

MARMo: An Evidence-Guided and Memory-Augmented Multi-Agent Repair Framework for Accurate Text-to-NoSQL Parsing

Anonymous ACL submission

Abstract

Text-to-NoSQL aims to extend natural language interfaces (NLIs) beyond traditional Text-to-SQL, enabling intuitive querying of NoSQL systems such as MongoDB—widely adopted in modern data-intensive applications. While recent advances in Text-to-NoSQL have yielded benchmark datasets, existing systems still suffer from erroneous query generation and a lack of robust mechanisms for error detection, explanation and iterative correction. In this paper, we propose MARMo: a novel evidence-guided, memory-augmented multi-agent repair framework for accurate Text-to-NoSQL parsing (focusing on MongoDB). The framework employs a modular, plug-and-play architecture that coordinates multiple specialized agents to perform iterative error diagnosis and adaptive refinement. Extensive experiments on the TEND benchmark validate that MARMo, as a modular and interpretable framework, achieves overall improvement on repair accuracy and generalization for Text-to-NoSQL. Our work not only enhances the reliability of natural language-driven NoSQL query systems but also paves the way for future research into robust, human-in-the-loop NLIs for unstructured and semi-structured data management.

1 Introduction

Natural language interfaces, by lowering the barrier to data access for non-technical users, have catalyzed significant research interest in NLP and databases, particularly in the Text-to-SQL task of translating plain-language queries into executable statements over structured relational data (Qin et al., 2022; Iacob et al., 2020; Katsogiannis-Meimarakis et al., 2023). However, despite the remarkable progress in Text-to-SQL modeling, existing efforts are largely confined to relational databases, which presume rigid schemas and normalized tables (Zhong et al., 2017; Yu et al., 2018). In contrast, NoSQL databases such as MongoDB

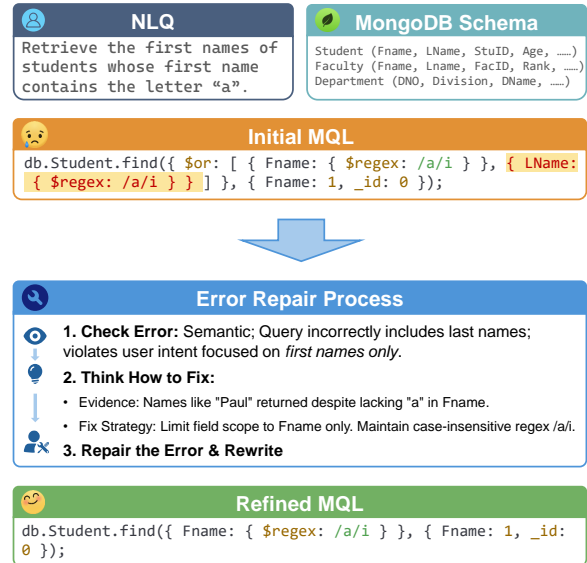


Figure 1: An illustrative example of the MQL repair task that simulates a realistic human correction process, including check errors in the initial query, formulating repair strategies based on contextual evidence, and rewriting the query to better align with the user’s intent.

are increasingly adopted in modern applications for their schema-less design, hierarchical data representation, and flexibility in modeling diverse, dynamic data. Although several studies have proposed benchmarks for the Text-to-NoSQL task, the interaction gap between natural language and NoSQL systems remains significantly underexplored, as current methods predominantly adopt few-shot generation paradigms with limited diversity, hindering the broader vision of democratized access to heterogeneous data platforms (Zhang et al., 2024a,b; Lu et al., 2025; Qin et al., 2025).

Recent advances in large language models (LLMs) have demonstrated remarkable capabilities in understanding complex instructions and generating structured outputs, including executable database queries (Achiam et al., 2023; Xie et al., 2024b; Zhu et al., 2024; Chung et al., 2025). Mean-

043
044
045
046
047
048
049
050
051
052
053
054
055
056
057
058
059
060

061 while, multi-agent systems have emerged as a
062 promising paradigm for decomposing and coordin-
063 ating complex tasks through the collaboration
064 of specialized agents each responsible for subtasks
065 such as planning, reasoning, or correction. Within
066 the context of NLIs for NoSQL databases, the syn-
067 ergistic combination of LLMs and multi-agent ar-
068 chitectures offers a promising solution to the limi-
069 tations of one-shot generation (Wang et al., 2025;
070 Cen et al., 2025). This hybrid approach not only
071 enhances interpretability and modularity, but also
072 enables dynamic feedback integration and adap-
073 tive repair, paving the way for more robust and
074 extensible query generation pipelines.

075 In this paper, we propose **MARMo**, a novel
076 evidence-guided and **Memory-Augmented Multi-**
077 **Agent Repair** framework for accurate Text-to-
078 NoSQL parsing (MongoDB), which operates in
079 a feedback-driven refinement loop, consisting of
080 three key agents: the Error Checker Agent, the
081 Feedback Manager Agent and the Error Repair
082 Agent. As shown in Figure 1, this repair workflow
083 is inspired by the realistic human correction process
084 when people find errors. Given a natural language
085 query (NLQ), MARMo refines the corresponding
086 MQL (MongoDB Querying Language) output, the
087 Error Checker Agent first detects different errors by
088 leveraging analysis from an MQL execution tool.
089 In the setting of MQL application, we categorize
090 MQL errors into three main types: the syntax error,
091 schema error, and semantic error. The Feedback
092 Manager Agent invokes different sub-agents based
093 on the error type and generate the feedback based
094 on the aggregated information. The Error Repair
095 Agent finally leverages this aggregated feedback
096 to retrieve and adapt similar error-repair patterns
097 from memory to repair and regenerate the refined
098 MQL, which is then re-executed for validation, al-
099 lowing for iterative correction if necessary. In all,
100 MARMo is designed to be pluggable, extensible,
101 and generalizable, making it applicable across di-
102 verse scenarios for Text-to-NoSQL models and a
103 promising step toward robust natural language in-
104 terfaces for data system diversity.

105 In our experiments, we evaluate MARMo on the
106 Text-to-MongoDB benchmark TEND (Lu et al.,
107 2025), which contains diverse natural language
108 queries paired with corresponding MQL anno-
109 tations. We assess the overall performance of
110 MARMo and compare it against a set of competi-
111 tive baseline methods. Experimental results show
112 that our framework consistently outperforms prior

113 approaches, highlighting the benefits of multi-agent
114 collaboration and error-aware iterative repair in the
115 context of Text-to-NoSQL generation. The modu-
116 lar design of MARMo enables flexible and robust
117 handling of diverse error types. In particular, the
118 incorporation of database schema and query ex-
119 ecution feedback facilitates more accurate query
120 refinement, contributing to improved structural and
121 semantic alignment between generated queries and
122 gold references.

