

SPARQ: A Cost-Efficient Framework for Offline Table Question Answering via Adaptive Routing

Yang Liu¹, Mengyi Yan², Jiao Xue⁴, Weilong Ren³, Yutong Ye¹, Haoyi Zhou¹, Jianxin Li¹, Zhumin Chen²

¹Beihang University

²Shandong University

³Shenzhen Institute of Computing Sciences

⁴Inspur Cloud Information Technology Co., Ltd.

{ly_act, yutongye, haoyi, lijx}@buaa.edu.cn, {yanmy, chenzhumin}@sdu.edu.cn, renweilong@sics.ac.cn, xuejiao02@inspur.com

Abstract—Table Question Answering (TQA), which aims to answer natural language questions over tabular data, has recently attracted growing interest in the database community. While state-of-the-art (SOTA) methods based on online Large Language Models (LLMs) achieve remarkable accuracy, they suffer from several drawbacks, including data privacy risks, high latency, high cost, and overthinking due to excessively long reasoning chains. To address these challenges, we propose SPARQ (Sufficient Precise Adaptive Routing for TableQA), a cost-efficient TQA framework designed for robust offline deployment. SPARQ extends the operator pool for table reasoning and introduces an adaptive query routing mechanism that dynamically selects optimal operators, assisted by a verifier with rollback/fallback strategies. Through extensive evaluations, we demonstrate that SPARQ achieves remarkable performance in an offline setting: it improves accuracy by over 5% on WikiTQ and Tab Fact datasets, while reducing the average end-to-end latency by up to $10.19\times$ on consumer-grade hardware (e.g., RTX 4090). To the best of our knowledge, this is the first systematic framework that deploys offline LLMs to achieve SOTA performance for TQA under realistic hardware constraints, balancing both effectiveness and efficiency. Code, full version and artifacts are provided.¹

I. INTRODUCTION

Table Question Answering (TQA), a core task in table reasoning, bridges natural language understanding and structured data analysis [1], [2]. It aims to derive concise factual answers, detailed explanations, or fact-checking results from a user’s textual query over schema-free tables [3], e.g., *which country was the first to reach new heights with a skyscraper?* TQA is inherently challenging, as it requires translating ambiguous natural language into executable reasoning steps that combine logical, numerical, and textual inference over unstructured or semi-structured tables [4], [5]. Recently, the database community has increasingly studied TQA as a natural interface for querying heterogeneous data without predefined schemas [6]–[12]. Despite these difficulties, TQA holds great promise: it empowers non-experts to query and reason over large-scale data using plain language, producing interpretable and verifiable analytical results [2].

Research in TQA has evolved rapidly across three major stages. (a) *Supervised methods* rely on annotated datasets to train sequence-to-sequence models that map table–query pairs to answers [13]–[18]. (b) *LLM-based approaches* leverage in-context learning, often enhanced by retrieval mechanisms such as RAG [19] and table decomposition [20], [21]. Within this paradigm, diverse strategies have emerged, including Chain-of-Thought prompting [22], multi-agent collaboration [9], and

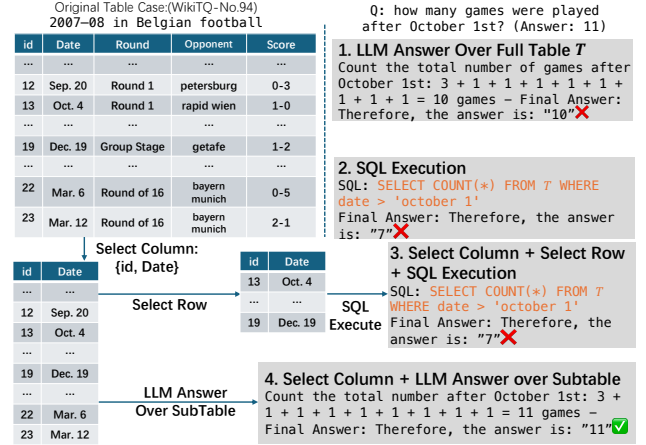


Fig. 1: Comparison of TQA strategies on WikiTQ: (a) offline LLMs (Answer 1), (b) SQL-based reasoning (Answer 2), (c) hybrid method (Answer 3), and (d) a concise yet sufficient Select Column + LLM Answer strategy (Answer 4).

program generation via SQL/Python [23]. (c) *Hybrid frameworks*, e.g., H-STAR [24] and TableRAG [19], integrate symbolic and semantic reasoning via hierarchical operators that repeatedly invoke LLMs and SQL engines for information aggregation. Most of these systems, however, rely on online APIs (e.g., GPT series [25], [26] and PaLM [27]), offering strong interpretability but at the expense of privacy, latency, and cost.

Despite their success, existing methods face three fundamental challenges. (a) *Data and resource constraints*. Training-based models require large-scale annotation, heavy computation, and centralized data that are both costly and often infeasible due to human labor and data privacy concerns [7], [28]. (b) *Reasoning and privacy limitations of LLMs*. Even advanced LLMs struggle with structured or long-form inputs and exhibit weak numerical reasoning, leading to hallucinations and error propagation [6], [29], [30]. Querying online LLMs further exposes sensitive user or table information to external servers, risking privacy leakage [31]. (c) *Efficiency and cost*. Online LLMs incur high financial and latency costs, while hybrid pipelines amplify computation overhead (e.g., invoking multiple LLM calls per query [9], [11]) and often “overthink”, performing redundant reasoning that degrades accuracy and increases response time [7], [10], [21].

Example 1: Fig. 9 illustrates the pitfalls of existing TQA approaches for the query “How many games were played after October 1st?” on a football match table from the WikiTQ

¹<https://github.com/authorlord/SPARQ>

dataset [5], where the correct answer is “11”. (a) Feeding the entire table to an offline LLM (*i.e.*, Answer 1) overwhelms its context window and numerical reasoning ability, producing a hallucinated answer of “10”. This demonstrates the *long-context collapse* problem [6], [30], where LLMs fail to reason effectively over long and structured inputs. (b) Symbolic methods relying on auto-generated SQL (*i.e.*, Answer 2) also fail because the Date column is semi-structured (*e.g.*, “Oct. 4”) and spans multiple years, which string-based comparisons misinterpret, leading to an incorrect result of “7”. (c) Even hybrid pipelines that attempt to filter rows and columns before execution (*i.e.*, Answer 3) simply propagate these underlying errors and returns a wrong answer “7”, since they omit next-year games in rows 20–23. This indicates that increasing reasoning complexity does not ensure correctness; it often adds computational overhead while failing to address fundamental issues such as handling semi-structured data.

These show that the best TQA performance to date is achieved by hybrid approaches powered by online LLMs, yet at the expense of high monetary cost, increased latency, and limited data privacy. Worse still, more complex reasoning pipelines do not necessarily yield better results. These observations naturally raise three key questions: (1) Can we accurately answer TQA queries using offline LLMs while fully preserving user privacy? (2) Can such offline models achieve accuracy comparable to or even surpassing state-of-the-art online systems despite hardware and model-size constraints? (3) Can we accelerate table reasoning by adaptively selecting the optimal strategy for each query, balancing computational efficiency with practical deployability?

Contributions & Organization. To address the challenges discussed above, this paper presents SPARQ (Sufficient and Precise Adaptive Routing for TQA), a cost-efficient framework that enables accurate, privacy-preserving, and efficient offline table reasoning on consumer-level hardware.

The main contributions are summarized as follows:

(1) *Preliminary analysis* (Section III). We systematically analyze the failure modes of existing TQA methods and identify common pitfalls on hard samples. Guided by these insights, we construct a unified library of multi-view table operators that combine symbolic reasoning (*e.g.*, SQL-based retrieval) and semantic reasoning (*e.g.*, dense retrieval), forming a comprehensive and extensible operator pool for table manipulation.

(2) *Framework* (Section IV). We propose SPARQ, a cost-efficient and adaptive framework for offline TQA that integrates a dynamic query router with a mutual-information verifier. It features an extended pool of semantic/symbolic operators for flexible reasoning and employs coordinated CPU–GPU scheduling to ensure efficient and scalable execution.

(3) *Query routing mechanism* (Section V-A). We develop a lightweight and semantic-aware router E that dynamically selects the most suitable operator set for each query based on its predicted complexity. This adaptive mechanism minimizes redundant computation, mitigates LLM overthinking on simple

queries, and prevents information loss from overly aggressive pruning on complex ones, balancing accuracy and efficiency.

(4) *Verification model* (Section V-B). To ensure reliable routing, we introduce a verification model Q that learns a unified representation space for natural-language queries, symbolic code, and tabular data through contrastive learning [32], [33]. By maximizing cross-modal mutual information, Q evaluates in real time whether the pruned subtable and selected operations contain sufficient evidence to answer the query, enabling automatic rollback or supplementation when necessary.

(5) *Scheduling strategy* (Section V-C). We accelerate SPARQ’s execution with a complexity-aware dynamic batching scheduler optimized for offline TQA. Leveraging complexity estimates from E , the scheduler groups queries of similar difficulty and overlaps CPU-based retrieval with GPU-based inference. This joint optimization minimizes idle time and achieves on average $12.19\times$ speedup on 4B–30B models, making SPARQ efficient and practical for on-premise deployment.

(6) *Experimental evaluation* (Section VI). Extensive experiments show that SPARQ, fully deployed offline on two RTX 4090 GPUs, beats existing methods by over 5% in accuracy and achieves average speedup of $10.19\times$ in end-to-end latency.

Finally, we review related work in Section II and conclude the paper in Section VII.

II. PRELIMINARY AND RELATED WORK

Given a question (also called query) q and a corresponding table T , the TQA task aims to generate the answer a , which can be a short-form answer or a long-form natural language description. In recent years, TQA relies heavily on LLMs for understanding natural language and reasoning over structured data. Despite remarkable progress, “out-of-the-box” LLMs still struggle with multi-step reasoning, symbolic computation, and interaction with external knowledge sources [34], [35]. We classify related work into the following categories.

Large Language Models for TQA. LLMs have revolutionized natural language processing and human–computer interaction [7], [36]. They can be categorized into online and offline models. Online models (*e.g.*, GPT [25] and PaLM [27]) offer strong general capabilities but rely on remote APIs, leading to high latency, limited deployability, and privacy risks. Offline models (*e.g.*, Qwen [37], LLaMA [38]) enable local deployment and full data control, but their smaller size results in less stable performance on complex reasoning tasks.

Despite their broad linguistic ability, transformer-based LLMs perform poorly on TQA due to weak numerical reasoning [29], hallucinations [39], and limited understanding of relational structures [8]. Although some research has focused on pre-training table-tailored LLMs to mitigate these issues [8], [40], but they remain costly and underperform in practice. Thus, neither generic nor specialized training-based LLMs provide an ideal solution for TQA.

Natural Language to SQL. Translating natural language questions into executable SQL (NL2SQL) has long been

studied. Early works relied on handcrafted grammars [41], [42] but struggled with linguistic ambiguity. This led to a paradigm shift towards model-based methods, which learn complex mapping between natural language and code from large datasets [43]–[45]. With the rise of LLMs, few-shot prompting enabled direct SQL/code generation [46], [47], enhanced by constrained decoding and syntax verification [48], [49]. Benchmarks like Spider [50], CoSQL [51], WikiSQL [52], and BIRD [53] further accelerated such progress.

Recent TQA systems extend this paradigm through iterative SQL reasoning. ReAcTable [7] integrates step-by-step code execution; LEVER [54] adds verification of execution results; and TabSQLify [23] performs multi-step SQL-based table extraction, feeding subtables back to LLMs for final answers. However, such approaches often suffer from information loss and low SQL parsing success rate with offline LLMs, making purely programmatic pipelines unreliable for robust TQA.

Chain-of-Thought Reasoning. Chain-of-Thought (CoT) enhances LLMs’ multi-step reasoning by decomposing complex problems into sequential steps [35]. Multi-agent frameworks extend this idea by coordinating specialized agents that collaboratively solve complex tasks [9], [10], [55]. However, as a test-time scaling solution [56], CoT-based methods often suffer from *LLM overthinking* [57], [58]: the tendency to apply overly complex reasoning to simple problems. This incurs significant computational overhead without improving final accuracy [11], an issue widely observed in domains like mathematical QA [59], [60].

Methods like DATER [20], CoT [22], H-STAR [24] and RoT [21] adopt CoT-style step-by-step reasoning, while AutoTQA [9], [10] employs multi-agent decomposition for multi-table queries. However, these methods rely heavily on online LLMs; when transferred to smaller offline models, their performance drops sharply and hallucinations increase due to limited context retention. Unlike these paradigms, SPARQ employs a lightweight router to assign queries to pre-defined operators and a verifier to assess sufficiency, balancing effectiveness, efficiency, and stability in offline settings.

III. MOTIVATION AND PRELIMINARY ANALYSIS

Test setting. In practice, when directly migrating online LLM-based solutions to offline LLM-based environment, we observe the following failure mode: (1) offline LLMs still struggle with structured or long-form inputs, and exhibit weak numerical reasoning/code generation ability, leading to hallucinations and error propagation when applying CoT-based or NL2SQL-based solutions; and (2) most online LLM services incur substantial latency costs when migrating into consumer-level hardware, with the tendency of overthinking. We validate these observations through a controlled test where all methods are evaluated under a unified offline setting, as described below.

Datasets. We conducted our tests on the WikiTQ dataset [5], a widely used benchmark for TQA. Following H-STAR [24], we selected 476 challenging samples from its validation set whose

table token count exceeds 4,000, representing complex and long-context reasoning tasks. On average, these tables contain 78.5 (resp. 7.73) rows (resp. columns), up to 753 (resp. 21).

Methods. We tested ten methods: six representative single-operator approaches, two state-of-the-art (SOTA) hybrid methods for TQA, and two hand-drafted hybrid methods.

- 1) Base: a baseline that provides LLM with full table T and three few-shot demonstrations. Following [35], it uses chain-of-thought, prompting LLM to divide queries into several sub-tasks. The demonstrations and CoT prompt are shared across all methods below.
- 2) Select_Column: prompts the LLM to select relevant columns, forming a subtable $T' \subset T$ for the final QA.
- 3) Select_Row: prompts the LLM to generate two SQL queries that retrieve relevant rows; their intersection serves as additional QA context.
- 4) Execute_SQL: directly generates SQL queries to answer the question; three SQL candidates are executed on the full table, and their union result is provided as additional context for the final answer.
- 5) RAG: a retrieval-augmented approach [19] that treats table rows and columns as independent documents and retains up to 50 rows and 10 columns by their semantic similarity with the question.
- 6) RAG +Rewrite: A hybrid retrieval approach that first leverages LLM for multi-level query rewriting, and then combines dense retrieval (semantic similarity) with sparse retrieval (BM25), merged by reciprocal rank fusion. A subsequent re-ranker model weights and reorders these results to select the final rows and columns, retaining up to 50 rows and 10 columns.
- 7) H-STAR [24]: a SOTA CoT-based approach that performs three sequential steps: column selection, row selection, and SQL execution.
- 8) ReAcTable [7]: a SOTA NL2SQL-based approach that iteratively performs 5 rounds of subtable selection and SQL execution, followed by majority voting(*s-vote*).
- 9) Route/Route +Check: following [61], it first prompts LLM to generate a combination of table operators from the operators mentioned above, and then executes them in order. For Route +Check, we manually verify the LLM-selected operators and remove the unnecessary operators.

