

DMRAL: Decomposition-Guided Retrieval and Reasoning for Multi-Table QA

Feng Luo¹, Hai Lan², Hui Luo³, Zhifeng Bao², Xiaoli Wang⁴, J. Shane Culpepper², Shazia Sadiq²

¹RMIT University, ²The University of Queensland, ³University of Wollongong, ⁴Xiamen University

ABSTRACT

In this paper, we study the problem of numerical multi-table question answering (MTQA) using tables in the wild (e.g., online data repositories). This task is essential in many analytical applications. Existing MTQA solutions, such as text-to-SQL or open-domain MTQA methods, are designed for databases and struggle when applied to the siloed tables in the wild. The key limitations include: (1) Limited support for incomplete metadata and complex table relationships; (2) Ineffective retrieval of relevant tables at scale; (3) Inaccurate answer generation. To overcome these limitations, we propose DMRAL, a Decomposition-guided Multi-table Retrieval and Answering framework for MTQA using tables in the wild, consisting of: (1) Preprocessing, which constructs a graph to capture the table relationships and infers missing metadata; (2) Table-Aligned Question Decomposer and Coverage-Aware Retriever, which jointly enable the effective identification of relevant tables from large-scale corpora by enhancing the question decomposition quality and maximizing the question coverage of retrieved tables; (3) Sub-question Guided Reasoner, which produces correct answers by progressively generating and refining the reasoning program based on sub-questions. Experiments on two MTQA datasets demonstrate that DMRAL significantly outperforms existing state-of-the-art MTQA methods, with an average improvement of 24% in table retrieval and 55% in answer accuracy.

PVLDB Reference Format:

Feng Luo¹, Hai Lan², Hui Luo³, Zhifeng Bao², Xiaoli Wang⁴, J. Shane Culpepper², Shazia Sadiq². DMRAL: Decomposition-Guided Retrieval and Reasoning for Multi-Table QA. PVLDB, 14(1): XXX-XXX, 2020. doi:XX.XX/XXX.XX

PVLDB Artifact Availability:

The source code, data, and/or other artifacts have been made available at https://github.com/JrjessyLuo/multitab_qa_dmrал.

1 INTRODUCTION

Multi-table question answering (MTQA) is a well-known task that requires identifying and integrating information from multiple tables to answer a single question [41]. To better support various analytical applications, there is a growing need for MTQA systems to also support numerical questions [27]. At the same time, the growing prevalence of *tables in the wild*—such as collections of tables found on the web, in public data lakes, or available

through data markets—offers rich opportunities for MTQA, where the valuable data for answering the numerical questions is often dispersed across those soiled tables. While promising, working with such tables presents unique challenges due to their *large scale* (e.g., tens of thousands of tables), *incomplete metadata* (e.g., missing column headers) [15, 39], and *complex inter-table relationships*, namely unionability (tables that can be unioned on similar column headers) and joinability (tables that can be joined based on matching columns) [9, 10]. In this paper, we study MTQA using tables in the wild, a challenging but increasingly realistic scenario driven by real-world data acquisition and usage trends.

Current MTQA approaches are primarily designed for relational databases and can be categorized into two paradigms: Text-to-SQL and Open-domain MTQA. As shown in Figure 1, Text-to-SQL methods [20, 26] were originally developed for a single database, typically involving fewer than ten tables and relying on a well-defined database schema with complete table metadata and PK-FK constraints [6]. More recently, state-of-the-art methods increasingly leverage large language models (LLMs) to select the most relevant tables and generate executable SQL queries [24]. In contrast, open-domain MTQA methods [7, 41] use moderately sized collections (i.e., hundreds of tables) that are aggregated from multiple databases. These approaches typically begin by decomposing a question into a set of sub-questions using the internal knowledge of LLMs, then retrieving relevant tables based on a relevance score between sub-questions and tables using table joinability, and finally invoking an LLM to generate SQL over the retrieved tables.

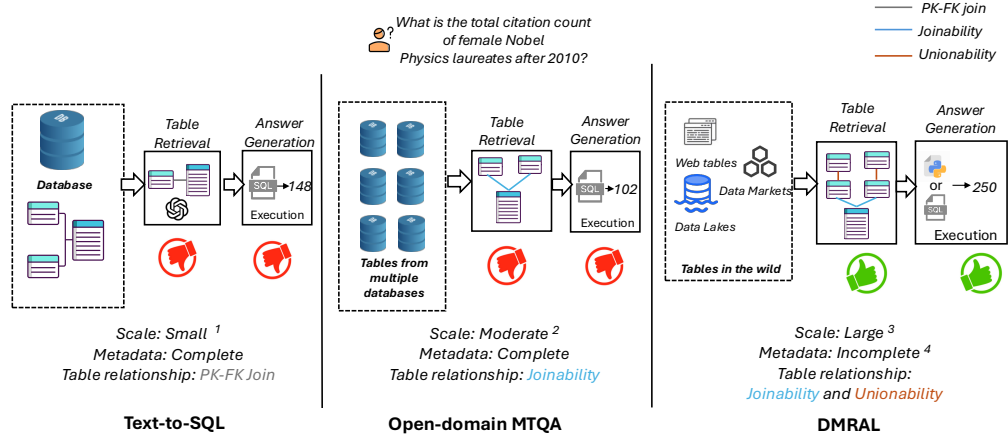
However, applying these methods to the tables in the wild presents three major limitations. **(L1) Limited support for incomplete metadata and complex table relationships:** Both approaches assume that complete metadata exists and overlook complex table relationships with unionability, which limits their robustness in handling such table corpora. **(L2) Ineffective retrieval of relevant tables at scale:** With token limits, state-of-the-art LLM-based Text-to-SQL methods cannot directly select relevant tables from large-scale corpora. Meanwhile, Open-domain MTQA depends on LLMs to decompose sub-questions for table retrieval, while the quality of these decompositions is often low, leading to suboptimal retrieval effectiveness. **(L3) Inaccurate answer generation:** While generating programs such as SQL and running them to derive answers is promising for answering numerical questions, existing approaches rarely produce fully correct programs (e.g., contain incorrect joins) for execution, which results in low answer accuracy [35].

Our Contributions. To resolve these limitations, we propose DMRAL, a Decomposition-guided Multi-table Retrieval and Answering framework for MTQA using tables in the wild. Our key contributions are as follows:

- We develop a *Preprocessing* pipeline (§2.3) that builds a *Table Relationship Graph* to capture complex table relationships, and

This work is licensed under the Creative Commons BY-NC-ND 4.0 International License. Visit <https://creativecommons.org/licenses/by-nc-nd/4.0/> to view a copy of this license. For any use beyond those covered by this license, obtain permission by emailing info@vldb.org. Copyright is held by the owner/author(s). Publication rights licensed to the VLDB Endowment.

Proceedings of the VLDB Endowment, Vol. 14, No. 1 ISSN 2150-8097.
doi:XX.XX/XXX.XX



¹ We measure the scale based on the number of tables. For example, the Text-to-SQL benchmark Bird averages ~ 7 tables per database.

² The datasets utilized in Open-domain MTQA study contains at most ~ 500 tables.

³ We collected at least 73,688 tables for evaluation.

⁴ Incomplete metadata is common. For example, $\sim 23\%$ tables in NQ-Tables (extracted from Wikipedia) have at least one missing column header.

Figure 1: A comparison of problem settings and system capabilities across different MTQA approaches: Text-to-SQL, Open-domain MTQA, and our proposed DMRAL.

performs metadata inference to recover missing column headers (i.e., addressing L1).

- We propose an effective decomposition-guided multi-table retrieval strategy to address L2, consisting of: (1) *Table-Aligned Question Decomposer* (§3), which enhances decomposition quality by identifying distinct information needs and aligning them with the underlying table structures; (2) *Coverage-Aware Retriever* (§4), which effectively retrieves relevant tables from a large-scale table corpus by maximizing their coverage of the question through coverage scoring and verification mechanisms.
- We propose a *Sub-question Guided Reasoner* (§5) to improve answer accuracy by guiding LLMs to progressively generate and refine the program based on the decomposed sub-questions (i.e., addressing L3).
- To comprehensively evaluate this task, we prepared two large-scale datasets, SpiderWild and BirdWild, consisting of 73,688 and 109,949 tables curated from real-world table sources [8, 22, 47] (§6). Extensive experiments demonstrate the effectiveness of DMRAL, achieving an average improvement of 24% in identifying relevant tables and 55% in producing accurate answers (§7).

In summary, DMRAL is a robust, scalable, and traceable framework designed for MTQA over tables in the wild. For each answer, DMRAL enables fine-grained tracing and verification (e.g., whether the retrieved tables are appropriate, whether the reasoning program is correct, and whether the sub-questions are well-decomposed). This traceability is crucial for MTQA, as it not only ensures transparency in how answers are derived but also provides actionable insights to diagnose and improve each component of our framework.

2 PROBLEM FORMULATION AND SOLUTION OVERVIEW

2.1 Problem Definition

A table T consists of a set of columns $T_c = \{c_1, \dots, c_m\}$ and rows T_r , along with associated metadata. The metadata includes a table

title T_p and column headers $T_h = \{h_1, \dots, h_m\}$, which may be partially missing. We refer to *tables in the wild* as a table collection $\mathcal{T} = \{T^1, T^2, \dots, T^{|\mathcal{T}|}\}$. These tables are often related through two representative inter-table relationships:

- **Joinability:** Two tables T^i and T^j are considered joinable if there exists at least one column from T_c^i and one from T_c^j such that the two columns share overlapping or semantically similar values.
- **Unionability:** Two tables T^i and T^j are considered unionable if their column headers T_h^i and T_h^j are sufficiently similar to allow a one-to-one alignment across the two tables.

DEFINITION 1 (NUMERICAL MULTI-TABLE QUESTION ANSWERING). Given a table collection \mathcal{T} and a numerical question q that requires a numeric answer, it aims to compute the correct answer ans by learning the function $f : (\mathcal{T}, q) \rightarrow ans$, where the answer ans is derived by identifying a subset of relevant tables $\mathcal{T}_q \subseteq \mathcal{T}$ and performing reasoning over \mathcal{T}_q .

2.2 Related Work

Text-to-SQL. Text-to-SQL is a long-standing task that aims to convert natural language (NL) questions into executable SQL queries over relational databases. Existing solutions for Text-to-SQL fall into three paradigms: Rule-based methods, Neural network-based methods, and LLM-based methods. *Rule-based methods* [11, 21, 46] use hand-crafted rules and semantic parsers to map NL queries to SQL. These methods struggle to generalize to unseen domains and require extensive manual engineering [24]. *Neural network-based methods* [3, 43, 45] adopt deep learning models, such as encoder-decoder architectures [43], pretrained language models [45], and graph neural networks [3], to enhance schema understanding and SQL generation, thus improving the generalization across databases. For example, GNN-Parser [3] employs graph neural networks to represent the database schema structure and integrates this representation into an encoder-decoder framework for generating executable SQL queries. While more robust, they are

still limited by fixed model capacity and the availability of high-quality supervised training data. *LLM-based methods* [5, 20, 49] have recently gained attention for leveraging LLMs to generate SQL via prompt engineering [31, 44], supervised fine-tuning [19], and tool-augmented reasoning (e.g., LLM agents) [20, 36], achieving state-of-the-art performance [24].