123 To summarize, the main contributions of this
124 work are as follows:

- 125 • We introduce MARMo, a modular multi-agent
126 repair framework for Text-to-NoSQL query
127 generation. MARMo systematically inte-
128 grates error diagnosis, feedback, and iterative
129 repair, marking the first framework in this do-
130 main that empowers agents to not only gener-
131 ate but also correct NoSQL queries.
- 132 • We design a structured error taxonomy and a
133 corresponding feedback mechanism that dis-
134 tinguish syntax, schema, and semantic errors
135 in Text-to-NoSQL tasks. To address schema-
136 related issues, we develop an Evidence Agent
137 for precise error localization and informed re-
138 pair, and a global repair memory tailored to
139 the Text-to-NoSQL setting, which supports
140 iterative corrections and robust MQL regener-
141 ation.
- 142 • Through extensive experiments on the TEND
143 benchmark, we demonstrate the effectiveness
144 of MARMo and provide empirical insights
145 that point toward promising future directions
146 for robust and accurate Text-to-NoSQL sys-
147 tems.

148 2 Related Works

149 **Text-to-SQL** has long been a core research prob-
150 lem at the intersection of natural language process-
151 ing and databases (Hong et al., 2024). Early sys-
152 tems were limited to narrow domains and simple
153 schemas, often relying on rule-based grammars or
154 templates (Li and Jagadish, 2014; Kim et al., 2020;
155 Gkini et al., 2021). The introduction of large-scale
156 benchmarks, such as WikiSQL (Zhong et al., 2017),
157 Spider and its variants (Yu et al., 2018; Gan et al.,
158 2021a,b), and BIRD (Li et al., 2023), has driven
159 the development of more robust and generalizable
160 models. Pre-trained language models have further
161 advanced the field, enabling improved contextual

reasoning, schema linking, and handling of complex SQL constructs, including in conversational settings (Cao et al., 2021; Li et al., 2021; Scholak et al., 2021; Fan et al., 2023). More recently, LLMs and multi-agent paradigms have been investigated to decompose Text-to-SQL into subtasks such as planning, reasoning, and verification, enhancing both robustness and interpretability (Pourreza and Rafiei, 2023; Zhang et al., 2023; Wang et al., 2025; Li et al., 2024a; Ren et al., 2024; Mao et al., 2024; Pourreza et al., 2024; Yun and Lee, 2025; Cen et al., 2025; Shkapenyuk et al., 2025; Dönder et al., 2025). In parallel, growing attention has been given to error analysis and repair, with new benchmarks and methods targeting error classification and correction in SQL generation (Liu and Tan, 2024; Xu et al., 2025; Liu et al., 2025; Shen et al., 2025).

NoSQL Databases play an indispensable role in modern infrastructures due to their scalability, flexibility, and support for semi-structured and heterogeneous data (Dritsas and Trigka, 2025). Unlike relational databases, NoSQL systems adopt dynamic, hierarchical representations, spanning document stores, key-value stores, and graph databases (Zhang et al., 2024b). Recent efforts have increasingly addressed Text-to-NoSQL translation, most focusing on graph-based systems such as Neo4j and Cypher (Zhou et al., 2024; Tran et al., 2024; Lin et al., 2024; Gao et al., 2024). A key milestone is TEND (Lu et al., 2025), which provides the first end-to-end pipeline for constructing Text-to-NoSQL datasets and introduces a MongoDB-specific benchmark. Its multilingual extension, MultiTEND (Qin et al., 2025), underscores the significance of cross-lingual generalization. Building on these foundations, our work expands the solution space by introducing novel architectural components and refinement strategies that improve the accuracy and robustness of Text-to-NoSQL parsing in MongoDB.

LLM-based Agents have recently become a dominant paradigm for complex, multi-step reasoning tasks, from programming to document understanding (Yang et al., 2023; Nunez et al., 2024; Ishibashi and Nishimura, 2024; Xiao et al., 2024; Li et al., 2024b, 2025; Han et al., 2025). In Text-to-SQL, agent-based approaches decompose parsing into sub-modules such as intent detection, schema linking, and query generation (Wang et al., 2025; Xie et al., 2024a; Cen et al., 2025; Deng et al., 2025), improving modularity and interpretability. Specifically, recent work in Text-to-SQL and code gen-

eration has demonstrated that many agent-based approaches combine error detection and self-repair mechanisms with tool-use strategies to achieve further performance improvements. For instance, CRITIC (Gou et al., 2024) introduces a framework that enables LLMs to iteratively validate and refine their own outputs using external tools, thereby improving reliability and correctness in tasks such as code generation. MAC-SQL (Wang et al., 2025) introduces a multi-agent, tool-assisted framework that boosts Text-to-SQL performance on large databases, achieving new SOTA results on BIRD. SQLFixAgent (Cen et al., 2025) further enhances Text-to-SQL reliability by focusing specifically on detecting and repairing semantic errors in generated SQL.

However, prior work has largely focused on query construction, devoting limited attention to systematic error detection and repair. We address this gap by introducing a multi-agent repair framework that integrates query generation with robust error analysis and iterative refinement for Text-to-NoSQL tasks. This framework not only advances Text-to-NoSQL research but also provides a foundation for systematic exploration of natural language interaction with semi-structured data.

3 Overview

In this section, we provide an overview of the problem formulation and our proposed framework. We first formalize the Text-to-NoSQL and MQL Repair tasks, highlighting the unique challenges posed by MongoDB’s semi-structured data model. We then briefly introduce the overall architecture of MARMo, a modular, agent-based system designed to iteratively detect, diagnose, and repair errors in generated MQL queries.

3.1 Task Formulation

3.1.1 Text-to-NoSQL

The Text-to-NoSQL task aims to translate a natural language question into an executable NoSQL query, conditioned on a target database and its schema. Formally, let the natural language input be denoted as a token sequence $x = \{x_1, x_2, \dots, x_n\}$, $x \in \mathcal{N}$, where \mathcal{N} is NLQ sets. Let the NoSQL database be represented as $\mathcal{D} = \{C_1, C_2, \dots, C_k\}$, where each C_i is a collection and consists of a set of fields $F_i = \{f_{i1}, f_{i2}, \dots, f_{im_i}\}$. The database schema \mathcal{S} can thus be viewed as the structural metadata $\mathcal{S} = \{(C_i, F_i)\}_{i=1}^k$. The goal of the task is to generate

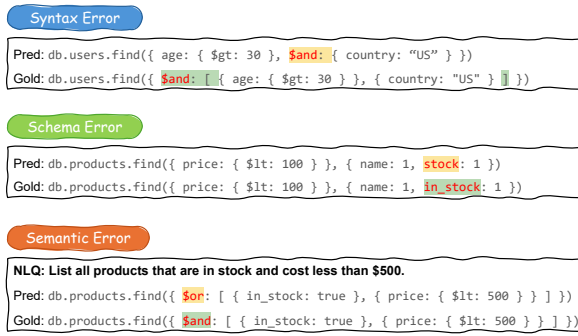


Figure 2: Three error types: Syntax Error: misuse of `$and` operator corrected by wrapping conditions in an array; Schema Error: invalid field `stock`, replaced with the correct schema `in_stock`; Semantic Error: misuse of `$or` instead of `$and`, leading to a mismatch with the user intent.

a MQL $y = \{y_1, y_2, \dots, y_m\}$, where $y \in \mathcal{Y}$ is a sequence of tokens from the target MQLs, such that y faithfully captures the semantics of x and is executable against \mathcal{D} .

