# CHESS: CONTEXTUAL HARNESSING FOR EFFICIENT SQL SYNTHESIS

**Shayan Talaei** [1]  **Mohammadreza Pourreza** [2]  **Yu-Chen Chang** [1]  **Azalia Mirhoseini** [1] [*]  **Amin Saberi** [1] [*]

## ABSTRACT

Translating natural language questions into SQL queries, known as text-to-SQL, is a long-standing research problem. Effective text-to-SQL synthesis can become very challenging due to (i) the extensive size of database catalogs (descriptions of tables and their columns) and database values, (ii) reasoning over large database schemas, (iii) ensuring the functional validity of the generated queries, and (iv) navigating the ambiguities of natural language questions. We introduce CHESS, a Large Language Model (LLM) based multi-agent framework for efficient and scalable SQL synthesis, comprising four specialized agents, each targeting one of the aforementioned challenges: the Information Retriever (IR) extracts relevant data, the Schema Selector (SS) prunes large schemas, the Candidate Generator (CG) generates high-quality candidates and refines queries iteratively, and the Unit Tester (UT) validates queries through LLM-based natural language unit tests. Our framework offers configurable features that adapt to various deployment constraints, including 1) Supporting industrial-scale databases: leveraging the Schema Selector agent, CHESS efficiently narrows down very large database schemas into manageable sub-schemas, boosting system accuracy by approximately 2% and reducing the number of LLM tokens by ×5. 2) State-of-the-Art privacy-preserving performance: Among the methods using open-source models, CHESS achieves state-of-the-art performance, resulting in a high-performing, privacy-preserving system suitable for industrial deployment. 3) Scalablity with additional compute budget: In settings with high computational budgets, CHESS achieves 71.10% accuracy on the BIRD test set, within 2% of the leading proprietary method, while requiring approximately 83% fewer LLM calls.

## 1 INTRODUCTION

Generating SQL queries from natural language questions, commonly known as text-to-SQL, is a persistent and high-impact research problem. The growing complexity of databases in recent years has exacerbated this problem, due to the increasing sizes of schemas (sets of columns and tables), values (content), and catalogs (metadata describing schemas and values) stored within them. Even some of the largest proprietary models, such as GPT-4, lag significantly behind human performance on text-to-SQL benchmarks, with a notable accuracy gap of 40% (Li et al., 2024b). Beyond the complexity of writing SQL queries, this substantial gap is primarily caused by the need to effectively retrieve and integrate multiple sources of information, including database values, catalogs, and schema, each in different formats, which complicates the process. Considering all these complexities, the lack of a robust verification method during the inference stage of text-to-SQL systems can undermine the reliability of the generated SQL queries.

Figure 1 illustrates some of the challenges facing modern text-to-SQL systems. For instance, users' questions might not directly match the stored values in the database (Gan et al., 2021), making it crucial to accurately identify the value format for effective SQL query formulation. Additionally, real-world database schemas often contain ambiguous column names, table names, and messy data, further complicating the SQL translation process and necessitating a robust retrieval system to identify relevant information (Pourreza and Rafiei, 2023). Moreover, some questions may be inherently ambiguous, leading large language models (LLMs) to make subtle mistakes when generating SQL queries. For example, in the question illustrated on the right side of Figure 1, one approach might use `ORDER BY` and `LIMIT 1` to identify the highest average score, thereby returning only a single result. Another approach could involve a subquery with the `MAX()` function, which would return all entries that match the maximum `AvgScrRead` value. This distinction in methods can lead to different outputs.

**CHESS Framework.** To address these challenges, we introduce **CHESS: Contextual Harnessing for Efficient SQL Synthesis**, a multi-agent framework designed to improve SQL query generation in complex, real-world databases. Each agent within CHESS is specifically designed to tackle one or more of the aforementioned challenges:

---

[1]Stanford University, Stanford, CA, USA [2]University of Alberta, Alberta, Canada. Correspondence to: Shayan Talaei <stalaei@stanford.edu>.
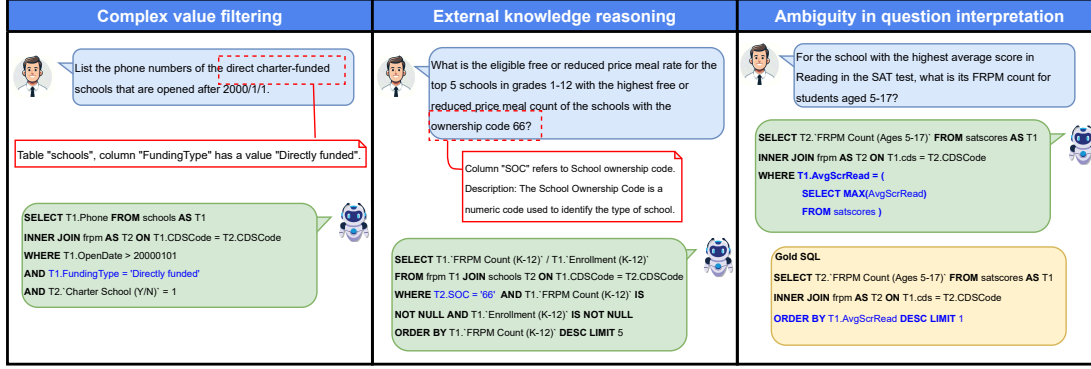
[*]Equal senior authorship.

Figure 1. Examples of challenges in text-to-SQL translation. (1) Users' questions might not have the exact database value. (2) Column names might not accurately represent their content, making database catalogs essential. (3) Ambiguity in question interpretation can result in different, seemingly valid SQL queries.

The **Information Retriever (IR)** agent aims to retrieve the relevant entities (values in the database) and context (schema descriptions provided in the database catalog). To achieve this, we present scalable and efficient methods using locality-sensitive hashing to retrieve database values from millions of rows, leveraging keyword detection, and vector databases to extract contextual information from database catalogs. Our approach utilizes both semantic and syntactic similarities between the database content and the user's query to enhance the retrieval quality.

As previous works (Guerreiro et al., 2022) and our ablation studies (Figure 3) confirm, overloading the LLM context with redundant information can degrade its performance on downstream task. To address this, the **Schema Selector (SS)** agent is specifically designed to prune very large database schemas, allowing the candidate generator agent (CG) to perform a more effective reasoning when synthesizing SQL queries. This SS agent is equipped with three tools, *filter_column*, *select_tables*, and *select_columns*, which allow it to maintain a scalable performance even with complex, real-world schemas with more than 4000 columns.

The **Candidate Generator (CG)** agent generates SQL queries based on the question posed to the database. It then executes and revises these queries as needed based on the results.

The **Unit Tester (UT)** agent assesses the quality of the final query by generating multiple natural language unit tests and evaluating the candidate queries against them. UT assigns a score to each candidate based on passed tests, selecting the highest-scoring query, thus ensuring robust query generation.

In summary, the CHESS framework leverages a collaborative multi-agent approach to enhance SQL generation for complex databases. Each agent contributes a specialized function, from retrieving relevant data and pruning schema

information to generating and validating SQL queries, collectively ensuring accurate, efficient, and context-aware query synthesis.

**CHESS's Adaptivity for Deployment Constraints.** Real-world deployment of text-to-SQL systems must consider various constraints, such as the reasoning capabilities and sampling budgets of large language models (LLMs), as well as data privacy concerns. CHESS is designed as a configurable framework that can adapt to different deployment scenarios while maintaining high performance. In subsection 4.2, we discuss various deployment settings and how CHESS can be customized to meet performance requirements while adhering to these constraints.

To address the challenge of complex industrial database schemas, CHESS incorporates a Schema Selector (SS) agent that efficiently narrows down large schemas into manageable sub-schemas, ensuring sufficient information to accurately answer queries. For environments with limited computational resources, CHESS optimizes performance by minimizing reliance on frequent and costly LLM calls, achieving comparable results using less resource-intensive models. Regarding data privacy concerns, CHESS supports the use of fully open-source models that can be deployed on-premise, eliminating the need to send private data to third-party providers. These features collectively enable CHESS to adapt to various deployment constraints, offering a flexible and high-performing solution for industrial text-to-SQL applications.

**Experiments.** To validate the effectiveness of CHESS, we conducted comprehensive experiments on challenging benchmarks, including BIRD and Spider, across various deployment scenarios. Our evaluations included testing on large-scale schemas, operating under limited computational resources, and ensuring data privacy by using fully open-source models. We also performed ablation studies to assess

the contributions of each component within CHESS.

The experimental results demonstrate that CHESS outperforms many existing state-of-the-art methods across diverse settings. Specifically, CHESS achieved 71.10% accuracy on the BIRD test set, within 2% of the leading proprietary method (Pourreza et al., 2024), while requiring approximately 83% fewer LLM calls. Among the methods using open-source models, CHESS sets a new benchmark with 61.5% accuracy on the BIRD development set. Moreover, we tested CHESS on very large industrial-scale database schema with over 4000 columns. Leveraging the Schema Selector agent, CHESS effectively prunes irrelevant columns, resulting in a 2% increase in accuracy and ×5 reduction in token usage. The ablation studies confirm the critical role of components like the Information Retriever and Unit Tester in enhancing overall performance, with each contributing substantially to the accuracy and efficiency of the system. Finally, it is shown that CHESS maintains high accuracy comparable to existing methods even under computational constraints, making it suitable for resource-limited deployments while addressing data privacy concerns.

**Contributions.** Concretely, here our contributions:

- **Novel Multi-Agent Framework**: We propose CHESS, a multi-agent framework that significantly improves text-to-SQL performance and efficiency on complex, real-world databases, utilizing four specialized agents: Information Retriever (IR), Schema Selector (SS), Candidate Generator (CG), and Unit Tester (UT).

- **Innovative Unit Test Generation**: We introduce a novel natural language unit test generation method for text-to-SQL, enabling the effective evaluation of generated SQL queries during inference.

- **Computationally Efficient Scalable Retrieval Approach**: We develop a highly efficient hierarchical retrieval method for extracting relevant entities and context from vast datasets, significantly enhancing SQL prediction accuracy.

- **Schema Selection Protocol applicable to Industrial-Scale Database Schemas**: Our schema selection protocol effectively prunes very large industrial-scale schemas, including those with over 4,000 columns. To our knowledge, this is the first method capable of handling such large schemas at this level of complexity.

- **State-of-the-Art Privacy-Preserving Performance**: Among the methods using open-source models, we achieve state-of-the-art performance, resulting in a high-performing, privacy-preserving system suitable for industrial deployment.

The remainder of this paper is organized as follows: Section 2 reviews related work, Section 3 details the CHESS framework, Section 4 presents experimental results, and Section 5 concludes with discussions on future work. All codes to reproduce the results reported in this paper are available on our GitHub repository.*

## 2 RELATED WORK

Generating accurate SQL queries from natural language questions, known as text-to-SQL, is an active research area in NLP and database fields. Initial progress involved custom templates (Zelle and Mooney, 1996), requiring significant manual effort. More recent approaches leverage transformer-based, sequence-to-sequence models (Vaswani et al., 2017; Sutskever et al., 2014), which are well-suited for tasks involving sequence generation, including text-to-SQL (Qin et al., 2022). Initial sequence-to-sequence models, such as IRNet (Guo et al., 2019), used bidirectional LSTM architectures with self-attention for encoding queries and database schemas. To better integrate schema information, models like RAT-SQL (Wang et al., 2019) and RASAT (Qi et al., 2022) incorporated relation-aware self-attention, while SADGA (Cai et al., 2021) and LGESQL (Cao et al., 2021) employed graph neural networks for schema-query relationships. Despite these advancements, sequence-to-sequence models still fall short of human-level performance, with none achieving over 80% execution accuracy on the Spider hold-out test set (Yu et al., 2018).