Despite their success, these methods are not directly applicable to MTQA over tables in the wild, due to three key limitations. First, they assume the existence of explicit table relationships via PK-FK joins [12], which are unavailable in real-world table collections. Second, they require complete metadata (i.e., table titles and column headers) for schema linking [37], while such metadata can be incomplete for tables in the wild. Third, the SOTA LLM-based methods are limited by the context length of LLMs, and do not scale for large-scale table corpora [42]. Given these limitations, we do not consider Text-to-SQL methods as primary baselines, while we provide an experimental evaluation to compare our effectiveness with Text-to-SQL methods in identifying relevant tables for SQL generation (§7.5).

Open-domain MTQA. Open-domain MTQA operates over table collections by first retrieving the relevant tables using the decomposed sub-questions, and then generating SQL queries directly using LLMs [50]. To identify the relevant tables, MMQA [41] formulates the problem as a multi-hop retrieval problem. In contrast, JAR [7] formulates the problem as an optimization problem by jointly considering sub-question-table relevance and overall table relevance. However, both approaches are designed for relational database tables, which limits their ability to support incomplete metadata and table unionability for tables in the wild. Moreover, MMQA is vulnerable to cascading errors introduced by early retrieval mistakes [18], reducing its *effectiveness*. Conversely, JAR suffers from a high computational overhead due to the optimization algorithm used, limiting *efficiency*.

2.3 Solution Overview

DMRAL begins with a *Preprocessing* pipeline that captures complex table relationships to support efficient identification of multiple relevant tables. To achieve this, we construct a *Table Relationship Graph* $\mathcal{G} = (\mathcal{V}, \mathcal{E})$, where each node $v \in \mathcal{V}$ represents a cluster of unionable tables whose unionability [34] exceeds a threshold τ_u . Edges $e \in \mathcal{E}$ connect these clusters if any pair of tables in the clusters is joinable, as measured by a joinability [7] exceeding a threshold τ_j . By default, we set the unionability threshold τ_u to 0.9 and the joinability threshold τ_j to 0.5, based on their strong empirical effectiveness (§7.2.2). As incomplete metadata (i.e., missing column headers) may hinder accurate unionability computation and negatively affect the answer accuracy, we introduce a *metadata inference* module that leverages LLMs to infer and complete the missing headers. A naive approach here would sample a few rows (e.g., 10) from the table as the whole context and feed them to an LLM to infer missing headers. However, this may include irrelevant or loosely correlated columns, which dilute the context and reduce effectiveness [48]. To mitigate this, we apply the column semantics discovery algorithm [15] to partition each table into semantically coherent column groups. We extract the columns in each group

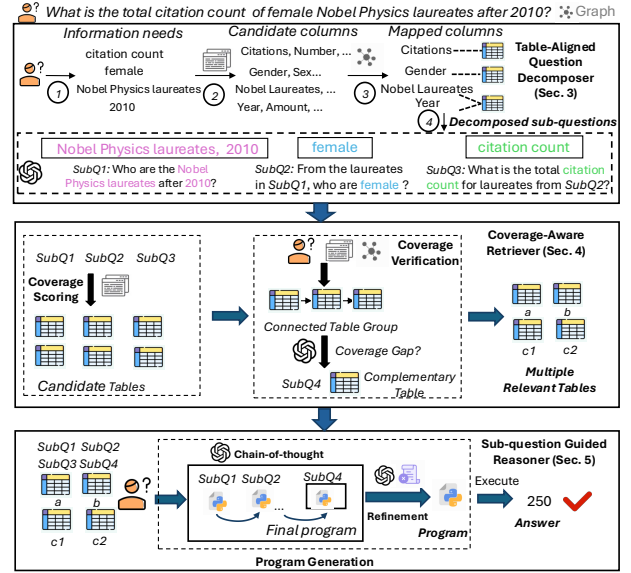


Figure 2: MTQA processing with DMRAL, in which intermediate results at each step can be traced and refined.

as a subtable and prompt the LLM using only the subcontext (i.e., sampled rows from the partitioned subtable) for inference.

Built upon this preprocessed graph, DMRAL answers questions using three core modules, as illustrated in Figure 2:

Table-Aligned Question Decomposer (§3) decomposes the input question into multiple sub-questions. To improve decomposition quality, it proposes a four-step decomposition approach: (1) extract the information needs from the question, (2) align each information need to the columns in the tables, (3) select the most promising mapping between the information needs and columns, and (4) group contextually coherent information needs to generate subquestions using LLM-based decomposition.

Coverage-Aware Retriever (§4) retrieves multiple relevant tables based on the decomposed subquestions. To support the effective retrieval over large-scale table corpora, it first introduces a learning-based coverage scoring function to retrieve the candidate tables for each subquestion. Then it proposes a coverage verification function to first construct the connected table group to collectively answer the question, and then identify the complementary tables to fill the potential coverage gap.

Sub-question Guided Reasoner (§5) generates executable programs (e.g., SQL or Python) over the retrieved tables to derive the final answer. To improve answer reliability, it first uses chain-of-thought prompting to incrementally generate intermediate programs using a sequence of decomposed sub-questions. It then applies an execution-guided refinement mechanism to verify and iteratively revise the generated program.

3 TABLE-ALIGNED QUESTION DECOMPOSER

3.1 Principles of Effective Decomposition

Correct decomposition of complex user questions is essential for MTQA, as it enables precise retrieval. Given a question q , we define

the information needs contained as a set $I(q)$, where each element $I_q \in I(q)$ represents a specific concept, entity, or condition included in q that is required to derive the correct answer.

We argue that high-quality question decomposition should satisfy these three guiding principles: (1) **Completeness**: The generated sub-questions must collectively cover all of the information needs contained in q . Missing any aspect of a question may result in an incomplete answer that fails to fully meet the user’s intent; (2) **Non-redundancy**: There should be no overlap between the information needs in each sub-question. Redundant information requires unnecessary computation, increases retrieval latency, and can create misleading input for the reasoning; (3) **Table-specificity**: Each sub-question targets a single table or a group of unionable tables. This reduces the complexity of multi-table joins and ensures that each sub-question is answerable within a coherent context.

While LLMs have shown promise when decomposing questions [30, 41, 51] and improved decomposition accuracy compared to previous approaches [14, 25], our empirical analysis (§7.2) reveals that using LLMs directly to decompose questions often fails to meet the above principles, leading to incorrect answers.

3.2 From Question to Table-Specific Sub-questions

To address these limitations and satisfy the principles outlined above, we propose a four-step decomposition approach that consistently aligns sub-questions with the structural information of the table collection.

3.2.1 Identifying Information Needs. To extract the most salient information needs from each question, we first parse a question into a constituency tree¹. Then, we extract all the noun phrases, which represent the core concepts and entities required, and include adjective phrases and verb phrases that may represent the conditions, based on the syntactic labels [17]. This produces a comprehensive set of “information needs” to ensure **Completeness**.

3.2.2 Hybrid Column Matching. To ensure alignment between information needs and columns, we construct a text snippet to represent each column c_j in table T by concatenating the table title T_p , column header h_j , and distinct cell values from c_j , separated by spaces. Since information needs can differ from table contents both lexically and semantically, we use an M3-Embedding (M3) [4], a unified embedding model designed to encode both lexical and semantic features of text.

We then encode all column snippets in the table collection \mathcal{T} using M3 and index them with FAISS to support efficient similarity search. For each information need I_q , we encode an M3 embedding and retrieve the top-30 column snippets based on their similarity scores. A retrieval depth of 30 was chosen empirically to balance the trade-off between relevance and noise (number of irrelevant columns for the information need).

3.2.3 Context-Aware Column Disambiguation. For each information need I_q , we identify the candidate columns C_{I_q} , along with their similarity scores $\text{sim}(I_q, c_{I_q})$ for each $c_{I_q} \in C_{I_q}$. To guide sub-question generation, we want to select the most promising mapping

¹We use the Stanford Stanza toolkit. See <https://stanfordnlp.github.io/stanza/>.

You are an expert in multi-hop question decomposition for table-based question answering.

Task: Decompose a complex question into a sequence of simpler, minimal sub-questions. Each sub-question must satisfy the following guidelines:

- Use the key phrases grouped together in each entry from a provided list called “information needs”.
- Preserve the full meaning and intent of the original question without omitting important details.
- Ensure the sub-questions are minimal, non-redundant, and natural.

Response Format:

{"Sub-questions": [...]}

Example:

[Question] Among the schools with the average Math score over 560 in the SAT test, how many schools are directly charter-funded?
[Information Needs] [['schools', 'SAT test', 'charter funded'], ['Math score over 560']]

Output:

```
{"Sub-questions": [
  "Which schools have an average Math score over 560?",
  "How many of the schools from #1 are directly charter-funded in the SAT test?"
]}
```

Now decompose the following question:

- [Question] {question}
- [Information Needs] {table-aligned groups}

Figure 3: Question Decomposition Prompt.

between the information needs and the columns and ensure that the information needs are contextually coherent.

Since the information needs for a question q may be interrelated, they should refer to columns from a single table or from a set of joinable tables. To achieve this, we use a table cluster graph \mathcal{G} that contains any required contextual knowledge, and define a “contextual relevance score” that jointly scores all candidate columns for each information need.

Contextual Relevance Score. Given a specific mapping M that assigns a single column from C_{I_q} to I_q for each $I_q \in \mathcal{I}_q$, the *context relevance score* $R(M, \mathcal{I}_q)$ is defined as $\sum_{(I_q, c_{I_q}) \in M} \text{sim}(I_q, c_{I_q})$ if $\mathbb{I}\left(\{T(c_{I_q}) \mid I_q \in \mathcal{I}(q)\}, \mathcal{G}\right) = 1$. Otherwise, $R(M, \mathcal{I}_q) = 0$. Here, $T(c_{I_q})$ denotes the table containing column c_{I_q} , and $\mathbb{I}(\cdot, \mathcal{G})$ is an indicator function returning 1 if all selected tables belong to nodes forming a connected component in \mathcal{G} .

Our objective is to identify a mapping M that achieves the highest contextual relevance score for all possible mappings. However, the number of candidate mappings grows exponentially with the number of information needs, which can be large. So, exhaustive search is not computationally efficient. Therefore, we apply a greedy strategy on \mathcal{G} . The strategy consists of the following steps.

Step 1: Ranking Information Needs. We rank all information needs in a descending order based on the highest similarity score between each information need and its candidate columns.

Step 2: Initializing the Mapping. We initialize the mapping M by selecting the candidate column with the highest similarity score for the top-ranked information need.

Step 3: Expanding the Mapping Progressively. We iterate through the remaining sorted information needs one by one. For each information need, we select the candidate column that (1) belongs to a table that is from the same connected component in \mathcal{G} with the tables already included in the current mapping M , and (2) with the highest similarity score among these connected candidates. If no candidate satisfies these requirements, the current expansion is considered unsuccessful. In such cases, we backtrack to Step 2 to retry the process by selecting the next-best candidate column based on the similarity scores.