Configuration. All experiments were conducted using Qwen3-4B-Instruct-2507 [37] as the core offline LLM for TQA, generating SQL/Code and all intermediate results. For retrieval, we adopted a hybrid configuration that combines BGE-m3 for dense embedding, BGE-v2-reranker-m3 for re-ranking, and BM25 for sparse matching [62], [63]. All tests run on a workstation with two NVIDIA RTX 4090 GPUs.

Metric. We used two metrics: (a) accuracy, measured as the ratio of correctly generated answers; and (b) computation cost, quantified by the average number of LLM calls per query, which reflects both efficiency and resource consumption.

Test findings. Fig. 2 shows the performance of different

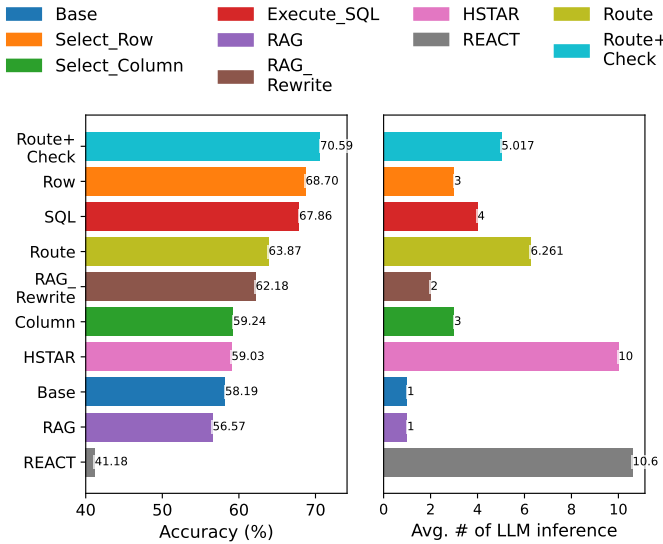


Fig. 2: Preliminary test result.

methods. Our key observations are as follows.

(1) *Offline LLMs struggle with long-form tables.* Base performs poorly, with an accuracy of only 0.582, since most samples contain tables exceeding 4,000 tokens, beyond the effective reasoning range of offline models. This can be mitigated by (a) providing smaller and more informative inputs, *e.g.*, Select_Row outperforms Base by 18.07% by querying the LLM with pruned tables; or (b) leveraging SQL-based computation, *e.g.*, Execute_SQL beats Base by 16.78% by executing LLM-generated SQL queries on the full tables.

(2) *SQL is powerful but brittle.* Select_Row, Select_Column, and Execute_SQL outperform Base by 18.07%, 1.77%, and 16.78%, respectively. This shows that SQL-based methods retain clear advantages over purely LLM-based answering approaches. However, an inspection of their internal execution reveals that SQL statements frequently fail due to formatting errors (*e.g.*, in the WHERE clause); the failure rates of the three methods above caused by such errors are 26.47%, 2.52%, and 21.63%, respectively, which significantly hinder further improvements of SQL-based approaches. That is because although SQL supports precise numerical computation, its intolerance to “dirty” tables severely limits its overall effectiveness.

(3) *Overthinking wastes resources without accuracy gains.* Complex hybrid approaches often consume substantially more resources yet yield limited or even degraded performance, *e.g.*, H-STAR and ReAcTable invoke LLM 3.3 and 3.6 times more frequently than Select_Row, but achieve 13.93% and 40.05% lower accuracy than Select_Row respectively; and worse still, the accuracy of ReAcTable is 29.23% worse than Base despite requiring 9 more LLM inferences per query. This degradation arises because H-STAR’s multi-round pruning discards critical information, while ReAcTable’s majority-voting scheme suffers from inconsistent SQL results across rounds. Although SQL can enhance performance, higher complexity does not necessarily lead to better outcomes,

i.e., simpler SQL queries often yield superior results.

(4) Conventional RAG paradigms are unsuitable for TQA.

RAG methods relying solely on semantic similarity perform poorly on tabular data, *e.g.*, the simple RAG operator scores 2.80% lower than Base. This is because it treats rows and columns as independent documents, which breaks the table’s relational structure and prevents models from capturing dependencies across attributes and tuples. Consequently, essential numerical and categorical relationships are lost, making standard RAG fundamentally incompatible with TQA. However, when properly integrated, retrieval can enhance performance: RAG + Rewrite improves accuracy by 6.85% over Base and 9.91% over simple RAG, by combining query rewriting, dense-sparse retrieval, and re-ranking; this shows its effectiveness for handling ambiguous or complex questions when used properly (see Section VI for more).

(5) TQA benefits from minimal yet sufficient information.

We evaluated two prototype methods, Route and Route + Check, designed to test adaptive operator selection. (a) Route prompts the LLM to generate a set of operators for table pruning, where the final subtable is formed by intersecting the results of all selected operators. It outperforms Base by 9.8% while reducing the average table size in the query context by 73.3%, using an average of 2.08 operators per query. (b) Route + Check further refines the LLM-generated operator sets through manual verification to remove unreasonable operators. This adjustment performs the best among all tested methods, *e.g.*, it improves accuracy by 21.31% over Base and 10.52% over Route, reduces LLM inference time by 19.86%, and simplifies the average operator set size to 1.26. The resulting average table reduction stabilizes at a more moderate 56.42%. These results indicate that (a) LLM-based router tends to over-prune by producing overly complex operator sets, which may lead to insufficient context for reasoning; and (b) identifying a smaller yet sufficient operator set is key to achieving better TQA performance.

Our findings. These results show that no single operator is universally optimal. The best strategy depends on the query’s complexity and table characteristics: simple lookups benefit from lightweight pruning, while analytical questions require SQL-based computation. Improper strategies, however, lead to a trade-off between excessive computation (*overthinking*) and insufficient evidence (*underthinking*).

Among all methods, Route + Check performs best by combining LLM-based operator selection with manual inspection to remove unreasonable operators. Yet, this approach has two key drawbacks: (a) the LLM may hallucinate or generate sub-optimal operator sets; and (b) manual checking is impractical for real-world deployment. Therefore, a practical TQA system must automatically determine suitable operator sets for each query and verify their sufficiency without human involvement.

IV. FRAMEWORK DESIGN

This section introduces SPARQ, an adaptive framework for offline TQA that integrates a dynamic query router with a mutual-information verifier. We first present an extended

TABLE I: Notation Table

Notation	Description
$o^i \in \mathcal{O}$	A pool of table operators, each is denoted as o^i .
$\mathcal{S} \in \mathcal{S}$	A set of operators from \mathcal{O} .
(q, T)	Query(question)-Table pair for TQA task.
E, Q	Router model E and verify model Q .
LLM	A pre-trained offline Large Language Model(LLM).
Context(\mathcal{S})	LLM-generated result after conducting a operator set \mathcal{S}

operator pool used in SPARQ (Section IV-A), followed by an overview of the overall framework (Section IV-B). The detailed design of its components will be described in Section V. Table I summarizes the notations in the paper.

A. An Extended Pool of Semantic and Symbolic Operators

Table operators. Effective table reasoning requires both *semantic* and *symbolic* processing. Semantic operators capture contextual relevance, while symbolic ones ensure precise computation. SPARQ defines a modular operator pool \mathcal{O} that integrates both types, including a new *retrieval operator* that bridges dense retrieval and structured reasoning. Each operator transforms an intermediate subtable T_i into T_{i+1} , enabling compositional reasoning. This modular design avoids the brittleness of one-shot SQL generation in prior end-to-end prompting methods [23], [64], which often fail on offline LLMs due to syntax or execution errors.

Table preprocessing. The preprocessing operator o^{prep} cleans and normalizes the input table T , serving as the first step in TQA. It standardizes text and numbers (lowercasing, merging redundant columns, normalizing malformed entries such as $32k \rightarrow 32000$). For large tables, it also extracts key statistics (Min/Max/Avg) or top- k frequent items to form a compact, noise-reduced input for LLM reasoning.

Column extraction. The column extraction operator o^{col} identifies query-relevant columns through two complementary paths: (1) direct LLM-based enumeration $t^{\text{col}}(\cdot)$ and (2) SQL-based projection $S^{\text{col}}(\cdot)$. Their union $o^{\text{col}}(q, T) = t^{\text{col}}(\cdot) \cup S^{\text{col}}(\cdot)$ yields a subtable T^{col} , reducing hallucination and syntax errors from single-method extraction.

Row extraction. A row extraction operator, o^{row} , is to retrieve query-relevant rows via two complementary paths: (1) direct enumeration based on LLM reasoning $t^{\text{row}}(\cdot)$, and (2) SQL-based selection $S^{\text{row}}(\cdot)$. The union $o^{\text{row}}(q, T) = t^{\text{row}}(\cdot) \cup S^{\text{row}}(\cdot)$ forms the merged subtable T^{row} for later reasoning.

Retrieval operator. A retrieval operator, o^{Ret} , introduces hybrid semantic retrieval that combines dense and sparse search. Given (q, T) , LLM rewrites q into multiple granular forms (*i.e.*, general, balanced and specific) to fit for different retrieval methods [65]. Rows and columns are embedded with SBERT [66] (dense retrieval) and scored by BM25 [67] (sparse retrieval), which are then merged by weighted reciprocal rank fusion [68]; finally the top- M rows and top- N columns are retrieved as subtable. This operator bridges symbolic reasoning with RAG-style retrieval, allowing offline LLMs to efficiently access contextual evidence.

SQL execute operator. An SQL execution operator, o^{SQL} , is to perform symbolic computation by generating and executing SQL from (q, T) . It supports filtering, comparison, and aggregation functions (*e.g.*, COUNT and AVG), and moreover enhances robustness by returning [UNKNOWN] when execution is infeasible, such that useless operations are filtered out.

Remark. Unlike conventional TQA pipelines that rely on fixed symbolic or neural reasoning paths, SPARQ adopts a modular operator design that enables flexible composition and self-correction. In particular, the newly introduced retrieval operator o^{Ret} bridges semantic retrieval with symbolic computation (Section III), allowing offline LLMs to access contextual evidence efficiently, an ability limited to online RAG systems.

B. Overview

We first define the notion of *minimal sufficiency* and then present the overall workflow of SPARQ.

Minimally sufficient. Given a query-table pair (q, T) and a pretrained LLM, a subtable $T^* \subset T$ is *minimally sufficient* if it (a) contains the information necessary to answer q correctly (sufficiency), and (b) minimizes the token length required by the LLM (minimality). However, identifying such a subtable is extremely challenging, as it requires enumerating all possible row-column combinations of T and evaluating each candidate with the LLM, leading to an exponential complexity of $O(2^{m+n} \cdot \text{COST}(\text{LLM}))$, where m and n denote the number of rows and columns in T , respectively, and $\text{COST}(\text{LLM})$ represents the cost of LLM inference on a single subtable. Ambiguous queries or overly long or incomplete tables usually incur a high $\text{COST}(\text{LLM})$ and often lead to error propagation.

Framework of SPARQ. Given a query-table pair (q, T) and a pretrained LLM, SPARQ aims to retrieve a subtable from T that approximates the *minimally sufficient* one; that is, it contains just enough information for the LLM to accurately answer q while minimizing inference cost. In the sequel, we first introduce the core components of SPARQ and then describe its overall workflow.

Component. SPARQ comprises four cooperative components: (1) an operator pool \mathcal{O} that provides a unified interface for both semantic and symbolic table manipulation (Section IV-A); (2) a query router E that dynamically selects an operator set $\mathcal{S} \subseteq \mathcal{O}$ most suitable for a given query-table pair (q, T) , where each operator $o_i \in \mathcal{S}$ is assigned a fitness probability $o_i.p$ to guide the TQA process (Section V-A); (3) a mutual-information verifier Q that evaluates whether the resulting subtable retains sufficient information to answer q , triggering rollback when necessary (Section V-B); and (4) an execution optimizer that orchestrates CPU-based database operations and GPU-based inference in parallel to reduce latency and resource consumption (Section V-C). Together, these components enable SPARQ to approximate a minimally sufficient subtable for accurate TQA without reliance on online LLMs.

Workflow. As shown in Fig. 4, SPARQ takes as input a

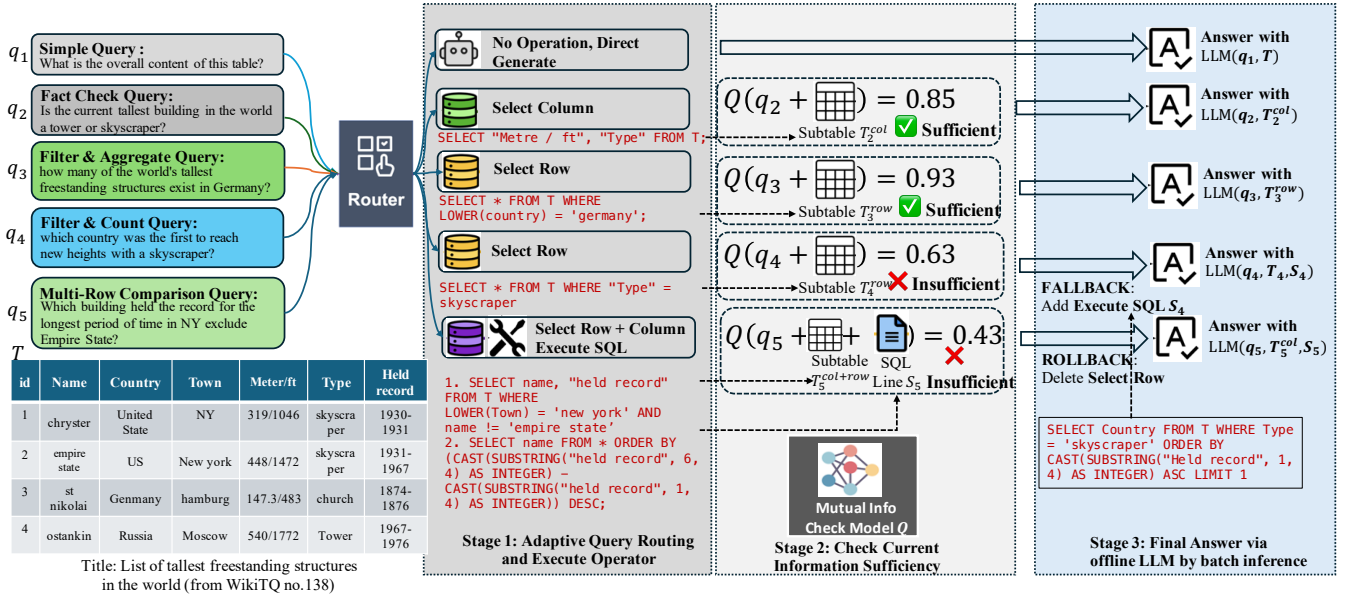


Fig. 3: Overview of SPARQ. Given a query-table pair (q, T) , the router E selects a cost-efficient operator set \mathcal{O} . Operators are executed with parallel batching on CPU/GPU. After each step the verifier Q checks information sufficiency; failures trigger ROLLBACK, and persistent insufficiency triggers a FALLBACK before final QA.