Alongside the widespread adoption of LLMs across various NLP domains, the text-to-SQL field has similarly benefited from recent methodological innovations with LLMs to enhance performance. Early approaches (Rajkumar et al., 2022), leveraged the zero-shot in-context learning capabilities of LLMs for SQL generation. Building on this, subsequent models including DIN-SQL (Pourreza and Rafiei, 2024a), DAIL-SQL (Gao et al., 2023), MAC-SQL (Wang et al., 2023), and C3 (Dong et al., 2023) have enhanced LLM performance through task decomposition. In addition to in-context learning, proposals in DAIL-SQL (Gao et al., 2023), DTS-SQL (Pourreza and Rafiei, 2024b), and CodeS (Li et al., 2024a) have sought to elevate the capabilities of open-source LLMs through supervised fine-tuning, aiming to rival or exceed their larger, proprietary counterparts. However, the most notable performance gains have been observed in proprietary LLMs utilizing in-context learning methods (Li et al., 2024b). Unlike previous efforts, this paper introduces a hybrid approach that combines both in-context learning and supervised fine-tuning to further enhance performance. Moreover, we propose novel methods to integrate contextual data such as database values and database catalog into the text-to-SQL pipeline, leveraging a rich yet often overlooked

---

*https://github.com/ShayanTalaei/CHESS

source of information. Moreover we also propose natural language generation for SQL query generation to enhance the reliability of the text-to-SQL systems. In contrast to most previous works, the Distillery method (Maamari et al., 2024) demonstrates that the latest versions of LLMs can effectively handle database schema information with up to 200 columns within their prompt, eliminating the need for a separate schema linking step that could introduce errors into the pipeline. In our study, we confirm that for benchmarks like BIRD (Li et al., 2024b), where the schema contains approximately 100 columns, schema linking becomes unnecessary. However, for larger schemas with around 4000 columns, performance degradation still occurs due to the increased prompt size, making a separate schema linking step necessary.

Independently, but concurrently with our work, Large Language Monkeys (Brown et al., 2024), Archon (Saad-Falcon et al., 2024), MCS-SQL (Lee et al., 2024), and CHASE-SQL (Pourreza et al., 2024) introduced methods that rely on generating a large number of candidate responses for given questions at inference time. The latter two methods, designed specifically for the text-to-SQL domain, use a selection algorithm to compare and select the best candidate responses. However, unlike these approaches that require numerous LLM calls, we propose a framework of multiple agents that can be used adaptively based on the available computational budget. This framework can either generate multiple candidates and filter them using unit tests, or, in cases where minimizing LLM calls is essential, it allows the combination of different agents to achieve accurate results while staying within budget constraints.

## 3 METHODOLOGY

"CHESS: Contextual Harnessing for Efficient SQL Synthesis" is a multi-agent framework consisting of four agents, each equipped with specialized tools to achieve specific objectives. Below, we provide a detailed explanation of the functionality and tools of each agent. See Figure 2 for an illustration of overall structure of CHESS.

Real-world deployment of text-to-SQL systems must consider various constraints, such as the reasoning capabilities and sampling budgets of large language models (LLMs), as well as data privacy concerns. CHESS is designed as a configurable framework that can adapt to different deployment scenarios while maintaining high performance. In subsection 4.2, we discuss various deployment settings and how CHESS can be customized to meet performance requirements while adhering to these constraints.

### 3.1 Agents: Descriptions and Tools

In this subsection, we introduce the core agents of the CHESS framework and their respective tools, each specifically designed to tackle the unique challenges of the text-to-SQL task. The four agents, Information Retriever (IR), Schema Selector (SS), Candidate Generator (CG), and Unit Tester (UT), play essential roles in the framework, leveraging specialized tools to perform their assigned functions. Below, we provide a detailed explanation of each agent's role and the tools they employ to enable efficient and high-performance SQL query generation. Implementation details and prompt templates are provided in Appendix A and C.

**Information Retriever (IR)** The information retriever agent equipped with the tools, *extract_keywords*, *entity_retriever*, and *context_retriever*, gathers relevant information related to the input, including entities mentioned in the question and the contextual information provided in the database catalog. The tools are described as follows:

***extract_keywords***. To search for the similar values in the database and schema description, the agent needs to extract the main keywords from the natural language question. This tool uses a few-shot LLM call to extract the primary keywords and key phrases from the input.

***retrieve_entity***. From the list of keywords extracted from the question, some may correspond to entities present in the database. This tool allows the agent to search for similar values in the database and return the most relevant matches, along with their corresponding columns, for each keyword. To assess syntactic similarity between the keywords and the database values, we employ the edit distance similarity metric. Additionally, to enhance retrieval efficiency, we propose a hierarchical retrieval strategy based on Locality Sensitive Hashing (LSH) and semantic (embedding) similarity measures, which we explain in detail in Appendix A. This approach enables the agent to retrieve values that exhibit both syntactic and semantic similarity to the keywords effectively.

***retrieve_context***. In addition to retrieving values, the IR agent can access the database catalog, which often includes schema metadata, such as column descriptions, extended column names (to resolve abbreviations), and value descriptions. As shown in Figure 1, this information can be useful, and not providing it to the model can lead to suboptimal performance. The IR agent uses this tool to retrieve the most relevant descriptions from the database catalog by querying a vector database of descriptions, constructed during the preprocessing. Retrieval is based on semantic (embedding) similarity, ensuring that the most relevant context is provided to the model.

**Schema Selector (SS)** The goal of the Schema Selector (SS) agent is to reduce the schema size by selecting only the
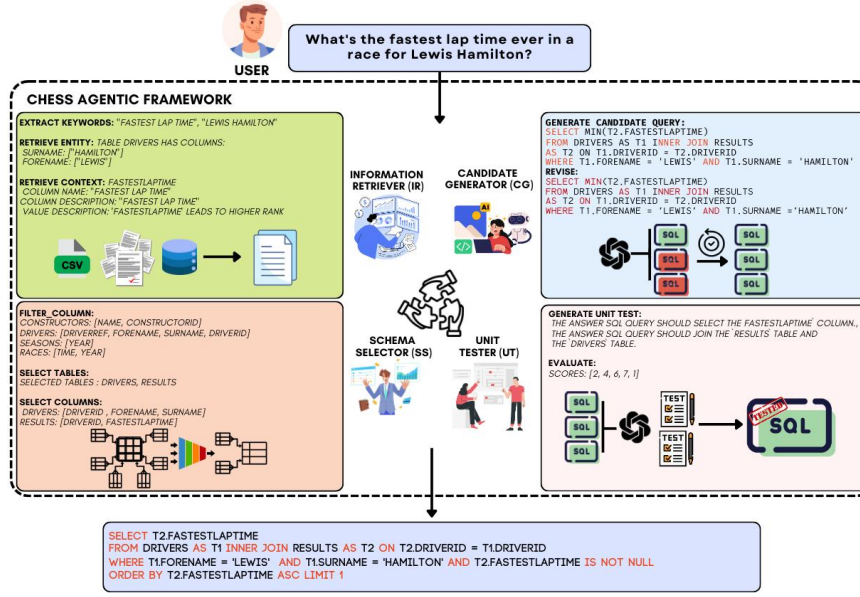
*Figure 2.* CHESS multi-agent framework for Text-to-SQL task including four agents, Information Retriever (IR), Schema Selector (SS), Candidate Generation (CG), and Unit Tester (UT).

necessary tables and columns required for generating the SQL query. To achieve this, the SS agent is equipped with three tools, *filter_column*, *select_tables*, and *select_columns*.

*filter_column* Databases often contain hundreds of columns, many of which may be semantically irrelevant to the user's query. Inspired by this observation, we design this tool which takes a column name and the question as input and determines whether the column is relevant. While this task could be perform using embedding-based similarity methods, we opt for a relatively inexpensive LLM to ensure high accuracy.

Although this tool is effective at filtering out irrelevant columns, identifying relevant schema elements accurately often requires processing multiple schema items together. The following tools complement this process by handling such scenarios.

*select_tables* This tool takes a sub-schema of the database and the question as input, and through LLM prompting, returns the tables that are necessary for answering the query.

*select_columns* To further narrow down the relevant schema items, the *select_columns* tool can be used. By inputting a sub-schema and the question, the agent retrieves only the necessary columns.

The SS agent can utilize these tools either individually or in combination, depending on its configuration. However, its actions involve a trade-off between precision and recall when selecting relevant schema items. Notably, when deal-

ing with extremely large databases—containing over 4,000 columns—this agent can significantly reduce irrelevant information. A detailed discussion of the importance of SS agent is presented in Section 4.2.2.

**Candidate Generator (CG)** The Candidate Generator (CG) is responsible for synthesizing SQL query that answers the question asked from the database. To accomplish this, the CG agent first calls the *generate_candidate_query* tool to generate candidate queries. It then executes these candidates on the database, checks the result, and identifies any faulty queries (those containing syntactic errors or empty result). For each faulty candidate, the agent repeatedly attempts to *revise* the query, until the issue is resolved or a maximum number of allowed revisions is reached.

*generate_candidate_query* This tool generates a single candidate query that answers the question. It takes the question, the schema, and the context (entities and descriptions) and prompts an LLM to follow a multi-step reasoning guideline to write a candidate SQL query.

*revise* Sometimes, the generated candidate queries may contain syntax errors or produce empty result. In such cases, the agent identifies the issue by executing the queries on the database and then uses this tool to fix them. This tool takes the question, schema, context, the faulty query and a description of the issue as input. It then prompts an LLM with this information, asking it revise the query. Note that the issue description (or execution log) is critical in guiding the model, as it provides a direct signal of the query's failure

point.

**Unit Tester (UT)** The Unit Tester (UT) agent is responsible for selecting the most accurate SQL query from the pool of candidates generated by the CG agent. UT identifies the best candidate by 1) generating multiple unit tests that highlight differences between the candidate queries and 2) evaluating the candidates against these unit tests. It then assigns a score to each query based on the number of unit tests it passes, selecting the top-scoring candidate as the final SQL query for the given question.

*generate_unit_test* This tool prompts an LLM to generate $k$ unit tests, where $k$ is an input parameter, designed such that only the correct SQL query can pass each of them. The prompt is carefully constructed to produce high-quality unit tests that highlight the semantic differences between the candidate queries. Detailed prompt used to generate unit tests is provided in Appendix C.

*evaluate* After generating the unit tests, the UT agent evaluates the candidate queries against them. This tool takes multiple candidate queries and a single unit test as input, prompting an LLM to reason through each candidate and determine whether it passes the unit test. While this tool could be implemented to evaluate a single candidate against multiple unit tests simultaneously, our experiments in Section 4.2.1 have shown that the current approach, evaluating multiple candidates against one unit test at a time, yields better results.

## 3.2 Preprocessing

To accelerate the IR agent's tools, *retrieve_entity* and *retrieve_context*, and improve their efficiency, we preprocess database values and catalogs before running the system. For database values, we perform a syntactic search by creating a Locality-Sensitive Hashing (LSH) index, as explained in *retrieve_entity*. For database catalogs, which contain longer texts that require semantic understanding, we use a vector database retrieval method to measure semantic similarity.

**Locality Sensitive Hashing Indexing of Values.** To optimize the entity retrieval process, we employ a method capable of efficiently searching through large databases, which may contain millions of rows, to retrieve the most similar values. This process doesn't require perfect accuracy but should return a reasonably small set of similar values, such as around a hundred elements. Locality Sensitive Hashing (LSH) is an effective technique for approximate nearest-neighbor searches, allowing us to retrieve database values most similar to a given keyword. During preprocessing, we index unique database values using LSH. Then, in the *retrieve_entity* tool, we query this index to quickly find the top similar values for a keyword.

