-	code, text	Reasoning Difficulty
MMQA (2024)	Spider (2018)	3,313	3,312		JSON	sub-table	Input Complexity

Table 1: Summary of benchmarks for Table-based Question Answering. **Sizes** shows the number of questions and tables. **Passage** indicates if an input passage is included. **Directions** categories each benchmark’s primary focus compare to previous ones.

due to their adaptability, accessibility, and ability to preserve structural information (Wydmanski et al., 2024). Specifically, Zheng et al. (2024) achieved superior results using images with a fine-tuned LLaVA model (Liu et al., 2023b), outperforming models with OCR and serialization settings. They found that additional training focused on table structure understanding—such as cell extraction and cell location—enhance the model’s ability to accurately interpret tables.

**Image resolution.** While images offer the advantage of preserving the original table layout, they face constraints similar to serialization: the amount of data they can present is limited by image size and resolution, which can significantly impact model performance (Li et al., 2024). As tables grow larger, the information becomes blurred at a fixed resolution, leading to deteriorated performance. One potential approach is to use images as supplemental input alongside serialized text or data schema (Luo et al., 2023). This combined input strategy could potentially allow the model to receive structural information directly from the image while accessing detailed content from the text-based format. However, to the best of our knowledge, systematic evaluations of this approach remain lacking.

### 3.4 Table Encoder

Specific table encoder designs have been employed in smaller-scale language models to handle table-

related tasks, utilizing various embeddings such as column-based (Iida et al., 2021), row-based (Herzig et al., 2020), tree-structured (Wang et al., 2021c), and graph-based embeddings (Wang et al., 2021a). Building on these approaches, recent work has demonstrated a trend toward employing specialized encoders in larger base models, effectively creating table foundation models (van Breugel and van der Schaar, 2024; Su et al., 2024; Ma et al., 2024). In particular, TableGPT2 leverages a specialized table encoder—with column- and row-wise attention—to integrate tabular data during the pretraining and fine-tuning stages of 7B and 72B base models (Su et al., 2024), outperforming other table generalist models across a range of tasks while remaining competitive with task-specific methods.

## 4 Table-Related Tasks

In this section, we introduce key table-related tasks such as Table Question Answering (TQA), Table-to-Text, and Table Fact Verification (TFV), along with other intriguing applications like leaderboard construction that actively utilize tables.

### 4.1 Table Question Answering

TQA<sup>1</sup> is one of the most common and well-studied table tasks, with various benchmarks developed as shown in Table 1. It typically involves a free-form question and a single table, sometimes accompanied by an optional passage or passage links, and

Benchmark	Sources / Domain	# Q	# T	Table Format	Focus	Text Input	Directions
Rotowire (2017)	NBA	-	4,853	JSON	N/A		Domain Knowledge
ToTTo (2020)	Wikipedia	134,161	83,141	Index Mapping	Highlight Span	Caption	-
Logic2Text (2020d)	WikiTable	10,800	5,600	Row-Separated	N/A		Logic Summarization
LogicNLG (2020a)	TabFact (2020b)	37,000	7,300	Data Matrices	N/A		Logic Comparison
SciGen (2021)	Scientific Paper	53,000	-	Row-Separated	N/A	Caption	Domain Knowledge
NumericNLG (2021)	Scientific Paper	1,300	1,300	JSON	N/A	Caption	Domain Knowledge
FetaQA (2021a)	ToTTo (2020)	-	10,330	Matrices	Text Query		Input Complexity
	E2E (2020), WTQ						
DART (2021b)	WikiTable (2023)	82,191	5,623	XML, JSON	N/A	Table Title	Table Structure
	WebNLG (2019)						
QTSUMM (2023)	Wikipedia	7,111	2,934	Data Matrices	Text Query		Input Complexity
FindSUM (2023c)	Company Report	-	21,125	Data Matrices	N/A	Long Text	Input Complexity

Table 2: Summary of benchmarks for Table-to-Text and Table Summarization. **Focus** specifies the subset of table content intended for natural language generation, while N/A indicates the entire table should be transformed to natural language.

Benchmark	Sources / Domain	# Q	# T	Table Format	Output	Directions
TabFact (2020b)	Wikipedia	117,843	18,000	Row-Separated	S, R	-
InfoTabs (2020)	Wikipedia	23,738	2,540	HTML, JSON	S, R, N	Output Format
FEVEROUS (2021)	Wikipedia	87,062	-	JSON / Mapping	S, R, N	Output Format
SEM-TAB-FACTS (2021b)	Science	5,715	2,961	XML	S, R, N, EC	Domain Knowledge
XInfoTabs (2022)	InfoTabs	23,738	2,540	JSON	S, R, N	Multi-Language
El-InfoTabs (2022)	InfoTabs	23,738	2,540	JSON	S, R, N	Indic-Language
SciTab (2023b)	SciGen(Moosavi et al., 2021)	1,255	-	JSON / Mapping	S, R, N	Domain Knowledge

Table 3: Summary of benchmarks for Table-based Fact Verification. *S* in the output denotes Supported, *R* represents Refuted, *N* stands for Neither or Not Enough Evidence, and *EC* refers to Evidence Cells.

the output is expected to be information derived from the table or passage, generally presented as cell spans, calculated values, or minimal text spans.

TQA benchmarks have expanded significantly over the past two years, inspiring future work across multiple directions, including domain knowledge, answer format, input complexity, and reasoning difficulty. Domain-specific benchmarks now better reflect real-world scenarios in fields such as airlines (Katsis et al., 2022) and finance (Zhu et al., 2021; Chen et al., 2022). Answer formats have also diversified, with benchmarks requiring free-form responses (Nan et al., 2021a; Zhao et al., 2023; Wang et al., 2024a) and SQL queries (Kweon et al., 2023), beyond traditional cell values or text spans. Input complexity has increased through multi-table datasets (Pal et al., 2023; Zhao et al., 2022), hierarchical tables (Cheng et al., 2022), and semi-structured tables (Lu et al., 2023a), which challenge models to navigate intricate structures. Reasoning requirements have similarly intensified, incorporating hypothetical questions (Li et al., 2023b), implicit time-based inference (Gupta et al., 2023), and sequential or conversational queries (Iyyer et al., 2017; Li et al., 2022; Liu et al., 2023a). Overall, recent benchmarks generally demand more complex reasoning steps and operations to yield accurate answers.