Input: A query-table pair (q, T) , a pretrained offline LLM, and an operator pool \mathcal{O} .

Output: The answer a of query q on T .

/* Offline */

1. $T^{\text{train}} :=$ On a small set of hard samples, collect LLM's predictions from a set \mathcal{S} of operator sets $\mathcal{S} \subset \mathcal{O}$;
2. $E, Q :=$ Train small router and verifier models on T^{train} ;
3. $T :=$ Preprocessing T with table operator o^{prep} ;
4. $\mathcal{S} := E(q, T)$; // Router E selects the suitable \mathcal{S} in \mathcal{S} .
5. $(T^{\mathcal{S}}, \text{Context}(\mathcal{S})) := Q(\mathcal{S})$; // Verifier Q produces a subtable $T^{\mathcal{S}}$ with optional execution result $\text{Context}(\mathcal{S})$.
6. $a :=$ querying LLM with $(q, T^{\mathcal{S}}, \text{Context}(\mathcal{S}))$.

Fig. 4: Workflow of SPARQ.

query-table pair (q, T) and a predefined operator pool \mathcal{O} . It operates in two phases: an *offline phase* and an *online phase*.

In the offline phase, SPARQ selects a small set of hard samples with ground-truth labels (e.g., a subset of the WikiTQ validation set) and collects LLM predictions by executing various operator sets $\mathcal{S} \subset \mathcal{O}$ in a set \mathcal{S} , thereby generating positive and negative supervision signals following a self-distillation paradigm (line 1). It then trains the lightweight router E and verifier Q on this data (line 2).

In the online phase, SPARQ first preprocesses T with o^{prep} for data cleaning and normalization (line 3). The router E selects an operator set $\mathcal{S} \in \mathcal{S}$ for (q, T) , and the LLM executes it to produce intermediate results $\text{Context}(\mathcal{S})$ (e.g., SQL commands from o^{SQL} or rewritten queries from o^{Ret}) (line 4). The verifier Q then assesses \mathcal{S} with the generated context, finalizes execution, and outputs a refined subtable $T^{\mathcal{S}}$ (line 5). Finally, SPARQ queries the LLM with $(q, T^{\mathcal{S}}, \text{Context}(\mathcal{S}))$ to produce the final answer a (line 6).

Remark. Unlike conventional TQA frameworks with fixed symbolic or neural pipelines, SPARQ adopts a modular and adaptive architecture that dynamically balances efficiency and information sufficiency. By coupling the router E with the verifier Q , SPARQ automatically selects and refines operator sets to approximate the minimally sufficient subtable, while the execution optimizer ensures hardware-level efficiency. This closed-loop design enables SPARQ to achieve competitive accuracy, using only offline LLMs with a low cost.

V. MODEL AND ALGORITHM

This section details the core components of SPARQ, including (a) an adaptive query router that dynamically selects a cost-efficient operator set approximating the minimally sufficient subtable (Section V-A); (b) a mutual-information verifier that ensures information sufficiency through rollback control (Section V-B); and (c) system-level optimizations that enable efficient execution on commodity hardware (Section V-C).

A. Adaptive Query Router

Challenge. For a given query-table pair (q, T) and an offline LLM, multiple operator sets may produce the same correct answer. For instance, both a single operator o^{row} and a composite set $\{o^{\text{row}}, o^{\text{SQL}}\}$ can succeed. However, different operators incur varying computational costs: o^{SQL} requires three LLM generations, o^{row} and o^{col} involve two calls, and o^{Ret} performs about $3(m+n)$ similarity computations for a table with m rows and n columns plus one LLM generation for query rewrite (e.g., as general, balanced or specific forms). As the operator set expands, computational cost typically increases, creating an inherent trade-off between reasoning depth and efficiency.

Mechanism. Given (q, T) , LLM, and an operator pool \mathcal{O} ,

SPARQ aims to select an operator set $\mathcal{S} \subseteq \mathcal{O}$ that minimizes both the token length of the final context, $\text{TOKEN}(T^\mathcal{O})$, and the total computational cost, $\text{COST}(\mathcal{O})$, while ensuring that the extracted information remains sufficient to answer the query. To achieve this, SPARQ employs a cost-aware router E that adaptively selects the optimal operator set \mathcal{S} for each query q , balancing reasoning accuracy with computational efficiency. For each operator $o_i \in \mathcal{S}$, the router assigns a fitness probability $o_i.p = P(o_i | q, T, \mathcal{S})$, indicating how likely o_i is to effectively assist the LLM in answering q over T . Once \mathcal{S} is determined, the verifier Q evaluates its information sufficiency and triggers rollback when necessary (Section V-B).

Example 2: As illustrated in Fig. 3, given a table T , the router E dispatches queries of varying complexity (q_1 – q_5) to different operator sets \mathcal{S} based on a cost–benefit analysis:

- $q_1 \rightarrow \{o^{\text{Base}}\}$: q_1 only needs general table description, which can be answered directly without pruning/transformation.
- $q_2 \rightarrow \{o^{\text{col}}\}$: q_2 focuses on specific attributes (e.g., “Meter/ft”, “Type”), requiring column extraction before inference.
- $q_3 \rightarrow \{o^{\text{row}}\}$: q_3 involves filtering rows by “country”, followed by aggregation by the LLM.
- $q_4 \rightarrow \{o^{\text{row}}\}$: q_4 requires filtering/ranking rows to identify the first country, with temporal or ranking reasoning via LLM.
- $q_5 \rightarrow \{o^{\text{col}}, o^{\text{row}}, o^{\text{SQL}}\}$: q_5 demands multi-step reasoning: filtering by location, computing date differences, and comparing durations; accordingly, column extraction, row extraction, and SQL-based arithmetic are applied simultaneously.

Router design. We introduce two designs for the router E .

Training-free router. The training-free router leverages the intrinsic reasoning ability of a pretrained LLM to select the most suitable operator set from \mathcal{O} without additional training. Given a query–table pair (q, T) , SPARQ prompts the LLM with routing instructions and a few in-context examples to predict the optimal operator set from \mathcal{O} . This approach offers high flexibility and zero training cost, making it particularly suitable when data or computational resources are limited.

Training-based router. We introduce a lightweight training-based router as a complementary option when the training-free approach lacks sufficient accuracy. Given the operator pool \mathcal{O} , we define a set of operator sets $\mathcal{S} = \{\mathcal{S}_0, \mathcal{S}_1, \dots, \mathcal{S}_k\}$, where each $\mathcal{S} \in \mathcal{S}$ represents a unique combination of operators from \mathcal{O} , sorted by increasing complexity, i.e., \mathcal{S}_0 denotes the base case (no operator from \mathcal{O} is actually applied), and \mathcal{S}_k corresponds to the most costly configuration.

To construct the training data, we sample a subset of tables T^{train} with ground-truth labels (answers) l from the validation set. For each sampled table $T \in T^{\text{train}}$, we apply every operator set $\mathcal{S} \in \mathcal{S}$ to generate a reduced sub-table $T^\mathcal{S}$ by intersecting the sub-tables produced by executing all operators in \mathcal{S} on T . Given a set of queries q , SPARQ queries the offline LLM with triplets $(q, \text{Context}(\mathcal{S}), T^\mathcal{S})$ to obtain predictions p . Each triplet is assigned a reward score $\mathcal{R}(q, \mathcal{S}, T^\mathcal{S})$, defined as the product of three components: a prediction evaluation

score, an operator-cost penalty, and a context-length penalty:

$$\mathcal{R}(q, \mathcal{S}, T^\mathcal{S}) = \text{acc}(p, l) \times \rho_{\text{cost}}(\mathcal{S}) \times \rho_{\text{len}}(\mathcal{S}), \quad (1)$$

where $\text{acc}(p, l)$ measures the answer accuracy between prediction p and ground-truth label l , $\rho_{\text{cost}}(\mathcal{S}) = \exp(-\lambda_c \cdot \text{COST}(\mathcal{S}))$ penalizes computationally expensive operator sets, and $\rho_{\text{len}}(\mathcal{S})$ penalizes overly long LLM contexts.

Using this training data, the router E is trained as a multi-class classifier to predict the optimal operator-set distribution $P(\mathcal{S} | q, T)$. We formulate this as a knowledge distillation task, where the pre-computed reward \mathcal{R} serves as soft labels. In addition, we use the InfoNCE [69] loss to contrast high-reward and low-reward operator sets, guiding E toward selecting cost-efficient operator paths for queries of varying complexity.

Remark. An operator set with too few operators may produce an overly large context for the LLM, while an excessively long set risks information loss. Both situations can degrade answer accuracy. Verifier Q is proposed to address this (see below).

B. Mutual-Information Verifier Q

Motivation. As discussed in Section V-A, (a) an over-aggressive operator set \mathcal{S} may cause information loss by filtering out critical evidence; and (b) each operator $o_i \in \mathcal{S}$ is associated with a fitness probability $o_i.p$ that o_i can effectively assist the LLM in answering a query q on T . To avoid an operator set causes information loss due to excessive execution, SPARQ introduces a mutual-information verifier Q as a real-time information sufficiency checker.

Inspired by the retrieve–rerank paradigm in RAG systems, Q assesses whether the current subtable $T^\mathcal{S}$, produced by executing the operator set \mathcal{S} on T , contains sufficient information to answer the query q , while additionally considering the intermediate execution result $\text{Context}(\mathcal{S})$ generated by the LLM, which serves as an auxiliary reranker.

Mechanism. Given a query–table pair (q, T) and a subtable $T^\mathcal{S}$ with $\text{Context}(\mathcal{S})$ generated by applying operator set \mathcal{S} on T , Q estimates a sufficiency score $Q(q, T^\mathcal{S}, \text{Context}(\mathcal{S})) \in [0, 1]$ that quantifies the mutual information (MI) between the query and the selected evidence (see below). Here a higher Q indicates that the subtable preserves the essential information required for reasoning, while a lower score suggests potential information loss. If the score falls below a predefined threshold (τ), SPARQ initiates a rollback by reverting to a larger sub-table $T^{\mathcal{S}'}$ derived from an operator subset $\mathcal{S}' \subset \mathcal{S}$, where the operator with the lowest fitness probability $o_i.p$ is excluded. The verifier Q then re-evaluates $T^{\mathcal{S}'}$ iteratively until a sufficient context is found or the base case (\mathcal{S}_0) is reached, i.e., no operator from \mathcal{O} is applied. Such design aims to avoid potential error accumulation via multi-step generation, without invoking LLM re-generation. Only table extraction operator e.g., $\{o^{\text{row}}, o^{\text{col}}, o^{\text{ret}}\}$ will trigger rollback.

Fallback strategy. If the base case still fails to provide sufficient information to answer the query q , SPARQ invokes a fallback mechanism that prompts the LLM to generate an SQL

statement over the full table T to retrieve additional evidence. The retrieved results are then incorporated into the context for a final LLM inference to produce the answer. We treat this *Fallback* as a safeguard, since SQL parsing failures in fallback cases are typically more complex than in other scenarios.

Example 3: Continue with Example 2. We consider two failure cases, q_4 and q_5 , where the router dispatches them wrongly.

Rollback. For q_5 , the operator o^{row} mistakenly filters out the first row by misinterpreting the keyword NY, leading to $T_5^S = \emptyset$. The rollback mechanism then reverts the operator sequence from $\mathcal{S} = \{o^{\text{col}}, o^{\text{row}}, o^{\text{SQL}}\}$ to $\{o^{\text{col}}, o^{\text{SQL}}\}$, restoring the missing evidence. Here, o^{row} is removed because it has the lowest fitness probability $o^{\text{row}}.p$ among all operators in the set.

Fallback. For q_4 , neither o^{row} nor o^{Base} provides sufficient information. Therefore, SPARQ triggers the fallback strategy, prompting the LLM to generate an SQL query (shown in the bottom-right corner of Fig. 3) to retrieve supplementary evidence. The resulting subtable $S_4 = o^{\text{SQL}}(q, T)$, combined with q and T , is then used to produce the final answer.

Model training. We first outline the challenges of training the verifier Q and then describe its training procedure.

Challenges. Training verifier Q is non-trivial because TableQA involves three heterogeneous modalities: (a) the unstructured natural-language query q , (b) the semi-structured table T^S (generated by an operator set \mathcal{S}) containing textual, numerical, and categorical data, and (c) the structured SQL or execution result $\text{Context}(\mathcal{S})$, which encodes logical and numerical evidence. These modalities differ in both representation and granularity, making it difficult to consistently assess whether the information in T^S is sufficient to answer q . A unified representation space is therefore required to align their semantics and enable consistent quantification of information sufficiency.

To tackle above challenge, we define Q as a learnable function that holistically evaluates the alignment between q , T^S , and $\text{Context}(\mathcal{S})$. It is calculated as a weighted aggregation of the semantic alignment scores across all three modality pairs, thereby assessing both the relevance of the evidence to the query and the internal consistency of the evidence itself:

$$Q(q, T^S, \text{Context}(\mathcal{S})) = \sigma(\alpha \cdot f_\theta(q, T^S) + \beta \cdot f_\theta(q, \text{Context}(\mathcal{S})) + \gamma \cdot f_\theta(T^S, \text{Context}(\mathcal{S}))), \quad (2)$$

where the function f_θ produces an alignment score refactored to $[-1, 1]$, and the non-negative coefficients α, β, γ are balancing weights defined in the training objective (see Equation 3).

Training via cross-modal contrastive alignment. Q is trained using a cross-modal contrastive learning framework that maximizes the mutual information between semantically aligned modality pairs: (q, T^S) , (q, S) , and (T^S, S) . Each pair is encoded by a Transformer-based CrossEncoder f_θ , which takes two linearized inputs, $L_T(\cdot)$ for tables and $L_S(\cdot)$ for SQL, and outputs a scalar similarity score $s = f_\theta(\text{text}_1, \text{text}_2)$. We optimize f_θ using a listwise ranking loss that encourages positive pairs to score higher than negatives, effectively max-

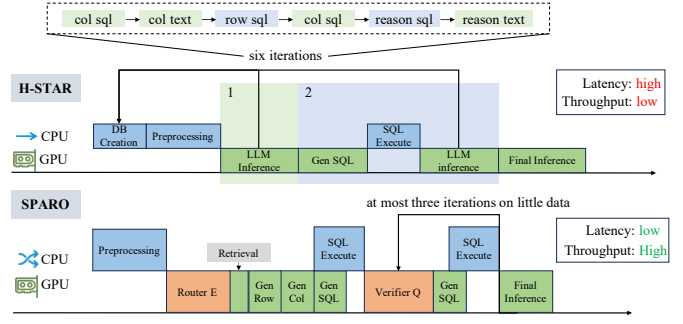


Fig. 5: Scheduling Pipeline of baseline H-STAR and SPARQ. imizing their mutual information [69]. For a generic modality pair (X, Y) , the InfoNCE is defined as:

$$\mathcal{L}_{XY} = -\log \frac{\exp(s^{XY+})}{\exp(s^{XY+}) + \sum_{j=1}^k \exp(s_j^{XY-})}, \quad (3)$$

where s^{XY+} denotes the score of a positive pair and s_j^{XY-} represents the scores of k negative pairs. The total objective aggregates all modality pairs:

$$\mathcal{L}_{\text{total}}(\theta) = \alpha \cdot \mathcal{L}_{qT} + \beta \cdot \mathcal{L}_{qS} + \gamma \cdot \mathcal{L}_{TS}, \quad (4)$$

where α, β , and γ are balancing coefficients controlling the contributions of the query-table, query-SQL, and table-SQL alignment losses, respectively. By default, we set $\alpha = \beta = 1$ and $\gamma = 0.2$, due to the inherent weak ability of f_θ in understanding structural and symbolic context. The training set for Q is T^{train} , augmented with hard negative samples within each batch, e.g., alternative queries, SQL statements, or subtables extracted from the same table T , to enhance discrimination without relying on router E .