The resulting mapping M is obtained until all information needs have been successfully processed.

3.2.4 Question Decomposition. Once the optimal column selection is obtained, we group the information needs based on the tables containing the columns selected. These table-aligned groups are then sent to an LLM to generate a single sub-question for each group, using the prompt as shown in Figure 3. This step **reduces redundancy** by ensuring that each sub-question corresponds to a disjoint set of targeted needs, and leads to more consistent **Table-specificity** by limiting the scope of each sub-question.

4 COVERAGE-AWARE RETRIEVER

Retrieving complete and relevant sets of tables is a key requirement for accurate multi-table question answering. However, existing approaches either incur high computational costs [7] or include overly rigid multi-hop pipelines that increase the error rate [41]. To address this, we introduce Coverage-Aware Retriever to support more efficient multi-hop retrieval, which relies on two key innovations: (1) A coverage scoring function that prioritizes question semantic coverage to improve the effectiveness of early-stage candidate selection (§4.1); (2) A coverage verification function that detects and corrects any missing information using *residual sub-questions* (§4.2).

4.1 Maximizing Question Coverage via Learning-based Scoring

Our approach uses a two-stage pipeline. First, we perform candidate retrieval using FAISS and M3 table embeddings. Then, we apply a more fine-grained reranking model to increase the precision.

Coarse Retrieval. We represent each table cluster in \mathcal{G} as a document by concatenating the table metadata and encoding it using M3. The embeddings produced are then indexed using FAISS. For each sub-question sq_i , we retrieve all of the candidate clusters by querying the index.

Learned Scoring. Coarse retrieval often introduces irrelevant documents due to superficial semantic similarity, which can propagate errors in later stages. To increase precision during candidate selection, we train a scoring function $f_\theta(q, T_q)$ that estimates the *semantic coverage* of a candidate table T_q w.r.t a question q . We use

ColBERTv2 [33] to create f_θ due to its effectiveness in capturing contextual interactions.

To construct the training data, we include single-table QA datasets [22, 47]. Each training instance consists of a question q , a positive table T^+ (i.e., the ground-truth table containing the correct answer), and a hard negative table T^- created by removing the answer-bearing column from T^+ to simulate partial relevance. We train f_θ using the following margin-based ranking loss:

$$\mathcal{L} = \sum_{(q, T^+, T^-)} \max(0, 1 - f_\theta(q, T^+) + f_\theta(q, T^-)).$$

At inference time, f_θ is used to rerank the retrieved candidates based on semantic coverage.

4.2 Ensuring Completeness using Coverage Verification

Although reranking improves the candidate quality, non-relevant information in the retrieval stage may still lead to only partial coverage of the information needs. To address this problem, we introduce a coverage verification function using a gap detection and refinement algorithm.

Connected Table Group Construction. We construct connected table groups by selecting a single table for each sub-question such that the corresponding clusters form a connected component in \mathcal{G} . Each group represents a candidate set of tables that collectively cover the full question. Each group is assigned a score by concatenating all of the tables involved and applying f_θ on the concatenated content and question.

Gap Detection and Refinement. If the score of the top-ranked group is less than a predefined threshold, we assume the coverage is incomplete. Next, we use an LLM to generate a residual sub-question following [52]. We then retrieve and rerank candidate tables using this new sub-question and select an alternative table—connected to the current group—that maximizes the joint coverage score.

We finally assign each table a score based on the best-performing group containing it, and select the top- k tables for answer generation. This design enables higher-precision retrieval that is effective and scalable for large-scale table corpora.

5 SUB-QUESTION GUIDED REASONER

Once the relevant tables are retrieved, a common *answer generation* strategy is to perform the reasoning over tables directly using a sequence-to-sequence model [29]. However, this strategy performs poorly on numerical reasoning tasks, often failing to model necessary arithmetic operations. Program-based reasoning [54] offers a structured alternative – by first generating an executable program (e.g., SQL or Python) that is then executed over multiple tables to produce the final answer. However, existing program generators are unreliable when applied to a collection of soiled tables as they frequently select non-relevant tables, fail to infer which tables can be joined, or contain incorrect reasoning steps, leading to low-quality results [35].

To address these challenges, we propose **Sub-question Guided Reasoner**. Instead of generating the entire program in one shot, our approach incrementally constructs the program using a sequence

of decomposed sub-questions. Each sub-question can be reliably answered using a single table or a unionable group, and inter-sub-question dependencies determine which intermediate results should be joined to get a complete reasoning program. The resulting program is then executed to derive the final answer. This design leads to a more modular, interpretable, and accurate reasoning pipeline, capable of supporting diverse types of reasoning programs over siloed retrieved tables.

5.1 Program Generation

Given the original question, a set of decomposed sub-questions, and the retrieved tables, our reasoner generates a program in two stages:

CoT-Guided Multi-step Program Generation. Following recent advances in multi-hop reasoning [40], we integrate chain-of-thought (CoT) prompting to generate each program step by step based on the sequence of sub-questions. The process begins by generating an initial sub-program that uses a unionable group of tables relevant to the initial sub-question. Then, for each subsequent sub-question, the program is incrementally improved by joining the current intermediate program with a new sub-program generated for that sub-question. The final executable program is obtained when all sub-questions have been processed.

This step-by-step process provides two benefits: (1) It explicitly encodes reasoning to solve high complexity tasks using table dependencies as input, and (2) It enables more robust program construction by constructing each reasoning step based on smaller, coherent subsets of the table collection. To handle column-level inconsistencies when joining unionable tables, we also include a fuzzy-join operator [16] when necessary.

Execution-guided Refinement. Despite the use of carefully constructed prompts, the generated program may still contain errors (e.g., syntax errors). To improve robustness, we introduce an *execution-guided refinement* step in our solution. We first execute the program using the retrieved tables and check the output for any failures. If errors are detected, we re-prompt the LLM and include the error message to help refine the program. This process is repeated until a valid program is generated or a maximum retry limit is reached.

Our sub-question guided approach ensures that the reasoning program aligns with the table structure and retrieved content, enabling more reliable answers in complex multi-table scenarios.

6 DATA PREPARATION FOR EVALUATION

In this section, we first examine the limitations of existing benchmarks in evaluating our problem setting. Then, we introduce our designed solutions to prepare the datasets for our evaluation.

6.1 Motivation

To the best of our knowledge, there are no existing benchmarks that support the evaluation of numerical MTQA over tables in the wild. Two main gaps exist. First, datasets commonly adapted for MTQA evaluation [7, 41] are directly sourced from text-to-SQL benchmarks, such as Spider [47], and Bird [22]. These datasets typically consist of tables from one or multiple databases, with limited scale, complete metadata, and focus only on explicit PK-FK

joins. Thus, they fail to reflect the *scalability*, *incomplete metadata*, and *complex table relationships* for tables in the wild. Second, recent large-scale table benchmarks that aggregate tables from diverse sources (e.g., LakeBench [8]) primarily focus on table discovery tasks (i.e., finding joinable or unionable tables) given a query table. Thus, they lack the real-world question-answer pairs to evaluate our MTQA problem.

6.2 Preparation Process

To address these gaps, our data preparation follows two steps. First, we collect real-world question-answer pairs that are explicitly grounded in multiple relevant tables. Second, we expand these grounded tables into a large-scale table repository, ensuring the *scalability*, *incomplete metadata*, and *complex table relationships* characteristic of tables in the wild.

To realize this, we repurpose the question-answer pairs and their grounded tables from existing text-to-SQL benchmarks, Spider [47] and Bird [22]. The resulting datasets are referred to as SpiderWild and BirdWild. In the next, we first introduce how we prepare tables in the wild by repurposing grounded tables and introducing external tables, and then make the question annotation.

6.2.1 Repurposing Grounded Tables. For each text-to-SQL benchmark, we begin by aggregating all tables from the original databases as the ground tables into a centralized repository. To reflect the characteristics of tables in the wild, we then apply a three-stage table transformation pipeline:

Stage 1: Table Decomposition. To overcome the limited scale and inter-table relationship diversity of the original datasets, we adopt the idea of table decomposition inspired by [8]. Different from their random decomposition, our decomposition is designed to produce more semantically meaningful tables. Specifically, we decompose large tables (i.e., those with more than 5 columns and 50 rows) into multiple disjoint subtables using both column-wise and row-wise splitting strategies. This design is motivated by common organizational patterns observed in data lakes, where tables are often constructed using semantically related columns (e.g., Year, Month, Day) or partitioned by value ranges or categories [28, 34].

For the column-wise splitting strategy, we first identify key columns (i.e., columns with all distinct values) and use the LLM² to group the remaining non-key columns into semantically related column subsets. To achieve this, we design a custom prompt to instruct the model to cluster columns by topic or theme based on their column headers. Each column group is then combined with a randomly selected key column to form a new subtable. To ensure the decomposed subtables can be joined together to reconstruct the original table (i.e., values from the same row in the original table remain correctly aligned), we verify whether any subtable contains all key columns. If no such subtable exists, an additional subtable containing all key columns is created. Since the original database tables only provide the PK-FK joins, we aim to introduce additional joinability of semantic joins. Thus, we choose subtables derived from the same key column that does not participate in PK-FK joins to create joinability, which we will describe later.

²The LLM utilized here is GPT-4o-mini.

Row-wise splitting strategy is applied after the column-wise step, which further partitions each subtable based on the distribution of a randomly selected non-key column. If a numerical column is chosen, we partition the rows into a random number of buckets (between 5 and 20) based on value ranges (e.g., splitting a sales table by ranges of sale amounts). If a categorical column is chosen, we randomly divide the table into 2 to 20 tables, each containing rows for a specific group of categories (e.g., splitting a sales table by region might yield separate tables for “North America Sales”, “Europe Sales” etc.). This row-wise splitting produces unionable tables that share the same column headers but cover different table content.

Stage 2: Metadata Incompleteness Simulation. To mimic the incomplete metadata, we randomly select 20% of the decomposed tables and mask 50% of their column headers and table titles, using the placeholder MASK based on the incomplete metadata statistics reported in [15, 39]. For example, an employee table with columns emp_id, department, salary, and role may be transformed into MASK table with columns MASK, department, MASK, and role.

Stage 3: Joinability Simulation. Among various types of joinability, fuzzy joins (where values differ slightly in spelling or format) are particularly common in real-world data lakes [13, 23, 53]. To simulate this, we inject value-level variations into the textual key columns of decomposed tables. Specifically, following [13], we introduce perturbations such as typographical errors and character deletions into 20% of the cell values. This transformation creates realistic join scenarios where approximate string matching is required (e.g., “New York” vs. “Nw York”).

6.2.2 Incorporating External Tables. While our table decomposition strategy increases the number of tables from hundreds to thousands, it remains insufficient to reflect the large-scale settings typically encountered in practice, such as the massive collections of web tables available online. To further enhance scalability, we incorporate additional tables from the widely used table corpora WebTables and OpenData [8], considering two dimensions: (1) tables relevant to the questions—to ensure that the tables are useful for answering the question, and (2) tables from a similar domain as existing tables, where domain relevance is estimated using the relevance of table titles—to simulate realistic distractor tables. Specifically, for each collected question or decomposed table, we retrieve the top- N candidate tables from these corpora using BM25 [32] over both the question and the table names. This enrichment significantly enlarges the table repository. For example, with $N=100$, the number of tables increases from 2,210 to 73,688 on SpiderWild, and from 5,136 to 109,949 on BirdWild.