<sup>1</sup>For a more comprehensive understanding of TQA, see this curated list of relevant papers: <https://github.com/lfy79001/Awesome-Table-QA>

## 4.2 Table-to-Text and Table Summarization

Table-to-Text and Table Summarization are table tasks initially developed to evaluate whether models could accurately interpret and describe table content. In these tasks, the input typically includes a table, sometimes with specified cell spans as shown in the *Focus* column in Table 2. If a span or region is provided, the model generates a textual description or summary of that specific area; if not, it summarizes the entire table. With advances in models’ table understanding, this task has become less prominent, as the number of related publications has steadily decreased since 2021.

**Query Focused Summarization.** A recent, noteworthy benchmark in this area is QTSUMM (Zhao et al., 2023), which requires models to generate text-based summaries of specific table regions in response to questions. By integrating the aspect of table search based on textual queries from TQA with the descriptive demands of Table-to-Text, QTSUMM introduces new complexities that push models to move beyond simple fact retrieval. Notably, QTSUMM includes “why” questions, prompting models to reason about underlying causes or explanations—a shift that aligns more closely with human interests and highlights the importance of generating responses that incorporate causal understanding and contextual depth.

**Lack of Multilingual Benchmarks.** A notable gap in current research is the absence of multilin-



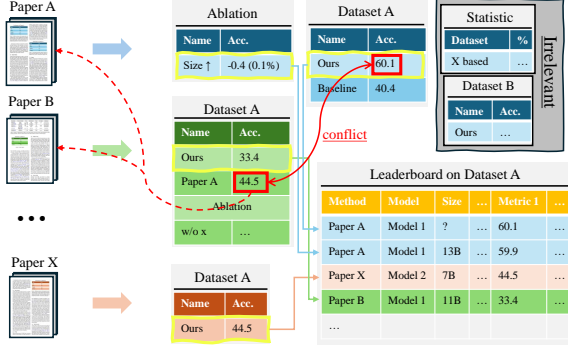


Figure 6: Illustration of automatic leaderboard construction pipeline. Results are extracted from ablation and performance tables in each paper. The red line highlights inconsistency across paper that may require examination across texts.

gual benchmarks for table-to-text tasks. As highlighted in (Osuji et al., 2024), to the best of our knowledge, no table-to-text benchmarks exist in languages other than English, significantly limiting the applicability and inclusivity of this task.

### 4.3 Table Fact Verification

Table Fact Verification (also referred to as Table Reasoning or Table Natural Language Inference) is a task designed to assess fact-searching and logic inference capabilities within tables. In this task, the input typically consists of a statement or claim alongside a reference table. The model’s output is a verification label—such as “Supported,” “Refuted,” or “Not Enough Information”—indicating whether the claim aligns with the table content. Some benchmarks also require a justification for the answer, as shown in Table 3. Recent methods have enabled models to achieve over 80% accuracy on widely used benchmarks like TabFact and FEVEROUS (Sui et al., 2024; Ye et al., 2023; Wang et al., 2024b), demonstrating substantial progress in fact-checking within tabular data. However, scenarios involving longer contexts, multiple tables, or complex table structures remain unassessed.

### 4.4 Leaderboard Construction

Beyond the widely studied tasks, an intriguing direction proposed by Kardas et al. (2020) is leaderboard construction. This task aims to streamline the comparison of experimental results within a research domain through scientific papers, offering a concise and structured view of progress.

Existing methods, such as those proposed in (Kardas et al., 2020; Yang et al., 2022), have made notable strides in automating this process. These approaches typically employ pipelines that classify and extract data from performance and ab-

lation tables in scientific papers, leveraging techniques like Named Entity Recognition (NER) or string matching to form tuples (Task, Dataset, Metric) or quadruples (Task, Dataset, Metric, Score). Such methods provide a foundational framework for building leaderboards and have proven effective in capturing basic performance comparisons across different methods and datasets. However, as scientific tasks and methodologies grow increasingly complex, these pipelines face limitations. Tasks often require varying schemas to account for unique aspects, and surface-level extraction may not fully capture the nuances of more intricate experiments or analyses. For instance, discrepancies in reported results between papers, as illustrated in Figure 6, often necessitate a deeper comparison and reasoning over both tables and textual content to resolve.

### 4.5 Other Tasks

Emerging new table-related tasks include innovations such as tabular synchronization across languages (Khinchin et al., 2023) and column name abbreviation expansion (Zhang et al., 2023b). Among these, Text-to-Table has gained increasing attention in 2024 (Ramu et al., 2024; Jiang et al., 2024; Deng et al., 2024). The task was first formalized by Wu et al. (2022) as a sequence-to-sequence task by inversely applying table-to-text datasets. Recent studies have explored various methods, such as incorporating knowledge graphs (Jiang et al., 2024), to enhance its utility as a data integration task for field like finance, medicine, and law.

## 5 Further Reading

For readers seeking deeper insights into table-related research areas, several survey papers offer valuable perspectives. For methodologies aimed at improving table reasoning with LLMs, work by Zhang et al. (2024b) provides a detailed taxonomy and an analysis of emerging trends. Lu et al. (2024) explores prompting and training techniques for table-related tasks in the context of LLMs and VLMs. Meanwhile, Badaro et al. (2023); Ren et al. (2025) presents a focused analysis of transformer-based, smaller-scale models designed for tabular data. For an in-depth perspective, the comprehensive 30-page survey by Fang et al. (2024) provides an extensive overview of table understanding tasks, datasets, and corresponding fundamental methods.



## 534 Limitations

535 This study presents a comprehensive survey of  
536 table-related tasks with LLMs and MLLMs, high-  
537 lighting key trends and emerging opportunities.  
538 While we have made our best effort to provide  
539 a thorough review, certain limitations remain. Due  
540 to space constraints, we focus on summarizing  
541 the main trends rather than providing exhaustive  
542 technical details for each approach. Our selection  
543 of works primarily draws from major NLP con-  
544 ferences, including ACL, EMNLP, NAACL, and  
545 ICLR, along with relevant studies from other do-  
546 mains and preprints. While we strive to incorporate  
547 the latest research, many new works continue to  
548 emerge during our submission of this paper. Given  
549 the rapid evolution of this field, our survey offers a  
550 snapshot of current progress rather than a definitive  
551 account. We will continue to track developments  
552 and refine our analysis in future updates.

## 553 References

554 Chaitanya Agarwal, Vivek Gupta, Anoop Kunchukut-  
555 tan, and Manish Shrivastava. 2022. [Bilingual tabular  
556 inference: A case study on Indic languages](#). In *Pro-  
557 ceedings of the 2022 Conference of the North Amer-  
558 ican Chapter of the Association for Computational  
559 Linguistics: Human Language Technologies*, pages  
560 4018–4037, Seattle, United States. Association for  
561 Computational Linguistics.