Unlike SQL, NoSQL queries typically operate over semi-structured or hierarchical data without fixed relational schemas. For MongoDB, especially, its nested document-oriented data model introduces unique challenges such as inferring collection names, nested field paths, and handling heterogeneous document structures. As a result, solving this task requires not only understanding the natural language semantics but also adapting to the dynamic, schema-less nature of NoSQL databases.

3.1.2 MQL Repair Task

The MQL Repair task focuses on identifying and correcting errors in initially generated MongoDB queries, with the ultimate goal of producing executable MQL statements that accurately reflect the user’s intent.

Within the context of our framework MARMo, the MQL Repair task functions as a plug-and-play module that can be integrated with existing Text-to-NoSQL systems. By post-editing their generated queries, it improves the overall generation accuracy and enhances system usability. This modular repair mechanism ensures robustness in practical deployment, enabling seamless and reliable end-to-end workflows.

3.1.3 Error Types

To address the challenges in robust MQL repair, we design a fine-grained error categorization scheme along with corresponding identification and repair

strategies. Specifically, we classify MQL errors into three primary types: syntax errors, schema errors, and semantic error, as illustrated in Figure 2.

(1) **Syntax Error** refers to cases where the generated MQL query violates the formal grammar rules of the query language, rendering it unparseable or unexecutable.

(2) **Schema Error** arises when the structural form of the query does not match the expected template, such as incorrect nesting, misplaced clauses, or missing components, regardless of specific value content.

(3) **Semantic Error** occurs when the query structure is correct but the intent it conveys deviates from the user’s original question, often due to incorrect predicates, misused operators, or improper value conditions.

3.2 Framework Overview

MARMo operates in a feedback-driven refinement loop comprising three core agents: the *Error Checker Agent*, *Feedback Manager Agent*, and *Error Repair Agent* shown in Figure 3. First, error information is leveraged by the Error Checker Agent to detect and classify MQL failures (e.g., syntactic, schema-related, or semantic). Based on the identified error type, the Feedback Manager Agent selects an appropriate reasoning path to activate different sub-agents and generate targeted repair suggestions, selectively incorporating evidence such as schema metadata, domain constraints, or self-reflection strategies. Finally, the Error Repair Agent synthesizes corrected MQL queries guided by this feedback. The algorithm is shown as Algorithm 1.

This iterative, agent-based design enables progressive query refinement and modular error correction. As a plug-and-play architecture, MARMo enhances MQL generation across diverse upstream systems by enabling fine-grained error resolution and offering transparent diagnostic traces for each repair step.

4 Methodology

In this section, a detailed description of MARMo is presented by systematically introducing its three main agents, the Error Checker Agent, the Feedback Manager Agent, and the Error Repair Agent, in the order they operate. Each component is designed to emulate distinct stages of a realistic hu-

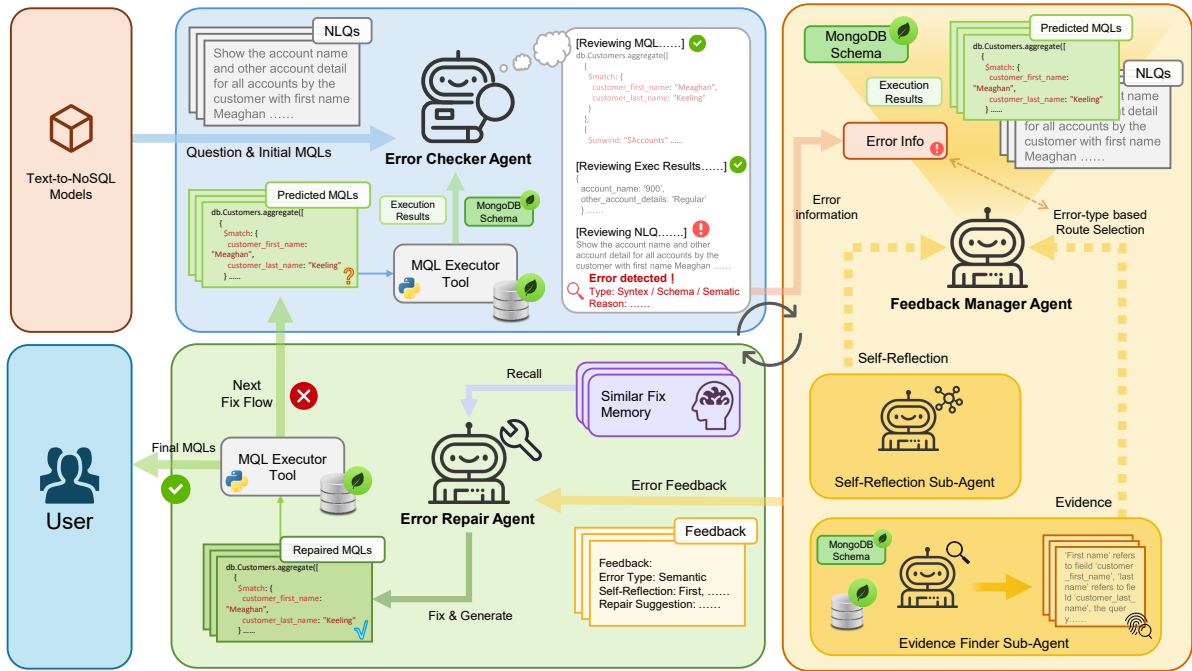


Figure 3: An overview of **MARMo** (MongoDB-focused, **M**emory-Augmented, **M**ulti-Agent **R**epair), an evidence-guided framework for accurate Text-to-NoSQL parsing. It features an iterative workflow with three specialized agents: the *Error Checker Agent* detects errors via execution traces and semantic cues; the *Feedback Manager Agent* orchestrates evidence-aware reasoning to produce targeted correction signals; and the *Error Repair Agent* integrates feedback and contextual memory to iteratively refine queries toward syntactic validity, schema compliance, and intent alignment.

man correction process, from error detection to guided repair.

4.1 Error Checker Agent

The Error Checker Agent is responsible for identifying and classifying errors in the predicted MQL query based on its execution behavior and contextual information. Given a natural language query and a corresponding MQL query generated by an external model, the system attempts to determine whether the MQL is correct or contains errors.

Specifically, the Error Checker Agent first invokes a Python-based executor of MQL to run the query against the MongoDB instance. The execution result including returned documents or error messages is then parsed and analyzed. Upon detecting a failure, the Error Detector uses a hybrid strategy: it combines rule-based inspection of execution logs with LLM-based reasoning over both the original input and query to determine the error type. Errors are categorized into syntactic (e.g., invalid query structure), schema-related (e.g., incorrect collection or field names), or semantic (e.g., mismatched intent). The Error Checker Agent performs comprehensive analysis and validation of

the generated MQL query against the input question and schema, ultimately identifying the error type by detecting syntax error, schema mismatch or logical misalignment with the user intent and providing a rationale for the classification. Finally, the agent outputs a structured error report including the detected error type and contextual clues for downstream repair of the Feedback Manager Agent.