6.2.3 Question Annotation. Finally, we curate a set of numerical questions from the original questions that require reasoning over multiple tables. For each selected question, we then construct annotations consisting of both the ground-truth answer and the set of relevant tables to answer this question. Specifically, we first execute the original SQL query on their provided databases to obtain the ground-truth answer. Next, to determine the relevant tables, we analyze how the original SQL query draws information from multiple tables. To achieve this, we manually rewrite the query to extract the values of the primary key and other relevant columns (e.g., those mentioned in WHERE or GROUP BY clauses) for each

Property	SpiderWild	BirdWild
# Tables	73,688	109,949
% Joinability	9%	13%
Avg. # Columns	4.4	5.0
Avg. # Rows	1,384	8,629
# Numerical Questions	274	461
% Easy/Moderate/Hard	52%/43%/5%	10%/84%/6%
% Relevant Table (2 / >=3)	67% / 33%	64% / 36%
% Incomplete Metadata	42%	48%
% Requiring Unionability	13%	56%

Table 1: Dataset statistics.

relevant table. We then identify the specific table rows and columns containing these values. These records are then mapped back to the decomposed tables in the table repository by locating the tables that contain the corresponding rows and columns. The subset of decomposed tables that jointly cover all these records is regarded as the set of relevant tables for that question.

7 EXPERIMENTAL EVALUATION

Evaluation Goals.

Primary Evaluation Goals. We aim to assess the effectiveness, efficiency, scalability, and robustness of our proposed method in realistic settings of MTQA over tables in the wild. Specifically, we answer the following four questions:

- **Q1: How well does our method retrieve multiple relevant tables and answer numerical questions?** Unlike prior work in open-domain MTQA, which struggles to retrieve relevant tables at scale and generate reliable answers (§2.2), our framework aims to address such limitations. We evaluate the effectiveness and efficiency, in both table retrieval and answer accuracy, against strong baselines. (§ 7.1)
- **Q2: What is the contribution of our key design choices and parameter settings to overall performance?** We perform ablation studies to assess the impact of core components (i.e., preprocessing, decomposer, retriever, and reasoner) and analyze the sensitivity of key parameters such as table joinability and unionability thresholds. (§ 7.2)
- **Q3: How well does our method scale with the size of the table corpus?** We evaluate whether DMRAL maintains high effectiveness and low latency as the number of tables in the table corpora increases. (§ 7.3)
- **Q4: How robust is our method under different challenging scenarios?** In particular, we evaluate its robustness w.r.t. the varying number of involved tables, their relationships (i.e., unionability), and degrees of completeness of metadata. (§ 7.4)

Secondary Evaluation Goal: Comparison with Text-to-SQL in

Multi-Table Retrieval. As discussed in §2.2, text-to-SQL methods are designed for relational databases, where complete metadata and explicitly defined table relationships are available. This setting provides advantages that are not present in our scenario. Nevertheless, both our approach and text-to-SQL methods involve a common step, namely, the identification of relevant tables. Accordingly, we assess the effectiveness of this table identification process.

Top-k	Method	SpiderWild									BirdWild								
		Easy			Moderate			Hard			Easy			Moderate			Hard		
		P	R	F1	P	R	F1	P	R	F1	P	R	F1	P	R	F1	P	R	F1
Top-3	JARUnion	46.5	69.7	55.8	45.2	60.6	51.3	44.8	59.8	50.7	39.7	55.0	45.8	39.2	53.7	45.7	45.0	53.7	44.8
	MMQAUnion	47.2	70.8	65.6	37.9	52.4	43.7	37.9	52.1	43.5	41.0	56.6	49.1	36.7	49.7	41.8	35.9	48.8	41.0
	DMRAL	54.9	82.4	65.9	52.9	71.3	60.2	52.3	70.3	59.4	52.2	71.0	59.6	52.7	70.7	59.8	51.5	69.5	58.6
Top-5	JARUnion	31.3	78.2	44.7	32.3	71.5	44.1	31.9	70.4	43.5	28.6	65.6	39.5	28.4	64.4	39.1	28.1	64.1	38.8
	MMQAUnion	31.1	77.8	44.5	26.4	60.3	36.5	26.2	59.7	36.2	26.7	60.9	36.9	26.9	60.7	37.0	26.3	59.6	36.2
	DMRAL	34.2	85.6	48.9	33.1	74.4	45.4	32.7	73.5	45.0	33.3	74.8	45.4	33.3	74.7	45.7	32.8	73.7	45.0

Table 2: Comparison of table retrieval effectiveness on SpiderWild and BirdWild datasets.

Top-k	Method	SpiderWild			BirdWild		
		Easy	Moderate	Hard	Easy	Moderate	Hard
Top-3	JARUnion	40.8	28.9	10.0	28.8	24.9	1.2
	MMQAUnion	40.8	25.6	14.0	26.7	24.9	3.8
	DMRAL	57.0	43.8	20.0	45.7	43.3	26.9
Top-5	JARUnion	46.5	38.0	30.0	41.7	36.0	3.8
	MMQAUnion	46.5	33.9	28.0	35.1	31.9	11.5
	DMRAL	57.7	41.3	35.0	47.5	44.6	19.2

Table 3: Answer accuracy comparison using top-k retrieved tables .

7.1 Evaluation on Table Retrieval and Answer Generation

7.1.1 Experimental Setup. We evaluate DMRAL on the prepared SpiderWild and BirdWild datasets, with the construction details in §6. There are three primary factors that directly impact MTQA effectiveness: (1) the number of table joins, (2) the presence of union operations, and (3) whether the relevant tables contain incomplete metadata. To facilitate a more fine-grained evaluation over these factors, we categorize the questions in both datasets into three levels of complexity—*Easy*, *Moderate*, and *Hard*—based on the number of table joins and unions required to derive the correct answer, as well as whether incomplete metadata is involved in the associated tables.

Specifically, *Easy* questions involve only two table joins, require no union, and rely on complete metadata. By contrast, *Hard* questions require more than two joins, at least one union, and involve incomplete metadata in their associated tables. Questions that do not fully meet the criteria for either the *Easy* or *Hard* category are classified as *Moderate*. Table 1 summarizes the dataset statistics.

Evaluation Metrics. Following previous work [7, 41], we use Precision@k (P@k), Recall@k (R@k), and F1@k to evaluate the effectiveness for multi-table retrieval. To measure answer accuracy on numerical questions, we adopt Arithmetic Exact Match (EM@k) from [38], which measures the percentage of answers obtained by executing the generated program using top-k retrieved tables that match the ground truth. We choose 3 and 5 for k following the previous study [41].

Implementation & Hardware. We use the BGE-M3³ to encode the questions and tables, due to its effectiveness in previous retrieval

tasks [55]. We include GPT-4.1 mini as the primary LLM model. The prompts used by our framework are provided in §A of our technical report [2].

All experiments were conducted on a server running Red Hat 7.9, equipped with an Intel(R) Xeon(R) E5-2690 CPU, 512GB RAM, and an NVIDIA Tesla P100 GPU with 16GB of memory.

Competitors. We evaluate our approach against JARUnion and MMQAUnion, which are adapted from the state-of-the-art open-domain MTQA systems JAR [7] and MMQA [41], respectively. Both methods were originally designed for databases with complete metadata and do not consider table unionability. To ensure a fair comparison, we apply their original table retrieval strategies on our processed metadata-complete tables and augment the retrieved tables by incorporating unionability. Specifically, for each retrieved table, we expand it by merging all tables in its cluster from our constructed graph to form the final retrieved tables. We then use their methods on these retrieved tables to obtain answers.

7.1.2 Effectiveness for Multi-table Retrieval. Table 2 presents the table retrieval effectiveness results across varying levels of question complexity. We highlight two key observations: (1) DMRAL consistently outperforms all baseline methods across both datasets and all complexity levels. For example, in terms of R@3, it achieves an average relative improvement of 17.8% on SpiderWild and 30.1% on BirdWild compared to the second-best baseline JARUnion. (2) Almost all methods exhibit a noticeable performance drop as question complexity increases. For example, on SpiderWild, the R@3 recall for Easy questions is up to 17% higher than for Hard ones, and on BirdWild the gap is also around 2%. These trends highlight an increased difficulty of retrieving relevant tables when questions require more joinable and unionable tables, especially in the presence of incomplete metadata.

7.1.3 Effectiveness for Answer Generation. Table 3 presents the answer accuracy using the top-k retrieved results. We report EM scores under both Top-3 and Top-5 retrieval settings, stratified by question complexity. We make several key observations: (1) DMRAL consistently outperforms all baselines across both datasets and all complexity levels. For example, under the Top-3 setting, DMRAL achieves relative EM improvements of 40%, 51%, and 100% over the second-best baseline JARUnion on Easy, Moderate, and Hard questions in the SpiderWild dataset, respectively. (2) As question complexity increases, answer accuracy also drops significantly on both datasets. This suggests that more complex questions—those involving more joins, unions, and incomplete metadata present

³https://bge-model.com/bge/bge_m3.html

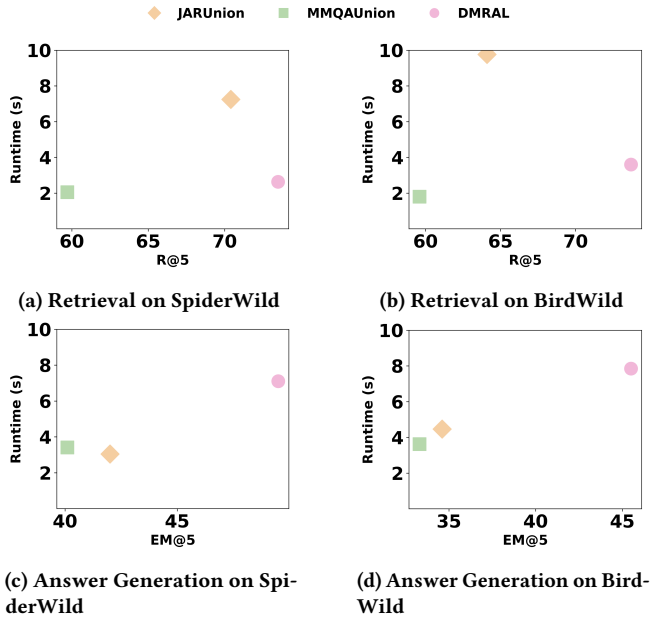


Figure 4: Efficiency breakdown of multi-table retrieval and answer generation, measured by the average runtime per question⁴.

greater challenges not only in identifying the correct tables but also in generating correct answers via reasoning. (3) Interestingly, the drop in answer accuracy across complexity levels is much larger than that of multi-table retrieval. For example, on BirdWild, DMRAL shows only a 2% drop in table retrieval performance R@3 from Easy to Hard questions, but EM@3 drops by 70%. This indicates that even with accurate table retrieval, reasoning over multiple tables becomes significantly more challenging as the complexity of joins, unions, and incomplete metadata increases.