Remark. (1) Unlike conventional discriminators that rely on heuristic thresholds or binary feedback [70], [71], Q measures sufficiency in a continuous and differentiable manner, enabling adaptive self-correction during reasoning. (2) We train the verifier Q by aligning three heterogeneous modalities, rather than focusing on pairwise alignments such as query-to-table as in [72], query-to-SQL as in [73] or table-to-SQL/Code in [20].

C. Optimization Strategies

In TQA systems, the end-to-end inference pipeline typically comprises two distinct stages: (1) CPU-intensive tabular preprocessing and database operations; and (2) GPU-accelerated LLMs inference. A naïve implementation executes these two stages sequentially in a single thread, processing one query at a time adapted by most online model systems. While simple, this design leads to severe GPU underutilization due to frequent idle periods, commonly referred to as CPU-GPU bubbles, during which the GPU stalls while waiting for the CPU to complete preprocessing for the next query.

To address this inefficiency, we introduce a batched pipeline execution strategy that decouples and parallelizes the CPU-bound and GPU-bound stages of the workflow. Our design consists of two core components as depicted in Fig. 5.

Multi-threaded preprocessing. We parallelize tabular preprocessing using multithread and optimize the database operations

through an asynchronous scheduling pipeline that enables streaming management of the SQL execution. To enhance robustness, we introduce an optimized parsing algorithm that automatically detects and filters out malformed SQL queries, thereby preventing erroneous queries from blocking the execution pipeline. This approach maximizes CPU throughput and reduces per-query preprocessing latency, especially in database-augmented TQA pipelines where SQL execution and robustness to malformed queries are critical.

Offline batching for GPU inference. Rather than dispatching each preprocessed query immediately, we batch structurally and semantically similar queries to maximize KV cache reuse and throughput. A resource-aware dynamic loader manages multiple GPU models, temporarily loading lightweight components (e.g., routing and verification) and releasing them upon completion to free memory for LLM inference. This offline batching exploits LLM parallelism, substantially improving GPU utilization and efficiency.

Batching combines multiple preprocessed queries into a single GPU operation, enabling concurrent processing that reduces CPU–GPU bubbles. For the offline TQA, this approach is far more effective than multithread (validated in Section VI-C). Notably, our operator set is mutually independent which can execute in parallel, with results combined via intersection. However, the current scheduler does not leverage this property, which we leave for future work.

VI. EXPERIMENT

In this section, we aim to answer the following questions:

- **Q1: Effectiveness v.s. Online Methods.** Can SPARQ achieve competitive or even superior accuracy *w.r.t.* online TQA systems while being fully deployed offline to ensure data privacy?
- **Q2: Effectiveness v.s. Offline Methods.** Can SPARQ consistently outperform other methods when all models are deployed in the same offline environment, particularly across LLMs of different parameter scales?
- **Q3: Efficiency.** How does the scheduling system of SPARQ utilize hardware resources to minimize end-to-end latency and computational overhead during offline inference?
- **Q4: Adaptive.** Can SPARQ dynamically adjust its operator selection and retrieval strategy to handle queries and tables of varying complexity while maintaining information sufficiency?
- **Q5: Benefits.** How do SPARQ’s components (*i.e.*, adaptive router, verification model, and symbolic–semantic operators) contribute to its overall performance and robustness?

A. Experimental Settings

Datasets. We evaluated SPARQ on three widely used TQA benchmarks: (a) *WikiTQ* [5] contains complex questions over Wikipedia tables that require compositional reasoning to produce short and factual answers. (b) *FeTaQA* [3] is a dataset with questions that demand reasoning over tables to generate long-form and descriptive answers. (c) *TabFact* [15] is a fact-checking benchmark where a textual statement must be veri-

fied as either *True* or *False* based on a given table. We fine-tune E and Q with 487 (resp. 476 and 376) labeled hard samples from validation set on WikiTQ (resp. TabFact and FeTaQA).

Baselines. We compared the performance of our SPARQ approach with various strong baselines, including both training-based and training-free paradigms. Note that, considering the nature of offline TQA, we only consider methods with publicly released code and evaluation artifacts.

Training-Based Methods. For WikiTQ and TabFact, we consider LEVER [54], TaPas [13] and TAPEX [14] as Pre-Trained model (PLM)-based baselines. For the LLM-based training method, we adopt Table-GPT (7B) [8] and TableGPT2 (70B) [40]. Following [23], we further evaluate PLM-based T5 [74] on the open-form FeTaQA benchmark.

LLM-Based Training-Free Method. We considered (a) program-aided reasoning methods: DATER [20], TabSQLify [23] and ReAcTable [7]; and (b) hybrid symbolic and semantic approaches: H-STAR [24] and Chain-of-Table [22]. By default, we report DATER and ReAcTable under OpenAI Codex [64], TabSQLify and HSTAR under gpt-3.5-turbo [25], and Chain-of-Tables under PaLM2 [27]. For H-STAR, we additionally report their performance under GPT-4o-mini and llama-3.1-70B [38], denoted as HSTAR_{gpt4} and HSTAR_{70B}, respectively. Since the above methods are built upon different PLMs and LLMs, and some of them are not publicly available (*e.g.*, PaLM-2 [27]), we reported and compared with the best performance in their papers in Table II and III.

Environment. All experiments are conducted on 2 RTX 4090 GPUs powered by 64GB RAM and 32 processors with Intel(R) Xeon(R) Platinum 8432C @2.00GHz. For large-scale baselines (*i.e.*, TABLEGPT2_{70B} and HSTAR_{70B}), we additionally use two A800 GPUs with 80G VRAM. To deploy the QA-LLM, we use vLLM [75] as the backend.

LLM Models and Configurations. For offline LLM backbones of our SPARQ approach, we employ the dense model Qwen3-4B [76](SPARQ_{4B}) and the sparse MoE model Qwen3-30B [77] (SPARQ_{30B}) with FP8 quantization, respectively. SPARQ_{30B} uses the MoE model that activates only 3B parameters during inference, which holds computational cost and resource usage comparable to SPARQ_{4B}. Moreover, the above two models are also used as training-free routers, denoted as E_{4B} and E_{30B} , respectively. For the default training-based router E , we adopt BGE-m3 [78] as the base model, and for the verification module Q , we use BGE-Reranker-v2-m3 [79]. For hyper-parameter setting, we set the sampling number of o^{col} , o^{row} as 2, and o^{SQL} as 3. The acceptance threshold for τ of Q is set to 0.8. The maximum length of selected \mathcal{S} is set to 3 to avoid aggressive pruning. For constructing T^{train} , we set operator pool size as $|\mathcal{O}| = 4$ excluding o^{Base} , and all possible solutions number $k = |\mathcal{S}| = 15$.

Evaluation Metrics. Following prior work [5], [7], [24], we use *denotation accuracy* [5] for WikiTQ, *binary classification accuracy* for TabFact, and *ROUGE-1*, *ROUGE-2* and *ROUGE-*

TABLE II: Results on WikiTQ and TabFact. The best result is marked in **bold**, second best in underline. Entries marked with / were not reported in the paper and are irreproducible with unavailable artifacts.

Method	WikiTQ	TabFact
Approaches requiring training		
TaPas (ACL’20)	48.8	83.9
TAPEX (ICLR’22)	57.5	86.7
LEVER (ICML’23)	62.9	/
Table-GPT (SIGMOD’24)	52.8	/
TABLEGPT2 _{7B}	61.4	77.8
TABLEGPT2 _{70B}	71.4	85.4
Approaches without training		
DATER(SIGIR’23)	65.9	85.6
TabSQLify (NAACL’24)	64.7	79.5
ReAcTable (VLDB’24)	68.0	86.1
Chain-Of-Tables (ICLR’24)	67.3	86.6
HSTAR _{gpt3.5} (NAACL’25)	69.56	86.5
HSTAR _{gpt4}	74.93	89.42
HSTAR _{70B}	75.76	89.23
Ours		
SPARQ _{4B}	<u>77.03</u>	91.30
SPARQ _{30B}	79.79	92.19

TABLE III: Experimental results on FeTaQA.

Model	Rouge-1	Rouge-2	Rouge-L
Approaches requiring training			
T5-small	0.55	0.33	0.47
T5-base	0.61	0.39	0.51
T5-large	0.63	0.41	0.53
Approaches without training			
DATER	0.66	0.45	0.56
ReAcTable	0.71	0.46	<u>0.61</u>
TabSQLify	0.58	0.35	0.48
HSTAR _{gpt3.5}	0.62	0.39	0.52
Chain-of-Tables	0.66	0.44	0.56
Ours			
SPARQ _{4B}	0.60	0.36	0.49
SPARQ _{30B}	<u>0.69</u>	0.51	0.65

L [80] for FeTaQA to evaluate long-form answer quality. Each experiment was run 3 times, and the average is reported.

B. Effectiveness Analysis (Q1&Q2)

SPARQ Outperforms Online Methods (Q1). To answer Q1, we compare our SPARQ approach with various training-based and training-free online-LLM based solutions in Tables II and III, except TABLEGPT2_{7B}, TABLEGPT2_{70B} and HSTAR_{70B} that are offline deployed. From these tables, we can find that SPARQ consistently outperforms all baselines with online models. For example, on the WikiTQ dataset shown in Table II, SPARQ_{4B} surpasses the performance of baselines using the coding-specialized OpenAI Codex, such as DATER and ReAcTable, by approximately 16.89% and 13.28%, respectively. Even compared to more general-purpose online LLMs, SPARQ_{4B} also exhibits significant improvements, leading H-STAR (with GPT-3.5 and GPT-4o-mini backends) by over 10.73% and 2.8%, respectively. Note that our SPARQ approach achieves the above performance by using a model with fewer than 5% of the parameters than these online LLMs. These results validate SPARQ achieves better performance

than online methods with a large number of parameters while achieving low-cost offline deployment to protect data privacy.

SPARQ Outperforms Offline Methods (Q2). To answer Q2, we offline deploy some baseline methods, including large offline LLMs (*i.e.*, TABLEGPT2_{7B}, TABLEGPT2_{70B}, and HSTAR_{70B} in Table III) and small models (*i.e.*, all baselines compared in Table IV).

Offline Large Models. The primary challenge in offline TQA is the scarcity of large amounts of annotated data and the computational resources required to train task-specific models. To address these issues, SPARQ uses an efficient way that fine-tunes only its lightweight router E and verifier Q models on a small set of hard negative samples. However, as shown in Table II compared with pre-training and fine-tuning based offline methods TABLEGPT2_{7B} and TABLEGPT2_{70B}, SPARQ_{4B} can still achieve a better performance by 25.45% and 7.88%, respectively. Even when compared to HSTAR_{70B}, which leverages a 70B offline model, our SPARQ_{4B} still achieves a 1.67% performance improvement with only 5.7% parameter size. This demonstrates that for TQA, our concise and effective SPARQ framework that maintains minimal and sufficient context is a more efficient way to achieve high performance than simply scaling up computational resources and model parameters.

Offline Small Models. Table IV presents a unified comparison where all baseline methods and SPARQ are benchmarked on the same offline models, served via a vLLM-based, OpenAI-compatible server. We observe that baseline methods suffer a significant performance degradation when migrated to offline models, particularly those with fewer parameters (*e.g.*, 4B). As noted in recent literature, smaller models are highly susceptible to hallucinations, formatting errors, and repetitive outputs when confronted with complex tasks like SQL generation or insufficient context [64], [81]. On the Qwen3-4B model, for example, SPARQ outperforms all baselines by over 14% on both TabFact and WikiTQ. This is because baseline methods are prone to overthinking, which relies on long reasoning chains where errors from small models accumulate across consecutive steps (discussed in Section III above). In contrast, SPARQ employs a straightforward yet effective route-check mechanism that detects and rolls back erroneous steps, ensuring the context provided to the LLM remains sufficient. As a direct comparison, H-STAR’s multi-step filtering process results in a SQL execution success rate of only 20.22%, with the accuracy drastically dropping to 8.98% due to information loss. While SPARQ_{4B} maintains a 77.79% SQL execution success rate and exhibits no significant performance degradation compared to SPARQ_{30B}.

SPARQ Demonstrates Stability Across Various LLMs (Q2). Designed with a model-ability-aware scheduling strategy, SPARQ holds consistently high performance across both small and large models with similar latency in table IV. For smaller models, SPARQ employs a form of test-time scaling [82] by increasing the number of samples and enforcing stricter verification thresholds, thereby trading additional LLM infer-

TABLE IV: Comparison of accuracy and end-to-end latency on TabFact and WikiTQ datasets using the same offline backend LLM with different training-free baseline methods. All offline model is deployed with 2 RTX 4090 GPUs. Latency means the average latency per query in seconds measured using 16 threads for all methods if supported.

Method	Qwen2.5-7B				Qwen3-4B				Qwen3-30B-MoE			
	TabFact		WikiTQ		TabFact		WikiTQ		TabFact		WikiTQ	
	Acc	Latency	Acc	Latency	Acc	Latency	Acc	Latency	Acc	Latency	Acc	Latency
ReAcTable	75.84	<u>0.50</u>	47.44	<u>0.70</u>	25.84	<u>0.40</u>	0.58	<u>0.59</u>	66.90	<u>0.56</u>	33.47	<u>0.63</u>
TabSQLify	31.27	4.46	<u>69.27</u>	6.40	6.23	3.14	<u>63.58</u>	5.08	45.06	2.96	<u>72.26</u>	4.46
H-STAR	84.04	1.53	61.28	1.79	<u>60.33</u>	1.66	8.98	2.06	<u>88.24</u>	2.26	70.33	2.48
Chain-of-Tables	21.44	1.04	44.04	1.96	17.74	1.20	49.15	2.02	83.75	2.12	58.31	1.96
SPARQ (Ours)	<u>82.76</u>	0.1931	72.79	0.3920	91.30	0.3112	77.03	0.4949	92.19	0.3295	79.79	0.5527

ences for improved performance. Conversely, for larger models like Qwen3-30B, which possess stronger native capabilities for long-context reasoning and SQL generation, SPARQ uses fewer samples and a more relaxed verification threshold. This allows the LLM to handle simpler problems directly without intermediate steps, reducing latency. Consequently, while baseline methods show performance gains ranging from 10-40% as model size increases, SPARQ maintains a consistently high level of performance across all scales. This validates SPARQ’s ability to strike a dynamic balance between effectiveness and efficiency in different conditions.