The motivation behind this design is to integrate execution-based error detection with the semantic reasoning capabilities of large language models, thereby enhancing the accuracy of query correctness assessment. Beyond mere error detection, the agent aims to precisely identify the type and underlying cause of each error, offering structured feedback to support downstream repair mechanisms.

4.2 Feedback Manager Agent

The Feedback Manager Agent is a composite agent composed of two collaborative sub-agents that work in tandem to produce structured and informative feedback for query repair. Taking the original natural language query, the initial MQL and the error type identified by the Error Checker Agent as

Algorithm 1 The Algorithm of MARMo

In: Natural Language Question N , Initial MQL Query \hat{Q} , MongoDB database D

Out: Correct MQL \hat{Q}_c

```
1: Let iteration  $i = 0$ .
2: while  $i < Max$  do
3:   ExecResult  $r = ExecTool(D, \hat{Q})$ .
4:   ErrorType  $e = ErrorChecker(N, r, \hat{Q})$ .
5:   if  $e = correct$  then
6:      $\hat{Q}_c = \hat{Q}$ .
7:     Break.
8:   end if
9:   Feedback  $f = Feedback(N, e, D, \hat{Q})$ .
10:  MQL  $\hat{Q}_r = ErrRefiner(N, f, D, \hat{Q})$ .
11:  ExecResult  $r' = ExecTool(D, \hat{Q}_r)$ .
12:  if  $r' = PASS$  then
13:     $\hat{Q}_c = \hat{Q}_r$ .
14:    Break.
15:  end if,
16:   $\hat{Q} = \hat{Q}_r$  ( $\rightarrow$  Next Iteration)
17: end while
18: return final MQL  $\hat{Q}_c$ 
```

input, it initiates a reasoning process to analyze and contextualize the failure. The Feedback Manager Agent selectively determines the processing route based on the identified error type, invoking specialized sub-agents to generate type-specific feedback, which is subsequently aggregated and integrated into a unified correction feedback response.

(1) **Self-Reflection Sub-Agent** performs a detailed, step-by-step inference process on the faulty query by leveraging chain-of-thought prompting. This approach enables the agent to systematically analyze each component of the query and reason through potential inconsistencies. The primary goal is to accurately localize the source of the error, whether it stems from semantic misunderstandings or structural flaws within the query.

(2) **Evidence-Finder Sub-Agent** is to extract schema-related information, such as field types, collection names, nesting hierarchies, and field frequency—from the MongoDB context to facilitate schema-aware reasoning. Inspired by the effectiveness of evidence provided by the BIRD (Li et al., 2023; Yun and Lee, 2025) in supporting query generation, we introduce this agent to similarly enhance the generation and optimization of MQL queries. By grounding reasoning in concrete schema evidence, the agent helps improve both accuracy and robustness of the query construction

process.

Together, these collaborative sub-agents enable the Feedback Manager to deliver precise, context-aware, and actionable feedback, thereby enhancing the effectiveness of downstream query repair and optimization.

4.3 Error Repair Agent

The Error Repair Agent is tasked with fixing and synthesizing corrected MQL queries based on structured feedback generated by upstream components in the pipeline, primarily the Feedback Manager Agent. Its core objective is to transform faulty queries into executable forms that align with both the database schema and the user intent.

To this end, this agent employs a similar fix memory mechanism and execution-guided judgment. It selects past cases that closely resemble the current query failure, and these retrieved examples serve as in-context demonstrations, providing inductive signals to guide the repair process. In parallel, the agent integrates the structured feedback to further constrain and inform the generation of candidate MQL queries, ensuring that the revisions are both semantically appropriate and structurally valid. Finally, each repaired MQL query is evaluated by the MQL executor to determine whether it satisfies the execution criteria as the final output or requires further refinement in the next repair iteration. This design allows the Error Repair Agent to effectively integrate retrieved patterns with contextual reasoning, enabling accurate and adaptable MQL query correction.

5 Experiments

In this section, we present the experimental setup and main results of MARMo, our proposed multi-agent repair framework for Text-to-NoSQL query generation. We conduct comprehensive experiments on the benchmark TEND to assess the framework’s performance in handling erroneous predictions, including its ability to detect, classify, and refine faulty queries. The results clearly demonstrate that MARMo improves the accuracy and robustness of the query generation process compared to existing baselines.

For implementation of MARMo, the backbone LLM of the whole framework and different agent roles is “DeepSeek-V3” with the parameter setting ‘temperature = 0.0’ for relatively more precise and stable generation of MQLs. The number of repair

Method	Query-based Metric Results			Execution-based Metric Results		
	EM	QSM	QFC	EX	EFM	EVM
Instructing LLM	5.91%	50.77%	64.32%	35.06%	53.05%	58.23%
Few-shot LLM	10.41%	48.40%	55.46%	35.82%	66.99%	68.25%
RAG for LLM	16.32%	63.06%	70.05%	53.26%	73.95%	62.02%
Fine-tuned LLaMA	20.54%	56.50%	67.68%	53.12%	84.36%	68.11%
SMART	21.77%	60.65%	72.86%	59.57%	84.90%	69.73%
MARMo (Ours)	22.09%	62.81%	74.52%	62.34%	83.71%	66.88%

Table 1: The overall performance comparison of different methods on the TEND benchmark, including both query-based and execution-based evaluation metrics.

iterations is set to be 3. The memory mechanism is constructed from the training set. Similarity is computed via vector-based semantic similarity by text-embedding-ada-002.

5.1 Dataset

We evaluate our approach on the TEND benchmark (Lu et al., 2025), a large-scale benchmark for complex NoSQL query understanding especially on MongoDB. As shown in Table 2, it comprises 154 MongoDB databases across 105 domains, with 347 collections and 5,960 unique fields. The dataset includes 17,020 natural language and NoSQL query pairs, with each NoSQL query aligned to five NLQs. This diversity in both schema and query structure makes TEND a challenging and comprehensive testbed for evaluating Text-to-NoSQL models on tasks that demand nuanced understanding of both natural language and structured query generation.

Property	Value
Databases	154
Domains	105
Collections	347
Unique fields	5,960
NL-MQL pairs	17,020

Table 2: Basic statistics of the TEND dataset.

5.2 Baselines and Metrics

We compare our proposed framework MARMo against a range of strong baselines which are aligned with the experimental configurations in (Lu et al., 2025), whose details are available in Appendix A.1, ensuring fair and consistent evaluation. Our main experiments compare different paradigms of closed-source LLM utilization, including zero-shot prompting, few-shot in-

context learning, and retrieval-augmented generation (RAG), as well as domain-adapted open-source LLMs via fine-tuning, alongside the most recent SMART framework. Collectively, these baselines provide a comprehensive spectrum of methods, enabling us to rigorously position MARMo against both prior neural solutions and the latest advances in LLM-based Text-to-NoSQL generation.