562 Rami Aly, Zhijiang Guo, Michael Schlichtkrull, James  
563 Thorne, Andreas Vlachos, Christos Christodoulopou-  
564 los, Oana Cocarascu, and Arpit Mittal. 2021. [Fever-  
565 ous: Fact extraction and verification over un-  
566 structured and structured information](#). *Preprint*,  
567 arXiv:2106.05707.

568 Anonymous. 2024. [MMQA: Evaluating LLMs with  
569 multi-table multi-hop complex questions](#). In *Submit-  
570 ted to The Thirteenth International Conference on  
571 Learning Representations*. Under review.

572 Akari Asai, Jacqueline He, Rulin Shao, Weijia Shi,  
573 Amanpreet Singh, Joseph Chee Chang, Kyle Lo,  
574 Luca Soldaini, Sergey Feldman, Mike D’arcy,  
575 David Wadden, Matt Latzke, Minyang Tian, Pan Ji,  
576 Shengyan Liu, Hao Tong, Bohao Wu, Yanyu Xiong,  
577 Luke Zettlemoyer, Graham Neubig, Dan Weld, Doug  
578 Downey, Wen tau Yih, Pang Wei Koh, and Han-  
579 naneh Hajishirzi. 2024. [Openscholar: Synthesiz-  
580 ing scientific literature with retrieval-augmented lms](#).  
581 *Preprint*, arXiv:2411.14199.

582 Gilbert Badaro, Mohammed Saeed, and Paolo Papotti.  
583 2023. [Transformers for tabular data representation:  
584 A survey of models and applications](#). *Transactions  
585 of the Association for Computational Linguistics*,  
586 11:227–249.

Fan Bai, Junmo Kang, Gabriel Stanovsky, Dayne Fre-  
itag, Mark Dredze, and Alan Ritter. 2024. [Schema-  
driven information extraction from heterogeneous  
tables](#). *Preprint*, arXiv:2305.14336.

Ruisheng Cao, Fangyu Lei, Haoyuan Wu, Jixuan Chen,  
Yeqiao Fu, Hongcheng Gao, Xinzhuang Xiong, Han-  
chong Zhang, Yuchen Mao, Wenjing Hu, Tianbao  
Xie, Hongshen Xu, Danyang Zhang, Sida Wang,  
Ruoxi Sun, Pengcheng Yin, Caiming Xiong, Ansong  
Ni, Qian Liu, Victor Zhong, Lu Chen, Kai Yu, and  
Tao Yu. 2024. [Spider2-v: How far are multimodal  
agents from automating data science and engineering  
workflows?](#) *Preprint*, arXiv:2407.10956.

Shuaichen Chang, Jun Wang, Mingwen Dong, Lin  
Pan, Henghui Zhu, Alexander Hanbo Li, Wuwei  
Lan, Sheng Zhang, Jiarong Jiang, Joseph Lilien,  
Steve Ash, William Yang Wang, Zhiguo Wang,  
Vittorio Castelli, Patrick Ng, and Bing Xiang.  
2023. [Dr.spider: A diagnostic evaluation bench-  
mark towards text-to-sql robustness](#). *Preprint*,  
arXiv:2301.08881.

Peter Baile Chen, Yi Zhang, and Dan Roth. 2024. Is  
table retrieval a solved problem? exploring join-  
aware multi-table retrieval. In *Proceedings of the  
62nd Annual Meeting of the Association for Computa-  
tional Linguistics (ACL)*. ArXiv:2404.09889 [cs.LR],  
<https://doi.org/10.48550/arXiv.2404.09889>.

Wenhu Chen, Ming-Wei Chang, Eva Schlinger, William  
Wang, and William W. Cohen. 2021. [Open ques-  
tion answering over tables and text](#). *Preprint*,  
arXiv:2010.10439.

Wenhu Chen, Jianshu Chen, Yu Su, Zhiyu Chen, and  
William Yang Wang. 2020a. [Logical natural lan-  
guage generation from open-domain tables](#). In *Pro-  
ceedings of the 58th Annual Meeting of the Asso-  
ciation for Computational Linguistics*, pages 7929–  
7942, Online. Association for Computational Lin-  
guistics.

Wenhu Chen, Hongmin Wang, Jianshu Chen, Yunkai  
Zhang, Hong Wang, Shiyang Li, Xiyu Zhou, and  
William Yang Wang. 2020b. [Tabfact: A large-scale  
dataset for table-based fact verification](#). *Preprint*,  
arXiv:1909.02164.

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

Zhiyu Chen, Wenhu Chen, Charese Smiley, Sameena  
Shah, Iana Borova, Dylan Langdon, Reema Moussa,  
Matt Beane, Ting-Hao Huang, Bryan Routledge, and  
William Yang Wang. 2022. [Finqa: A dataset of  
numerical reasoning over financial data](#). *Preprint*,  
arXiv:2109.00122.