**Vector Database for Descriptions.** Extracting the most semantically relevant pieces of information from database catalogs is crucial for generating accurate SQL queries. These catalogs can be extensive, with hundreds of pages explaining the entities and their relationships within the database, necessitating an efficient retrieval method. To enable a high-efficiency semantic similarity searches, we preprocess the database catalogs into a vector database. During the context retrieval step, we query this vector database to find the most relevant information for the question at hand.

For a more detailed description of our framework and the preprocessing phase, please refer to Appendix A.

## 3.3 Adapting Text-to-SQL Systems for Real-World Deployment

We designed our multi-agent framework to allow for flexible combinations of agents, enabling it to adapt to various deployments, database settings, and constraints. In setups where higher computational budgets and more powerful LLMs are available, CHESS can be configured to maximize accuracy by leveraging these resources. Specifically, the Candidate Generator (CG) and Unit Tester (UT) agents can be programmed to generate multiple candidate queries and unit tests, leading to more accurate final SQL queries (Brown et al., 2024). This increased sampling allows the system to refine its query generation process, resulting in better performance, as demonstrated in Section 4.2.1.

Interestingly, similar to the findings in (Maamari et al., 2024), we observed that when schema sizes are relatively small and the available LLMs are highly capable, the inclusion of the Schema Selector (SS) agent may introduce a precision-recall trade-off that can negatively impact accuracy. In these cases, omitting the SS agent can lead to better results.

However, as we will demonstrate in Section 4.2.2, for larger schemas, even powerful LLMs struggle to maintain performance. In such cases, the collaboration of all four agents is essential for achieving higher accuracy, as the SS agent becomes crucial for managing the complexity of the schema.

Finally, when computational resources are limited or when only weaker LLMs, such as small open-source models, are available to ensure data privacy, CHESS can be reconfigured for efficiency. In these scenarios, we streamline the system by excluding the UT agent and limiting the CG agent to generate only one candidate query. This approach minimizes computational overhead while still synthesizing a reliable SQL query, making CHESS suitable for resource-constrained environments.

# 4 EXPERIMENTS AND RESULTS

## 4.1 Datasets and Metrics

### 4.1.1 Datasets

**Spider.** The Spider dataset (Yu et al., 2018) includes 200 database schemas, with 160 schemas available for training and development, and 40 schemas reserved for testing. Notably, the databases used in the training, development, and test sets are distinct and do not overlap.

**Bird.** The recently introduced BIRD dataset (Li et al., 2024b) features 12,751 unique question-SQL pairs, covering 95 large databases with a combined size of 33.4 GB. This dataset spans 37 professional fields, including sectors such as blockchain, hockey, healthcare, and education. BIRD enhances SQL query generation by incorporating external knowledge and providing a detailed database catalog that includes column and database descriptions, thereby clarifying potential ambiguities. The SQL queries in BIRD are generally more complex than those found in the Spider dataset.

**Subsampled Development Set (SDS).** To facilitate ablation studies, reduce costs, and maintain the distribution of the BIRD development set, we subsampled 10% of each database in the development set, resulting in the Subsampled Development Set which we call SDS. This SDS consists of 147 samples: 81 simple, 54 moderate, and 12 challenging questions. For reproducibility, we included the SDS in our GitHub repository.

**Synthetic Industrial-scale Database Schema.** Production-level text-to-SQL systems often operate on real-world databases with extremely large schemas containing thousands of columns, far exceeding the scope of academic benchmarks like BIRD (Li et al., 2024b). As highlighted in the Lost-in-the-Middle paper (Liu et al., 2024), large language models (LLMs) can struggle with overly long prompts, leading to confusion. However, recent advancements, such as Gemini (Team et al., 2023) and GPT-4o (OpenAI, 2024b), have demonstrated the capability to generate accurate SQL queries without the need to filter prompts or exclude any part of the database schema (Maamari et al., 2024). To assess the scalability of these models on massive databases, we synthesized extremely large datasets to mimic real-world conditions. This involved merging multiple databases from the BIRD training and development sets, resulting in schemas with up to 4,337 columns, the largest in our experiments. This method enabled us to replicate the complexity and challenges of large, production-level databases.

**Metrics.** In this paper, we primarily used Execution accuracy as the metric for comparison following the previous works (Pourreza et al., 2024; Maamari et al., 2024; Pourreza

and Rafiei, 2024a). Execution Accuracy (EX) measures the correctness of the SQL output by comparing the results of the predicted query with those of the reference query when executed on specific database instances. This metric offers a nuanced evaluation by accounting for variations in valid SQL queries that can yield the same results for a given question. Additionally, given that the number of LLM calls and the tokens used in prompts can be substantial in some scenarios, we also report both the total number of tokens and the number of LLM calls required by our approach.

## 4.2 Different Configurations of CHESS per Deployment Constraints

In this section, we discuss various deployment constraints, detailing how CHESS is configured for each scenario, along with the corresponding experiments and results. First, in section 1, we consider a scenario with high computational budget that allows multiple LLM calls during inference using strong LLMs. For this case, we evaluate the performance of CHESS using the IR, CG, and UT agents, referred to as CHESS$_{(IR,CG,UT)}$, on the BIRD dataset. Next, in section 4.2.2, we address the challenge of industrial-scale database schemas beyond those in the BIRD dataset, where the reasoning capabilities of powerful proprietary models, such as Gemini-1.5-pro, may be insufficient. To tackle this issue, we introduce the SS agent to the team, demonstrating its effectiveness in handling very large schemas. Then, to accommodate scenarios with limited computational budgets, we exclude the UT agent from the team and restrict the number of candidates generated by CG to one. This configuration, labeled as CHESS$_{(IR,SS,CG)}$, is evaluated using both less powerful proprietary models and open-source models. Finally, in section 4.2.3, we perform ablation studies to assess the contribution of each component to the end-to-end performance of our system.

### 4.2.1 CHESS$_{(IR,CG,UT)}$ for High Computational Budget

Recent studies, such as MCS-SQL (Lee et al., 2024) and CHASE-SQL (Pourreza et al., 2024), have achieved state-of-the-art performance on text-to-SQL benchmarks by scaling LLM calls, generating a large set of candidate SQL queries, and selecting the best candidate among them. Following a similar strategy, we evaluate the performance of our text-to-SQL multi-agent framework by generating 20 candidate SQL queries using the Gemini-1.5-pro model for each sample question. We then create ten natural language unit tests to effectively differentiate correct candidates from incorrect ones. Given that schema linking is unnecessary for small databases within the BIRD benchmark (Maamari et al., 2024), we excluded the SS agent and retained only the IR, CG, and UT agents for this experiment. The performance comparison between our framework with Gemini-1.5-pro model and prior methods is presented in Table 1. As

shown in the table, our approach outperforms all previous methods except CHASE-SQL, which utilizes a fine-tuned model specifically designed to select the best candidate SQL query for the BIRD benchmark. In contrast, our framework, without any fine-tuning, surpasses all other methods, demonstrating its robustness across benchmarks without the need for task-specific optimization. Furthermore, we evaluated the performance of the UT agent in two distinct settings. In the first setting, a single LLM call receives one unit test and assigns scores to all candidate queries. In the second setting, a single LLM call receives all unit tests and assigns a score to one candidate query. The results for the first scenario are presented in Table 1, where we achieved a score of 68.31. In contrast, the second scenario yielded a lower score of 66.78, highlighting the significance of allowing the LLM to compare candidate queries when evaluating unit tests.

*Table 1.* Performance of $CHESS_{(IR,CG,UT)}$ on BIRD dataset, comparing to the methods with high computation budget.

| Method | dev EX | test EX |
|---|---|---|
| CHASE-SQL | 73.01 | 73.00 |
| **CHESS**$_{(IR,CG,UT)}$ | **68.31** | **71.10** |
| Distillery + GPT-4o | 67.21 | 71.83 |
| MCS-SQL + GPT-4 | 63.36 | 65.45 |
| SFT CodeS-15B | 58.47 | 60.37 |
| DTS-SQL + DeepSeek 7B | 55.8 | 60.31 |
| MAC-SQL + GPT-4 | 59.59 | 59.59 |

### 4.2.2 Introducing the Schema Selection Agent for Large Database Schema

In this section, we use our synthetically generated dataset, the Synthetic Industrial-scale Database Schema, to evaluate the performance of our multi-agent framework on extremely large schemas and analyze the impact of large prompts on the performance of recent LLMs. We assessed the performance of Gemini model using Pass@1 and Pass@5 metrics on BIRD development set using the synthetic dataset as target database. For database sizes smaller than the maximum size (4337), we randomly select columns while ensuring that the ground truth schema is retained. Pass@1 measures the success rate when the model is allowed to generate a single SQL query with execution accuracy of one, while Pass@5 measures the rate of generating at least one correct query from five attempts. We tested the Gemini-1.5-pro model and presented the results in Figure 3, which compares the Pass@1 and Pass@5 performances across varying schema sizes. The results reveal an 11% performance gap between scenarios where the exact ground truth schema (i.e., the correct columns and tables) is provided and those where a large number of columns are included in the prompt. As the number of columns decreases, the model's performance improves, underscoring the importance of schema linking to filter out irrelevant columns and tables.

**Effectiveness of the Schema Selection Agent.** To measure the effectiveness of our schema selection agent (SS agent) in dealing with very large schema, we repeated the previous experiment this time including the SS agent. After applying the SS agent to the largest schema with 4,337 columns, the Pass@1 and Pass@5 scores improved to 61% and 63%, respectively, representing a 2% increase compared to the performance without schema linking. This demonstrates the critical role of schema linking in enhancing query generation for very large databases.

Note that the Pass@1 of the system including SS agent has a slightly better performance of the Pass@5 of the system without any schema linking, which shows a higher performance while being more than $\times 5$ efficient in terms of used tokens. This suggests that for very large schema the schema selection agent can increase the accuracy as well as the efficiency of the system.
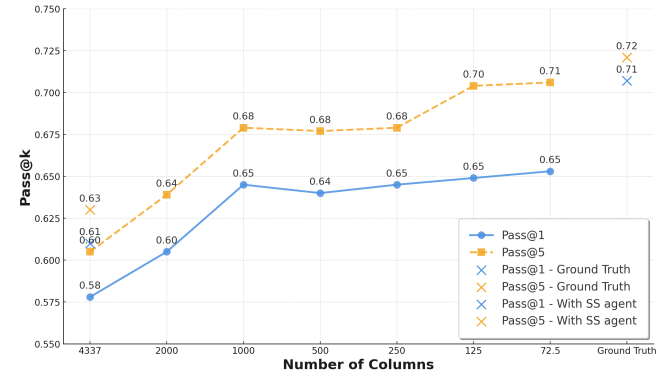


*Figure 3.* Text-to-SQL Pass@1 and Pass@5 performance of Gemini-1.5-pro model vs. the number of columns in the database.

### 4.2.3 $CHESS_{(IR,SS,CG)}$ for Limited Computational Budget

In real-world deployment scenarios, various constraints such as low computational power, limited access to powerful models due to budget restrictions, or data privacy concerns can make the use of high-compute methods with powerful proprietary models (e.g., (Lee et al., 2024; Pourreza et al., 2024)) infeasible. To accommodate such cases, we propose a configuration of our multi-agent system, $CHESS_{(IR,SS,CG)}$, that achieves high performance while adhering to these constraints.

To limit the number of LLM calls, we exclude the UT agent from the framework and restrict the number of candidate generations and revisions by GC to only one and three times, respectively. Because we cannot rely on multiple LLM calls to generate candidate queries, it becomes crucial to be more efficient with the schema passed to the candidate generator. Therefore, the Schema Selector (SS) plays a vital role in enhancing efficiency by narrowing down the schema to the most relevant subset before candidate generation.

Moreover, to demonstrate the performance of the system using less powerful models, we conduct experiments using older proprietary models such as GPT-3.5/4-turbo, as well as open-source models like Llama-3-70B and a fine-tuned DeepSeek model.