C. Efficiency Analysis (Q3)

To answer Q3, in Tables IV and V, we compare the efficiency of our SPARQ with various baseline methods, where all methods are evaluated within an identical offline environment and utilizing the same underlying models. The runtime of most training-free, LLM-based methods can be divided into two primary components: *prediction time*, which encompasses all queries to the LLM (e.g., for CoT generation, SQL generation, and final answer synthesis), and *execution time*, which includes the execution of external tools such as SQL, Python code, or retrieval operators. Table V provides a detailed breakdown of these metrics, where “avg. I/O token throughput” serves as a measure of the LLM’s responsiveness. To minimize end-to-end latency, for Table IV we set multi-thread number to 16 for all methods except TabSQLify which only supports single-thread, while for Table V, we set thread to 1 for all methods, towards a clearer breakdown of prediction/execution steps. Our findings are reported below.

SPARQ Achieves Superior Latency via Coordinated CPU-GPU Scheduling. For the experiments in Table IV, we compare SPARQ with the baselines using our batched pipeline. SPARQ achieves an average speedup of $1.49\times$ (resp. $1.37\times$) over the second-fastest baseline ReAcTable on TabFact (resp. WikiTQ), while consistently outperforming ReAcTable in accuracy. In the experiments reported in Table V, by decoupling and interleaving these heterogeneous CPU-bound and GPU-bound tasks as detailed in Fig. 5, SPARQ achieves an average speedup of over $10.19\times$ in end-to-end latency compared to baseline methods, with improvements reaching up to $17.93\times$ in high-contention scenarios. For LLM inference, SPARQ attains a $5.26\times$ average speedup in prediction time and $8.82\times$ higher output tokens throughput than the baselines,

TABLE V: Runtime and throughput analysis with single thread on WikiTQ dataset using the same offline model Qwen3-4B on 2 RTX 4090 GPUs. Time is in second(s), throughput is in tokens/s.

Method	Avg. prediction time	Avg. Execution time	Avg. I/O token throughput	Avg. End-to-End Latency
TabSQLify	1.96	1.49	2870 / 125	3.45
ReAcTable	3.52	<u>0.07</u>	11601 / 114	3.59
H-STAR	6.53	8.66	1709 / <u>179</u>	15.19
SPARQ _{4B}	0.7612	0.0860	<u>5260</u> / 1184	0.8472

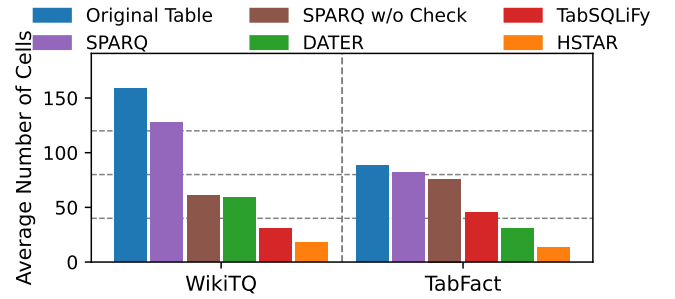


Fig. 6: Comparison of average table cells retrieved for final QA for SPARQ_{4B} and baselines with online LLM baselines.

demonstrating the benefits of batch inference.

D. Retrieval Quality Analysis (Q4)

To answer Q4, in Table VI and Figure 6, we provide a quantitative and qualitative analysis of SPARQ’s retrieval quality under different complexity configurations, where the evaluation metrics follow existing works [7], [23], [24].

SPARQ Achieves Robust Performance Across All Sample Complexities via its Adaptive Strategy. Table VI presents SPARQ’s performance on WikiTQ, stratified by table size. We define tables with fewer than 2000 tokens as *Small*, those with 2000-4000 tokens as *Medium*, and those with over 4000 tokens as *Large*. SPARQ’s adaptive routing, guided by a length-penalty mechanism, effectively allocates operators based on complexity. For *Small* tables, it preserves context, while for *Medium* and *Large* tables, it increasingly uses operators like `Execute_SQL` to condense information and mitigate risks from long inputs. Consequently, SPARQ achieves performance gains of over 10% on *Small* and *Large* tables and 12% on *Medium*, validating the efficacy of routing methodology.

The Verification Module Minimizes Information Loss from Aggressive Routing. Fig. 6 illustrates the average number

TABLE VI: Accuracy performance of various methods on different table sizes on WikiTQ.

Method	Small	Medium	Large
DATER	62.50	42.34	34.62
Chain-of-Table	68.13	52.25	44.87
TabSQLify	68.15	57.91	52.34
H-STAR	71.64	65.20	64.84
SPARQ _{4B}	79.01	74.79	70.34
SPARQ _{30B}	82.26	77.94	74.68

TABLE VII: Ablation study of the accuracy performance for SPARQ_{4B} on TabFact and WikiTQ datasets.

Method	TabFact	WikiTQ
SPARQ _{4B}	91.30	77.04
w. Training-Free Router E_{4B}	86.86	75.53
w/o Verification Model Q	87.20	74.06
w/o Table Extraction	87.10	74.38
w/o SQL Operator	87.70	74.24
w/o Retrieval Operator	87.55	76.26

of table cells retrieved per query by SPARQ compared to baseline methods. We observe that a router model, due to its limited input and representational capacity, is naturally biased towards selecting more complex and aggressive extraction paths, which introduces a significant risk of information loss. The verification model, Q , counteracts this tendency by triggering a rollback mechanism for erroneous selections that result in an insufficient context. This corrective action was activated in 51.9% of table size on the WikiTQ dataset and 7.4% on TabFact. Although this makes SPARQ’s retrieval strategy appear more conservative than the baselines (i.e., it retains more cells on average), this is a direct consequence of prioritizing information sufficiency.

E. Ablation Study (Q5)

To answer Q5, Fig. 7 and Table VII show the ablation study of all SPARQ’s components and an analysis of the operator set S , respectively, demonstrating the effectiveness of each component within the SPARQ framework.

A Training-Free Router Is Not a Plug-and-Play Solution.

As discussed in existing work [61], we find that a training-free router cannot accurately assess the quality and success rate of operators for offline LLMs. Its routing decisions are often overly aggressive (Fig. 7), favoring excessively long operator sequences. This leads to a significant performance degradation on both the WikiTQ and TabFact datasets (Table VII). Similarly, removing the check model Q results in a comparable performance drop, which underscores the critical importance of an information-aware verification step.

The Complementary Roles of Table Extraction and SQL Execution.

The removal of either the table extraction operators or the Execute_SQL operator leads to a notable decline in performance, though these ablations impact different subsets of samples. Specifically, removing table extraction severely impairs the LLM’s ability to parse and reason over long tables. In contrast, removing the Execute_SQL operator deprives the LLM of a mechanism to verify its own calculations, resulting in frequent numerical reasoning errors.

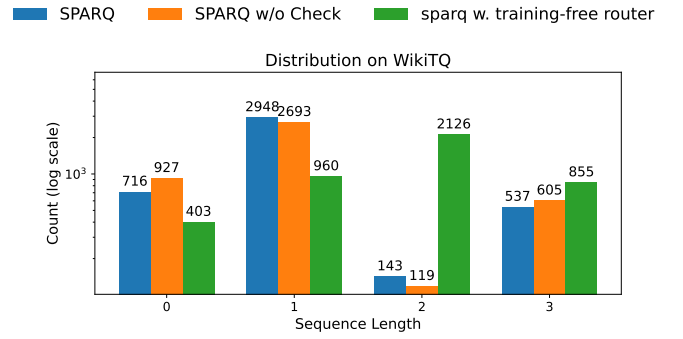


Fig. 7: Distribution of the length of S on WikiTQ by SPARQ_{4B}.

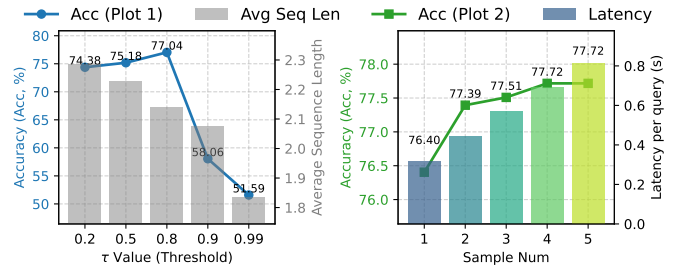


Fig. 8: Hyperparameter Analysis for SPARQ_{4B}.

F. Hyper-parameter Setting

Threshold τ for verify model. As shown in Fig. 8 (left), the value of τ presents a trade-off between TQA accuracy and operational complexity. A permissive threshold (e.g., $\tau = 0.2$) allows Q to approve most operator sets S , which increases the average number of operators by 6.5% but degrades accuracy by 2.66% due to a higher risk of error propagation. A tight threshold (e.g., $\tau = 0.99$) causes Q to reject most operator sets, forcing frequent rollbacks to the Base model or fallbacks to SQL execution, leads to a substantial performance drop of over 25%, where we observe that more than 20% of answers are incorrect, primarily due to errors from processing improperly long contexts or faulty SQL results.

Sampling Number. As shown in Fig. 8 (right), increasing the sample number for o^{row} , o^{col} , o^{SQL} can correct formatting errors, yielding a modest performance gain of 1.32%. However, this benefit comes at a significant latency cost. For complex task of SQL generation, increasing the sample number from 1 to 5 inflates the end-to-end latency by 2.56 \times . To balance accuracy and efficiency, we select a sample number of 2 for o^{row} and o^{col} , and 3 for the more computationally intensive o^{SQL} .

VII. CONCLUSION

In this paper, we presented SPARQ, an offline framework for TableQA that addresses the key challenges of data privacy, latency, and cost in online systems. By integrating an adaptive query router with a mutual-information verifier, SPARQ dynamically selects the most efficient reasoning path for each query while ensuring information sufficiency. We verify the effectiveness and efficiency of SPARQ via extensive tests.

The authors confirm that no part of the paper’s main content, including text, figures, tables, or code, was generated by artificial intelligence (AI) tools. Minor assistance was limited to grammar checking and reference formatting using standard productivity software. No generative AI system was used in the conception, writing, or analysis of the scientific content. This statement is included in accordance with the ICDE 2026 policy on the disclosure of AI-generated content.