643	Zhiyu Chen, Wenhu Chen, Hanwen Zha, Xiyu Zhou,	sheng Huang. 2021a. <a href="#">Towards robustness of text-to-</a>	700
644	Yunkai Zhang, Sairam Sundaresan, and William Yang	<a href="#">sql models against synonym substitution</a> . <i>Preprint</i> ,	701
645	Wang. 2020d. <a href="#">Logic2text: High-fidelity natural</a>	arXiv:2106.01065.	702
646	<a href="#">language generation from logical forms</a> . <i>Preprint</i> ,		
647	arXiv:2004.14579.		
648	Zhoujun Cheng, Haoyu Dong, Zhiruo Wang, Ran Jia,	Yujian Gan, Xinyun Chen, and Matthew Purver.	703
649	Jiaqi Guo, Yan Gao, Shi Han, Jian-Guang Lou, and	2021b. <a href="#">Exploring underexplored limitations of</a>	704
650	Dongmei Zhang. 2022. <a href="#">HiTab: A hierarchical table</a>	<a href="#">cross-domain text-to-sql generalization</a> . <i>Preprint</i> ,	705
651	<a href="#">dataset for question answering and natural language</a>	arXiv:2109.05157.	706
652	<a href="#">generation</a> . In <i>Proceedings of the 60th Annual Meet-</i>		
653	<i>ing of the Association for Computational Linguistics</i>	Dawei Gao, Haibin Wang, Yaliang Li, Xiuyu Sun,	707
654	(Volume 1: Long Papers), pages 1094–1110, Dublin,	Yichen Qian, Bolin Ding, and Jingren Zhou. 2023.	708
655	Ireland. Association for Computational Linguistics.	<a href="#">Text-to-sql empowered by large language models: A</a>	709
		<a href="#">benchmark evaluation</a> . <i>Preprint</i> , arXiv:2308.15363.	710
656	Xiang Deng, Ahmed Hassan Awadallah, Christopher	Deepak Gupta, Surabhi Kumari, Asif Ekbal, and Push-	711
657	Meek, Oleksandr Polozov, Huan Sun, and Matthew	pak Bhattacharyya. 2018. <a href="#">MMQA: A multi-domain</a>	712
658	Richardson. 2021. <a href="#">Structure-grounded pretraining</a>	<a href="#">multi-lingual question-answering framework for En-</a>	713
659	<a href="#">for text-to-sql</a> . In <i>Proceedings of the 2021 Confer-</i>	<a href="#">glish and Hindi</a> . In <i>Proceedings of the Eleventh In-</i>	714
660	<i>ence of the North American Chapter of the Associ-</i>	<i>ternational Conference on Language Resources and</i>	715
661	<i>ation for Computational Linguistics: Human Lan-</i>	<i>Evaluation (LREC 2018)</i> , Miyazaki, Japan. European	716
662	<i>guage Technologies</i> . Association for Computational	Language Resources Association (ELRA).	717
663	Linguistics.		
664	Zheye Deng, Chunkit Chan, Weiqi Wang, Yuxi Sun,	Vivek Gupta, Pranshu Kandoi, Mahek Bhavesh Vora,	718
665	Wei Fan, Tianshi Zheng, Yauwai Yim, and Yangqiu	Shuo Zhang, Yujie He, Ridho Reinanda, and Vivek	719
666	Song. 2024. <a href="#">Text-tuple-table: Towards information</a>	Srikumar. 2023. <a href="#">Temptabqa: Temporal question</a>	720
667	<a href="#">integration in text-to-table generation via global tuple</a>	<a href="#">answering for semi-structured tables</a> . <i>Preprint</i> ,	721
668	<a href="#">extraction</a> . In <i>Proceedings of the 2024 Conference on</i>	arXiv:2311.08002.	722
669	<i>Empirical Methods in Natural Language Processing</i> ,		
670	pages 9300–9322, Miami, Florida, USA. Association	Vivek Gupta, Maitrey Mehta, Pegah Nokhiz, and Vivek	723
671	for Computational Linguistics.	Srikumar. 2020. <a href="#">INFOTABS: Inference on tables</a>	724
672	Mingwen Dong, Nischal Ashok Kumar, Yiqun Hu, Anuj	<a href="#">as semi-structured data</a> . In <i>Proceedings of the 58th</i>	725
673	Chauhan, Chung-Wei Hang, Shuaichen Chang, Lin	<i>Annual Meeting of the Association for Computational</i>	726
674	Pan, Wuwei Lan, Henghui Zhu, Jiarong Jiang, Patrick	<i>Linguistics</i> , pages 2309–2324, Online. Association	727
675	Ng, and Zhiguo Wang. 2025. <a href="#">PRACTIQ: A practical</a>	for Computational Linguistics.	728
676	<a href="#">conversational text-to-SQL dataset with ambigu-</a>		
677	<a href="#">ous and unanswerable queries</a> . In <i>Proceedings of</i>	Moshe Hazoom, Vibhor Malik, and Ben Bogin. 2021.	729
678	<i>the 2025 Conference of the Nations of the Americas</i>	<a href="#">Text-to-sql in the wild: A naturally-occurring</a>	730
679	<i>Chapter of the Association for Computational Lin-</i>	<a href="#">dataset based on stack exchange data</a> . <i>Preprint</i> ,	731
680	<i>guistics: Human Language Technologies (Volume 1:</i>	arXiv:2106.05006.	732
681	<i>Long Papers)</i> , pages 255–273, Albuquerque, New		
682	Mexico. Association for Computational Linguistics.	Xinyi He, Mengyu Zhou, Xinrun Xu, Xiaojun Ma,	733
683	Longxu Dou, Yan Gao, Mingyang Pan, Dingzirui	Rui Ding, Lun Du, Yan Gao, Ran Jia, Xu Chen,	734
684	Wang, Wanxiang Che, Dechen Zhan, and Jian-Guang	Shi Han, Zejian Yuan, and Dongmei Zhang. 2024.	735
685	Lou. 2022. <a href="#">Multispider: Towards benchmarking</a>	<a href="#">Text2analysis: A benchmark of table question</a>	736
686	<a href="#">multilingual text-to-sql semantic parsing</a> . <i>Preprint</i> ,	<a href="#">answering with advanced data analysis and unclear</a>	737
687	arXiv:2212.13492.	<a href="#">queries</a> . <i>arXiv preprint arXiv:2312.13671</i> .	738
688	Ondřej Dušek, Jekaterina Novikova, and Verena Rieser.	Jonathan Herzig, Pawel Krzysztof Nowak, Thomas	739
689	2020. <a href="#">Evaluating the state-of-the-art of end-to-end</a>	Müller, Francesco Piccinno, and Julian Eisenschlos.	740
690	<a href="#">natural language generation: The e2e nlg challenge</a> .	2020. <a href="#">Tapas: Weakly supervised table parsing via</a>	741
691	<i>Computer Speech and Language</i> , 59:123–156.	<a href="#">pre-training</a> . In <i>Proceedings of the 58th Annual Meet-</i>	742
692	Xi Fang, Weijie Xu, Fiona Anting Tan, Jiani Zhang,	<i>ing of the Association for Computational Linguistics</i> ,	743
693	Ziqing Hu, Yanjun Qi, Scott Nickleach, Diego Socol-	pages 4320–4333, Online. Association for Computa-	744
694	insky, Srinivasan Sengamedu, and Christos Faloutsos.	tional Linguistics.	745
695	2024. <a href="#">Large language models(lms) on tabular data:</a>	Atin Sakkeer Hussain. 2025. <a href="#">Artemis-da: An advanced</a>	746
696	<a href="#">Prediction, generation, and understanding – a survey</a> .	<a href="#">reasoning and transformation engine for multi-</a>	747
697	<i>Preprint</i> , arXiv:2402.17944.	<a href="#">step insight synthesis in data analytics</a> . <i>Preprint</i> ,	748
698	Yujian Gan, Xinyun Chen, Qiuping Huang, Matthew	arXiv:2412.14146.	749
699	Purver, John R. Woodward, Jinxia Xie, and Peng-	Hiroshi Iida, Dung Thai, Varun Manjunatha, and Mohit	750
		Iyyer. 2021. <a href="#">TABBIE: Pretrained representations of</a>	751
		<a href="#">tabular data</a> . In <i>Proceedings of the 2021 Conference</i>	752
		<i>of the North American Chapter of the Association</i>	753
		<i>for Computational Linguistics: Human Language</i>	754
		<i>Technologies</i> , pages 3446–3456, Online. Association	755
		for Computational Linguistics.	756