7.1.4 Efficiency for Multi-table Retrieval and Answer Generation.

To evaluate the efficiency of MTQA methods, we measure the average runtime per question for both multi-table retrieval and answer generation across all questions in each dataset. Note that for the runtime of multi-table retrieval of DMRAL, we include the time consumed by both the question decomposer and retriever modules. Figure 4 presents a comparative analysis of efficiency, along with the corresponding effectiveness. From the results, we observe: (1) For table retrieval, DMRAL achieves a strong balance between effectiveness and latency. Compared to the most efficient baseline, MMQAUnion, it yields an average of 15% higher R@5 on both datasets, incurring only a 1.6× increase in runtime. (2) For answer generation, DMRAL delivers substantial accuracy gains—achieving 18% and 31% higher EM@5 scores on SpiderWild and BirdWild, respectively, compared to the second-best method JARUnion, at the cost of an average 3.7 seconds per question. This additional

⁴The method offering the optimal trade-off appears toward the bottom-right in each plot.

Method	SpiderWild	BirdWild
Direct LLM	0.602	0.602
DMRAL	0.697	0.652

Table 4: Comparison of metadata inference quality using Bert-F1 score.

Method	SpiderWild				BirdWild			
	P	R	F1	EM	P	R	F1	EM
Top-3								
No-fill	42.5	55.8	46.0	42.6	40.9	57.0	47.6	28.4
Direct LLM	50.3	67.0	57.4	44.5	48.6	66.2	56.0	40.3
DMRAL	52.3	70.3	59.4	49.0	51.5	69.5	58.6	44.3
GT	59.1	78.0	67.3	51.9	57.2	76.5	65.4	45.2
Top-5								
No-fill	25.3	59.8	35.7	43.0	26.7	58.5	36.6	29.2
Direct LLM	31.4	70.0	43.4	45.0	30.9	70.2	43.0	41.4
DMRAL	32.7	73.5	45.0	49.5	32.8	73.7	45.0	45.5
GT	37.0	81.6	50.9	52.5	36.4	81.1	50.2	46.4

Table 5: The impact of metadata quality to the effectiveness of table retrieval and answer generation.

time cost stems from two reasoning-oriented mechanisms incorporated in the reasoning component of DMRAL to enhance program robustness: (i) a step-by-step CoT prompting strategy to generate the reasoning program incrementally, and (ii) a refinement mechanism that verifies and improves the generated program before final execution.

7.2 Ablation Study and Parameter Study

7.2.1 Ablation Study. We conduct a comprehensive ablation study to evaluate the impact of our key designs – namely the preprocessing (§2.3), decomposer (§3), retriever (§4), and reasoner (§5) components, for the effectiveness.

Preprocessing. We begin by evaluating the inference quality of our metadata inference strategy against a naive baseline, *Direct LLM*, which directly prompts an LLM to complete missing metadata using the whole context. Following prior work [50], we use *BERT-F1* score to assess inference quality: for each missing column header, we compute the F1 score between the predicted and ground-truth header based on contextual embedding similarity, and report the average across all the missing column headers. As reported in Table 4, our proposed strategy yields significantly better scores on both datasets, highlighting the effectiveness of using column-semantics-based subcontexts to enhance metadata inference accuracy.

We further evaluate how metadata quality affects downstream retrieval and answer accuracy by comparing four configurations: (1) *No-fill*, which leaves incomplete metadata unchanged; (2) *Direct LLM*, which fills missing metadata by directly prompting an LLM on the whole context; (3) DMRAL, which infers missing metadata based on subcontext using our proposed inference strategy; (4) *GT*, which uses the corresponding ground-truth metadata to replace all incomplete metadata. The results are presented in Table 5, covering both Top-3 and Top-5 retrieval settings. From the results, we observe that DMRAL consistently improves over the

Method	SpiderWild			BirdWild		
	IRR (↑)	SR (↓)	SAR (↑)	IRR (↑)	SR (↓)	SAR (↑)
<i>Direct LLM</i>	93%	0.589	57%	94%	0.600	59%
DMRAL	96%	0.521	71%	96%	0.522	70%

Table 6: Comparison of decomposition quality. ↑: higher is better, ↓: lower is better.

Method	SpiderWild				BirdWild			
	P	R	F1	EM	P	R	F1	EM
Top-3								
<i>Direct LLM</i>	50.9	68.3	57.8	45.4	49.5	67.2	56.5	38.9
DMRAL	52.3	70.3	59.4	49.0	51.5	69.5	58.6	44.3
Top-5								
<i>Direct LLM</i>	31.5	70.7	43.2	49.1	32.0	72.1	43.9	39.6
DMRAL	32.7	73.5	45.0	49.5	32.8	73.7	45.0	45.5

Table 7: Effectiveness comparison using our decomposer vs. direct LLM-decomposer.

No-fill and *Direct LLM* baselines across all evaluation metrics on both SpiderWild and BirdWild datasets. These results highlight the importance of completing missing metadata and the effectiveness of using relevant subtable context for metadata recovery. Moreover, the performance gap between DMRAL and *GT* setting is relatively narrow, as compared to *No-fill*. This demonstrates that our inference strategy closely approximates gold metadata and significantly contributes to both table retrieval and answer generation in MTQA. **Table-Aligned Question Decomposer.** We first compare our decomposition strategy against a *Direct LLM* generation approach, which uses an LLM to generate the sub-questions without any structural alignment information. To assess the quality of decomposition, we use three metrics corresponding to the criteria defined in §3 – (i) *Information Retention Rate (IRR)* to measure completeness: The proportion of questions whose decompositions successfully preserve all information needs required; (ii) *Subquestion Redundancy (SR)* to measure redundancy: The average pairwise semantic similarity⁵; (iii) *Subquestion-Table Alignment Rate (SAR)* to measure table-specificity: Determine if the number of sub-questions matches the number of relevant tables required to answer the question.

From Table 6, we observe that our table-aligned question decomposer consistently improves all three metrics using both datasets. These results demonstrate the effectiveness of our decomposition strategy which leverages the structure of the table corpus. An additional case study illustrating the improvements possible is shown in §B of our technical report [2].

Next, we deepen our study of the impact of decomposition on both retrieval effectiveness and final answer accuracy in Table 7. Observe that our decomposition consistently achieves better performance on Top-3 and Top-5 retrieval settings. These results demonstrate that DMRAL improves both retrieval effectiveness and final answer accuracy by improving the quality of decomposed sub-questions.

Coverage-Aware Retriever. Now, we evaluate the retrieval effectiveness of the coverage scoring function and coverage verification

⁵The similarity is computed using embeddings generated from a pretrained Sentence-BERT model.

Method	SpiderWild		BirdWild	
	R@3	R@5	R@3	R@5
DMRAL	70.3	73.5	69.5	73.7
<i>NaiveScoring</i>	62.6 (−11)	66.7 (−9)	66.4 (−4)	71.3 (−3)
<i>w/o Verification</i>	56.3 (−20)	60.7 (−17)	55.9 (−20)	64.7 (−12)

Table 8: Ablation study of Coverage-Aware Retriever⁶.

Method	SpiderWild		BirdWild	
	EM@3	EM@5	EM@3	EM@5
DMRAL	49.0	49.5	44.3	45.5
<i>w/o CoT</i>	44.8 (−10)	45.2 (−10)	34.1 (−23)	37.0 (−19)
<i>w/o Refinement</i>	48.3 (−3)	49.7 (−1)	37.8 (−10)	42.0 (−8)

Table 9: Ablation study of Sub-question Guided Reasoner.

submodule using two ablations: (1) *NaiveScoring* is used to replace our trained coverage scoring function using a simple baseline which uses the sum of the individual embedding similarities computed between sub-questions and candidate tables, obtained using our coarse retrieval method described previously. (2) *w/o Verification* disables our residual sub-question generation that is used to retrieve complementary tables. In Table 8, we find: (1) Our trained coverage scoring function provides more effective retrieval performance by identifying subsets of tables that cover the question intent more fully. (2) Our coverage verification submodule enhances retrieval performance by filling potential coverage gaps introduced by non-relevant tables –those that appear to be relevant to the question but cannot provide a complete answer to the question.

Sub-question Guided Reasoner. Assessing the benefit of our reasoning module is achieved by comparing it to two other variants: (1) A *w/o CoT* baseline that prompts the LLM to generate a program in a single shot using the original question and the most similar tables retrieved; (2) A *w/o Refinement* executes the first program generated with no further refinement. The answer accuracy is shown in Table 9, which illustrates that: (1) Our chain-of-thought generation substantially improves answer quality, achieving up to 19% better EM scores. This result demonstrates the benefits of modeling the reasoning process using sub-questions, which enables the LLM to correctly infer table relationships for generating an accurate program, thus producing more reliable answers. (2) Our execution-guided refinement further improves answer quality by ensuring the correctness of the generated program.

7.2.2 Parameter Study. We conduct a parameter study to investigate the impact of different thresholds for *joinability* and *unionability* when modeling table relationships (§ 2.3). We set the threshold ranges—[0.1, 0.3, 0.5, 0.7, 0.9] for joinability and [0.5, 0.6, 0.7, 0.8, 0.9] for unionability. The lower bounds are selected based on their strong pruning effects observed across all table pairs, where they effectively filter out a large portion of low-quality or spurious relationships across tables. Figure 5 illustrates the retrieval and answering effectiveness under varying threshold settings.

⁶Values in parentheses indicate relative performance drops compared to our full method.

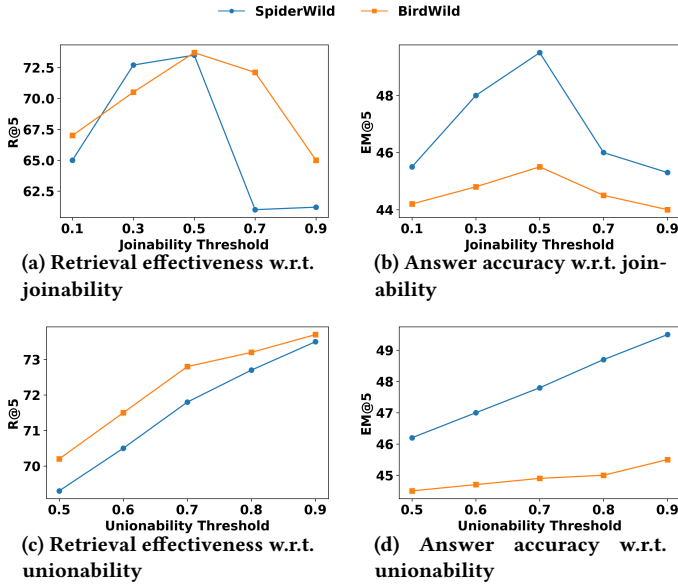


Figure 5: Impact of joinability and unionability thresholds on retrieval and answer accuracy.