**BIRD Results.** Since the test set of the BIRD benchmark is not available, we conducted our ablations and performance evaluations on the development set. We assessed our proposed method using both 1) proprietary and 2) open-source models. In the first case, we utilized a fine-tuned DeepSeek Coder model for candidate generation, GPT-3.5-turbo for column filtering, and GPT-4-turbo for the remaining LLM calls. In the second case, we used our fine-tuned DeepSeek Coder model for candidate generation, with all other LLM calls handled by Llama-3-70B.

As reported in Table 2, our approach using proprietary models achieved a high execution accuracy on both the development and test sets of BIRD. We want to highlight the efficiency of our method in this configuration where we only call GPT-4-turbo model 6 times where some other methods such as MCS-SQL performs around 100 LLM calls.

**CHESS using Open-source Models.** In most of the industrial use cases, the privacy of the data is at the highest priority, which limits the use of proprietary models as in (Pourreza et al., 2024), (Maamari et al., 2024), and (Lee et al., 2024). For such cases, companies would prefer to deploy a text-to-SQL system using only open-source models completely on-premis. Our method with open-source LLMs attained the highest performance among all open-source methods, making it very suitable for these scnearios.

**Table 2.** Performance of CHESS$_{(IR,SS,CG)}$ on the BIRD dataset, comparing to the methods with low computation budget.

| Method | test EX | dev EX |
|---|---|---|
| **CHESS**$_{(IR,SS,CG)}$ + proprietary | **66.69** | **65.00** |
| MCS-SQL + GPT-4 | 65.45 | 63.36 |
| **CHESS**$_{(IR,SS,CG)}$ + Open LLMs | – | **61.5** |
| SFT CodeS-15B | 60.37 | 58.47 |
| DTS-SQL + DeepSeek 7B | 60.31 | 55.8 |
| MAC-SQL + GPT-4 | 57.56 | 59.59 |

**Table 3.** Performance of CHESS$_{(IR,SS,CG)}$ on the Spider test set, comparing to all published methods.

| Method | EX |
|---|---|
| MCS-SQL+GPT-4 | 89.6 |
| **CHESS**$_{(IR,SS,CG)}$ | **87.2** |
| DAIL-SQL + GPT-4 | 86.6 |
| DIN-SQL + GPT-4 | 85.3 |
| C3 + ChatGPT | 82.3 |
| RESDSQL-3B | 79.9 |

**Spider Results.** To evaluate the generalizability of our proposed method beyond the BIRD benchmark, we tested it on the Spider test set without specifically fine-tuning a new

model for candidate generation or modifying the in-context learning samples. We followed our default engine setup. The only adjustment we made to our setting was the removal of the *retrieve_context* tool since the Spider test set lacks column or table descriptions, which are integral to our method. As shown in Table 3, our approach achieved an execution accuracy of 87.2% on the test set, ranking it as the second-highest performing method among those published. This underscores the robustness of our method across different databases without any modifications. Notably, the best propriety (and undisclosed) method on the Spider test set leaderboard is Miniseek with an accuracy of 91.2%.

**Table 4.** The execution accuracy (EX) of the framework by removing each tool on the (subsampled) development set.

| Tools | EX | ΔEX |
|---|---|---|
| All tools | 64.62 | – |
| w/o *retrieve_entity & _context* | 59.86 | -4.76 |
| w/o *filter_column* | 61.90 | -2.72 |
| w/o *select_tables* | 58.50 | -6.12 |
| w/o *select_columns* | 59.18 | -5.44 |
| w/o *revise* | 57.82 | -6.80 |
| with 1-time *revise* | 61.22 | -3.40 |

### 4.3 Ablation Studies

Table 4 presents the execution accuracy (EX), where different modules or components are omitted. In the configuration without entity and context retrieval, we retrieved a random example and included column descriptions for all columns. This approach highlights the significant impact of our selective retrieval, which outperforms naive context augmentation by 4.76% in execution accuracy. Additionally, we evaluated the effect of removing each tool of the SS agent, revealing that the *select_tables* tool is the most critical, contributing a 6.12% increase in performance. The table also illustrates the significant influence of the revision tool, with a 6.80% improvement. Increasing the number of revision samples for self-consistency led to higher performance gains, aligning with findings from (Lee et al., 2024).

**Unit Tests Improve Performance.** To assess the effectiveness of the UT agent in generating high-quality unit tests for distinguishing between correct and incorrect SQL candidates, we evaluated the performance of our agent-based framework by varying the number of unit tests and select the candidate with the highest score from the 20 candidates. Figure 4 illustrates the performance across different numbers of unit tests. As shown in the figure, performance improves as the number of unit tests increases, peaking at 10 tests. Beyond this point, the performance plateaus, indicating that additional tests provide diminishing returns.

**Evaluation of the Schema Selection Agent** In addition to conducting ablation studies to assess the impacts of context retrieval, schema selection, and revision methods, we
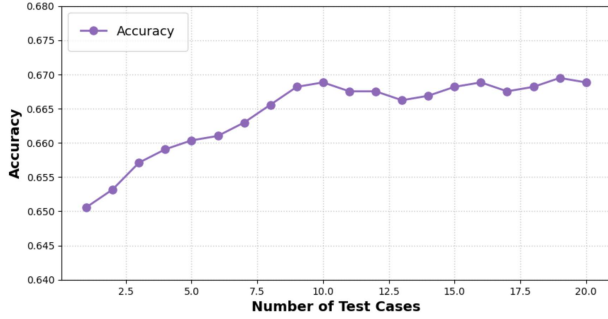
*Figure 4.* Text-to-SQL performance of CHESS with respect to different number of unit tests on BIRD development set.



*Figure 5.* The average number of columns remaining after each stage of the schema selection process across 11 BIRD dev databases. The initial count shows the total columns in each database, and the schema selection agent consistently reduces this number to around 10 columns for candidate generation, regardless of database complexity.

evaluate their specific influence on the precision and recall of tables and columns identified by the model for generating final SQL queries. Precision and recall are calculated by comparing the tables and columns used in correct SQL queries as ground truth. As shown in Table 5, the schema selection agent improves the precision of selected tables and columns, with only a slight reduction in recall.

This balance of high precision and recall enables the model to work with the most relevant information in a compact context window, enhancing the accuracy of SQL generation. Figure 16 illustrates this process, where the schema selection (SS) agent successfully narrows down the tables and columns, selecting two tables and five columns, including the two columns used in the correct SQL query. For a comprehensive example of the schema selection process, please refer to Appendix E.

Furthermore, Figure 5 illustrates how the schema selection agent effectively narrows down the number of columns at each step, adapting the remaining columns based on the complexity of the question rather than the database schema itself. This approach consistently reduces the columns to around 10, making our method scalable and applicable across databases of varying complexity by focusing only on the most relevant information for SQL generation.

*Table 5.* Precision and Recall of schema items for tables and columns used in correct SQL, evaluated after applying each tool by the schema selection agent: *filter_columns*, *select_tables*, and *select_columns*.

|  | Table | | Column | |
| --- | --- | --- | --- | --- |
|  | Recall | Precision | Recall | Precision |
| No Filtering or Selection | 1.0 | 0.33 | 1.0 | 0.11 |
| *filter_columns* | 1.0 | 0.33 | 0.98 | 0.21 |
| *select_tables* | 0.97 | 0.89 | 0.96 | 0.45 |
| *select_columns* | 0.96 | 0.90 | 0.94 | 0.71 |

**Performance Evaluation Across Queries with Varying Complexity** The BIRD benchmark categorizes questions and SQL pairs based on the number and type of SQL keywords used into three classes: easy, moderate, and chal-
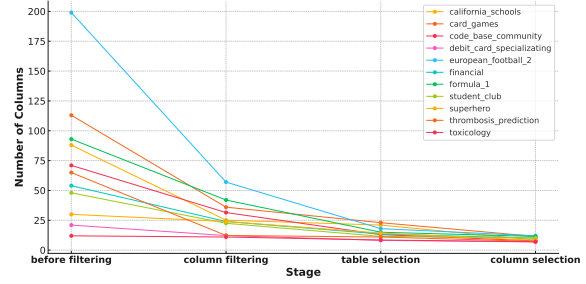
lenging. In this section, we evaluate the performance of our method 1) without computation constraints by using the IR, CG and UT agents, and 2) with limited computational budget by using the IR, SS, and CG agents. We compare the results to the original GPT-4 baseline as in the BIRD paper, where the question, evidence, and the complete schema with all tables and columns are presented to GPT-4 along with chain-of-thought reasoning prompts.

The analysis is conducted on the SDS dataset, and the results are detailed in Table 6. Our proposed method significantly improves performance across all classes in both settings. This demonstrates the effectiveness of our multi-agent framework in enhancing performance across varying difficulty levels and constraints.

*Table 6.* Comparing the performance of CHESS in two different settings with a naive GPT-4 baseline where we pass all the information in its context, across different question difficulty levels.

|  | Easy | Moderate | Challenging | Overall |
| --- | --- | --- | --- | --- |
| **CHESS**$_{(IR,CG,UT)}$ | 73.48 | 62.99 | 59.31 | 68.31 |
| **CHESS**$_{(IR,SS,CG)}$ | 65.43 | 64.81 | 58.33 | 64.62 |
| GPT-4-turbo (baseline) | 54.32 | 35.18 | 41.66 | 46.25 |

## 5 CONCLUSION AND LIMITATIONS

In this paper, we introduced CHESS, a multi-agent framework capable of handling complex, industrial-scale databases. CHESS devises 4 specialized agents, namely, Information Retriever, Schema Selector, Candidate Generator, and Unit Tester, to tackle key challenges such as efficient data retrieval, schema pruning, SQL query generation, and validation. We provide extensive ablation studies to show the effectiveness of each of these agents and their tools. We demonstrate that for large industrial-scale schema sizes, methods solely relying on scaling sampling at inference time become less effective, whereas CHESS's schema selection

agent can deliver a high accuracy while processing less than ×5 tokens. We also show that CHESS achieves state-of-the-art performance when controlled for computational budget. CHESS achieves 71.10% accuracy on the challenging BIRD test set, within 2% of the leading method (a proprietary approach), while requiring approximately 83% fewer LLM calls.

As discussed in the experiments section, there is a significant need for an extensive text-to-SQL benchmark that accurately reflects the challenges posed by large-scale databases. In our paper, we attempted to achieve this by combining different databases; however, having access to even larger real-world databases can meaningfully advance research in this domain. Additionally, previous studies (Pourreza et al., 2024) have shown that finetuning can improve LLM performance in query selection. As future work, we plan to finetune models specifically for test-case generation and evaluation to further improve the performance.

## REFERENCES

Meta AI. Exploring meta llama-3. `https://ai.meta.com/blog/meta-llama-3/`. Accessed: 2024-04-18.

Bradley Brown, Jordan Juravsky, Ryan Ehrlich, Ronald Clark, Quoc V Le, Christopher Ré, and Azalia Mirhoseini. Large language monkeys: Scaling inference compute with repeated sampling. *arXiv preprint arXiv:2407.21787*, 2024.

Ruichu Cai, Jinjie Yuan, Boyan Xu, and Zhifeng Hao. Sadga: Structure-aware dual graph aggregation network for text-to-sql. *Advances in Neural Information Processing Systems*, 34:7664–7676, 2021.

Ruisheng Cao, Lu Chen, Zhi Chen, Yanbin Zhao, Su Zhu, and Kai Yu. Lgesql: line graph enhanced text-to-sql model with mixed local and non-local relations. *arXiv preprint arXiv:2106.01093*, 2021.

Tim Dettmers, Artidoro Pagnoni, Ari Holtzman, and Luke Zettlemoyer. Qlora: Efficient finetuning of quantized llms. *Advances in Neural Information Processing Systems*, 36, 2024.