- 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.
- Peiwen Jiang, Xinbo Lin, Zibo Zhao, Ruhui Ma, Yvonne Jie Chen, and Jinhua Cheng. 2024. [TKGT: Redefinition and a new way of text-to-table tasks based on real world demands and knowledge graphs augmented LLMs](#). In *Proceedings of the 2024 Conference on Empirical Methods in Natural Language Processing*, pages 16112–16126, Miami, Florida, USA. Association for Computational Linguistics.
- Marcin Kardas, Piotr Czapla, Pontus Stenetorp, Sebastian Ruder, Sebastian Riedel, Ross Taylor, and Robert Stojnic. 2020. [AxCell: Automatic extraction of results from machine learning papers](#). In *Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing (EMNLP)*, pages 8580–8594, 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.
- Siddharth Khincha, Chelsi Jain, Vivek Gupta, Tushar Kataria, and Shuo Zhang. 2023. [InfoSync: Information synchronization across multilingual semi-structured tables](#). In *Findings of the Association for Computational Linguistics: ACL 2023*, pages 2536–2559, Toronto, Canada. Association for Computational Linguistics.
- 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](#). *Preprint*, arXiv:2305.07288.
- 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.
- Gyubok Lee, Woosog Chay, Seonhee Cho, and Edward Choi. 2024. [Trustsql: Benchmarking text-to-sql reliability with penalty-based scoring](#). *Preprint*, arXiv:2403.15879.
- Gyubok Lee, Hyeonji Hwang, Seongsu Bae, Yeonsu Kwon, Woncheol Shin, Seongjun Yang, Minjoon Seo, Jong-Yeup Kim, and Edward Choi. 2023. [Ehsql: A practical text-to-sql benchmark for electronic health records](#). *Preprint*, arXiv:2301.07695.
- Fangyu Lei, Jixuan Chen, Yuxiao Ye, Ruisheng Cao, Dongchan Shin, Hongjin Su, Zhaoqing Suo, Hongcheng Gao, Wenjing Hu, Pengcheng Yin, Victor Zhong, Caiming Xiong, Ruoxi Sun, Qian Liu, Sida Wang, and Tao Yu. 2024. [Spider 2.0: Evaluating language models on real-world enterprise text-to-sql workflows](#). *Preprint*, arXiv:2411.07763.
- Fangyu Lei, Tongxu Luo, Pengqi Yang, Weihao Liu, Hanwen Liu, Jiahe Lei, Yiming Huang, Yifan Wei, Shizhu He, Jun Zhao, and Kang Liu. 2023. [Tableqakit: A comprehensive and practical toolkit for table-based question answering](#). *Preprint*, arXiv:2310.15075.
- Jinyang Li, Binyuan Hui, Ge Qu, Jiayi Yang, Binhua Li, Bowen Li, Bailin Wang, Bowen Qin, Rongyu Cao, Ruiying Geng, Nan Huo, Xuanhe Zhou, Chenhao Ma, Guoliang Li, Kevin C. C. Chang, Fei Huang, Reynold Cheng, and Yongbin Li. 2023a. [Can llm already serve as a database interface? a big bench for large-scale database grounded text-to-sqls](#). *Preprint*, arXiv:2305.03111.
- Moxin Li, Wenjie Wang, Fuli Feng, Hanwang Zhang, Qifan Wang, and Tat-Seng Chua. 2023b. [Hypothetical training for robust machine reading comprehension of tabular context](#). In *Findings of the Association for Computational Linguistics: ACL 2023*, pages 1220–1236, Toronto, Canada. Association for Computational Linguistics.
- 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.
- Zhang Li, Biao Yang, Qiang Liu, Zhiyin Ma, Shuo Zhang, Jingxu Yang, Yabo Sun, Yuliang Liu, and Xiang Bai. 2024. [Monkey: Image resolution and text label are important things for large multi-modal models](#). In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*. ArXiv:2311.06607 [cs.CV], <https://doi.org/10.48550/arXiv.2311.06607>.
- Chuang Liu, Junzhuo Li, and Deyi Xiong. 2023a. [Tab-CQA: 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.
- Haotian Liu, Chunyuan Li, Qingyang Wu, and Yong Jae Lee. 2023b. [Visual instruction tuning](#). *Preprint*, arXiv:2304.08485.