From the result, we observe that: (1) For joinability, both the retrieval and answering accuracy improve from 0.1 to 0.5 but decline beyond that point, as the threshold increases. This is because very low thresholds (e.g., 0.1) admit many noisy or spurious links between unrelated tables, which dilute the quality of retrieval and introduce irrelevant candidates for reasoning. On the other hand, very high thresholds (e.g., 0.9) are overly restrictive, excluding moderately joinable but still relevant tables. This reduces coverage and prevents the model from retrieving all necessary tables for reasoning, ultimately harming answer accuracy. (2) For unionability, both performance consistently improves as the threshold increases. This is because stricter thresholds ensure that only highly unionable tables—those with strong semantic similarity—are grouped together. This reduces the inclusion of irrelevant or semantically misaligned tables in union groups, thereby enhancing both retrieval precision and the grounding quality for reasoning.

7.3 Scalability Study

7.3.1 Experimental Setup. To evaluate the scalability of DMRAL under varying sizes of table corpus, we vary the number of top- N external tables retrieved per question/table (§ 6.2.2) with $N \in \{100, 150, 200, 250, 300\}$. These settings yield table corpora containing approximately 109K, 149K, 183K, 214K, and 243K tables, respectively. Since the scale of the table corpus largely influences the table retrieval effectiveness and efficiency, we are mainly focusing on the retrieval effectiveness measured by R@5, and the retrieval efficiency measured by the retrieval time cost per question.

7.3.2 Main results. Figure 6 presents the results. As the size of the table corpus increases, we observe that DMRAL remains robust in retrieval effectiveness—exhibiting only a modest 5% drop in R@5 despite more than doubling the number of tables. In terms of efficiency, the query time grows gradually with scale, indicating

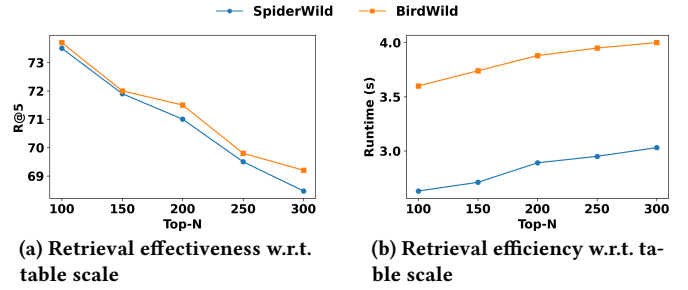


Figure 6: Retrieval effectiveness and efficiency as the number of tables increases for DMRAL.

that the method remains computationally efficient even under a larger table corpus.

7.4 Robustness Cross Challenge Scenarios

Robustness on Questions Requiring a Varying Number of Relevant Tables. To evaluate robustness across questions involving different numbers of tables, we group questions by the number of relevant tables retrieved (2 vs. ≥ 3) and assess both retrieval effectiveness and answer accuracy. As shown in Figure 7, DMRAL consistently outperforms all other baselines for both groupings, demonstrating the robustness of our solution on questions of higher cognitive complexity. This is because: (1) Our retriever explicitly optimizes for question coverage, ensuring that the necessary tables are selected even as the required set grows; (2) Our reasoner generates programs that are guided using decomposed sub-questions, which enables accurate inference of table relationships for questions of increasing difficulty.

Robustness of Questions Involving Incomplete Metadata. We also consider the robustness of our method by comparing the performance on questions with relevant tables, which include incomplete metadata, with the complete metadata case. The results are shown in Figure 8. Observe that: (1) All methods have degraded performance when introducing questions that require tables that have incomplete metadata. This highlights the challenges that incomplete metadata can introduce. (2) DMRAL is consistently the most effective approach for retrieval effectiveness and answer accuracy across both datasets, with notably less performance degradation when compared to the baselines (e.g., an average EM@5 drop of 26% for DMRAL vs. 62% for the baselines). These results demonstrate DMRAL is more robust when incomplete metadata exists.

Robustness on Questions that Require Unionable Tables. We now compare the performance of our method using questions that require unionable tables to answer versus those that do not. The results are shown in Figure 9. Once again, DMRAL is consistently the most effective retrieval model and produces higher-quality answers for both question types and datasets. While existing baselines have notable performance degradation when answering questions that require one or more union operations, DMRAL is consistently good in both scenarios. The advantage stems from our design decisions, which explicitly model table unionability using a graph and groups unionable tables into clusters during the retrieval and reasoning

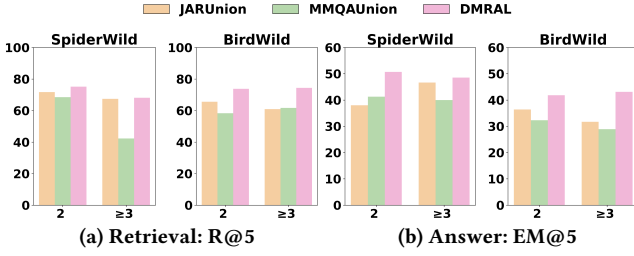


Figure 7: Effect of the number of relevant tables (2 vs. ≥ 3) on effectiveness.

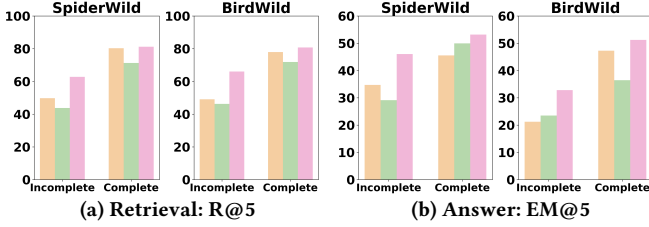


Figure 8: Evaluation on questions involving tables with incomplete vs. complete metadata.

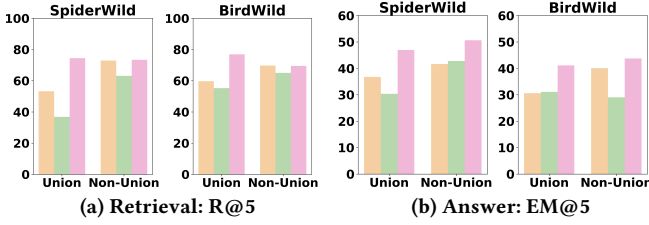


Figure 9: Evaluation on questions that require integrating unionable tables (Union) versus those that do not (Non-Union).

stage. In contrast, the extended baselines assume that each table is independent and fail to incorporate unionability into the retrieval model, reducing their ability to locate multiple joinable tables and effectively use them for reasoning in such questions.

7.5 Comparison with Text-to-SQL Methods

7.5.1 Experimental Setup. Text-to-SQL methods typically assume a simplified MTQA setting, where all tables reside within a relational database, with complete metadata and explicit table relationships (i.e., PK-FK constraints). We evaluate these methods under their standard setup, where all the tables and such structured information are both available for SQL generation.

In contrast, DMRAL is designed for the scenario involving tables in the wild, where explicit inter-table relationships and complete metadata are unavailable. To enable the comparison using DMRAL, we construct a localized table corpus for each question by aggregating all tables from the corresponding relational database. We then apply our full pipeline, treating SQL as the target reasoning program to derive the final SQL query.

Method	Table Selection Accuracy	Execution Accuracy
CHESS	0.733	0.733
OpenSearch-SQL	0.880	0.867
DMRAL	0.867	0.778

Table 10: Comparison with Text-to-SQL baselines.

Dataset. We use the Bird development set [22], and focus on the curated subset of 45 numerical questions that require integration and reasoning over multiple tables.

Competitors. We compare against two resource-efficient and high-performing Text-to-SQL systems on Bird benchmark [1]:

- **CHESS** [36]: A multi-agent LLM-based framework that decomposes SQL generation into four stages—table retrieval, schema pruning, candidate generation, and query validation.
- **OpenSearch-SQL** [44]: A lightweight and modular pipeline that improves SQL generation via structured CoT prompting, SQL-Like intermediate representation, and consistency alignment.

Evaluation Metrics. We adopt two standard metrics commonly used in Text-to-SQL evaluation [22, 47], with table selection accuracy as the primary evaluation metric for table identification:

- **Table Selection Accuracy.** This metric is computed as the percentage of the predicted SQL query references exactly the ground-truth relevant tables, without missing or spurious ones.
- **Execution Accuracy.** This metric is computed as the percentage of predicted SQL query yields the correct answer when executed on the database, matching the executed result using the ground-truth SQL query.

7.5.2 Main Results. Table 10 reports the results on table selection and SQL execution accuracy. From the results, we observe that DMRAL achieves competitive table selection accuracy to the best-performing baseline, OpenSearch-SQL, and clearly outperforms CHESS. This demonstrates its effectiveness in identifying the correct table subset, despite lacking access to explicit PK-FK constraints. In terms of execution accuracy, DMRAL achieves reasonable performance, with most failures (57%) stemming from value errors (i.e., the values used in the predicted SQL query do not align with the values in the database) due to our reasoner is not specifically optimized for SQL generation.

8 CONCLUSION

In this paper, we propose DMRAL, a novel decomposition-guided multi-table retrieval and answering framework designed for MTQA over tables in the wild. DMRAL operates with four modules: Preprocessing, which infers any missing metadata and constructs a graph to capture table relationships; Table-Aligned Question Decomposer and Coverage-Aware Retriever, which jointly improve retrieval by enhancing the quality of decomposition and maximizing the coverage; Sub-question Guided Reasoner improves the answer quality using guided LLMs to generate an accurate executable program based on decomposed sub-questions. Experiments using two pre-pared datasets demonstrate that DMRAL significantly outperforms existing state-of-the-art MTQA methods, achieving an average improvement of 24% in table retrieval and 55% in answer accuracy.