Meanwhile, to comprehensively evaluate model performance on the Text-to-NoSQL task, we adopt six complementary metrics, which can be broadly categorized into query-based and execution-based evaluations, following the settings of TEND, and details are elaborated in the Appendix A.2. Specifically, Query-based metrics emphasize structural fidelity, measuring how well the generated NoSQL queries align with the gold-standard annotations at the syntactic and representational levels. In contrast, execution-based metrics assess functional correctness, focusing on whether the generated queries return the intended results when executed against the underlying databases. This dual perspective allows us to jointly capture surface-level accuracy and end-task utility.

5.3 Experimental Results

Table 1 summarizes the performance of MARMo against representative baselines on the TEND benchmark. Evaluation is conducted using six complementary metrics that jointly capture query-level fidelity and execution-level correctness.

MARMo the method achieves the overall balance across all evaluated metrics, highlighting the effectiveness of its repair-oriented design. In terms of query-based metrics, MARMo achieves an Exact Match (EM) of 22.09%, which, despite the strictness of surface-form alignment, demonstrates the framework’s ability to generate queries that are structurally close to gold references. For structural

evaluation, MARMo attains 62.81% Query Stages Match (QSM) and 74.52% Query Fields Coverage (QFC), reflecting its strength in reconstructing logical steps and correctly referencing schema elements. On execution-based metrics, MARMo reaches 62.34% Execution Accuracy (EX), showing robust end-to-end performance in producing executable queries. However, the proposed method exhibits a certain degradation on the EFM and EVM metrics compared to SMART. By analyzing the execution results of the golden queries and the queries predicted by MARMo, we observe that this issue may stem from inconsistencies in field naming introduced during the repair process. Specifically, different models may adopt non-unified naming conventions when representing the same semantic fields. For example, for the question Give me the names of schools and their associated drivers for all school bus routes., the golden query uses the fields School and Name, whereas the repaired predicted query employs School and Driver_Name. By explicitly using Driver_Name, the repaired query avoids the semantic ambiguity associated with the generic field name Name. Taken together, these results indicate that MARMo achieves improvements primarily at the symbolic and structural levels of query formulation, while generally preserving correctness in execution-based evaluation. The overall performance gains validate the utility of its modular multi-agent repair strategy, underscoring its potential as a foundation for advancing Text-to-NoSQL query generation.

5.4 Case Study

To illustrate the effectiveness of MARMo in repairing incorrect Text-to-NoSQL predictions, we present a representative case study from our experimental dataset, as summarized in Table 3. The natural language query (NLQ) posed was 'What is the predominant age of editors in the dataset?'. The initially predicted query computed the average age of all editors rather than identifying the most frequent age. This mismatch demonstrates a typical semantic error where the model captures a general statistical aggregation rather than the mode, reflecting a subtle misalignment between the NLQ and the predicted aggregation logic. MARMo successfully detected this semantic discrepancy through its error classification and feedback generation modules. This case highlights MARMo's capability to not only identify incorrect aggregation functions

NLQ	What is the predominant age of editors in the dataset?
Target NoSQL	<pre>db.editor.aggregate([{ \$group: { _id: "\$Age", count: { \$sum: 1 } } }, { \$sort: { count: -1 } }, { \$limit: 1 }, { \$project: { _id: 0, Age: "\$_id" } }]]);</pre> <p>Target: Success to find the most frequent Age among editors, returns 1 result.</p> <p>Result: [{"Age": 20}]</p>
Before Fix	<pre>db.editor.aggregate([{ \$group: { _id: null, avg_Age: { \$avg: "\$Age" } } }, { \$project: { _id: 0, avg_Age: 1 } }]]);</pre> <p>Error: Wrong computing the average age instead of finding the most frequent age (mode), leading to unexpected results.</p> <p>Result: [{"avg_Age": 32.66666}]</p>
After Fix by MARMo (Ours)	<pre>db.editor.aggregate([{ \$group: { _id: "\$Age", count: { \$sum: 1 } } }, { \$sort: { count: -1 } }, { \$limit: 1 }, { \$project: { _id: 0, Age: "\$_id" } }]]);</pre> <p>Correction: Find the semantic mismatch of NLQ and query, fix and generate the correct query.</p> <p>Result: [{"Age": 20}]</p>

Table 3: The shown example is a natural language query, the original incorrect NoSQL query, and the corrected query by MARMo, along with their execution results.

but also systematically transform the query to align with the user intent expressed in natural language.

6 Conclusions

In this paper, we propose MARMo, a novel multi-agent repair framework designed to address Text-to-NoSQL query fix generation, with a focus on MongoDB. MARMo systematically handles prediction errors by incorporating dedicated agents for error detection, classification, feedback generation, and query refinement. Extensive experiments demonstrate that MARMo enhances the accuracy and reliability of query generation overall, outperforming existing approaches. This work highlights the potential of collaborative agent-based architectures in improving the interpretability and robustness of natural language interfaces to NoSQL databases. Future work will explore MARMo to broader NoSQL paradigms, such as document-graph hybrids and multi-modal interfaces.

606 Limitations

607 We propose MARMo, a novel multi-agent repair
608 framework for enhancing the accuracy of Text-to-
609 NoSQL parsing, inspired by the human process
610 of error correction. MARMo employs an efficient
611 error detection mechanism coupled with a type-
612 aware repair strategy implemented in a plug-and-
613 play fashion. Through extensive comparisons with
614 competitive baselines and detailed case studies, we
615 demonstrate that MARMo consistently achieves a
616 high level of repair performance.

617 Although the empirical results are promising,
618 the current evaluation is primarily conducted on
619 the TEND dataset, with a relatively limited set
620 of baselines. Due to data availability and incom-
621 plete literature preparation, we were unable to
622 obtain the complete MultiTEND and DocSpider
623 datasets. To ensure reproducibility and consistency,
624 these datasets were not included in the main ex-
625 periments. MARMo was originally developed and
626 implemented with the TEND Benchmark in mind,
627 incorporating optimizations specific to MongoDB
628 semantics, error types, and repair patterns. Ex-
629 tending the framework to other datasets would re-
630 quire additional engineering effort. As one of the
631 major benchmarks in the Text-to-NoSQL domain,
632 TEND’s complexity and diversity allow MARMo’s
633 capabilities to be effectively demonstrated. These
634 choices allow us to maintain a controlled and fo-
635 cused experimental setting, while also emphasizing
636 the extensibility of MARMo to Text-to-NoSQL
637 parsing tasks and the novelty of integrating multi-
638 agent collaboration with error repair in this context.
639 Meanwhile, to ensure the accuracy of our evalu-
640 ation stability, the current implementation of the
641 fixed memory in MARMo relies solely on the train-
642 ing set and does not incorporate any data from the
643 evaluation set and real-time fix process.