981	Xinyu Pi, Bing Wang, Yan Gao, Jiaqi Guo, Zhou-	Barkavi Sundararajan, Somayajulu Sripada, and Ehud	1037
982	jun Li, and Jian-Guang Lou. 2022. <a href="#">Towards ro-</a>	Reiter. 2024. <a href="#">Improving factual accuracy of neural</a>	1038
983	<a href="#">bustness of text-to-sql models against natural and</a>	<a href="#">table-to-text output by addressing input problems in</a>	1039
984	<a href="#">realistic adversarial table perturbation.</a> <i>Preprint</i> ,	<a href="#">tutto.</a> <i>Preprint</i> , arXiv:2404.04103.	1040
985	arXiv:2212.09994.		
986	Mohammadreza Pourreza and Davood Rafiei. 2023a.	Anh Tuan Nguyen, Mai Hoang Dao, and Dat Quoc	1041
987	<a href="#">Din-sql: Decomposed in-context learning of text-to-</a>	Nguyen. 2020. <a href="#">A pilot study of text-to-SQL semantic</a>	1042
988	<a href="#">sql with self-correction.</a> <i>Preprint</i> , arXiv:2304.11015.	<a href="#">parsing for Vietnamese.</a> In <i>Findings of the Associa-</i>	1043
989		<i>tion for Computational Linguistics: EMNLP 2020</i> ,	1044
990	Mohammadreza Pourreza and Davood Rafiei. 2023b.	pages 4079–4085, Online. Association for Computa-	1045
991	<a href="#">Evaluating cross-domain text-to-SQL models and</a>	tional Linguistics.	1046
992	<a href="#">benchmarks.</a> In <i>Proceedings of the 2023 Conference</i>		
993	<i>on Empirical Methods in Natural Language Process-</i>	Boris van Breugel and Mihaela van der Schaar. 2024.	1047
994	<i>ing</i> , pages 1601–1611, Singapore. Association for	Why tabular foundation models should be a re-	1048
995	Computational Linguistics.	search priority. In <i>Proceedings of the 41st Inter-</i>	1049
996	Pritika Ramu, Aparna Garimella, and Sambaran Bandy-	<i>national Conference on Machine Learning (ICML).</i>	1050
997	<a href="#">opadhyay. 2024. Is this a bad table? a closer look at</a>	ArXiv:2405.01147 [cs.LG], <a href="https://doi.org/10.48550/arXiv.2405.01147">https://doi.org/10.</a>	1051
998	<a href="#">the evaluation of table generation from text.</a> <i>Preprint</i> ,	<a href="#">48550/arXiv.2405.01147.</a>	1052
999	arXiv:2406.14829.		
1000	Weijieying Ren, Tianxiang Zhao, Yuqing Huang, and	Fei Wang, Kexuan Sun, Muhao Chen, Jay Pujara, and	1053
1001	Vasant Honavar. 2025. <a href="#">Deep learning within tabular</a>	Pedro Szekely. 2021a. <a href="#">Retrieving complex tables</a>	1054
1002	<a href="#">data: Foundations, challenges, advances and future</a>	<a href="#">with multi-granular graph representation learning.</a> In	1055
1003	<a href="#">directions.</a> <i>Preprint</i> , arXiv:2501.03540.	<i>Proceedings of the 44th International ACM SIGIR</i>	1056
1004		<i>Conference on Research and Development in Infor-</i>	1057
1005	Aofeng Su, Aowen Wang, Chao Ye, Chen Zhou,	<i>mation Retrieval (SIGIR).</i> ArXiv:2105.01736 [cs.IR],	1058
1006	Ga Zhang, Gang Chen, Guangcheng Zhu, Haobo	<a href="https://doi.org/10.48550/arXiv.2105.01736">https://doi.org/10.48550/arXiv.2105.01736.</a>	1059
1007	Wang, Haokai Xu, Hao Chen, Haoze Li, Haoxuan		
1008	Lan, Jiaming Tian, Jing Yuan, Junbo Zhao, Jun-	Nancy X. R. Wang, Diwakar Mahajan, Marina	1060
1009	lin Zhou, Kaizhe Shou, Liangyu Zha, Lin Long,	Danilevsky, and Sara Rosenthal. 2021b. <a href="#">SemEval-</a>	1061
1010	Liyao Li, Pengzuo Wu, Qi Zhang, Qingyi Huang,	<a href="#">2021 task 9: Fact verification and evidence finding</a>	1062
1011	Saisai Yang, Tao Zhang, Wentao Ye, Wufang Zhu,	<a href="#">for tabular data in scientific documents (SEM-TAB-</a>	1063
1012	Xiaomeng Hu, Xijun Gu, Xinjie Sun, Xiang Li,	<a href="#">FACTS).</a> In <i>Proceedings of the 15th International</i>	1064
1013	Yuhang Yang, and Zhiqing Xiao. 2024. <a href="#">Tablegpt2: A</a>	<i>Workshop on Semantic Evaluation (SemEval-2021)</i> ,	1065
1014	<a href="#">large multimodal model with tabular data integration.</a>	pages 317–326, Online. Association for Computa-	1066
1015	<i>Preprint</i> , arXiv:2411.02059.	tional Linguistics.	1067
1016			
1017	Lya Hulliyyatus Suadaa, Hidetaka Kamigaito, Kotaro	Yuqi Wang, Lyuhao Chen, Songcheng Cai, Zhijian Xu,	1068
1018	Funakoshi, Manabu Okumura, and Hiroya Takamura.	and Yilun Zhao. 2024a. <a href="#">Revisiting automated evalua-</a>	1069
1019	2021. <a href="#">Towards table-to-text generation with numer-</a>	<a href="#">tion for long-form table question answering.</a> In <i>Pro-</i>	1070
1020	<a href="#">ical reasoning.</a> In <i>Proceedings of the 59th Annual</i>	<i>ceedings of the 2024 Conference on Empirical Meth-</i>	1071
1021	<i>Meeting of the Association for Computational Lin-</i>	<i>ods in Natural Language Processing</i> , pages 14696–	1072
1022	<i>guistics and the 11th International Joint Conference</i>	14706, Miami, Florida, USA. Association for Com-	1073
1023	<i>on Natural Language Processing (Volume 1: Long</i>	putational Linguistics.	1074
1024	<i>Papers)</i> , pages 1451–1465, Online. Association for		
1025	Computational Linguistics.	Zhiruo Wang, Haoyu Dong, Ran Jia, Jia Li, Zhiyi Fu,	1075
1026	Yuan Sui, Mengyu Zhou, Mingjie Zhou, Shi Han,	Shi Han, and Dongmei Zhang. 2021c. <a href="#">TUTA: Tree-</a>	1076
1027	and Dongmei Zhang. 2024. Table meets llm: Can	<a href="#">based transformers for generally structured table pre-</a>	1077
1028	large language models understand structured table	<a href="#">training.</a> In <i>Proceedings of the 27th ACM SIGKDD</i>	1078
1029	data? a benchmark and empirical study. In <i>Pro-</i>	<i>Conference on Knowledge Discovery &amp; Data Mining</i>	1079
1030	<i>ceedings of the 17th ACM International Confer-</i>	<i>(KDD).</i> ArXiv:2010.12537 [cs.IR], <a href="https://doi.org/10.48550/arXiv.2010.12537">https://doi.org/10.48550/arXiv.2010.12537.</a>	1080
1031	<i>ence on Web Search and Data Mining (WSDM).</i>		1081
1032	ArXiv:2305.13062 [cs.CL], <a href="https://doi.org/10.48550/arXiv.2305.13062">https://doi.org/10.</a>	Zilong Wang, Hao Zhang, Chun-Liang Li, Julian Martin	1082
1033	<a href="#">48550/arXiv.2305.13062.</a>	Eisenschlos, Vincent Perot, Zifeng Wang, Lesly Mi-	1083
1034	Yuan Sui, Jiaru Zou, Mengyu Zhou, Xinyi He, Lun Du,	culicich, Yasuhisa Fujii, Jingbo Shang, Chen-Yu Lee,	1084
1035	Shi Han, and Dongmei Zhang. 2023. TAP4LLM:	and Tomas Pfister. 2024b. <a href="#">Chain-of-table: Evolving</a>	1085
1036	Table provider on sampling, augmenting, and pack-	<a href="#">tables in the reasoning chain for table understanding.</a>	1086
	ing semi-structured data for large language model	<i>Preprint</i> , arXiv:2401.04398.	1087
	reasoning. <i>arXiv preprint arXiv:2312.09039.</i> <a href="https://doi.org/10.48550/arXiv.2312.09039">https://doi.org/10.48550/arXiv.2312.09039.</a>		
		Sam Wiseman, Stuart M. Shieber, and Alexander M.	1088
		Rush. 2017. <a href="#">Challenges in data-to-document genera-</a>	1089
		<a href="#">tion.</a> <i>Preprint</i> , arXiv:1707.08052.	1090
		Xueqing Wu, Jiacheng Zhang, and Hang Li. 2022.	1091
		<a href="#">Text-to-table: A new way of information extraction.</a>	1092
		<i>Preprint</i> , arXiv:2109.02707.	1093