REFERENCES

- [1] N. Jin, J. Siebert, D. Li, and Q. Chen, “A survey on table question answering: Recent advances,” in *arXiv preprint arXiv:2207.05270*, 2022.
- [2] X. Zhang, D. Wang, L. Dou, Q. Zhu, and W. Che, “A survey of table reasoning with large language models,” *Frontiers of Computer Science*, vol. 19, no. 9, p. 199348, 2025.
- [3] L. Nan, C. Hsieh, Z. Mao, X. V. Lin, N. Verma, R. Zhang, W. Kryściński, N. Schoelkopf, R. Kong, X. Tang, M. Mutuma, B. Rosand, I. Trindade, R. Bandaru, J. Cunningham, C. Xiong, and D. Radev, “FeTaQA: Free-form table question answering,” *Transactions of the Association for Computational Linguistics*, vol. 10, pp. 35–49, 2022.
- [4] X. Deng, H. Sun, A. Lees, Y. Wu, and C. Yu, “Turl: Table understanding through representation learning,” *ACM SIGMOD Record*, vol. 51, no. 1, pp. 33–40, 2022.
- [5] P. Pasupat and P. Liang, “Compositional semantic parsing on semi-structured tables,” in *Proceedings of the 53rd Annual Meeting of the Association for Computational Linguistics and the 7th International Joint Conference on Natural Language Processing (Volume 1: Long Papers)*. Association for Computational Linguistics, 2015, pp. 1470–1480. [Online]. Available: <https://aclanthology.org/P15-1142/>
- [6] Y. Chung, G. T. Kakkar, Y. Gan, B. Milne, and F. Özcan, “Is long context all you need? leveraging llm’s extended context for nl2sql,” *VLDB*, vol. 18, no. 8, p. 2735–2747, Sep. 2025.
- [7] Y. Zhang, J. Henkel, A. Floratou, J. Cahoon, S. Deep, and J. M. Patel, “ReAcTable: Enhancing react for table question answering,” *Proc. VLDB Endow.*, vol. 17, no. 8, pp. 1981–1994, 2024. [Online]. Available: <https://db.cs.cmu.edu/publications/>
- [8] P. Li, Y. He, D. Yashar, W. Cui, S. Ge, H. Zhang, D. Rifinski Fainman, D. Zhang, and S. Chaudhuri, “Table-gpt: Table fine-tuned gpt for diverse table tasks,” *Proceedings of the ACM on Management of Data*, vol. 2, no. 3, pp. 1–28, 2024.
- [9] J.-P. Zhu, P. Cai, K. Xu, L. Li, Y. Sun, S. Zhou, H. Su, L. Tang, and Q. Liu, “Autotqa: Towards autonomous tabular question answering through multi-agent large language models,” *VLDB*, vol. 17, no. 12, pp. 3920–3933, 2024.
- [10] —, “Unitqa: A unified automated tabular question answering system with multi-agent large language models(demo),” in *SIGMOD*, 2025, pp. 279–282.
- [11] G. Xiao, D. He, J. Wang, and M. Balazinska, “Cents: A flexible and cost-effective framework for llm-based table understanding,” *Proceedings of the VLDB Endowment*, vol. 18, no. 11, pp. 4574–4587, 2025.
- [12] A. Biswal, L. Patel, S. Jha, A. Kamsetty, S. Liu, J. E. Gonzalez, C. Guestrin, and M. Zaharia, “Text2sql is not enough: Unifying ai and databases with tag,” *CIDR*, 2025.
- [13] J. Herzig, P. K. Nowak, T. Müller, F. Piccinno, and J. Eisenschlos, “TaPas: Weakly supervised table parsing via pre-training,” in *ACL*. Association for Computational Linguistics, 2020, pp. 4320–4333.
- [14] Q. Liu, B. Chen, J. Guo, M. Ziyadi, Z. Lin, W. Chen, and J.-G. Lou, “Tapex: Table pre-training via learning a neural sql executor,” in *ICLR*, 2022.
- [15] W. Chen, H. Wang, J. Chen, Y. Zhang, H. Wang, S. Li, X. Zhou, and W. Y. Wang, “Tabfact: A large-scale dataset for table-based fact verification,” in *ICLR*, 2020.
- [16] H. Zhang, Y. Wang, S. Wang, X. Cao, F. Zhang, and Z. Wang, “Table fact verification with structure-aware transformer,” in *EMNLP*, 2020, pp. 1624–1629.
- [17] T. Yu, C.-S. Wu, X. V. Lin, Y. C. Tan, X. Yang, D. Radev, C. Xiong *et al.*, “Grappa: Grammar-augmented pre-training for table semantic parsing,” in *ICLR*, 2021.
- [18] Z. Gu, J. Fan, N. Tang, P. Nakov, X. Zhao, and X. Du, “Pasta: Table-operations aware fact verification via sentence-table cloze pre-training,” in *EMNLP*, 2022, pp. 4971–4983.
- [19] S.-A. Chen, L. Miculicich, J. M. Eisenschlos, Z. Wang, Z. Wang, Y. Chen, Y. Fujii, H.-T. Lin, C.-Y. Lee, and T. Pfister, “TableRAG: Million-token table understanding with language models,” in *Thirty-eighth Conference on Neural Information Processing Systems*, 2024. [Online]. Available: <https://openreview.net/forum?id=41lovPOCo5>
- [20] Y. Ye, B. Hui, M. Yang, B. Li, F. Huang, and Y. Li, “Large language models are versatile decomposers: Decomposing evidence and questions for table-based reasoning,” in *SIGIR*, 2023, pp. 174–184.
- [21] X. Zhang, D. Wang, K. Xu, Q. Zhu, and W. Che, “Rot: Enhancing table reasoning with iterative row-wise traversals,” *arXiv preprint arXiv:2505.15110*, 2025.
- [22] Z. Wang, H. Zhang, C.-L. Li, J. M. Eisenschlos, V. Perot, Z. Wang, L. Miculicich, Y. Fujii, J. Shang, C.-Y. Lee, and T. Pfister, “Chain-of-table: Evolving tables in the reasoning chain for table understanding,” in *ICLR*, 2024. [Online]. Available: <https://openreview.net/forum?id=4L0xnS4GQM>
- [23] M. M. H. Nahid and D. Rafiei, “TabSQLify: Enhancing reasoning capabilities of LLMs through table decomposition,” in *NAACL*, 2024.
- [24] N. Abhyankar, V. Gupta, D. Roth, and C. K. Reddy, “H-STAR: LLM-driven hybrid SQL-Text adaptive reasoning on tables,” in *Proceedings of the 2025 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*. Association for Computational Linguistics, 2025. [Online]. Available: <https://aclanthology.org/2025.naacl-long.445/>
- [25] T. Brown, B. Mann, N. Ryder, M. Subbiah, J. D. Kaplan, P. Dhariwal, A. Neelakantan, P. Shyam, G. Sastry, A. Askell *et al.*, “Language models are few-shot learners,” *NIPS*, vol. 33, pp. 1877–1901, 2020.
- [26] M. Chen, J. Tworek *et al.*, “Evaluating large language models trained on code,” *arXiv preprint arXiv:2107.03374*, 2021.
- [27] R. Anil, A. M. Dai, O. Firat, M. Johnson, D. Lepikhin, A. Passos, S. Shakeri, E. Taropa, P. Bailey, Z. Chen *et al.*, “Palm 2 technical report,” *arXiv preprint arXiv:2305.10403*, 2023.
- [28] M. Scannapieco, I. Figotin, E. Bertino, and A. K. Elmagarmid, “Privacy preserving schema and data matching,” in *Proceedings of the 2007 ACM SIGMOD international conference on Management of data*, 2007, pp. 653–664.
- [29] J. Ahn, R. Verma, R. Lou, D. Liu, R. Zhang, and W. Yin, “Large language models for mathematical reasoning: Progresses and challenges,” *arXiv preprint arXiv:2402.00157*, 2024.
- [30] Q. Zhang, C. Hu, S. Upasani, B. Ma, F. Hong, V. Kamanuru, J. Rainton, C. Wu, M. Ji, H. Li *et al.*, “Agentic context engineering: Evolving contexts for self-improving language models,” *arXiv preprint arXiv:2510.04618*, 2025.
- [31] Y. Zhang, A. Floratou, J. Cahoon, S. Krishnan, A. C. Müller, D. Banda, F. Psallidas, and J. M. Patel, “Schema matching using pre-trained language models,” in *ICDE*. IEEE, 2023, pp. 1558–1571.
- [32] Y. Mao, X. Yan, Q. Guo, and Y. Ye, “Deep mutual information maximin for cross-modal clustering,” in *Proceedings of the AAAI Conference on Artificial Intelligence*, vol. 35, no. 10, 2021, pp. 8893–8901.
- [33] F. Zhang, D. Cai, W. Song, and C. You, “Connecting multi-modal contrastive representations,” in *Advances in Neural Information Processing Systems*, vol. 36, 2023.
- [34] Z. Cheng, T. Xie, P. Shi, C. Li, R. Nadkarni, Y. Hu, C. Xiong, D. Radev, M. Ostendorf, L. Zettlemoyer *et al.*, “Binding language models in symbolic languages,” in *ICLR*, 2023.
- [35] J. Wei, X. Wang, D. Schuurmans, M. Bosma, F. Xia, E. Chi, Q. V. Le, and D. Zhou, “Chain-of-thought prompting elicits reasoning in large language models,” in *Advances in Neural Information Processing Systems*, vol. 35. Curran Associates, Inc., 2022, pp. 24 824–24 837. [Online]. Available: https://proceedings.neurips.cc/paper_files/paper/2022/file/9d5609613524ecf4f15af0f7b31abca4-Paper-Conference.pdf
- [36] T. Yu, R. Zhang, K. Yang, M. Yasunaga, D. Wang, Z. Li, J. Ma, I. Li, Q. Yao, S. Roman, Z. Zhang, and D. Radev, “Spider: A large-scale human-labeled dataset for complex and cross-domain semantic parsing and text-to-SQL task,” in *Proceedings of the 2018 Conference on Empirical Methods in Natural Language Processing*, E. Riloff, D. Chiang, J. Hockenmaier, and J. Tsujii, Eds. Brussels, Belgium: Association for Computational Linguistics, Oct.-Nov. 2018, pp. 3911–3921. [Online]. Available: <https://aclanthology.org/D18-1425/>
- [37] The Qwen Team, “Qwen3 technical report,” *arXiv preprint arXiv:2507.06850*, 2025.

- [38] A. Grattafiori, A. Dubey, A. Jauhri, A. Pandey, A. Kadian, A. Al-Dahle, A. Letman, A. Mathur, A. Schelten, A. Vaughan *et al.*, “The llama 3 herd of models,” *arXiv preprint arXiv:2407.21783*, 2024.
- [39] P. Shojaei, I. Mirzadeh, K. Alizadeh, M. Horton, S. Bengio, and M. Farajtabar, “The illusion of thinking: Understanding the strengths and limitations of reasoning models via the lens of problem complexity,” *arXiv preprint arXiv:2506.06941*, 2025.
- [40] A. Su, A. Wang, C. Ye, C. Zhou, G. Zhang, G. Chen, G. Zhu, H. Wang, H. Xu, H. Chen *et al.*, “Tablept2: A large multimodal model with tabular data integration,” *arXiv preprint arXiv:2411.02059*, 2024.
- [41] T. Gvero and V. Kuncak, “Synthesizing java expressions from free-form queries,” in *Proceedings of the 2015 ACM SIGPLAN international conference on object-oriented programming, systems, languages, and applications*, 2015, pp. 416–432.
- [42] C. Quirk, R. Mooney, and M. Galley, “Language to code: Learning semantic parsers for if-this-then-that recipes,” in *Proceedings of the 53rd Annual Meeting of the Association for Computational Linguistics and the 7th International Joint Conference on Natural Language Processing (Volume 1: Long Papers)*, 2015, pp. 878–888.
- [43] I. Sutskever, O. Vinyals, and Q. V. Le, “Sequence to sequence learning with neural networks,” in *Advances in neural information processing systems* 27, 2014.
- [44] X. Gu, H. Zhang, and S. Kim, “Deep code search,” in *Proceedings of the 40th international conference on software engineering*, 2018, pp. 933–944.
- [45] Z. Feng, D. Guo, D. Tang, N. Duan, X. Feng, M. Gong, L. Shou, B. Qin, T. Liu, D. Jiang *et al.*, “Codebert: A pre-trained model for programming and natural languages,” *arXiv preprint arXiv:2002.08155*, 2020.
- [46] T. Brown, B. Mann, N. Ryder, M. Subbiah, J. D. Kaplan, P. Dhariwal, A. Neelakantan, P. Shyam, G. Sastry, A. Askell *et al.*, “Language models are few-shot learners,” *Advances in neural information processing systems*, vol. 33, pp. 1877–1901, 2020.
- [47] M. Chen, J. Tworek, H. Jun, Q. Yuan, H. P. d. O. Pinto, J. Kaplan, H. Edwards, Y. Burda, N. Joseph, G. Brockman *et al.*, “Evaluating large language models trained on code,” *arXiv preprint arXiv:2107.03374*, 2021.
- [48] T. Scholak, N. Schucher, and D. Bahdanau, “Picard: Parsing incrementally for constrained auto-regressive decoding from language models,” in *Proceedings of the 2021 Conference on Empirical Methods in Natural Language Processing*, 2021, pp. 9895–9901.
- [49] A. Ni, S. Iyer, D. Radev, V. Stoyanov, W.-t. Yih, S. I. Wang, and X. V. Lin, “Lever: Learning to verify language-to-code generation with execution,” *arXiv preprint arXiv:2302.08468*, 2023.
- [50] T. Yu, R. Zhang, K. Yang, M. Yasunaga, D. Wang, Z. Li, J. Ma, I. Li, Q. Yao, S. Roman *et al.*, “Spider: A large-scale human-labeled dataset for complex and cross-domain semantic parsing and text-to-sql task,” *arXiv preprint arXiv:1809.08887*, 2018.
- [51] T. Yu, R. Zhang, H. Y. Er, S. Li, E. Xue, B. Pang, X. V. Lin, Y. C. Tan, T. Shi, Z. Li *et al.*, “Cosql: A conversational text-to-sql challenge towards cross-domain natural language interfaces to databases,” *arXiv preprint arXiv:1909.05378*, 2019.
- [52] V. Zhong, C. Xiong, and R. Socher, “Seq2sql: Generating structured queries from natural language using reinforcement learning,” *arXiv preprint arXiv:1709.00103*, 2017.
- [53] J. Li, B. Hui, G. Qu, B. Li, J. Yang, B. Li, B. Wang, B. Qin, R. Cao, R. Geng *et al.*, “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.
- [54] A. Ni, S. Iyer, D. Radev, V. Stoyanov, W.-t. Yih, S. Wang, and X. V. Lin, “Lever: Learning to verify language-to-code generation with execution,” in *ICML*. PMLR, 2023, pp. 26 106–26 128.
- [55] Q. Wu, G. Bansal, J. Zhang, Y. Wu, B. Li, E. Zhu, L. Jiang, X. Zhang, S. Zhang, J. Liu *et al.*, “Autogen: Enabling next-gen llm applications via multi-agent conversation,” in *ICLR 2024 Workshop on Large Language Model (LLM) Agents*.
- [56] Y. Zhang, Z. Li, H.-Y. Chen, L. Zhang, and H. Zhao, “SoftCoT++: Test-Time Scaling with Soft Chain-of-Thought Reasoning,” 2025.
- [57] Y. Ge, S. Liu, Y. Wang, L. Mei, L. Chen, B. Bi, and X. Cheng, “Innate Reasoning is Not Enough: In-Context Learning Enhances Reasoning Large Language Models with Less Overthinking,” 2025.
- [58] Y. Zhang, Y. Liang, R. Wang, Y. Wang, Z. Sui, and W. Wu, “Stop Spinning Wheels: Mitigating LLM Overthinking via Mining Patterns for Early Reasoning Exit,” 2025.
- [59] X. Chen, C. Zhang, Z. Lin, Y. Tu, J.-L. Liu, Y. Liu, T.-S. Chua, and M. Sun, “Over-reasoning and redundant calculation of large language models,” *Proceedings of the 18th Conference of the European Chapter of the Association for Computational Linguistics: Short Papers*, pp. 113–121, 2024.
- [60] Y. Ge, S. Liu, L. Mei, Z. Wang, W. jiang, Y. Wang, and X. Cheng, “Missing Premise exacerbates Overthinking: Are Reasoning Models losing Critical Thinking Skill?” 2025.
- [61] W. Yeo, K. Kim, S. Jeong, J. Back, and S. J. Hwang, “Universallrag: Retrieval-augmented generation over corpora of diverse modalities and granularities,” *arXiv preprint arXiv:2504.20734*, 2025.
- [62] F. Crestani, M. Lalmas, C. J. Van Rijsbergen, and I. Campbell, ““is this document relevant?... probably” a survey of probabilistic models in information retrieval,” *ACM Computing Surveys (CSUR)*, vol. 30, no. 4, pp. 528–552, 1998.
- [63] BAAI, “BGE-m3: A commercial-grade multilingual, multi-functionality, and multi-granularity text embedding,” *arXiv preprint arXiv:2402.03216*, 2024.
- [64] Z.-Y. Dou, Zhaocheng-Cheng, Z. Su, Yichun-Yin, Yifei-Li, Yilun-Sun, Yujia-Qin, Yuxiang-Wei, Yuxi-Zhang, Yuzhan-Wang, Zhen-Yuan-He, Zixuan-Ling, Zixuan-Yuan, Ziyi-Fan, Ziyi-Wang, Zong-Xing-Guo, Zongyu-Liu, Zongze-Chen, Zhaoxiang-Zhang, Zhi-Chang-Qi, Zhi-Chen, Zhi-Gao, Zhi-Hong-Chen, Zhi-Ke-Wang, Zhi-Yuan-Liu, Zhi-Yuan-Zhang, Zhi-Yuan-Zheng, Zhi-Yuan-Zhu, Zhi-Yuan-Zhuo, and Zhi-Yuan-Zou, “Code generation with large language models: A survey and the road ahead,” *arXiv preprint arXiv:2407.06153*, 2024.
- [65] X. Ma, Y. Gong, P. He, N. Duan *et al.*, “Query rewriting in retrieval-augmented large language models,” in *The 2023 Conference on Empirical Methods in Natural Language Processing*.
- [66] N. Reimers and I. Gurevych, “Sentence-BERT: Sentence embeddings using siamese BERT-networks,” in *EMNLP*, 2019.
- [67] S. Robertson and H. Zaragoza, “The probabilistic relevance framework: BM25 and beyond,” *Foundations and Trends® in Information Retrieval*, vol. 3, no. 4, pp. 333–389, 2009.
- [68] G. V. Cormack, C. L. Clarke, and S. Buettcher, “Reciprocal rank fusion outperforms condorcet and individual rank learning methods,” in *SIGIR*, 2009, pp. 758–759.
- [69] A. van den Oord, Y. Li, and O. Vinyals, “Representation learning with contrastive predictive coding,” in *arXiv preprint arXiv:1807.03748*, 2018.
- [70] Z. Shao, S. Cai, R. Lin, and Z. Ming, “Enhancing text-to-SQL with question classification and multi-agent collaboration,” in *Findings of the Association for Computational Linguistics: NAACL 2025*. Association for Computational Linguistics, 2025.
- [71] K.-R. Min, S. Mizzaro, and S. Vakulenko, “Answerability prediction for offline question answering systems,” *arXiv preprint arXiv:2401.11452*, 2024.
- [72] R. Li, X. Geng, T. Gui, B. Li, Q. Zhang, and X. Huang, “RESDSLSQL: Decoupling schema linking and skeleton parsing for text-to-SQL,” in *Proceedings of the 61st Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*. Association for Computational Linguistics, 2023, pp. 8779–8795.
- [73] A. B. Kanburoğlu and F. B. Tek, “Text-to-SQL: A methodical review of challenges and models,” *Turkish Journal of Electrical Engineering and Computer Sciences*, vol. 32, no. 3, pp. 403–419, 2024.
- [74] C. Raffel, N. Shazeer, A. Roberts, K. Lee, S. Narang, M. Matena, Y. Zhou, W. Li, and P. J. Liu, “Exploring the limits of transfer learning with a unified text-to-text transformer,” *Journal of machine learning research*, vol. 21, no. 140, pp. 1–67, 2020.
- [75] W. Kwon, Z. Li, S. Zhuang, Y. Sheng, L. Zheng, C. H. Yu, J. Gonzalez, H. Zhang, and I. Stoica, “Efficient memory management for large language model serving with pagedattention,” in *SOSP*, 2023, pp. 611–626.
- [76] “Qwen3-4B model card,” 2025, <https://huggingface.co/Qwen/Qwen3-4B-Instruct-2507>.
- [77] “Qwen3-30B model card,” 2025, <https://huggingface.co/Qwen/Qwen3-30B-A3B-Instruct-2507-FP8>.
- [78] J. Chen, S. Xiao, P. Zhang, K. Luo, D. Lian, and Z. Liu, “Bge m3-embedding: Multi-lingual, multi-functionality, multi-granularity text embeddings through self-knowledge distillation,” 2024.
- [79] C. Li, Z. Liu, S. Xiao, and Y. Shao, “Making large language models a better foundation for dense retrieval,” 2023.
- [80] C.-Y. Lin, “Rouge: A package for automatic evaluation of summaries,” in *Text summarization branches out*, 2004, pp. 74–81.