Xuemei Dong, Chao Zhang, Yuhang Ge, Yuren Mao, Yunjun Gao, Jinshu Lin, Dongfang Lou, et al. C3: Zero-shot text-to-sql with chatgpt. *arXiv preprint arXiv:2307.07306*, 2023.

Yujian Gan, Xinyun Chen, Qiuping Huang, Matthew Purver, John R Woodward, Jinxia Xie, and Pengsheng Huang. Towards robustness of text-to-sql models against synonym substitution. *arXiv preprint arXiv:2106.01065*, 2021.

Dawei Gao, Haibin Wang, Yaliang Li, Xiuyu Sun, Yichen Qian, Bolin Ding, and Jingren Zhou. Text-to-sql empowered by large language models: A benchmark evaluation. *arXiv preprint arXiv:2308.15363*, 2023.

Nuno M Guerreiro, Elena Voita, and André FT Martins. Looking for a needle in a haystack: A comprehensive study of hallucinations in neural machine translation. *arXiv preprint arXiv:2208.05309*, 2022.

Daya Guo, Qihao Zhu, Dejian Yang, Zhenda Xie, Kai Dong, Wentao Zhang, Guanting Chen, Xiao Bi, Y Wu, YK Li, et al. Deepseek-coder: When the large language model meets programming–the rise of code intelligence. *arXiv preprint arXiv:2401.14196*, 2024.

Jiaqi Guo, Zecheng Zhan, Yan Gao, Yan Xiao, Jian-Guang Lou, Ting Liu, and Dongmei Zhang. Towards complex text-to-sql in cross-domain database with intermediate representation. *arXiv preprint arXiv:1905.08205*, 2019.

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

Dongjun Lee, Choongwon Park, Jaehyuk Kim, and Heesoo Park. Mcs-sql: Leveraging multiple prompts and multiple-choice selection for text-to-sql generation. *arXiv preprint arXiv:2405.07467*, 2024.

Haoyang Li, Jing Zhang, Hanbing Liu, Ju Fan, Xiaokang Zhang, Jun Zhu, Renjie Wei, Hongyan Pan, Cuiping Li, and Hong Chen. Codes: Towards building open-source language models for text-to-sql. *arXiv preprint arXiv:2402.16347*, 2024a.

Jinyang Li, Binyuan Hui, Ge Qu, Jiaxi Yang, Binhua Li, Bowen Li, Bailin Wang, Bowen Qin, Ruiying Geng, Nan Huo, et al. Can llm already serve as a database interface? a big bench for large-scale database grounded text-to-sqls. *Advances in Neural Information Processing Systems*, 36, 2024b.

Nelson F Liu, Kevin Lin, John Hewitt, Ashwin Paranjape, Michele Bevilacqua, Fabio Petroni, and Percy Liang. Lost in the middle: How language models use long contexts. *Transactions of the Association for Computational Linguistics*, 12:157–173, 2024.

Karime Maamari, Fadhil Abubaker, Daniel Jaroslawicz, and Amine Mhedhbi. The death of schema linking? text-to-sql in the age of well-reasoned language models. *arXiv preprint arXiv:2408.07702*, 2024.

OpenAI. Embeddings, 2024a. Retrieved May 15, 2024, from https://platform.openai.com/docs/guides/embeddings.

OpenAI. Hello gpt-4o. https://openai.com/index/hello-gpt-4o/, 2024b. Accessed: October 15, 2024.

Mohammadreza Pourreza and Davood Rafiei. Evaluating cross-domain text-to-sql models and benchmarks. *arXiv preprint arXiv:2310.18538*, 2023.

Mohammadreza Pourreza and Davood Rafiei. Din-sql: Decomposed in-context learning of text-to-sql with self-correction. *Advances in Neural Information Processing Systems*, 36, 2024a.

Mohammadreza Pourreza and Davood Rafiei. Dts-sql: Decomposed text-to-sql with small large language models. *arXiv preprint arXiv:2402.01117*, 2024b.

Mohammadreza Pourreza, Hailong Li, Ruoxi Sun, Yeounoh Chung, Shayan Talaei, Gaurav Tarlok Kakkar, Yu Gan,

Amin Saberi, Fatma Ozcan, and Sercan O Arik. Chase-sql: Multi-path reasoning and preference optimized candidate selection in text-to-sql. *arXiv preprint arXiv:2410.01943*, 2024.

Jiexing Qi, Jingyao Tang, Ziwei He, Xiangpeng Wan, Yu Cheng, Chenghu Zhou, Xinbing Wang, Quanshi Zhang, and Zhouhan Lin. Rasat: Integrating relational structures into pretrained seq2seq model for text-to-sql. *arXiv preprint arXiv:2205.06983*, 2022.

Bowen Qin, Binyuan Hui, Lihan Wang, Min Yang, Jinyang Li, Binhua Li, Ruiying Geng, Rongyu Cao, Jian Sun, Luo Si, et al. A survey on text-to-sql parsing: Concepts, methods, and future directions. *arXiv preprint arXiv:2208.13629*, 2022.

Nitarshan Rajkumar, Raymond Li, and Dzmitry Bahdanau. Evaluating the text-to-sql capabilities of large language models. *arXiv preprint arXiv:2204.00498*, 2022.

Jon Saad-Falcon, Adrian Gamarra Lafuente, Shlok Natarajan, Nahum Maru, Hristo Todorov, Etash Guha, E. Kelly Buchanan, Mayee Chen, Neel Guha, Christopher Ré, and Azalia Mirhoseini. Archon: An architecture search framework for inference-time techniques, 2024. URL https://arxiv.org/abs/2409.15254.

Ilya Sutskever, Oriol Vinyals, and Quoc V Le. Sequence to sequence learning with neural networks. *Advances in neural information processing systems*, 27, 2014.

Gemini Team, Rohan Anil, Sebastian Borgeaud, Yonghui Wu, Jean-Baptiste Alayrac, Jiahui Yu, Radu Soricut, Johan Schalkwyk, Andrew M Dai, Anja Hauth, et al. Gemini: a family of highly capable multimodal models. *arXiv preprint arXiv:2312.11805*, 2023.

Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Łukasz Kaiser, and Illia Polosukhin. Attention is all you need. *Advances in neural information processing systems*, 30, 2017.

Bailin Wang, Richard Shin, Xiaodong Liu, Oleksandr Polozov, and Matthew Richardson. Rat-sql: Relation-aware schema encoding and linking for text-to-sql parsers. *arXiv preprint arXiv:1911.04942*, 2019.

Bing Wang, Changyu Ren, Jian Yang, Xinnian Liang, Jiaqi Bai, Qian-Wen Zhang, Zhao Yan, and Zhoujun Li. Mac-sql: Multi-agent collaboration for text-to-sql. *arXiv preprint arXiv:2312.11242*, 2023.

Tao Yu, Rui Zhang, Kai Yang, Michihiro Yasunaga, Dongxu Wang, Zifan Li, James Ma, Irene Li, Qingning Yao, Shanelle 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.

John M Zelle and Raymond J Mooney. Learning to parse database queries using inductive logic programming. In *Proceedings of the national conference on artificial intelligence*, pages 1050–1055, 1996.

# A    IMPLEMENTATION DETAILS

## A.1    Locality Sensitive Hashing Indexing of Database Values

Our goal in the *retrieve_entity* tool is to retrieve database values that most closely match a set of keywords derived from the question. It is important to recognize that keywords from the question may not exactly correspond to the database values due to potential typos, variations in expression, or the common scenario where users are unaware of the precise format used to store data in the database. This reality demands a retrieval strategy that is both robust and adaptable, capable of accommodating such discrepancies. Relying solely on exact match retrieval, as suggested in prior studies (Li et al., 2024a), may not be sufficiently effective.

To address this, we employ string similarity measures, such as edit distance and semantic embedding, to retrieve the values most similar to the keywords. However, computing the edit distance and embedding similarity for every keyword against all values in the database is computationally expensive and time-consuming. To balance efficiency and accuracy, we utilize a hierarchical retrieval method.

Locality Sensitive Hashing (LSH) is an efficient technique for approximate nearest neighbor searches, which allows us to retrieve the most similar values to a keyword in the database. In the pre-processing stage, we index unique values in the database using LSH. Then, in the *retrieve_entity* tool of the IR agent, we query this index to rapidly find the top similar values to a keyword. Our approach involves using LSH queries to retrieve the top 10 similar values, after which we compute the edit distance and semantic similarity between the keyword and these values to further refine the results.

To simultaneously utilize edit distance and embedding similarity, we first identify the top 10 values closest to each keyword based on cosine similarity between their embedding vectors (obtained using OpenAI text-embedding-3-small (OpenAI, 2024a)) and the keyword's embedding vector. We then filter out values that fall below a specific threshold. Finally, for each keyword and column, we retain only the value that has the smallest edit distance.

We observed a significant reduction in time complexity, from 5 minutes to 5 seconds, using this method compared to a naive approach of computing the edit distance for all unique values in the database on the fly. While computing edit distance is proportional to the size of the database, significantly increasing the time complexity for processing a single question, using LSH allows us to index values in the pre-processing step and, during entity retrieval, rapidly query the index to find the most similar values to a keyword in a much more time-efficient manner.

## A.2    Vector database

Each database schema in the BIRD benchmark (Li et al., 2024b) includes detailed descriptions for columns, specifying the contents of each column and the values for categorical columns. Providing these descriptions to the model is essential for guiding the SQL query generation process. However, incorporating all descriptions in the prompt for the weaker open-source models can overwhelm the model, potentially leading to the generation of incorrect SQL queries, as observed in section 4.3. It is important to note that the database catalog in the BIRD benchmark provides a relatively limited view of database metadata. In contrast, real-world production-level databases often contain more diverse information, including value ranges, constraints, and usage instructions for each table. Our proposed method can effectively utilize this extensive metadata to enhance performance.

To evaluate the relevance of descriptions to a given question, we employ embedding similarity (OpenAI, 2024a), which quantifies the semantic similarity between the question and each description. To enhance the efficiency of the retrieval process, we pre-process the descriptions and created the embedding vectors for each of them and stored in a vector database, utilizing ChromaDB in our implementation. For the *retrieve_context* tool of the IR agent, we query this vector database to identify descriptions that are most semantically aligned with the question. This targeted approach ensures that only the most pertinent information is provided to the model, thereby improving the accuracy of the generated SQL queries.

## A.3    Local Column filtering

In the *filter_column* tool, decisions to retain a column for subsequent steps are made independently, without considering relative information. The data provided to the LLM for column filtering includes: 1) the table name, 2) the column name, 3) the data type, 4) descriptions, if retrieved by the IR agent, and 5) database values, if retrieved by *retrieve_entity* tool. To enhance the model performance and make the task definition clear to the model, we used few-shot samples for this tool.

Some key columns for SQL generation, which we call linking columns, such as those with foreign and primary key constraints, are crucial for writing SQL queries. For instance, questions about counting entities often requires primary keys, and joining tables necessitates foreign key columns. However, in the *filter_column*, *select_table select_column* tools, some of these essential columns may be initially rejected because they do not semantically relate to the given question. Despite this, these columns are indispensable for SQL generation. Therefore, in all of our sub-modules, we consistently retain foreign key and primary key columns, irrespective of the outputs from column selection and filtering processes.

## A.4 Revise Tool

Revising generated candidate SQL queries is a critical aspect for the Candidate Generator Agent. In addition to the database schema, the question, and the candidate SQL query, we also provide the execution result of the SQL query. This gives the LLM an opportunity to view the retrieved data and revise the SQL query accordingly. This process mirrors human behavior when writing complex SQL queries; typically, we start with a draft query and refine it based on the results of its execution. Furthermore, this method allows the LLM to make necessary adjustments to the SQL query in instances of execution syntax errors.