REFERENCES

- [1] [n.d.]. Bird-SQL Leaderboard. <https://bird-bench.github.io/>.
- [2] [n.d.]. Our Technical Report. https://github.com/JrjessyLuo/multitab_qa_dmr/blob/main/technical_report.pdf.
- [3] Ben Bogin, Jonathan Berant, and Matt Gardner. 2019. Representing Schema Structure with Graph Neural Networks for Text-to-SQL Parsing. In *Proceedings of the 57th Conference of the Association for Computational Linguistics, ACL 2019, Florence, Italy, July 28–August 2, 2019, Volume 1: Long Papers*, Anna Korhonen, David R. Traum, and Lluís Màrquez (Eds.). Association for Computational Linguistics, 4560–4565.
- [4] Jianyu Chen, Shitao Xiao, Peitian Zhang, Kun Luo, Defu Lian, and Zheng Liu. 2024. M3-embedding: Multi-linguality, multi-functionality, multi-granularity text embeddings through self-knowledge distillation. In *Findings of the Association for Computational Linguistics ACL 2024*. 2318–2335.
- [5] Kaiwen Chen, Yueting Chen, Nick Koudas, and Xiaohui Yu. 2025. Reliable Text-to-SQL with Adaptive Abstention. *Proc. ACM Manag. Data* 3, 1 (2025), 69:1–69:30.
- [6] Peter Baile Chen, Fabian Wenz, Yi Zhang, Devin Yang, Justin Choi, Nesime Tatbul, Michael Cafarella, Çağatay Demiralp, and Michael Stonebraker. 2024. BEAVER: an enterprise benchmark for text-to-sql. *arXiv preprint arXiv:2409.02038* (2024).
- [7] 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 Computational Linguistics (Volume 1: Long Papers)*. 2687–2699.
- [8] Yuhao Deng, Chengliang Chai, Lei Cao, Qin Yuan, Siyuan Chen, Yanrui Yu, Zhaoze Sun, Junyi Wang, Jiajun Li, Ziqi Cao, et al. 2024. LakeBench: A Benchmark for Discovering Joinable and Unionable Tables in Data Lakes. *Proceedings of the VLDB Endowment* 17, 8 (2024), 1925–1938.
- [9] Yuyang Dong, Chuan Xiao, Takuma Nozawa, Masafumi Enomoto, and Masafumi Oyamada. 2023. DeepJoin: Joinable Table Discovery with Pre-Trained Language Models. *Proceedings of the VLDB Endowment* 16, 10 (2023), 2458–2470.
- [10] Grace Fan, Jin Wang, Yuliang Li, Dan Zhang, and Renée J Miller. 2023. Semantics-Aware Dataset Discovery from Data Lakes with Contextualized Column-Based Representation Learning. *Proceedings of the VLDB Endowment* 16, 7 (2023), 1726–1739.
- [11] Han Fu, Chang Liu, Bin Wu, Feifei Li, Jian Tan, and Jianling Sun. 2023. CatSQL: Towards Real World Natural Language to SQL Applications. *Proc. VLDB Endow.* 16, 6 (2023), 1534–1547.
- [12] Jonathan Fürst, Catherine Kosten, Farhad Nooralahzadeh, Yi Zhang, and Kurt Stockinger. 2024. Evaluating the data model robustness of text-to-SQL systems based on real user queries. *arXiv preprint arXiv:2402.08349* (2024).
- [13] Xuming Hu, Chuan Lei, Xiao Qin, Asterios Katsifodimos, Christos Faloutsos, and Huzefa Rangwala. 2025. POLYJOIN: Semantic Multi-key Joinable Table Search in Data Lakes. In *Findings of the Association for Computational Linguistics: NAACL 2025*. 384–395.
- [14] Xiang Huang, Sitao Cheng, Yiheng Shu, Yuheng Bao, and Yuzhong Qu. 2023. Question decomposition tree for answering complex questions over knowledge bases. In *Proceedings of the AAAI Conference on Artificial Intelligence*, Vol. 37. 12924–12932.
- [15] Aamod Khatiwada, Grace Fan, Roece Shraga, Zixuan Chen, Wolfgang Gatterbauer, Renée J Miller, and Mirek Riedewald. 2023. SANTOS: Relationship-based Semantic Table Union Search. *Proceedings of the ACM on Management of Data* 1, 1 (2023), 1–25.
- [16] Aamod Khatiwada, Roece Shraga, and Renée J Miller. 2025. Fuzzy Integration of Data Lake Tables. *arXiv preprint arXiv:2501.09211* (2025).
- [17] Zhibo Lan and Shuangyin Li. 2024. PS-SQL: Phrase-based Schema-Linking with Pre-trained Language Models for Text-to-SQL Parsing. In *2024 6th International Conference on Natural Language Processing (ICNLP)*. IEEE, 31–35.
- [18] Hyunji Lee, Sohee Yang, Hanseok Oh, and Minjoon Seo. 2022. Generative Multi-hop Retrieval. In *2022 Conference on Empirical Methods in Natural Language Processing, EMNLP 2022*.
- [19] Boyan Li, Yuyu Luo, Chengliang Chai, Guoliang Li, and Nan Tang. 2024. The Dawn of Natural Language to SQL: Are We Fully Ready? [Experiment, Analysis & Benchmark]. *Proc. VLDB Endow.* 17, 11 (2024), 3318–3331.
- [20] Boyan Li, Jiayi Zhang, Ju Fan, Yanwei Xu, Chong Chen, Nan Tang, and Yuyu Luo. 2025. Alpha-SQL: Zero-Shot Text-to-SQL using Monte Carlo Tree Search. *arXiv preprint arXiv:2502.17248* (2025).
- [21] Fei Li and H. V. Jagadish. 2014. NaLIR: an interactive natural language interface for querying relational databases. In *International Conference on Management of Data, SIGMOD 2014, Snowbird, UT, USA, June 22–27, 2014*, Curtis E. Dyreson, Feifei Li, and M. Tamer Özsu (Eds.). ACM, 709–712.
- [22] Jinyang Li, Binyuan Hui, Ge Qu, Binhua Li, Jiayi Yang, Bowen Li, Bailin Wang, Bowen Qin, Rongyu Cao, Ruiying Geng, et al. 2023. Can Llm already serve as a database interface? a big bench for large-scale database grounded text-to-sqls. *arXiv preprint arXiv:2305.03111* (2023).
- [23] Peng Li, Xiang Cheng, Xu Chu, Yeye He, and Surajit Chaudhuri. 2021. Auto-fuzzyjoin: Auto-program fuzzy similarity joins without labeled examples. In *Proceedings of the 2021 international conference on management of data*. 1064–1076.
- [24] Xinyu Liu, Shuyu Shen, Boyan Li, Peixian Ma, Runzhi Jiang, Yuxin Zhang, Ju Fan, Guoliang Li, Nan Tang, and Yuyu Luo. 2025. A Survey of Text-to-SQL in the Era of LLMs: Where are we, and where are we going? *IEEE Transactions on Knowledge and Data Engineering* (2025).
- [25] Sewon Min, Victor Zhong, Luke Zettlemoyer, and Hannaneh Hajishirzi. 2019. Multi-hop Reading Comprehension through Question Decomposition and Rescoring. In *Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics*. 6097–6109.
- [26] Ali Mohammadjafari, Anthony S Maida, and Raju Gottumukkala. 2024. From natural language to sql: Review of llm-based text-to-sql systems. *arXiv preprint arXiv:2410.01066* (2024).
- [27] Rungsiman Nararatwong, Chung-Chi Chen, Natthawut Kertkeidkachorn, Hiroya Takamura, and Ryutaro Ichise. 2024. DBQR-QA: A Question Answering Dataset on a Hybrid of Database Querying and Reasoning. In *Findings of the Association for Computational Linguistics ACL 2024*. 15169–15182.
- [28] Fatemeh Nargesian, Ken Q Pu, Erkang Zhu, Bahar Ghadiri Bashardoost, and Renée J Miller. 2020. Organizing data lakes for navigation. In *Proceedings of the 2020 ACM SIGMOD International Conference on Management of Data*. 1939–1950.
- [29] 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)*. 6322–6334.
- [30] Yixing Peng, Quan Wang, Licheng Zhang, Yi Liu, and Zhendong Mao. 2024. Chain-of-Question: A Progressive Question Decomposition Approach for Complex Knowledge Base Question Answering. In *ACL (Findings)*.
- [31] Mohammadreza Pourreza and Davood Rafiei. 2023. Din-sql: Decomposed in-context learning of text-to-sql with self-correction. *Advances in Neural Information Processing Systems* 36 (2023), 36339–36348.
- [32] Stephen E. Robertson and Hugo Zaragoza. 2009. The Probabilistic Relevance Framework: BM25 and Beyond. *Found. Trends Inf. Retr.* 3, 4 (2009), 333–389.
- [33] Keshav Santhanam, Omar Khattab, Jon Saad-Falcon, Christopher Potts, and Matei Zaharia. 2022. CoBERTv2: Effective and Efficient Retrieval via Lightweight Late Interaction. In *Proceedings of the 2022 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*. 3715–3734.
- [34] Anish Das Sarma, Lujun Fang, Nitin Gupta, Alon Y Halevy, Hongrae Lee, Fei Wu, Reynold Xin, and Cong Yu. 2012. Finding related tables.. In *SIGMOD Conference*, Vol. 10. 2213836–2213962.
- [35] Jiawei Shen, Chengcheng Wan, Ruoyi Qiao, Jiazhen Zou, Hang Xu, Yuchen Shao, Yueling Zhang, Weikai Miao, and Geguang Pu. 2025. A Study of In-Context-Learning-Based Text-to-SQL Errors. *CoRR abs/2501.09310* (2025).
- [36] Shayan Talaei, Mohammadreza Pourreza, Yu-Chen Chang, Azalia Mirhoseini, and Amin Saberi. 2024. Chess: Contextual harnessing for efficient sql synthesis. *arXiv preprint arXiv:2405.16755* (2024).
- [37] Bailin Wang, Richard Shin, Xiaodong Liu, Aleksandr Polozov, and Matthew Richardson. 2020. RAT-SQL: Relation-Aware Schema Encoding and Linking for Text-to-SQL Parsers. In *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*. 7567–7578.
- [38] Dingziru Wang, Longxu Dou, Xuanliang Zhang, Qingfu Zhu, and Wanxiang Che. 2024. Enhancing Numerical Reasoning with the Guidance of Reliable Reasoning Processes. In *Proceedings of the 62nd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*. 10812–10828.
- [39] Qiming Wang and Raul Castro Fernandez. 2023. Solo: Data Discovery Using Natural Language Questions Via A Self-Supervised Approach. *Proc. ACM Manag. Data* 1, 4 (2023), 262:1–262:27.
- [40] Jason Wei, Xuezhi Wang, Dale Schuurmans, Maarten Bosma, Fei Xia, Ed Chi, Quoc V Le, and Denny Zhou. 2022. Chain-of-thought prompting elicits reasoning in large language models. *Advances in neural information processing systems* 35 (2022), 24824–24837.
- [41] Jian Wu, Linyi Yang, Dongyuan Li, Yuliang Ji, Manabu Okumura, and Yue Zhang. 2025. MMQA: Evaluating LLMs with Multi-Table Multi-Hop Complex Questions. In *The Thirteenth International Conference on Learning Representations*.
- [42] Xianjie Wu, Jian Yang, Linzheng Chai, Ge Zhang, Jiaheng Liu, Xinrun Du, Di Liang, Daixin Shu, Xianfu Cheng, Tianzhen Sun, et al. 2024. TableBench: A Comprehensive and Complex Benchmark for Table Question Answering. *CoRR* (2024).
- [43] Chunyang Xiao, Marc Dymetman, and Claire Gardent. 2016. Sequence-based Structured Prediction for Semantic Parsing. In *Proceedings of the 54th Annual Meeting of the Association for Computational Linguistics, ACL 2016, August 7–12, 2016, Berlin, Germany, Volume 1: Long Papers*. The Association for Computer Linguistics.
- [44] Xiangjin Xie, Guangwei Xu, Lingyan Zhao, and Ruijie Guo. 2025. OpenSearch-SQL: Enhancing Text-to-SQL with Dynamic Few-shot and Consistency Alignment. *arXiv:2502.14913 [cs.CL]* <https://arxiv.org/abs/2502.14913>
- [45] Pengcheng Yin, Graham Neubig, Wen-tau Yih, and Sebastian Riedel. 2020. TaBERT: Pretraining for Joint Understanding of Textual and Tabular Data. In

Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics, ACL 2020, Online, July 5-10, 2020, Dan Jurafsky, Joyce Chai, Natalie Schluter, and Joel R. Tetreault (Eds.). Association for Computational Linguistics, 8413–8426.