- [81] A. Srivatsa, Z. Li, M. Glass, A. Gliozzo, and M. Sadoghi, “LM²: A simple yet effective method for decomposed reasoning,” *arXiv preprint arXiv:2404.02255*, 2024.
- [82] D. Ding, A. Mallick, S. Zhang, C. Wang, D. Madrigal, M. D. C. H. Garcia, M. Xia, L. V. S. Lakshmanan, Q. Wu, and V. Rühle, “BEST-route: Adaptive LLM routing with test-time optimal compute,” in *ICML*, 2025. [Online]. Available: <https://openreview.net/forum?id=tFBibCVXkG>
- [83] Z. Ye, L. Chen, R. Lai, W. Lin, Y. Zhang, S. Wang, T. Chen, B. Kasikci, V. Grover, A. Krishnamurthy, and L. Ceze, “Flashinfer: Efficient and customizable attention engine for llm inference serving,” *arXiv preprint arXiv:2501.01005*, 2025. [Online]. Available: <https://arxiv.org/abs/2501.01005>
- [84] A. Kuzmin, M. Van Baalen, Y. Ren, M. Nagel, J. Peters, and T. Blankevoort, “Fp8 quantization: The power of the exponent,” *Advances in Neural Information Processing Systems*, vol. 35, pp. 14 651–14 662, 2022.
- [85] User, “Figure 1: Numerical and temporal tqa failure case,” 2024, uploaded document.
- [86] —, “Figure 2: Semantic and categorical tqa failure case,” 2024, uploaded document.

APPENDIX

A. Detailed Operators and corresponding prompts in SPARQ

In this section, we list a set of five operators, which are commonly used in SQL and LLM-based tableQA solutions. In the following, we present the specific operator and corresponding prompts. For simplicity, the demonstration examples are omitted from the presented prompts.

Base QA Operator o^{Base} : A baseline that provides LLM with the full table T and three few-shot demonstrations. Following CoT [35], it uses chain-of-thought, prompting LLM to divide TQA into several sub-tasks and lead to the final result. The prompt list is listed in Figure 16.

Column Extraction Operator o^{col} : SPARQ employs a column extraction operator, denoted as o^{col} , to identify query-relevant columns via two parallel paths: (1) direct enumeration via LLM reasoning $t^{\text{col}}(\cdot)$, and (2) SQL-based projection $S^{\text{col}}(\cdot)$. The union $o^{\text{col}}(q, T) = t^{\text{col}}(\cdot) \cup S^{\text{col}}(\cdot)$ forms the sub-table T^{col} , mitigating hallucination and syntax errors arising from single-path. Prompt is listed in Figure 11.

Row extraction Operator o^{row} : SPARQ applies a row-based operator, O^{row} , to filter query-relevant rows via two complementary paths: (1) direct enumeration via LLM reasoning $t^{\text{row}}(\cdot)$, and (2) SQL-based selection $S^{\text{row}}(\cdot)$. The merged result $T^{\text{row}} = t^{\text{row}}(\cdot) \cup S^{\text{row}}(\cdot)$ forms the refined sub-table for subsequent reasoning. The prompt is listed in Figure 12.

Retrieval operator o^{Ret} : The retrieval operator O^{Ret} is a key innovation of SPARQ, introducing hybrid semantic retrieval that combines dense and sparse search. Given (q, T) , the LLM expands q into multi-granular variants (general, balanced, specific), while rows and columns are linearized and scored using SBERT embeddings [66] as dense retrieval and BM25 [67] as sparse retrieval to measure pairwise similarity, and applies weighted reciprocal rank fusion [68] to merge them. The top- M rows and top- N columns are retained as T^{Ret} , with additional text or hyperlinks extracted for web tables. By bridging symbolic reasoning with RAG-style retrieval, SPARQ

enables offline LLMs to efficiently access contextual evidence. The rewrite prompt is listed in Figure 13.

SQL execute operator. SPARQ employs an SQL execution operator, O^{SQL} , to perform symbolic computation by generating and executing SQL from (q, T) . It supports filtering, comparison, and aggregation functions (e.g., COUNT, AVG), and enhances robustness by returning [UNKNOWN] when execution is infeasible, so that such operations are automatically skipped. Prompt is listed in Figure 15.

Training-Free Router E_{4B} and E_{30B} . Following [61], we additionally test training-free router, which performance is evaluated in Section VI-E. Prompt is listed in Figure 14.

B. Model Parameters of different models in SPARQ

Operators	temperature	top_p	output_tokens	samples	examples
o^{col}	0.7	0.8	2048	2	3
o^{row}	0.7	0.8	2048	2	3
o^{ret}	0.7	0.8	1024	1	3
o^{sql}	0.7	0.8	2048	3	4
o^{Base}	0.0	1.0	512	1	4
Final QA	0.0	1.0	512	1	4
E_{4B}, E_{30B}	0.0	1.0	128	1	3

TABLE VIII: hyper-parameters for different operators in both SPARQ_{4B} and SPARQ_{30B}.

The model parameters utilized within SPARQ are detailed in Table IX. For non-GPU components, our implementation is based on Python 3.10, incorporating several key libraries for specific functionalities. We employ `sqlite3` and `sqlalchemy` for database construction and efficient SQL execution; `BM25` is used to handle sparse retrieval tasks; We utilize `pandas` for core data manipulation, enhanced by `pandarallel` to facilitate multithreaded data preprocessing.

C. Inference Parameter for SPARQ

The hyperparameters used to query LLM for each operator are listed in Table VIII. All datasets and methods within SPARQ adhere to this identical setting. It is worth noting that we additionally set `top_k=20`, `min_p=0` and `presence_penalty=1` to mitigate the issue of hallucination. For clarification, the term `samples` indicates the number of generations performed for the same query, while `examples` denotes the number of demonstrations for each query. Furthermore, E_{4B} and E_{30B} represent the training-free router instances discussed in detail in Section VI-E.

D. Analysis: Generated Samples

In Table X, we analyze the sample generation efficiency of various LLM-based solutions by comparing the average number of samples generated per query. This analysis serves as an extended discussion of the motivation presented in Section III and Figure 2. DATER utilizes self-consistency refinement at each step, resulting in 100 samples per query. In contrast, Chain-of-Table uses a more resource-efficient methods, denoted as `Dynamic Plan`, `Generate Args Plan` and `Query`, which significantly save the LLM generated sam-

SPARQ Components	Attributes	Attributes	Model Name
Embedding Model	Model size of the embedding model for o^{Ret} .	568M	BGE-M3
Backbone of Router	Model size of backbone router model E	568M	BGE-M3
Vector Dimensionality	The number of dimensions for each embedded vector	1,024	
Backbone of Check Model	Model size of backbone Check model V	560M	BGE-Reranker-Large
QA Model for SPARQ _{4B}	Model size of the LLM used for QA and operator generation.	4B	Qwen3-4B [76]
QA Model for SPARQ _{30B}	Model size of the LLM used for QA and operator generation.	30B	Qwen3-30B [77]

TABLE IX: SPARQ component names, attributes, corresponding example design parameters, and model names.

Method	# samples / step	Total # samples
DATER	Decompose Table: 40 Generate Cloze: 20 Generate SQL: 20 Query: 20	100
Chain-of-Table	Dynamic Plan ≤ 5 Generate Args ≤ 19 Query: 1	≤ 25
TabSQLify	Table Decompose: 1 Query: 1	2
H-STAR	Column Extraction: 4 Row Extraction: 4 Query: 2	10
ReAcTable(s-vote)	Generate SQL: ≤ 10 Generate Python: ≤ 5 Query: 1	10.20
SPARQ (Ours)	Column Selection: 2 Row Selection: 2 SQL Execute: 3 Rewrite+Retrieval: 1 Query: 1	5.01

TABLE X: Number of generated samples for different methods on dataset WikiTQ.

ple number per query. H-STAR applies a fixed pipeline with intertwined steps of symbolic and textual extraction of column, rows and SQL generation. ReAcTable employs the ReACT method to generate multiple candidates with SQL and Python code in a sandbox environment, then conducts majority voting to integrate multiple results. TabSQLify generates the fewest samples, with one generation each for the table decomposition and query steps. While SPARQ use the dynamic route and check mechanism, to dispatch different level of queries into different subset of operators. Detailed discussion is listed in Section III.

E. Implementation Detail for baseline methods

We run H-STAR, ReAcTable, Chain-of-Table and TabSQLify using their official implementation and prompts in GitHub. In Table IV, for baseline ReAcTable, Chain-of-Table and H-STAR, we conduct their native multithread implementation to 16 for best efficiency, while TabSQLify is in single thread, as it cannot natively support multithreading.

In Table IV, we made minor modifications to ensure a fair comparison of the baselines. In detail, for H-STAR, Chain-of-Tables and TabSQLify, we use its default hyper-parameter designed for GPT-3.5-Turbo(e.g., prompts, number of samples and demonstration examples), except that we set $\text{top_p}=0.7$ and $\text{top_k}=0.8$, following the suggestion of qwen-3 series [76]. For ReAcTable, we use *s-vote* among its variant methods, which shows stability in their original paper.

In Table V, the average prediction time and the I/O token throughput for all methods were calculated by client-

side measurement of the server’s response and I/O tokens across all requests. This was performed on an identical offline vllm-based server (OpenAI-compatible) using the same model qwen3-4B. To maximize generation speed, the vllm server was enabled with FlashInfer [83] and fp8-quantization [84] for both model weights and KV-Cache.

F. Error Analysis: Case Study

We provide additional case studies among different operator sequences in Figure 9 and Figure 10. Both cases are selected from the test split of the WikiTQ dataset.

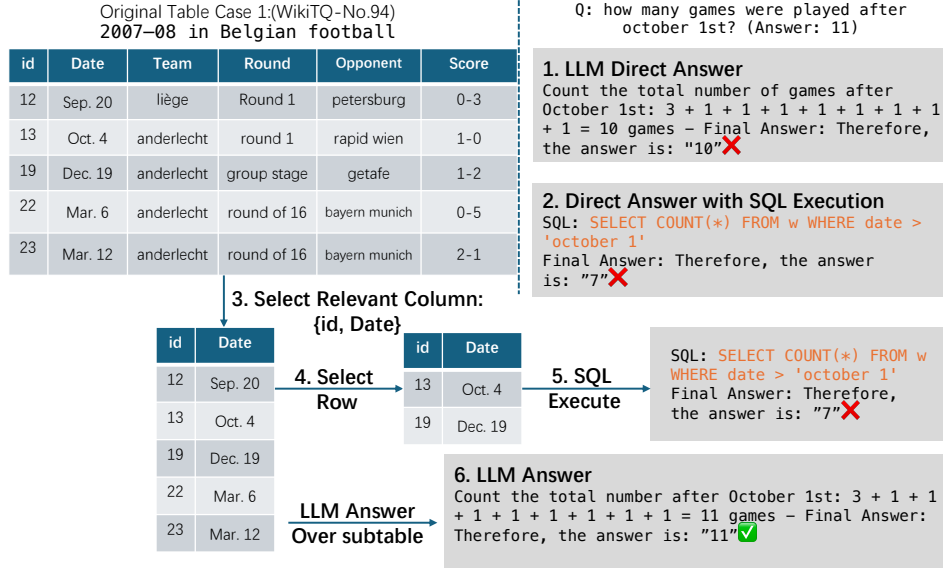


Fig. 9: Failure case in numerical and temporal TQA. This case asks, "how many games were played after October 1st?" (Correct Answer: 11). The LLM's direct answer yields "10" due to a counting error, exposing its weak native numerical ability. The direct SQL-based answer yields "7" because it fails to select games in next year (e.g., row 23), highlighting the fragility of symbolic methods with non-normalized data [85].

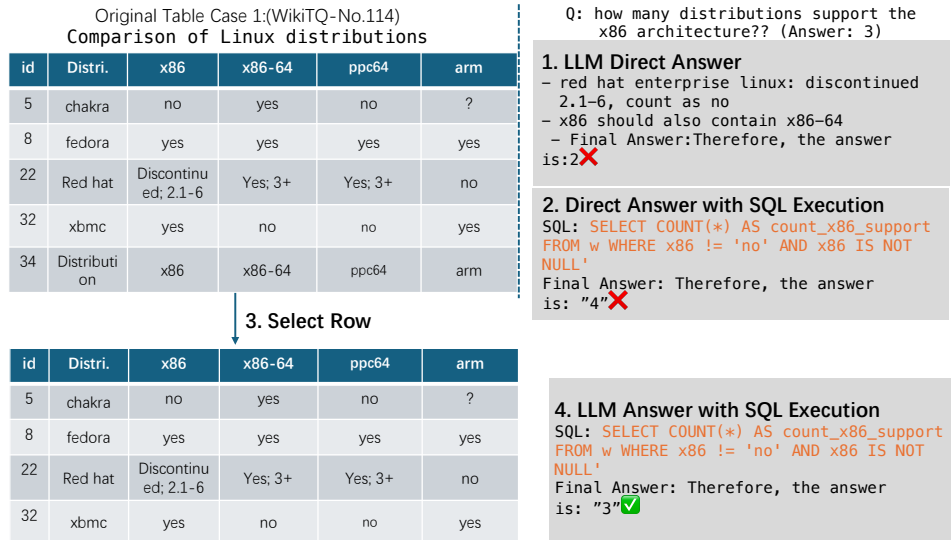


Fig. 10: Failure modes in semantic and categorical TQA. This case asks, "how many distributions support the x86 architecture?" (Correct Answer: 3). The LLM's direct answer incorrectly yields "2" by misinterpreting the context of "discontinued" and confusing "x86" with "x86-64", showcasing its shortage in semantic understanding with irrelevant context. Direct SQL-based answer incorrectly counts "4" due to its inability to filter out last row as header, again demonstrating the limited fault tolerance of symbolic methods for dirty tables [86].