In this step, we also incorporate instructions derived from our error analysis F to guide the model towards generating correct SQL queries. For instance, as shown in 6, we guided the model to ensure that all requested columns are included in the SQL query. In this specific example, the *revise* tool identified a missing column and successfully added it to the query.
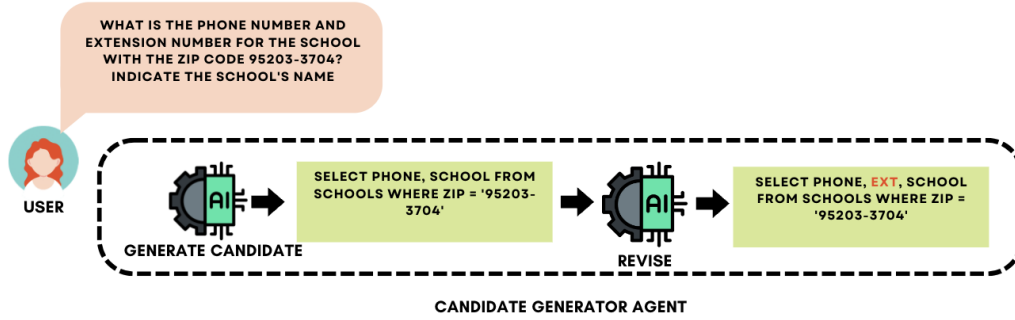


*Figure 6.* An example of the revise tool to fix missing columns in a candidate query.

# B  FINETUNING THE MODEL FOR GENERATE_CANDIDATE_QUERY

## B.1  Fine-tuning Dataset and Model

To enhance the generation of better candidate SQL queries for the Candidate Generator Agent, we fine-tuned the DeepSeek Coder 34B (Guo et al., 2024) on the training set of the BIRD benchmark, which comprises approximately 9,500 samples. In constructing the fine-tuning dataset, rather than solely using correct tables and columns like the work proposed in (Pourreza and Rafiei, 2024b), we developed a heuristic to address the error propagation issue. Recognizing that previous steps in the pipeline may not always pinpoint the most efficient schema with perfect accuracy, we intentionally introduced some noise into our dataset creation to train the model. Specifically, we included columns and tables that were incorrect but shared similar naming conventions and semantic attributes with the correct schema. We also utilized our keyword selection module to extract keywords from the questions, search for these keywords in the database, and incorporate them into the prompt.

## B.2  Hyperparameters

We fine-tuned the deepseek model for the *generate_candidate_query* tool using 4-bit quantization of the base model and LORA adapters (Hu et al., 2021), a technique formally referred to as QLORA (Dettmers et al., 2024). We configured the LORA rank parameter to 128 and set the LORA alpha parameter to 256. The fine-tuning process was conducted over two epochs on the constructed dataset, utilizing a batch size of 32 and a learning rate of 1e-4, along with a cosine scheduler, all on a single H100 GPU for 4 hours.

## C  PROMPT TEMPLATES

In this section we provide the exact prompts that have been used for the tools in our multi-agent setup. For easier parsing of the LLM's output, we instruct the model to generate its output in a structured format.

```
Objective:  Analyze the given question and hint to identify and extract
keywords, keyphrases, and named entities.  These elements are crucial for
understanding the core components of the inquiry and the guidance provided.
This process involves recognizing and isolating significant terms and phrases
that could be instrumental in formulating searches or queries related to the
posed question.

Instructions:
1.  Read the Question Carefully:  Understand the primary focus and specific
details of the question.  Look for any named entities (such as organizations,
locations, etc.), technical terms, and other phrases that encapsulate important
aspects of the inquiry.
2.  Analyze the Hint:  The hint is designed to direct attention toward certain
elements relevant to answering the question.  Extract any keywords, phrases, or
named entities that could provide further clarity or direction in formulating
an answer.
3.  List Keyphrases and Entities:  Combine your findings from both the question
and the hint into a single Python list.  This list should contain:
- Keywords:  Single words that capture essential aspects of the question or
hint.
- Keyphrases:  Short phrases or named entities that represent specific concepts,
locations, organizations, or other significant details.
Ensure to maintain the original phrasing or terminology used in the question
and hint.

{FEWSHOT_EXAMPLES}

Task:
Given the following question and hint, identify and list all relevant keywords,
keyphrases, and named entities.

Question:  {QUESTION}

Hint:  {HINT}

Please provide your findings as a Python list, capturing the essence of both
the question and hint through the identified terms and phrases.
Only output the Python list, no explanations needed.
```

*Figure 7.* Template for the *extract_keyword* tool

```
You are an expert and very smart data analyst.
Your task is to examine the provided database schema, understand the posed
question, and use the hint to pinpoint the specific columns within tables that
are essential for crafting a SQL query to answer the question.

Database Schema Overview:
{DATABASE_SCHEMA}

This schema offers an in-depth description of the database's architecture,
detailing tables, columns, primary keys, foreign keys, and any pertinent
information regarding relationships or constraints.  Special attention should
be given to the examples listed beside each column, as they directly hint at
which columns are relevant to our query.
For key phrases mentioned in the question, we have provided the most similar
values within the columns denoted by "-- examples" in front of the corresponding
column names.  This is a critical hint to identify the columns that will be used
in the SQL query.

Question:
{QUESTION}

Hint:
{HINT}

The hint aims to direct your focus towards the specific elements of the database
schema that are crucial for answering the question effectively.

Task:
Based on the database schema, question, and hint provided, your task is to
identify all and only the columns that are essential for crafting a SQL query
to answer the question.
For each of the selected columns, explain why exactly it is necessary
for answering the question.  Your reasoning should be concise and clear,
demonstrating a logical connection between the columns and the question asked.

Tip:  If you are choosing a column for filtering a value within that column,
make sure that column has the value as an example.

Please respond with a JSON object structured as follows:

{
  "chain_of_thought_reasoning":  "Your reasoning for selecting the columns, be
concise and clear.",
  "table_name1":  ["column1", "column2", ...],
  "table_name2":  ["column1", "column2", ...],
  ...
}

Make sure your response includes the table names as keys, each associated with
a list of column names that are necessary for writing a SQL query to answer the
question.
For each aspect of the question, provide a clear and concise explanation of your
reasoning behind selecting the columns.
Only output a json as your response.
```

*Figure 9.* Template for the *select_columns* tool.

```
You are an expert and very smart data analyst.
Your task is to analyze the provided database schema, comprehend the posed
question, and leverage the hint to identify which tables are needed to generate
a SQL query for answering the question.

Database Schema Overview:
{DATABASE_SCHEMA}

This schema provides a detailed definition of the database's structure,
including tables, their columns, primary keys, foreign keys, and any relevant
details about relationships or constraints.
For key phrases mentioned in the question, we have provided the most
similar values within the columns denoted by "-- examples" in front of the
corresponding column names.  This is a critical hint to identify the tables
that will be used in the SQL query.

Question:
{QUESTION}

Hint:
{HINT}

The hint aims to direct your focus towards the specific elements of the
database schema that are crucial for answering the question effectively.

Task:
Based on the database schema, question, and hint provided, your task is to
determine the tables that should be used in the SQL query formulation.
For each of the selected tables, explain why exactly it is necessary for
answering the question.  Your explanation should be logical and concise,
demonstrating a clear understanding of the database schema, the question, and
the hint.

Please respond with a JSON object structured as follows:

{
  "chain_of_thought_reasoning":  "Explanation of the logical analysis that led
to the selection of the tables.",
  "table_names":  ["Table1", "Table2", "Table3", ...]
}

Note that you should choose all and only the tables that are necessary to write
a SQL query that answers the question effectively.
Only output a json as your response.
```

*Figure 8.* Template for the *select_tables* tool.

```
You are a detail-oriented data scientist tasked with evaluating the relevance
of database column information for answering specific SQL query question based
on provided hint.
Your goal is to assess whether the given column details are pertinent to
constructing an SQL query to address the question informed by the hint.  Label
the column information as "relevant" if it aids in query formulation, or
"irrelevant" if it does not.

Procedure:
1.  Carefully examine the provided column details.
2.  Understand the question about the database and its associated hint.
3.  Decide if the column details are necessary for the SQL query based on your
analysis.

Here is an example of how to determine if the column information is relevant
or irrelevant to the question and the hint:

{FEWSHOT_EXAMPLES}

Now, it's your turn to determine whether the provided column information can
help formulate a SQL query to answer the given question, based on the provided
hint.

The following guidelines are VERY IMPORTANT to follow.  Make sure to check each
of them carefully before making your decision:
1.  You're given only one column's information, which alone isn't enough to
answer the full query.  Concentrate solely on this provided data and assess its
relevance to the question and hint without considering any missing information.
2.  Read the column information carefully and understand the description of
it, then see if the question or the hint is asking or referring to the same
information.  If yes then the column information is relevant, otherwise it is
irrelevant.
...

Column information:
{COLUMN_PROFILE}

Question:
{QUESTION}

HINT:
{HINT}

Take a deep breath and provide your answer in the following json format:

{   "chain_of_thought_reasoning":  "One line explanation of why or why not the
column information is relevant to the question and the hint.",
   "is_column_information_relevant":  "Yes" or "No"
}

Only output a json as your response.
```

*Figure 10.* Template for the *filter_column*.

```
You are a data science expert.
Below, you are presented with a database schema and a question.
Your task is to read the schema, understand the question, and generate a valid
SQLite query to answer the question.
Before generating the final SQL query think step by step on how to write the
query.

Database Schema:
{DATABASE_SCHEMA}

This schema offers an in-depth description of the database's architecture,
detailing tables, columns, primary keys, foreign keys, and any pertinent
information regarding relationships or constraints.  Special attention should
be given to the examples listed beside each column, as they directly hint at
which columns are relevant to our query.

Database admin instructions:
- Make sure you only output the information that is asked in the question.  If
the question asks for a specific column, make sure to only include that column
in the SELECT clause, nothing more.
- Predicted query should return all of the information asked in the question
without any missing or extra information.

Question:
{QUESTION}

Hint:
{HINT}

Please respond with a JSON object structured as follows:

{
  "chain_of_thought_reasoning":  "Your thought process on how you arrived at the
final SQL query.",
  "SQL":  "Your SQL query in a single string."
}

Priority should be given to columns that have been explicitly matched with
examples relevant to the question's context.

Take a deep breath and think step by step to find the correct SQLite SQL
query.
```

*Figure 11.* Template for the *generate_candidate_query* tool

```
Objective:  Your objective is to make sure a query follows the database admin
instructions and use the correct conditions.

Database Schema:
{DATABASE_SCHEMA}

Database admin instructions:
- Make sure you only output the information that is asked in the question.  If
the question asks for a specific column, make sure to only include that column
in the SELECT clause, nothing more.
- Predicted query should return all of the information asked in the question
without any missing or extra information.

{MISSING_ENTITIES}

Question:
{QUESTION}

Hint:
{EVIDENCE}

Predicted query:
{SQL}

Query result:
{QUERY_RESULT}

Please respond with a JSON object structured as follows (if the sql query is
correct, return the query as it is):

{
  "chain_of_thought_reasoning":  "Your thought process on how you arrived at the
solution.  You don't need to explain the instructions that are satisfied.",
  "revised_SQL":  "Your revised SQL query."
}

Take a deep breath and think step by step to find the correct SQLite SQL query.
```