644 For future work, we plan to broaden the empiri-
645 cal scope by applying MARMo to a wider variety
646 of datasets, different NoSQL databases and incor-
647 porating more diverse baseline systems, which will
648 further assess the generalizability and robustness
649 of our framework.

650 References

651 Josh Achiam, Steven Adler, Sandhini Agarwal, Lama
652 Ahmad, Ilge Akkaya, Florencia Leoni Aleman,
653 Diogo Almeida, Janko Altenschmidt, Sam Altman,
654 Shyamal Anadkat, and 1 others. 2023. Gpt-4 techni-
655 cal report. *arXiv preprint arXiv:2303.08774*.

- Ruisheng Cao, Lu Chen, Zhi Chen, Yanbin Zhao,
Su Zhu, and Kai Yu. 2021. [LGESQL: Line graph
enhanced text-to-SQL model with mixed local and
non-local relations](#). In *Proceedings of the 59th An-
nual Meeting of the Association for Computational
Linguistics and the 11th International Joint Confer-
ence on Natural Language Processing (Volume 1:
Long Papers)*, pages 2541–2555, Online. Association
for Computational Linguistics. 656–664
- Jipeng Cen, Jiaxin Liu, Zhixu Li, and Jingjing Wang.
2025. [Sqlfixagent: Towards semantic-accurate text-
to-sql parsing via consistency-enhanced multi-agent
collaboration](#). In *Proceedings of the AAAI Confer-
ence on Artificial Intelligence*, volume 39, pages 49–
57. 665–670
- Yeounoh Chung, Gaurav T. Kakkar, Yu Gan, Brenton
Milne, and Fatma Özcan. 2025. [Is long context all
you need? leveraging llm’s extended context for
nl2sql](#). *Proc. VLDB Endow.*, 18(8):2735–2747. 671–674
- Minghang Deng, Ashwin Ramachandran, Canwen Xu,
Lanxiang Hu, Zhewei Yao, Anupam Datta, and Hao
Zhang. 2025. [Reforce: A text-to-sql agent with self-
refinement, consensus enforcement, and column ex-
ploration](#). *arXiv preprint arXiv:2502.00675*. 675–679
- Yusuf Denizay Dönder, Derek Hommel, Andrea W
Wen-Yi, David Mimno, and Unso Eun Seo Jo. 2025.
[Cheaper, better, faster, stronger: Robust text-to-
sql without chain-of-thought or fine-tuning](#). *arXiv
preprint arXiv:2505.14174*. 680–684
- Elias Dritsas and Maria Trigka. 2025. [Database systems
in the big data era: Architectures, performance, and
open challenges](#). *IEEE Access*, 13:95068–95084. 685–687
- Yuankai Fan, Tonghui Ren, Zhenying He, X Sean Wang,
Ye Zhang, and Xingang Li. 2023. [Gensql: A genera-
tive natural language interface to database systems](#).
In *2023 IEEE 39th International Conference on Data
Engineering (ICDE)*, pages 3603–3606. IEEE. 688–692
- Yujian Gan, Xinyun Chen, Qiuping Huang, Matthew
Purver, John R. Woodward, Jinxia Xie, and Peng-
sheng Huang. 2021a. [Towards robustness of text-
to-SQL models against synonym substitution](#). In
*Proceedings of the 59th Annual Meeting of the Asso-
ciation for Computational Linguistics and the 11th
International Joint Conference on Natural Language
Processing (Volume 1: Long Papers)*, pages 2505–
2515, Online. Association for Computational Lin-
guistics. 693–702
- Yujian Gan, Xinyun Chen, and Matthew Purver. 2021b.
[Exploring underexplored limitations of cross-domain
text-to-SQL generalization](#). In *Proceedings of the
2021 Conference on Empirical Methods in Natural
Language Processing*, pages 8926–8931, Online and
Punta Cana, Dominican Republic. Association for
Computational Linguistics. 703–709
- Yingqi Gao, Yifu Liu, Xiaoxia Li, Xiaorong Shi, Yin
Zhu, Yiming Wang, Shiqi Li, Wei Li, Yuntao Hong,
710–711