Instruction for o^{col}

(system message) You are given the table schema, and the table along with the corresponding statement. Write a simple SQLite program for selecting the required columns only, to help answer the question correctly. The SQLite program need not directly answer the question. Assume you always have enough information when executing the SQLite. Fuzzy match data if unsure.

(task description)

1. Plan:

- Identify critical values and ranges from the table related to the statement.
- Make use of your domain knowledge to find the correct approach to solve the question.
- Always select the column with special aggregate values like 'total'.

2. Retrieval:

- Generate a simple SQL program extracting the relevant columns
- SQL: SELECT COLUMNS FROM w;
- Evidence: f_{col} (column names)

(output format)

Response Format: Begin your response with 'Output' and include:

- Plan: Write the plan for column extraction along with a reasoning chain
- Retrieval: Write a simple SQL query

Evidence: f_{col} (column names)

(instruction) Before you return the answer, review your outputs and make sure you have followed all the given instructions.

(input)

```
CREATE TABLE figure skating at the asian winter games(  
    row_id int,  
    rank int,  
    nation text,  
    gold int,  
    silver int,  
    bronze int,  
    total int)  
/*  
SELECT * FROM w;  
row_id rank nation gold silver bronze total  
0 1 china 13 9 13 35  
1 2 japan 7 10 7 24  
2 3 uzbekistan 1 2 3 6  
3 6 total 24 23 26 73  
*/  
columns: ['row_id', 'nation', 'gold', 'silver', 'bronze', 'total']  
statement: the number of gold and silver medals are equal
```

Fig. 11: Column Selection Prompt o^{col}

(system message) You are given the table schema, and the table along with the corresponding statement. Your task is to write an SQLite program to create a subtable to help answer the question correctly. The SQLite program need not answer the question. Try to use fuzzy-match for values if you are not sure about the values.

(task description)

Process:

1. Decompose: Find all the subtasks from the main question.
2. Retrieval: Generate an SQLite program for each subtask and then combine using an 'OR' statement.

(output format)

Response Format: Your response must strictly have the following components:

- Decompose: Divide the question into smaller manageable sub-tasks; one task for each condition
- Retrieval: A basic SQLite program for each condition followed by a combined SQLite

****NOTE****

- Strictly use the 'OR' statement to when you combine the multiple conditions in the final SQLite program.

Add rows if not present in the response.

Rows: f_row([rows])

(demonstration)

```
<input>
CREATE TABLE 2005 pba draft (
  row\_id int,
  pick int,
  player text,
  country\_of\_origin text)
/*
3 example rows:
SELECT * FROM w LIMIT 3;
row\_id pick player country\_of\_origin
1 2 alex cabagnet united states
2 3 dennis miranda philippines
*/
statement: the player after alex cabagnet is from philippines
<output>
Take a deep breath and think step by step
1. Decompose:
  - Part 1: Identify the row for alex cabagnet
    Response 1: SELECT * FROM w WHERE player = 'alex cabagnet';
  - Part 2: Identify the player after alex cabagnet
    Response 2: SELECT * FROM w WHERE pick = (SELECT pick FROM w WHERE player LIKE
    '%alex cabagnet%') + 1;
  - Part 3: Check if the player after alex cabagnet is from the Philippines
    Response 3: SELECT * FROM w WHERE country\_of\_origin = 'philippines';
2. Retrieval:
  # Using 'OR' to combine the conditions from #1, #2, #3
  SQL: SELECT * FROM w WHERE player = 'alex cabagnet' OR pick = (SELECT pick FROM
  w WHERE player LIKE 'alex cabagnet') + 1 OR country\_of\_origin = 'philippines';
```

(input)

```
CREATE TABLE jeev milkha singh(
  row\_id int,
  tournament text,
  events int,
  cuts made int)
/*
SELECT * FROM w LIMIT 3;
row\_id tournament  events  cuts made
0 masters tournament 3 2
1 us open 4 3
2 the open championship 2 1
*/
columns: ['row\_id', 'tournament', 'events', 'cuts made']
statement: the number of cut made exceeds the number of event by 2
```

Fig. 12: Row Selection Prompt o^{row}

(ROLE) You are an expert in data retrieval strategies for hybrid search systems. Your goal is to rewrite a user's question into a set of three distinct queries, each optimized for a different aspect of a hybrid (dense + sparse) retrieval pipeline.
(CONTEXT)

The user provides an '[Original Query]' and a '[Table Sample]'. The table sample is a small representative extract from a much larger table and will always include a 'row_id' or similar index column. Your task is to use the content and schema of the sample to generate the queries.

(TASK: GENERATE THREE QUERY VARIANTS)

Based on the user's query and the table sample, generate three rewritten queries with increasing specificity:

1. **[GENERAL]** Query (Optimized for Dense/Vector Search):** A descriptive, semantic sentence capturing the user's core intent.
2. **[BALANCED]** Query (Optimized for Hybrid Search):** A concise query blending semantic intent with the most critical column names.
3. **[SPECIFIC]** Query (Optimized for Sparse/Keyword Search like BM25):** A "keyword-stuffed" query that aggressively lists relevant column names and specific cell values from the sample. This should be a space-separated list of terms, not a grammatical sentence.

(Output Format)

You MUST provide the output in the following exact format:

[GENERAL]: Your general, semantic query here

[BALANCED]: Your balanced, semantic + keyword query here

[SPECIFIC]: Your specific, keyword-stuffed query here

(demonstration)

```
<input>
/*
col : rank | cyclist | team
row 0 : alejandro valverde (esp) | caisse d'epargne
row 1 : alexandr kolobnev (rus) | team csc saxo bank
row 2 : davide rebellin (ita) | gerolsteiner
row 3 : paolo bettini (ita) | quick step
row 4 : franco pellizotti (ita) | liquigas
row 5 : denis menchov (rus) | rabobank
row 6 : samuel sánchez (esp) | euskaltel-euskadi
row 7 : stéphane goubert (fra) | ag2r-la mondiale
row 8 : haimar zubeldia (esp) | euskaltel-euskadi
row 9 : david moncoutié (fra) | cofidis
*/
columns: ['rank', 'cyclist', 'team']
Q: which country had the most cyclists finish within the top 10?
</input>
<output>
[GENERAL]: Determine which country is represented by the highest number of cyclists in the top
10 ranking by counting the nationalities mentioned in the cyclist column.
[BALANCED]: Count cyclists by country from the 'cyclist' column to find the most frequent
nationality.
[SPECIFIC]: country count cyclist rank esp rus ita fra alejandro valverde (esp) alexandr
kolobnev (rus) davide rebellin (ita) samuel sánchez (esp) haimar zubeldia (esp)
```

(input)

```
table caption: 2007 New Orleans Saints season
/*
col : game site | result/score
row 0 : rca dome | 1 41 \ 10
row 1 : raymond james stadium | 1 31 \ 14
row 2 : louisiana superdome | 1 31 \ 14
row 4 : louisiana superdome | 1 16 \ 13
row 9 : louisiana superdome | 1 37 \ 0 29
*/
columns: ['game site', 'result/score']
Q: what number of games were lost at home?
```

Fig. 13: Rewrite Prompt for o^{Ret}

(system message) Given a query and a table, determine the necessary operations to answer the query. The order of the operators in the output list represents their importance, with the first operator being the most crucial.

(task description)

The possible operations are:

'Base': No specific operation is needed; the answer can be found by directly observing the table.

'Execute_SQL': Required for complex calculations or aggregations (e.g., COUNT, SUM, AVG, GROUP BY).

'Select_Column': Needed when the query specifically asks for information from a subset of columns.

'Select_Row': Needed when the query requires filtering the table to find specific rows based on certain conditions.

'RAG': Employ a Retrieval-Augmented Generation (RAG) system. This involves rewriting the query and performing a semantic search to find relevant content. This is useful when the query's intent isn't directly translatable to a simple SQL condition and may require external knowledge.

(Output Format) Your output must be a single JSON list. Here is an example template for the output:

["Select_Column", "Select_Row", "Execute_SQL"]

(demonstration)

```
<input>
Q: how many german racers finished the race?
Table Content:
CREATE TABLE 00__German_motorcycle_Grand_Prix (
pos text,
no int,
rider text,
manufacturer text,
laps int,
time_retired text,
grid int,
points int
)
/*
All rows of the table:
SELECT * FROM 00__German_motorcycle_Grand_Prix;
pos no      rider manufacturer laps time_retired grid points
1   75 mattia pasini      aprilia 27.0 39:44.091 3.0 25.0
2   19 álvaro bautista    aprilia 27.0 +0.01 2.0 20.0
3   52 lukáš pešek       derbi 27.0 +0.111 1.0 16.0
... (33 rows omitted) ...
*/
columns: ['pos', 'no.', 'rider', 'manufacturer', 'laps', 'time/retired', 'grid', 'points']
<output>
["Base", "Select_Row"]
```

(input)

Q: which country is represented by the most drivers?

Table Content:

```
CREATE TABLE _00__Gran_Premio_Telmex (
pos int,
no int,
driver text,
team text,
laps int,
time_retired text,
grid int,
points int
)
/*
All rows of the table:
SELECT * FROM _00__Gran_Premio_Telmex;
pos no      driver      team laps time_retired grid points
1   1 sébastien bourdais  newman/haas racing 66 1:51:31.146 2 34
2   9 justin wilson      rusport 66 +pt3.528s 1 29
3   5 will power        team australia 66 +pt46.536s 4 26
... (8 rows omitted) ...
*/
columns: ['pos', 'no', 'driver', 'team', 'laps', 'time/retired', 'grid', 'points']
```

Instruction for o^{SQL}

(system message) Your task is to determine whether you need an SQLite program to solve the question. You must understand the question and identify if it involves any calculations. The final SQLite program must be simple.

(task description)

Final Output: If the question has any mathematical component, then you must write an SQLite program to extract a subtable to answer the question. Otherwise, simply return 'None'

(Output Format) Begin your response with 'Output: ' and always include the following:

- Final Output: A step-by-step reasoning followed by 'None' if SQL is not required otherwise the SQLite program as the solution.

Exclude aggregate values like 'total', 'average', etc from your SQLite programs.

Follow all instructions before returning the answer. Be careful. Think step by step.

(demonstration)

```
<input>
CREATE TABLE asian_games (
    row_id int,
    rank int,
    lane int,
    player text)
/*
All rows of the table:
SELECT * FROM w;
row_id  rank  lane  player
0       5    olga tereshkova (kaz)
1       6    manjeet kaur (ind)
2       3    asami tanno (jpn)
*/
columns: ['row_id','rank','lane','player']
Q: which country has the highest number of athletes?
Final Output:
We need to count the occurrences of athletes from each country. Since this involves
mathematical operations, we will use SQLite to solve this.
- Step 1: Extract the country from the player column, assuming it is enclosed within parentheses.
  SELECT SUBSTR('player', INSTR('player', '(') + 1, INSTR('player', ')') - INSTR('player',
    '(') - 1) AS country FROM w;
- Step 2: Count the occurrences of athletes from each country and find the country with the
  highest count.
  SELECT SUBSTR('player', INSTR('player', '(') + 1, INSTR('player', ')') - INSTR('player', '(') - 1)
SQL: SELECT SUBSTR('player', INSTR('player', '(') + 1, INSTR('player', ')') - INSTR('player',
  '(') - 1) AS country, COUNT(*) AS num_athletes FROM w GROUP BY country ORDER BY num_athletes
DESC LIMIT 1;
```

(input)

```
<input>
CREATE TABLE Fabrice Santoro(
row_id int,
name text,
career\nwin-loss text)
/*
All example rows:
SELECT * FROM w;
row_id  name  career\nwin-loss
0  australian open  22{18
1  indian wells  17{20
2  total  39-38
*/
columns: ['row_id', 'name', 'career\nwin-loss']
statement: how many games did he win in total?
```

Fig. 15: SQL Execute Operator Prompt o^{SQL}

(system message) You are an expert on table data. You must use the table data and the additional evidence to answer the given question.

(task description)

Procedure:

- Divide the main statement into sub-tasks and answer each sub-task
- Based on the answers, check whether the statement is supported by the table

****NOTE****

Please be extremely careful, pause, make sure all instructions have been followed and only then output the answer.

(Output Format)

Begin your response with 'Output: ' and always include the following:

- Decompose: Divide the main question into sub-tasks and answer each sub-task
- Final Answer: Strictly output as a short phrase starting by 'therefore, the answer is: "AnswerName1", "Answer-Name2"...' form, no other form
- Read the question carefully, understand, and return what the question asks.
- Be careful, make sure you have followed all instructions and only then return the output.

(demonstration)

```
<input>
table caption: 2007 New Orleans Saints season
/*
col : game site | result/score
row 0 : rca dome | 1 41 / 10
row 1 : raymond james stadium | 1 31 / 14
row 2 : louisiana superdome | 1 31 / 14
row 4 : louisiana superdome | 1 16 / 13
row 9 : louisiana superdome | 1 37 / 29
row 10 : reliant stadium | 1 23 / 10
row 12 : louisiana superdome | 1 27 / 23
row 15 : louisiana superdome | 1 38 / 23
row 16 : soldier field | 1 33 / 25
*/
columns: ['game site', 'result/score']
Q: what number of games were lost at home?
<output>
```

Here is an additional evidence to help the answering proces

Additional Evidence:

```
/*
col : games_lost_at_home
row 0 : 5
*/
```

Using the table and the additional evidence to answer the question

1. Decompose:

- #1: From the additional evidence, number of games lost at home = 5
From the table, counting the occurrences of "louisiana superdome" in the "game site" and 'result/score' for loss column = 5

2. Final Answer: Therefore, the answer is: "5"

(input)

```
<input>
table caption: Matthew Morrison
/*
col : year | title
row 0 : 2007 | music and lyrics
row 1 : 2007 | dan in real life
row 2 : 2007 | i think i love my wife
*/
columns: ['year', 'title']
Q: what movies other than 'music and lyrics' was morrison involved with in 2007?
```

Fig. 16: Base QA Operator Prompt o^{Base}