*Figure 12.* Template for *revise* tool

```
** Instructions:  **

Given the following question database schema, and candidate responses, generate
a set of UNIT_TEST_CAP unit tests that would evaluate the correctness of SQL
queries that would answer the question.
Unit tests should be designed in a way that distinguish the candidate responses
from each other.

- The unit tests should cover various aspects of the question and ensure
comprehensive evaluation.
- Each unit test should be clearly stated and should include the expected
outcome.
- Each unit test should be designed in a way that it can distinguishes at lease
two candidate responses from each other.
- The unit test should be formatted like 'The answer SQL query should
mention...', 'The answer SQL query should state...', 'The answer SQL query
should use...', etc.  followed by the expected outcome.
- First think step by step how you can design the units tests to distinguish
the candidate responses using the <Thinking> tags.
- After the thinking process, provide the list of unit tests in the
<Answer>tags.

VERY IMPORTANT:
All of the unit tests should consider the logic of the SQL query do not
consider the formatting of the output or output values.

You are provided with different clusters of the candidate responses.  Each
cluster contains similar responses based on their results.
You MUST generate test cases that can distinguish between the candidate
responses in each cluster and the test case should promote the candidate
responses that you think are correct.

Example of the output format:
<Thinking> Your step-by-step reasoning here.  <Thinking>
<Answer>
['The answer SQL query should mention...', ...]
<Answer>
*** Database Schema:  **
{DATABASE_SCHEMA}

*** Candidate Clusters:  **
{CANDIDATE_QUERIES}

*** Question:  **
Question:  {QUESTION} (Hint:  {HINT})
```

*Figure 13.* Template for *generate_unit_tests* tool

```
 ** Instructions:  **
Given the following question, database schema, a candidate SQL query response,
and unit tests, evaluate whether or not the response passes each unit test.

- In your evaluation, you should consider how the responses align with the a
given unit test.
- Provide reasoning before you return your evaluation inside the
<Thinking>tags.
- At the end of your evaluation, you must finish with a list of verdicts
corresponding to each candidate responses in <Answer>tags.
- You must include a verdict with one of these formatted options:  'Passed' or
'Failed'
- Here is an example of the output format:
<Thinking>Your step by step reasoning here.
<Thinking>
<Answer>
Candidate Response #1:  Passed
Candidate Response #2:  Failed
Candidate Response #3:  Passed
....
<Answer>
- Each verdict should be on a new line and correspond to the candidate response
in the same order as they are provided.
- Here is the question, database schema, candidate responses, and the unit test
to evaluate the responses:

*** Database Schema:  **
{DATABASE_SCHEMA}

*** Candidate Clusters:  **
{CANDIDATE_QUERIES}

*** Question:  **
Question:  {QUESTION} (Hint:  {HINT})

*** Unit Test:  **
{UNIT_TEST}
```

*Figure 14.* Template for *evaluate* tool

```
CREATE TABLE satscores
(
  enroll12 INTEGER not null, -- Example Values: `(398,)`, `(62,)`, `(75,)` | Column Name Meaning: enrollment (1st-12nd grade) |
Column Description: enrollment (1st-12nd grade)

  dname TEXT null, -- Example Values: `('Alameda County Office of Education',)`, `('Alameda Unified',)`, `('Albany City
Unified',)` | Column Name Meaning: district name | Column Description: district segment

  NumGE1500 INTEGER null, -- Example Values: `(14,)`, `(9,)`, `(5,)` | Column Name Meaning: Number of Test Takers Whose
Total SAT Scores Are Greater or Equal to 1500 | Column Description: Number of Test Takers Whose Total SAT Scores Are
Greater or Equal to 1500 | Value Description: Number of Test Takers Whose Total SAT Scores Are Greater or Equal to 1500
Excellence Rate = NumGE1500 / NumTstTakr

  cname TEXT null, -- Example Values: `Alameda`

  rtype TEXT not null, -- Example Values: `D`, `S`

  NumTstTakr INTEGER not null, -- Example Values: `(88,)`, `(17,)`, `(71,)` | Column Name Meaning: Number of Test Takers |
Column Description: Number of Test Takers in this school | Value Description: number of test takers in each school \\
 foreign key (cds) references schools (CDSCode)

  AvgScrWrite INTEGER null, -- Example Values: `(417,)`, `(505,)`, `(395,)`
);
```

*Figure 15.* An example of how database schema, including the relevant values and descriptions, is formatted in the prompt.

# D EXTENDED EXPERIMENTS

## D.1 Query Generation with the Correct Context

To maintain high performance while ensuring privacy, we fine-tuned an open-source model for the *generate_candidate_query* tool of the Candidate Generator agent. This tuning incorporated the correct contextual information, including only relevant columns, tables, and their descriptions. As shown in 7, using precise contextual information aligned with the gold SQL improved performance to 72.4%. This result highlights the importance of retrieving efficient schema information, especially for open-source models, which struggle when handling excessively large input prompts.

*Table 7.* This table shows the maximum execution accuracy (EX) possible for our candidate SQL module generation by passing it the correct context for questions in the BIRD dataset.

| Engine | EX |
|---|---|
| CHESS | 64.62 |
| CHESS + correct context | 72.4 |

## D.2 Models Ablation

Thanks to our efficient retrieval process, which carefully controls the number of tokens passed to LLMs, we can utilize an open-source LLM with a small context window size, specifically Llama-3 with only 8K tokens (AI). This contrasts with previous works that predominantly use GPT-4 as their base model (Pourreza and Rafiei, 2024a; Lee et al., 2024; Wang et al., 2023). In Table 8, we present the results of our proposed framework using various LLMs from different families on a subsampled development set dataset. The results indicate that our fine-tuned model for candidate generation significantly enhances performance. Notably, Llama-3's performance surpasses that of GPT-3.5-turbo but does not yet reach the performance levels of GPT-4 in our analysis.

*Table 8.* This table shows the execution accuracy (EX) of different engine setups on the subsampled development set. Each engine setup is represented as a triplet (column filtering, candidate query generation, table/column selection + revision).

| Engine setups | EX |
|---|---|
| (GPT-3.5-turbo, Fine-tuned DeepSeek, GPT-4-turbo) | 64.62 |
| (GPT-3.5-turbo, GPT-4-turbo, GPT-4-turbo) | 55.78 |
| (GPT-3.5-turbo, GPT-3.5-turbo, GPT-3.5-turbo) | 49.65 |
| (Llama-3-70B, Llama-3-70B, Llama-3-70B) | 54.42 |
| (Llama-3-70B, Fine-tuned DeepSeek, Llama-3-70B) | 59.86 |

# E    SCHEMA SELECTION EXAMPLE

In this part, we use an example to showcase how we narrow down the initial schema throughout the *filter_column*, the *select_tables*, and the *select_columns*.

The example is chosen from *Formula 1* database from the BIRD benchmark, with a total of 13 tables and 96 columns. Here is the question and its evidence:

- **Question:** What's the fastest lap time ever in a race for Lewis Hamilton?

- **Evidence:** fastest lap time ever refers to min(fastestLapTime)

Figure 16 shows the number of tables and columns that are considered as the sub-selected schema after each tool. Starting from 13 tables and 96 columns, these numbers reduced to 36 columns in 13 tables after the *filter_column* tool. Subsequently, *select_tables* narrowed it down further to 2 tables and 7 columns. Finally, *select_columns* yielded the final schema with 2 tables and 5 columns, which is used for the SQL generation.
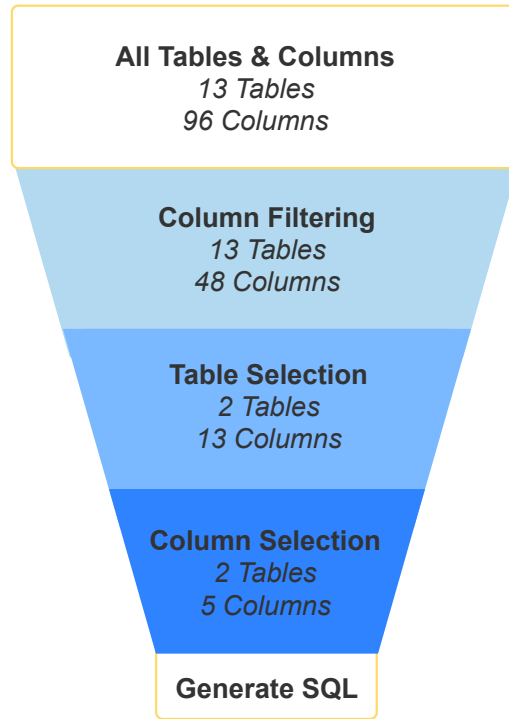


*Figure 16.* Funnel graph illustrating the progressive narrowing down of database schema with SS agent, leading to the final schema used for SQL generation.

To further illustrate the details of the schema selection process, we use the entity-relationship diagram (ERD) 17. In this figure, each table is represented as a block with its columns listed below it. Primary keys are underlined and the foreign keys are in italics connecting the corresponding columns. As shown in the legend, the columns that remained in the selected schema after each tool are colored with a gradient of white to dark blue; white represents the columns that were present in the initial schema and got filtered after *filter_column* while dark blue shows the columns that are selected after *select_columns* to be passed to the Candidate Generator agent.

There are some points worth emphasizing in the plot. First, the decision to filter linking columns (Primary and foreign keys) is a task requiring a global view of the schema and cannot be done in the local view of column filtering, hence we do not filter these columns and include all of them in the result of column filtering step, explaining why all of the primary and foreign keys are present after this step. Second, some columns such as "laps", "time", and "milliseconds" are all semantically related to the question because question asks for fastest lap time, which shows the *filter_column* tool successfully used the local information to find all relevant columns. However, all of these columns are not going to be used for crafting the SQL query, so we need to find the relevant columns with respect to their relative information, which is going

to be done in the *select_tables* and *select_columns* tools. In the *select_tables* tool, as it can be observed from the figure, the "lapTimes" table which has all information about time has been dropped by *select_tables* since there is a more relevant column, "fastestLapTime", which can be used to answer the question. This was a concrete example, which shows how local and global view of the columns and tables can help to pinpoint the correct schema.



*Figure 17.* Schema selection example formula_1_926

# F ERROR ANALYSIS

To analyze our failure cases, we subsampled 147 questions from the development set of BIRD (SDS) and processed these questions using our pipeline and a vanilla GPT-4 baseline. The vanilla GPT-4 baseline replicates the GPT-4 approach from BIRD, where the question, evidence, and the full schema with all tables and columns are provided to GPT-4 with chain-of-thought reasoning prompts. In this context, the evidence refers to the hint provided alongside some questions in the dataset.



(a) Vanilla GPT4　　　　　　　　　　　(b) CHESS

*Figure 18.* Distribution of Errors on Sampled Dev Set.

Figure 18 shows the categories of errors and their percentages for our approach and the baseline. "Incorrectly Predicted SQL" refers to failures in our pipeline that lead to incorrect final SQL, while ambiguous questions and incorrect golden SQL indicate problems with the dataset, where the data formats expected from the questions were ambiguous, and the golden SQL provided is incorrect.

For the vanilla baseline, 57.1% of the incorrectly predicted SQL results from incorrect schema linking issues, contributing 26.0% of the errors due to wrong columns in SELECT or JOIN. In comparison, our approach has 42.9% incorrect SQL, with only 5.4% of the errors attributed to incorrect column linking. These figures show that our approach distributes error types more evenly, indicating improvements in all categories of potential errors.