712	Zhiling Luo, and 1 others. 2024. Xiyan-sql: A multi-generator ensemble framework for text-to-sql. <i>arXiv preprint arXiv:2411.08599</i> .	<i>SIGMOD international conference on Management of data</i> , pages 709–712.	767
713			768
714			
715	Orest Gkini, Theofilos Belmpas, Georgia Koutrika, and Yannis Ioannidis. 2021. An in-depth benchmarking of text-to-sql systems. In <i>Proceedings of the 2021 International Conference on Management of Data</i> , pages 632–644.	Haoyang Li, Jing Zhang, Hanbing Liu, Ju Fan, Xiaokang Zhang, Jun Zhu, Renjie Wei, Hongyan Pan, Cuiping Li, and Hong Chen. 2024a. Codes: Towards building open-source language models for text-to-sql. <i>Proceedings of the ACM on Management of Data</i> , 2(3):1–28.	769
716			770
717			771
718			772
719			773
720	Zhibin Gou, Zhihong Shao, Yeyun Gong, Yelong Shen, Yujiu Yang, Nan Duan, and Weizhu Chen. 2024. CRITIC: Large language models can self-correct with tool-interactive critiquing . In <i>The Twelfth International Conference on Learning Representations</i> .	Jinyang Li, Binyuan Hui, Ge Qu, Jiayi Yang, Binhua Li, Bowen Li, Bailin Wang, Bowen Qin, Ruiying Geng, Nan Huo, and 1 others. 2023. Can llm already serve as a database interface? a big bench for large-scale database grounded text-to-sqls. <i>Advances in Neural Information Processing Systems</i> , 36:42330–42357.	774
721			775
722			776
723			777
724			778
725	Aaron Grattafiori, Abhimanyu Dubey, Abhinav Jauhri, Abhinav Pandey, Abhishek Kadian, Ahmad Al-Dahle, Aiesha Letman, Akhil Mathur, Alan Schelten, Alex Vaughan, and 1 others. 2024. The llama 3 herd of models. <i>arXiv preprint arXiv:2407.21783</i> .	Yuntao Li, Hanchu Zhang, Yutian Li, Sirui Wang, Wei Wu, and Yan Zhang. 2021. Pay more attention to history: A context modelling strategy for conversational text-to-sql. <i>arXiv preprint arXiv:2112.08735</i> .	779
726			780
727			781
728			782
729			783
730	Siwei Han, Peng Xia, Ruiyi Zhang, Tong Sun, Yun Li, Hongtu Zhu, and Huaxiu Yao. 2025. Mdocagent: A multi-modal multi-agent framework for document understanding. <i>arXiv preprint arXiv:2503.13964</i> .	Ziming Li, Qianbo Zang, David Ma, Jiawei Guo, Tuney Zheng, Minghao Liu, Xinyao Niu, Yue Wang, Jian Yang, Jiaheng Liu, and 1 others. 2024b. Autokaggle: A multi-agent framework for autonomous data science competitions. <i>arXiv preprint arXiv:2410.20424</i> .	784
731			785
732			786
733			787
734	Zijin Hong, Zheng Yuan, Qinggang Zhang, Hao Chen, Junnan Dong, Feiran Huang, and Xiao Huang. 2024. Next-generation database interfaces: A survey of llm-based text-to-sql. <i>arXiv preprint arXiv:2406.08426</i> .	Zhisheng Lin, Yifu Liu, Zhiling Luo, Jinyang Gao, and Yu Li. 2024. Momq: Mixture-of-experts enhances multi-dialect query generation across relational and non-relational databases. <i>arXiv preprint arXiv:2410.18406</i> .	788
735			789
736			790
737			791
738	Radu Cristian Alexandru Iacob, Florin Brad, Elena-Simona Apostol, Ciprian-Octavian Truică, Ionel Alexandru Hosu, and Traian Rebedea. 2020. Neural approaches for natural language interfaces to databases: A survey. In <i>proceedings of the 28th International Conference on Computational Linguistics</i> , pages 381–395.	Xinyu Liu, Shuyu Shen, Boyan Li, Nan Tang, and Yuyu Luo. 2025. Nl2sql-bugs: A benchmark for detecting semantic errors in nl2sql translation . In <i>Proceedings of the 31st ACM SIGKDD Conference on Knowledge Discovery and Data Mining V.2</i> , KDD '25, page 5662–5673, New York, NY, USA. Association for Computing Machinery.	792
739			793
740			794
741			795
742			796
743			797
744			798
745	Yoichi Ishibashi and Yoshimasa Nishimura. 2024. Self-organized agents: A llm multi-agent framework toward ultra large-scale code generation and optimization. <i>arXiv preprint arXiv:2404.02183</i> .	Xiping Liu and Zhao Tan. 2024. Epi-sql: Enhancing text-to-sql translation with error-prevention instructions. <i>arXiv preprint arXiv:2404.14453</i> .	799
746			800
747			801
748			802
749	George Katsogiannis-Meimarakis, Mike Xydas, and Georgia Koutrika. 2023. Natural language interfaces for databases with deep learning. <i>Proceedings of the VLDB Endowment</i> , 16(12):3878–3881.	Jinwei Lu, Yuanfeng Song, Zhiqian Qin, Haodi Zhang, Chen Zhang, and Raymond Chi-Wing Wong. 2025. Bridging the gap: Enabling natural language queries for nosql databases through text-to-nosql translation. <i>arXiv preprint arXiv:2502.11201</i> .	803
750			804
751			805
752			806
753	Hyeonji Kim, Byeong-Hoon So, Wook-Shin Han, and Hongrae Lee. 2020. Natural language to sql: Where are we today? <i>Proceedings of the VLDB Endowment</i> , 13(10):1737–1750.	Wenxin Mao, Ruiqi Wang, Jiyu Guo, Jichuan Zeng, Cuiyun Gao, Peiyi Han, and Chuanyi Liu. 2024. Enhancing text-to-sql parsing through question rewriting and execution-guided refinement. In <i>Findings of the Association for Computational Linguistics ACL 2024</i> , pages 2009–2024.	807
754			808
755			809
756			810
757	Bingxuan Li, Yiwei Wang, Jiuxiang Gu, Kai-Wei Chang, and Nanyun Peng. 2025. METAL: A multi-agent framework for chart generation with test-time scaling . In <i>Proceedings of the 63rd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)</i> , pages 30054–30069, Vienna, Austria. Association for Computational Linguistics.	Ana Nunez, Nafis Tanveer Islam, Sumit Kumar Jha, and Peyman Najafirad. 2024. Autosafecoder: A multi-agent framework for securing llm code generation through static analysis and fuzz testing. <i>arXiv preprint arXiv:2409.10737</i> .	811
758			812
759			813
760			814
761			815
762			816
763			817
764	Fei Li and Hosagrahar V Jagadish. 2014. Nalir: an interactive natural language interface for querying relational databases. In <i>Proceedings of the 2014 ACM</i>		818
765			819
766			820

821	Mohammadreza Pourreza, Hailong Li, Ruoxi Sun, Yeounoh Chung, Shayan Talaei, Gaurav Tarlok Kakkar, Yu Gan, Amin Saberi, Fatma Ozcan, and Sercan O Arik. 2024. Chase-sql: Multi-path reasoning and preference optimized candidate selection in text-to-sql. <i>arXiv preprint arXiv:2410.01943</i> .	878
822		879
823		880
824		881
825		882
826		
827	Mohammadreza Pourreza and Davood Rafiei. 2023. Din-sql: Decomposed in-context learning of text-to-sql with self-correction. <i>Advances in Neural Information Processing Systems</i> , 36:36339–36348.	883
828		884
829		885
830		886
831	Bowen Qin, Binyuan Hui, Lihan Wang, Min Yang, Jinyang Li, Binhua Li, Ruiying Geng, Rongyu Cao, Jian Sun, Luo Si, and 1 others. 2022. A survey on text-to-sql parsing: Concepts, methods, and future directions. <i>arXiv preprint arXiv:2208.13629</i> .	887
832		888
833		889
834		890
835		891
836	Zhiqian Qin, Yuanfeng Song, Jinwei Lu, Yuanwei Song, Shuaimin Li, and Chen Jason Zhang. 2025. Mul-tiTEND: A multilingual benchmark for natural language to NoSQL query translation . In <i>Findings of the Association for Computational Linguistics: ACL 2025</i> , pages 24632–24657, Vienna, Austria. Association for Computational Linguistics.	892
837		893
838		894
839		895
840		896
841		897
842		898
843		899
844	Tonghui Ren, Yuankai Fan, Zhenying He, Ren Huang, Jiaqi Dai, Can Huang, Yinan Jing, Kai Zhang, Yifan Yang, and X Sean Wang. 2024. Purple: Making a large language model a better sql writer. In <i>2024 IEEE 40th International Conference on Data Engineering (ICDE)</i> , pages 15–28. IEEE.	900
845		901
846		902
847		
848		
849	Torsten Scholak, Nathan Schucher, and Dzmitry Bahdanau. 2021. PICARD: Parsing incrementally for constrained auto-regressive decoding from language models . In <i>Proceedings of the 2021 Conference on Empirical Methods in Natural Language Processing</i> , pages 9895–9901, Online and Punta Cana, Dominican Republic. Association for Computational Linguistics.	903
850		904
851		905
852		906
853		907
854		908
855		909
856		910
857	Jiawei Shen, Chengcheng Wan, Ruoyi Qiao, Jiazhen Zou, Hang Xu, Yuchen Shao, Yueling Zhang, Weikai Miao, and Geguang Pu. 2025. A study of in-context-learning-based text-to-sql errors. <i>arXiv preprint arXiv:2501.09310</i> .	911
858		912
859		913
860		914
861		915
862	Vladislav Shkapenyuk, Divesh Srivastava, Theodore Johnson, and Parisa Ghane. 2025. Automatic metadata extraction for text-to-sql. <i>arXiv preprint arXiv:2505.19988</i> .	916
863		917
864		918
865		919
866		920
867	Quoc-Bao-Huy Tran, Aagha Abdul Waheed, and Sun-Tae Chung. 2024. Robust text-to-cypher using combination of bert, graphstage, and transformer (cobgt) model. <i>Applied Sciences</i> , 14(17):7881.	921
868		922
869		923
870	Bing Wang, Changyu Ren, Jian Yang, Xinnian Liang, Jiaqi Bai, LinZheng Chai, Zhao Yan, Qian-Wen Zhang, Di Yin, Xing Sun, and Zhoujun Li. 2025. MAC-SQL: A multi-agent collaborative framework for text-to-SQL . In <i>Proceedings of the 31st International Conference on Computational Linguistics</i> , pages 540–557, Abu Dhabi, UAE. Association for Computational Linguistics.	924
871		925
872		926
873		927
874		928
875		929
876		
877		
	Yihang Xiao, Jinyi Liu, Yan Zheng, Xiaohan Xie, Jianye Hao, Mingzhi Li, Ruitao Wang, Fei Ni, Yuxiao Li, Jintian Luo, and 1 others. 2024. Cellagent: An llm-driven multi-agent framework for automated single-cell data analysis. <i>arXiv preprint arXiv:2407.09811</i> .	930
		931
		932
		933
		934
	Wenxuan Xie, Gaochen Wu, and Bowen Zhou. 2024a. Mag-sql: Multi-agent generative approach with soft schema linking and iterative sub-sql refinement for text-to-sql. <i>arXiv preprint arXiv:2408.07930</i> .	935
		936
		937
		938
		939
	Yuanzhen Xie, Xinzhou Jin, Tao Xie, Matrixmxlin Matrixmxlin, Liang Chen, Chenyun Yu, Cheng Lei, Chengxiang Zhuo, Bo Hu, and Zang Li. 2024b. Decomposition for enhancing attention: Improving LLM-based text-to-SQL through workflow paradigm . In <i>Findings of the Association for Computational Linguistics: ACL 2024</i> , pages 10796–10816, Bangkok, Thailand. Association for Computational Linguistics.	940
		941
		942
		943
		944
		945
		946
		947
		948
		949
		950
		951
		952
		953
		954
		955
		956
		957
		958
		959
		960
		961
		962
		963
		964
		965
		966
		967
		968
		969
		970
		971
		972
		973
		974
		975
		976
		977
		978
		979
		980
		981
		982
		983
		984
		985
		986
		987
		988
		989
		990
		991
		992
		993
		994
		995
		996
		997
		998
		999
		1000