1094	Witold Wydmański, Ulvi Movsum-zada, Jacek Ta-	Yilun Zhao, Yunxiang Li, Chenying Li, and Rui Zhang.	1149
1095	bor, and Marek Śmieja. 2024. <a href="#">Vistabnet: Adapt-</a>	2022. <a href="#">MultiHiertt: Numerical reasoning over multi</a>	1150
1096	<a href="#">ing vision transformers for tabular data</a> . <i>Preprint</i> ,	<a href="#">hierarchical tabular and textual data</a> . In <i>Proceedings</i>	1151
1097	arXiv:2501.00057.	<i>of the 60th Annual Meeting of the Association for</i>	1152
		<i>Computational Linguistics (Volume 1: Long Papers)</i> ,	1153
1098	Sean Yang, Chris Tensmeyer, and Curtis Wigington.	pages 6588–6600, Dublin, Ireland. Association for	1154
1099	2022. <a href="#">TELIN: Table entity LINKer for extracting</a>	Computational Linguistics.	1155
1100	<a href="#">leaderboards from machine learning publications</a> . In		
1101	<i>Proceedings of the first Workshop on Information</i>	Yilun Zhao, Zhenting Qi, Linyong Nan, Boyu Mi, Yixin	1156
1102	<i>Extraction from Scientific Publications</i> , pages 20–25,	Liu, Weijin Zou, Simeng Han, Ruizhe Chen, Xiangru	1157
1103	Online. Association for Computational Linguistics.	Tang, Yumo Xu, Dragomir Radev, and Arman Co-	1158
		han. 2023. <a href="#">Qtsumm: Query-focused summarization</a>	1159
1104	Yunhu Ye, Binyuan Hui, Min Yang, Binhua Li, Fei	<a href="#">over tabular data</a> . In <i>Proceedings of the 2023 Con-</i>	1160
1105	Huang, and Yongbin Li. 2023. <a href="#">Large language mod-</a>	<i>ference on Empirical Methods in Natural Language</i>	1161
1106	<a href="#">els are versatile decomposers: Decompose evidence</a>	<i>Processing (EMNLP)</i> . Accepted at EMNLP 2023.	1162
1107	<a href="#">and questions for table-based reasoning</a> . <i>Preprint</i> ,		
1108	arXiv:2301.13808.	Mingyu Zheng, Xinwei Feng, Qingyi Si, Qiaoqiao	1163
		She, Zheng Lin, Wenbin Jiang, and Weiping	1164
1109	Tao Yu, Rui Zhang, Kai Yang, Michihiro Yasunaga,	Wang. 2024. Multimodal table understanding. In	1165
1110	Dongxu Wang, Zifan Li, James Ma, Irene Li, Qingn-	<i>Proceedings of the 62nd Annual Meeting of the</i>	1166
1111	ing Yao, Shanelle Roman, Zilin Zhang, and Dragomir	<i>Association for Computational Linguistics (ACL)</i> .	1167
1112	Radev. 2018. <a href="#">Spider: A large-scale human-labeled</a>	ArXiv:2406.08100 [cs.CL], <a href="https://doi.org/10.48550/arXiv.2406.08100">https://doi.org/10.</a>	1168
1113	<a href="#">dataset for complex and cross-domain semantic pars-</a>	48550/arXiv.2406.08100.	1169
1114	<a href="#">ing and text-to-SQL task</a> . In <i>Proceedings of the 2018</i>		
1115	<i>Conference on Empirical Methods in Natural Lan-</i>	Mingyu Zheng, Yang Hao, Wenbin Jiang, Zheng Lin,	1170
1116	<i>guage Processing</i> , pages 3911–3921, Brussels, Bel-	Yajuan Lyu, QiaoQiao She, and Weiping Wang. 2023.	1171
1117	gium. Association for Computational Linguistics.	<a href="#">IM-TQA: A Chinese table question answering dataset</a>	1172
		<a href="#">with implicit and multi-type table structures</a> . In <i>Pro-</i>	1173
1118	Hanchong Zhang, Ruisheng Cao, Lu Chen, Hongshen	<i>ceedings of the 61st Annual Meeting of the Associa-</i>	1174
1119	Xu, and Kai Yu. 2023a. Act-sql: In-context learning	<i>tion for Computational Linguistics (Volume 1: Long</i>	1175
1120	for text-to-sql with automatically-generated chain-of-	<i>Papers)</i> , pages 5074–5094, Toronto, Canada. Associ-	1176
1121	thought. <i>arXiv preprint arXiv:2310.17342</i> . <a href="https://doi.org/10.48550/arXiv.2310.17342">https:</a>	ation for Computational Linguistics.	1177
1122	<a href="https://doi.org/10.48550/arXiv.2310.17342">//doi.org/10.48550/arXiv.2310.17342</a> .		
		Victor Zhong, Caiming Xiong, and Richard Socher.	1178
1123	Jiani Zhang, Zhengyuan Shen, Balasubramaniam Srin-	2017. <a href="#">Seq2sql: Generating structured queries</a>	1179
1124	ivasan, Shen Wang, Huzefa Rangwala, and George	<a href="#">from natural language using reinforcement learning</a> .	1180
1125	Karypis. 2023b. <a href="#">NameGuess: Column name expan-</a>	<i>Preprint</i> , arXiv:1709.00103.	1181
1126	<a href="#">sion for tabular data</a> . In <i>Proceedings of the 2023</i>		
1127	<i>Conference on Empirical Methods in Natural Lan-</i>	Wei Zhou, Mohsen Mesgar, Heike Adel, and Annemarie	1182
1128	<i>guage Processing</i> , pages 13276–13290, Singapore.	Friedrich. 2024. <a href="#">FREB-TQA: A fine-grained robust-</a>	1183
1129	Association for Computational Linguistics.	<a href="#">ness evaluation benchmark for table question answer-</a>	1184
		<a href="#">ing</a> . In <i>Proceedings of the 2024 Conference of the</i>	1185
1130	Tianshu Zhang, Xiang Yue, Yifei Li, and Huan Sun.	<i>North American Chapter of the Association for Com-</i>	1186
1131	2024a. <a href="#">Tablellama: Towards open large generalist</a>	<i>putational Linguistics: Human Language Technolo-</i>	1187
1132	<a href="#">models for tables</a> . <i>Preprint</i> , arXiv:2311.09206.	<i>gies (Volume 1: Long Papers)</i> , pages 2479–2497,	1188
		Mexico City, Mexico. Association for Computational	1189
1133	Xuanliang Zhang, Dingzirui Wang, Longxu Dou,	Linguistics.	1190
1134	Qingfu Zhu, and Wanxiang Che. 2024b. <a href="#">A survey of</a>		
1135	<a href="#">table reasoning with large language models</a> . <i>Preprint</i> ,	Fengbin Zhu, Wenqiang Lei, Youcheng Huang, Chao	1191
1136	arXiv:2402.08259.	Wang, Shuo Zhang, Jiancheng Lv, Fuli Feng, and	1192
		Tat-Seng Chua. 2021. <a href="#">Tat-qa: A question answering</a>	1193
1137	Yi Zhang, Jan Deriu, George Katsogiannis-Meimarakis,	<a href="#">benchmark on a hybrid of tabular and textual content</a>	1194
1138	Catherine Kosten, Georgia Koutrika, and Kurt	<a href="#">in finance</a> . <i>Preprint</i> , arXiv:2105.07624.	1195
1139	Stockinger. 2023c. <a href="#">Sciencebenchmark: A complex</a>		
1140	<a href="#">real-world benchmark for evaluating natural language</a>		
1141	<a href="#">to sql systems</a> . <i>Preprint</i> , arXiv:2306.04743.		
		<b>A Text-to-SQL</b>	1196
1142	Zhehao Zhang, Xitao Li, Yan Gao, and Jian-Guang Lou.		
1143	2023d. <a href="#">CRT-QA: A dataset of complex reasoning</a>	Text-to-SQL is a semantic parsing task that is	1197
1144	<a href="#">question answering over tabular data</a> . In <i>Proceed-</i>	highly relevant to table-based applications: given a	1198
1145	<i>ings of the 2023 Conference on Empirical Methods</i>	natural language question, the model must gener-	1199
1146	<i>in Natural Language Processing</i> , pages 2131–2153,	erate a SQL query that accurately captures the	1200
1147	Singapore. Association for Computational Linguis-	intent of the query. Over time, these tasks have	1201
1148	tics.		