- [46] Tao Yu, Michihiro Yasunaga, Kai Yang, Rui Zhang, Dongxu Wang, Zifan Li, and Dragomir Radev. 2018. Syntaxsqlnet: Syntax tree networks for complex and cross-domain text-to-sql task. *arXiv preprint arXiv:1810.05237* (2018).
- [47] Tao Yu, Rui Zhang, Kai Yang, Michihiro Yasunaga, Dongxu Wang, Zifan Li, James Ma, Irene Li, Qingning Yao, Shanelle Roman, Zilin Zhang, and Dragomir R. Radev. 2018. Spider: A Large-Scale Human-Labeled Dataset for Complex and Cross-Domain Semantic Parsing and Text-to-SQL Task. In *EMNLP*. Association for Computational Linguistics, 3911–3921.
- [48] Xiaocan Zeng, Pengfei Wang, Yuren Mao, Lu Chen, Xiaozhe Liu, and Yunjun Gao. 2024. MultiEM: Efficient and Effective Unsupervised Multi-Table Entity Matching. In *2024 IEEE 40th International Conference on Data Engineering (ICDE)*. IEEE, 3421–3434.
- [49] Hanchong Zhang, Ruisheng Cao, Lu Chen, Hongshen Xu, and Kai Yu. 2023. ACT-SQL: In-Context Learning for Text-to-SQL with Automatically-Generated Chain-of-Thought. In *EMNLP (Findings)*. Association for Computational Linguistics, 3501–3532.
- [50] Jiani Zhang, Zhengyuan Shen, Balasubramaniam Srinivasan, Shen Wang, Huzefa Rangwala, and George Karypis. 2023. NameGuess: Column Name Expansion for Tabular Data. In *Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing*. 13276–13290.
- [51] Kun Zhang, Jiali Zeng, Fandong Meng, Yuanzhuo Wang, Shiqi Sun, Long Bai, Huawei Shen, and Jie Zhou. 2024. Tree-of-reasoning question decomposition for complex question answering with large language models. In *Proceedings of the AAAI Conference on artificial intelligence*, Vol. 38. 19560–19568.
- [52] Xuanliang Zhang, Dingzirui Wang, Longxu Dou, Qingfu Zhu, and Wanxiang Che. 2025. MURRE: Multi-Hop Table Retrieval with Removal for Open-Domain Text-to-SQL. In *Proceedings of the 31st International Conference on Computational Linguistics*. 5789–5806.
- [53] Erkang Zhu, Yeye He, and Surajit Chaudhuri. 2017. Auto-join: Joining tables by leveraging transformations. *Proceedings of the VLDB Endowment* 10, 10 (2017), 1034–1045.
- [54] Jun-Peng Zhu, Peng Cai, Kai Xu, Li Li, Yishen Sun, Shuai Zhou, Haihuang Su, Liu Tang, and Qi Liu. 2024. Autotqa: Towards autonomous tabular question answering through multi-agent large language models. *Proceedings of the VLDB Endowment* 17, 12 (2024), 3920–3933.
- [55] Rongzhi Zhu, Xiangyu Liu, Zequn Sun, Yiwei Wang, and Wei Hu. 2025. Mitigating Lost-in-Retrieval Problems in Retrieval Augmented Multi-Hop Question Answering. *arXiv preprint arXiv:2502.14245* (2025).

You are an expert in metadata inference for table-based data understanding.

Task: You are provided with a table containing partially missing metadata, such as incomplete or masked table titles and column headers. Your goal is to infer and recover the missing metadata based on the provided table structure and sampled rows.

Important Instructions:

- You must return exactly the same number of column headers as provided in the "Column Headers".
- Do not add, remove, or reorder any columns—only replace missing ones with your best inference.
- Maintain the original header order.

Response Format:

```
{
  "updated_title": "...",
  "updated_headers": ["...", "...", ...]
}
```

Now perform the metadata inference on the provided table:

- [Table Title] {table_title}
- [Column Headers] {headers}
- [Sample Rows (10 randomly sampled rows)] {sampled_rows}

Figure 10: Metadata Inference Prompt.

You are an expert in organizing tables by grouping related columns into smaller, meaningful subtables.

Task: You are provided with a cluttered table and asked to reorganize its columns into subtables.

Reorganization Guidelines:

- Group columns that naturally belong together.
- Do not rename, remove, or reorder the columns within each group.

Response Format:

```
{
  "Tables": [
    {
      "Table Title": "...",
      "Column Headers": ["...", "...", ...]
    },
    ...
  ]
}
```

Now perform the table reorganization based on the provided table context.

- Table Title: {table_title}
- Column Headers: {column_headers}

Figure 11: Semantic Column Grouping Prompt.

You are an expert in identifying semantic gaps in query decomposition.

Task: Identify any missing tables needed to fully answer the question as a residual sub-question, given a user question and a set of currently available tables. If all required tables are already provided, return None.

Response Format:

```
{"Residual Sub-question": ...}
```

Examples:

[Question] How many female clients opened their accounts in Jesenik branch?

[Provided Tables] financial.client(client_id, gender, birth date, district_id)
financial.disp(disposition_id, client_id, account_id, type)

Output:

```
{"Residual Sub-question": "What is the district name of the Jesenik branch?"}
```

[Question] Among the atoms that contain element carbon, which one does not contain compound carcinogenic?

[Provided Tables] toxicology.atom(atom_id, molecule_id, element)
toxicology.molecule(molecule_id, label)

Output:

```
{"Residual Sub-question": None}
```

Now perform the task on the following input:

- [Question] {question}
- [Provided Tables] {tables}

Figure 12: Residual Sub-question Generation Prompt.

A PROMPTS

In this section, we supplement the main prompts used in this paper. The metadata inference prompt (§2.3) is shown in Figure 10. The residual sub-question generation prompt (§4.2) is shown in Figure 12. The prompts for multi-step program generation and execution-guided refinement (§5.1) are shown in Figure 13 and Figure 14, respectively. The prompt used to prepare our evaluation datasets by grouping non-key columns into column subsets (§6.2) is shown in Figure 11.

B CASE ANALYSIS

We conduct a qualitative comparison between the subquestions produced by the *Direct LLM* approach and the table-aligned question decomposer in DMRAL. Table 11 presents representative examples highlighting the major decomposition issues commonly observed with the *Direct LLM* approach, alongside the improvements achieved by DMRAL. We categorize these issues into three types:

Missing Key Information: Critical elements required to fully specify the question are omitted. For example, in the first case, the *Direct LLM* decomposition misses the mention of "account opened", which is essential to correctly identify the relevant table account.

You are an expert in SQL program synthesis for multi-table question answering over structured tabular data.

Task: Your objective is to reason step by step to generate the correct SQL program. You are provided with:

- A natural language question.
- A set of retrieved tables with metadata.
- A set of decomposed sub-questions (which may be inaccurate or incomplete).
- Optional external knowledge (e.g., mappings between question phrases and columns).

Reasoning Process:

- **Step 1:** Carefully read the question and table metadata to understand the full requirements.
- **Step 2:** Evaluate the decomposed sub-questions. Revise or rewrite them if needed.
- **Step 3:** For each revised sub-question, reason step by step to write the corresponding SQL query.
- **Step 4:** After processing all sub-questions, check if the final sub-question's SQL fully answers the original question.

Response Format:

```
{
  "reasoning": "...",
  "Final SQL": "..."
}
```

Example:

[Question] List the names of female students who have enrolled in more than 3 courses.

[Retrieved Tables] enrollment(enrollment_id, student_id, course_id) student(student_id, name, age, gender)

[External Knowledge] "female" refers to gender = 'F'

[Decomposed Sub-questions] ['Which students have enrolled in more than 3 courses?', 'What are the names of female students from #1?']

Output:

```
{
  "reasoning": "Let's think step by step.
  Step 1: Revised sub-question 1: Which students have enrolled in more than 3 courses?
  SQL: SELECT student_id FROM enrollment GROUP BY student_id HAVING COUNT(course_id) > 3.

  Step 2: Revised sub-question 2: What are the names of female students from #1?
  SQL: SELECT name FROM student WHERE student_id IN (...) AND gender = 'F'.",
  "Final SQL": "SELECT name FROM student WHERE student_id IN (SELECT student_id FROM enrollment GROUP BY student_id HAVING COUNT(course_id) > 3) AND gender = 'F'"
}
```

Now perform the reasoning and SQL program generation on the provided input:

- [Question] {question}
- [Retrieved Tables] {tables}
- [External Knowledge] {external_knowledge}
- [Decomposed Sub-questions] {decomposed sub-questions}

You are an expert SQL correction assistant for multi-table question answering over structured tabular data.

Task: Given a natural language question, a set of related tables with metadata, a faulty SQL query that failed with a SQL execution error, and optional external knowledge (e.g., mappings between question phrases and columns), your task is to correct the SQL query by fixing the identified error.

Response Format:

```
{
  "SQL": "..."
}
```

Now perform the SQL correction based on the provided context:

- [Question] {question}
- [Tables] {tables}
- [Error SQL] {sql}
- [SQL Error] {sql_error}
- [External Knowledge] {external_knowledge}

Figure 14: Program Refinement Prompt.

In contrast, DMRAL explicitly retains this information in the second sub-question (i.e., better completeness).

Redundant Decomposition: The decomposition produces overlapping or repetitive sub-questions. In the second case, both sub-questions generated by *Direct LLM* redundantly reference "*superheroes*", resulting in duplicate query intent for table *superhero*. In contrast, DMRAL avoids such redundancy (i.e., better non-redundancy).

Entangled Sub-question: The decomposition produces a sub-question that entangles multiple distinct information needs into a single sub-question, increasing the answering complexity. In the third case, *Direct LLM* merges "*segment SME*" and "*year 2013*" into a single sub-question, even though both elements are less likely to co-occur in the same context. In contrast, DMRAL separates them into two sub-questions (i.e., better table-specificity).

Figure 13: Chain-of-Thought Guided Program Generation Prompt.

Issue	Question	Direct LLM	DMRAL
Missing Key Information	Among the <i>account opened</i> , how many female customers born before 1950 and stayed in Sokolov?	SubQ1: Which female customers were born before 1950? SubQ2: How many of #1 stayed in Sokolov?	SubQ1: Which female customers were born before 1950 and stayed in Sokolov? SubQ2: How many of #1 opened an account?
Redundant Decomposition	How many superheroes have brown eyes?	SubQ1: Which <i>superheroes</i> have brown eyes? SubQ2: How many of the <i>superheroes</i> from #1?	SubQ1: Which are brown eyes? SubQ2: How many superheroes have #1?
Entangled Sub-question	What was the average monthly consumption of customers in segment SME for the year 2013?	SubQ1: Which customers are in <i>segment SME</i> for the <i>year 2013</i> ? SubQ2: What was the average monthly consumption for #1?	SubQ1: Which customers are in segment SME? SubQ2: What was the average monthly consumption for #1 in year 2013?

Table 11: Examples of decomposition issues observed in *Direct LLM* vs. *DMRAL*.