- 935 Victor Zhong, Caiming Xiong, and Richard Socher.
936 2017. Seq2sql: Generating structured queries from
937 natural language using reinforcement learning. *arXiv*
938 *preprint arXiv:1709.00103*.
- 939 Yuhang Zhou, Yu He, Siyu Tian, Yuchen Ni, Zhangyue
940 Yin, Xiang Liu, Chuanjun Ji, Sen Liu, Xipeng
941 Qiu, Guangnan Ye, and Hongfeng Chai. 2024. *r³-*
942 *NL2GQL: A model coordination and knowledge*
943 *graph alignment approach for NL2GQL*. In *Find-*
944 *ings of the Association for Computational Linguistics:*
945 *EMNLP 2024*, pages 13679–13692, Miami, Florida,
946 USA. Association for Computational Linguistics.
- 947 Yizhang Zhu, Shiyin Du, Boyan Li, Yuyu Luo, and Nan
948 Tang. 2024. Are large language models good statisti-
949 cians? *Advances in Neural Information Processing*
950 *Systems*, 37:62697–62731.

951
952
953
954
955
956
957
958
959
960
961
962
963
964
965
966
967
968
969
970
971
972
973
974
975
976
977
978
979
980
981
982
983
984
985
986
987
988
989
990
991
992
993
994
995
996
997

A Appendix

A.1 Baselines

We compare our proposed framework against several strong baselines, including both traditional neural models and recent LLM-based approaches for Text-to-NoSQL generation. Baselines settings follow (Lu et al., 2025).

- **Zero-shot LLM** leverages the inherent zero-shot capabilities of LLMs to generate contextually appropriate responses without requiring task-specific examples or fine-tuning.
- **Few-shot LLM** employs in-context learning by incorporating a limited number of static task-relevant examples into the prompt, enabling LLMs to adapt to domain-specific tasks with context.
- **RAG for LLM** enhances LLM performance through Retrieval-Augmented Generation, which dynamically incorporates external knowledge into prompts, mitigating hallucinations caused by incomplete or ambiguous input.
- **Finetuned LLaMA (Grattafiori et al., 2024)** refers to the LLaMA model finetuned with TEND training set to predicted MQL. In practice we use LLaMA-3.2-1B-Instruct to implement the Text-to-NoSQL task.
- **SMART** is the first framework targeting Text-to-NoSQL generation, SMART integrates SLMs and RAG to enable schema prediction, query generation, refinement via retrieved examples, and optimization based on execution feedback.
- **MARMo** is the first solution that applies multi-agent collaboration architecture to Text-to-NoSQL. It introduces iterative refinement process and three dedicated agents responsible for error detection and classification, feedback synthesis, and query refinement, respectively. In addition, its modular and pluggable design allows it to be seamlessly integrated into a wide range of Text-to-NoSQL systems as an auxiliary repair component.

A.2 Metrics

To comprehensively evaluate model performance on the Text-to-NoSQL task, we adopt six complementary metrics, grouped into **query-based** and

execution-based evaluations, following TEND. Query-based metrics focus on structural fidelity: 998 999

- **Exact Match (EM)** requires the generated query to exactly match the gold query in both structure and content. It provides a strict measure of syntactic and semantic alignment. 1000 1001 1002 1003
- **Query Stages Match (QSM)** evaluates the correctness of the query pipeline by comparing stage types (e.g. match, group, lookup) and order with these of the gold query. 1004 1005 1006 1007
- **Query Fields Coverage (QFC)** measures whether all relevant fields—either from the schema or derived—are correctly included, reflecting schema comprehension. It evaluates the performance of methods on database schema coverage. 1008 1009 1010 1011 1012 1013

On the other side, execution-based metrics assess functional correctness: 1014 1015

- **Execution Accuracy (EX)** checks if the generated query yields fully identical results to the gold query when executed. EX is the most critical performance metric for evaluating Text-to-NoSQL models. 1016 1017 1018 1019 1020
- **Execution Fields Match (EFM)** evaluates whether the field set of output matches the correct set of field names, as the JSON-like query result of document-based database like MongoDB has the field part (key) and the value part (value). 1021 1022 1023 1024 1025 1026
- **Execution Value Match (EVM)** verifies the correctness of output values -. 1027 1028

Together, these metrics offer a holistic view of model performance, from structural precision to executable reliability. 1029 1030 1031