Benchmark	Sources / Domain	Sizes	Input Format	T / Q	Directions
WikiSQL (2017)	Wikipedia	80,654	Row Header, Row-Separated	1.0	-
Spider (2018)	Academic Databases, Online CSV, WikiSQL	10,181	Table(col), Type, PF	1.6	-
SEDE (2021)	Stack Exchange	12,023	Table(col), Type, PF	1.3	Noise Utterance
SpiderDK (2021b)	Spider	535	Table(col), Type, PF	> 1	Domain Knowledge
SpiderSyn (2021a)	Spider	8,034	Table(col), Type, PF	> 1	Query Perturbation
SpiderRealistic (2021)	Spider	508	Table(col), Type, PF	> 1	Query Perturbation
MIMICSQL (2021)	Electronic Medical Records	10,000	Row Header, Row-Separated	1.8	Domain Knowledge
KaggleDBQA (2021)	ATIS, GeoQuery, Restaurants, Yelp, Academic, IMDB, Scholar, Advising	272	Table(col), Type, PF, context	1.2	Domain Knowledge
ADVETA (2022)	Spider, WikiSQL, WTQ	-	Table(col), Type, PF	> 1	Table Perturbation
BIRD (2023a)	Kaggle, Machine Learning platform	12,751	Table(col), Type, PF, context	> 1	Table Size
Dr.Spider (2023)	Spider	15,000	Table(col), Type, PF	> 1	Table, Query Perturbation
EHRSQL (2023)	Electronic Medical Records	24,000	Table(col), Type, PF	2.4	Domain, Reasoning
ScienceBench (2023c)	CORDIS, SDSS, OncoMX	6,000	Table(col), Type, PF	> 1	Data Synthesis, Domain
TrustSQL (2024)	ATIS, Advising, EHRSQL, Spider	27,784	CREATE(EoT)	> 1	Reasoning
Spider2 (2024)	Cloud Data Warehouses	632	Table(col), PF	> 1	Reasoning, Table Size
Spider2V (2024)	Cloud Data Warehouses	494	Agent Workspace	> 1	Input Modality

Table 4: Summary of benchmarks for Text-to-SQL. **Sizes** refers to the number of SQL query pairs, and **T/Q** indicates the number of tables required to answer a single query.

evolved to incorporate additional contextual information—such as table schemas and optional sample rows—with the evaluation focus shifting from exact match (EM) to execution accuracy (EX) as the primary metric. A prominent benchmark in this area, Spider (Yu et al., 2018), significantly increased task complexity by introducing databases composed of multiple tables, foreign keys, and the requirement to employ a variety of functions.

Building on Spider, several adaptations and extensions have broadened the task’s scope and complexity. Multilingual adaptations (Min et al., 2019; Tuan Nguyen et al., 2020; Dou et al., 2022) expanded Text-to-SQL to cross-lingual and multilingual settings, enabling SQL generation across diverse languages. Other extensions include SpiderDK (Gan et al., 2021b), which incorporates domain knowledge, and Spider-Syn (Gan et al., 2021a) and Spider-Realistic (Deng et al., 2021), which obscure schema-related words or column names to simulate noisy utterances and more realistic queries.

Text-to-SQL has been well-studied with question decomposition pipelines (Gao et al., 2023; Ye et al., 2023; Wang et al., 2024b), with current models nearing saturation on some commonly used benchmarks.

**Effect of Noisy Input.** Beyond evaluation issues, Text-to-SQL faces inherent challenges, especially when handling ambiguity, or on very large tables. As noted in (Chen et al., 2024), performance drops significantly without PF keys, as variations in column names across tables and limited sample rows complicate element matching. Moreover, as highlighted in (Lei et al., 2024; Maamari et al., 2024), model performance deteriorates sharply when pro-

cessing extremely long database schema, a scenario prevalent in real-world industrial databases.

## B Responsible NLP Miscellanea

### B.1 AI Assistants

We acknowledge the use of GPT-4o and GPT-o3-mini for grammar checking and word polishing.