## F.1 Incorrectly Predicted SQL

### F.1.1 Column:

Description: The columns used in the predicted SQL are either missing or incorrect.

| Question | Which drivers born after 1975 have been ranked 2? Please give their forenames and surnames. |
|---|---|
| **Evidence** | born after 1975 refers to year(dob) >1975; |
| **Gold SQL** | SELECT T2.forename, T2.surname<br>FROM results AS T1<br>INNER JOIN drivers AS T2<br>ON T1.driverId = T2.driverId<br>WHERE STRFTIME('%Y', T2.dob) > '1975'<br>AND T1.rank = 2 |
| **Predicted SQL** | SELECT T1.forename, T1.surname<br>FROM drivers AS T1<br>INNER JOIN driverStandings AS T2<br>ON T1.driverId = T2.driverId<br>WHERE T2.position = 2<br>AND STRFTIME('%Y', T1.dob) > '1975' |
| **Error Explanation** | "rank" is the correct column. |
| **Question** | Find all cards illustrated by Stephen Daniel and describe the text of the ruling of these cards. State if these cards have missing or degraded properties and values. |
| **Evidence** | cards have missing or degraded properties and value refers to hasContentWarning = 1; 'Stephen Daniele' is artist; |
| **Gold SQL** | SELECT T1.id, T2.text, T1.hasContentWarning<br>FROM cards AS T1<br>INNER JOIN rulings AS T2<br>ON T1.uuid = T2.uuid<br>WHERE T1.artist = 'Stephen Daniele' |
| **Predicted SQL** | SELECT T1.uuid, T2.text, T1.hasContentWarning<br>FROM cards AS T1<br>INNER JOIN rulings AS T2<br>ON T1.uuid = T2.uuid<br>WHERE T1.artist = 'Stephen Daniele' |
| **Error Explanation** | "id" is the correct column. |

### F.1.2 Table:

Description: The tables used in the predicted SQL are either missing or incorrect.

| Question | Among the schools with the average score in Math over 560 in the SAT test, how many schools are directly charter-funded? |
|---|---|

| | |
|---|---|
| **Gold SQL** | ```SELECT COUNT(T2.School Code)```<br>```FROM satscores AS T1```<br>```INNER JOIN frpm AS T2```<br>```ON T1.cds = T2.CDSCode```<br>```WHERE T1.AvgScrMath > 560```<br>```AND T2.Charter Funding Type = 'Directly funded'``` |
| **Predicted SQL** | ```SELECT COUNT(T1.cds)```<br>```FROM satscores AS T1```<br>```INNER JOIN schools AS T2```<br>```ON T1.cds = T2.CDSCode```<br>```WHERE T1.AvgScrMath > 560```<br>```AND T2.FundingType = 'Charter'``` |
| **Error Explanation** | "frpm" is the correct table with the correct column. |
| **Question** | Who placed the order with the id 32423? |
| **Gold SQL** | ```SELECT T3.client_id```<br>```FROM order AS T1```<br>```INNER JOIN account AS T2```<br>```ON T1.account_id = T2.account_id```<br>```INNER JOIN client AS T3```<br>```ON T2.district_id = T3.district_id```<br>```WHERE T1.order_id = 32423``` |
| **Predicted SQL** | ```SELECT T2.client_id```<br>```FROM order AS T1```<br>```INNER JOIN disp AS T2```<br>```ON T1.account_id = T2.account_id```<br>```WHERE T1.order_id = 32423``` |
| **Error Explanation** | "account" and "client" are the correct tables with the correct columns. |

*F.1.3   Entity:*

Description: The entity used for filtering is incorrect in the predicted SQL.

| | |
|---|---|
| **Question** | For the customers who paid in the euro, what is their average total price of the transactions? |
| **Gold SQL** | ```SELECT AVG(T1.Price)```<br>```FROM transactions_1k AS T1```<br>```INNER JOIN gasstations AS T2```<br>```ON T1.GasStationID = T2.GasStationID```<br>```INNER JOIN customers AS T3```<br>```ON T1.CustomerID = T3.CustomerID```<br>```WHERE T3.Currency = 'EUR'``` |
| **Predicted SQL** | ```SELECT AVG(T1.Price)```<br>```FROM transactions_1k T1```<br>```JOIN customers T2```<br>```ON T1.CustomerID = T2.CustomerID```<br>```WHERE T2.Currency = 'Euro'``` |

| Error Explanation | "EUR" is the correct entity for filtering. |
|---|---|
| **Question** | List out the clients who choose statement of issuance after transaction are Disponent? |
| **Gold SQL** | ```SELECT T2.client_id FROM account AS T1 INNER JOIN disp AS T2 ON T1.account_id = T2.account_id WHERE T1.frequency = 'POPLATEK PO OBRATU' AND T2.type = 'DISPONENT'``` |
| **Predicted SQL** | ```SELECT T2.client_id FROM account AS T1 INNER JOIN disp AS T2 ON T1.account_id = T2.account_id WHERE T1.frequency = 'POPLATEK PO OBRATU' AND T2.type = 'Disponent'``` |
| **Error Explanation** | "DISPONENT" is the correct entity for filtering. |

*F.1.4   Revision:*

Description: The revision node incorrectly changes the correct predicted SQL into an incorrectly predicted SQL.

| **Question** | In Los Angeles how many schools have more than 500 free meals but less than 700 free or reduced price meals for K-12? |
|---|---|
| **Gold SQL** | ```SELECT COUNT(CDSCode) FROM frpm WHERE County Name = 'Los Angeles' AND Free Meal Count (K-12) > 500 AND 'Free Meal Count (K-12) < 700``` |
| **Predicted SQL** | ```SELECT COUNT(CDSCode) FROM frpm WHERE Free Meal Count (K-12) > 500 AND 'FRPM Count (K-12) < 700 AND County Name = 'Los Angeles'``` |
| **Error Explanation** | "Free Meal Count (K-12)" is the correct column for filtering. |
| **Question** | Find and list the names of sets which doesn't have Japanese translation but have Korean translation. |
| **Gold SQL** | ```SELECT name FROM sets WHERE code IN ( SELECT setCode FROM set_translations WHERE language = 'Korean' AND language NOT LIKE 'Japanese')``` |

| | |
|---|---|
| **Predicted SQL** | ```SELECT T1.name```<br>```FROM sets AS T1```<br>```WHERE EXISTS (SELECT 1```<br>```FROM set_translations AS T2```<br>```WHERE T1.code = T2.setCode```<br>```AND T2.language = 'Korean')```<br>```AND NOT EXISTS (SELECT 1```<br>```FROM set_translations AS T3```<br>```WHERE T1.code = T3.setCode```<br>```AND T3.language LIKE 'Japanese')``` |
| **Error Explanation** | Incorrect logic in subquery. |

### F.1.5   Malformed Date:

Description: The date format in the predicted SQL is incorrect.

| | |
|---|---|
| **Question** | Which year recorded the most consumption of gas paid in CZK? |
| **Evidence** | The first 4 strings of the values in the table yearmonth can represent year. |
| **Gold SQL** | ```SELECT SUBSTRING(T2.Date, 1, 4)```<br>```FROM customers AS T1```<br>```INNER JOIN yearmonth AS T2```<br>```ON T1.CustomerID = T2.CustomerID```<br>```WHERE T1.Currency = 'CZK'```<br>```GROUP BY SUBSTRING(T2.Date, 1, 4)```<br>```ORDER BY SUM(T2.Consumption) DESC```<br>```LIMIT 1``` |
| **Predicted SQL** | ```SELECT strftime('%Y', T2.Date) AS Year```<br>```FROM yearmonth T2```<br>```JOIN customers T1```<br>```ON T2.CustomerID = T1.CustomerID```<br>```WHERE T1.Currency = 'CZK'```<br>```GROUP BY Year```<br>```ORDER BY SUM(T2.Consumption) DESC```<br>```LIMIT 1``` |
| **Error Explanation** | Incorrect year formating. |

### F.1.6   Evidence:

Description: The evidence was not used correctly in the predicted SQL.

| | |
|---|---|
| **Question** | Name schools in Riverside which the average of average math score for SAT is greater than 400, what is the funding type of these schools? |
| **Evidence** | Average of average math = sum(average math scores) / count(schools). |

| Gold SQL | SELECT T1.sname, T2.`Charter Funding Type`<br>FROM satscores AS T1<br>INNER JOIN frpm AS T2<br>ON T1.cds = T2.CDSCode<br>WHERE T2.District Name LIKE 'Riverside%'<br>GROUP BY T1.sname, T2.Charter Funding Type`<br>HAVING CAST(SUM(T1.AvgScrMath) AS REAL) / COUNT(T1.cds) > 400 |
|---|---|
| Predicted SQL | SELECT DISTINCT T1.CDSCode, T1.FundingType<br>FROM schools AS T1<br>INNER JOIN satscores AS T2<br>ON T1.CDSCode = T2.cds<br>WHERE T1.County = 'Riverside'<br>AND T2.AvgScrMath > 400 |
| Error Explanation | Did not use the equation in the evidence. |
| Question | What is the detailed content of the comment of the post which was created on 7/19/2010 7:37:33 PM? |
| Evidence | detailed content of the comment refers to Text; created on 7/19/2010 7:37:33 PM CreationDate = 2010-07-19 19:37:33.0 |
| Gold SQL | SELECT T1.Text<br>FROM comments AS T1<br>INNER JOIN posts AS T2<br>ON T1.PostId = T2.Id<br>WHERE T1.CreationDate = '2010-07-19 19:37:33.0' |
| Predicted SQL | SELECT T2.Text<br>FROM posts AS T1<br>INNER JOIN comments AS T2<br>ON T1.Id = T2.PostId<br>WHERE T1.CreationDate = '2010-07-19 19:37:33' |
| Error Explanation | Did not use the time formatting in the evidence. |

### F.1.7   NULL/DISTINCT/GROUP BY:

Description: Missing NULL, DISTINCT, or GROUP BY in the predicted SQL.

| Question | When was the oldest mythic card released and what are its legal play formats? |
|---|---|
| Evidence | the oldest card refers to MIN(originalReleaseDate); mythic card refers to rarity = 'mythic'; legal play refers to status = 'legal'; play format refers to format |

| Gold SQL | `SELECT T1.originalReleaseDate, T2.format`<br>`FROM cards AS T1`<br>`INNER JOIN legalities AS T2`<br>`ON T1.uuid = T2.uuid`<br>`WHERE T1.rarity = 'mythic'`<br>`AND T1.originalReleaseDate IS NOT NULL`<br>`AND T2.status = 'Legal'`<br>`ORDER BY T1.originalReleaseDate`<br>`LIMIT 1` |
|---|---|
| **Predicted SQL** | `SELECT T1.originalReleaseDate, T2.format`<br>`FROM cards AS T1`<br>`INNER JOIN legalities AS T2`<br>`ON T1.uuid = T2.uuid`<br>`WHERE T1.rarity = 'mythic'`<br>`AND T2.status = 'Legal'`<br>`ORDER BY T1.originalReleaseDate ASC`<br>`LIMIT 1` |
| **Error Explanation** | Did not check if the selected column has NULL values. |
| **Question** | Name the foreign name of the card that has abzan watermark? List out the type of this card. |
| **Evidence** | |
| **Gold SQL** | `SELECT DISTINCT T1.name, T1.type`<br>`FROM cards AS T1`<br>`INNER JOIN foreign_data AS T2`<br>`ON T2.uuid = T1.uuid`<br>`WHERE T1.watermark = 'abzan'` |
| **Predicted SQL** | `SELECT T2.name, T2.type`<br>`FROM cards AS T1`<br>`INNER JOIN foreign_data AS T2`<br>`ON T1.uuid = T2.uuid`<br>`WHERE T1.watermark = 'abzan'` |
| **Error Explanation** | Did not use the DISTINCT keyword for the selected columns. |

*F.1.8   Operation:*

Description: Incorrect comparison or division operations in predicted SQL.

| Question | Among the male patients who have a normal level of white blood cells, how many of them have an abnormal fibrinogen level? |
|---|---|
| **Evidence** | male patients refers to Sex = 'M'; normal level of white blood cells refers to WBC between 3.5 and 9.0; abnormal fibrinogen level refers to FG <= 150 or FG >= 450; Don't compute repetitive ones. |

| Gold SQL | `SELECT COUNT(DISTINCT T1.ID)`<br>`FROM Patient AS T1`<br>`INNER JOIN Laboratory AS T2`<br>`ON T1.ID = T2.ID`<br>`WHERE T2.FG <= 150 OR T2.FG >= 450`<br>`AND T2.WBC > 3.5`<br>`AND T2.WBC < 9.0`<br>`AND T1.SEX = 'M'` |
|---|---|
| Predicted SQL | `SELECT COUNT(DISTINCT T1.ID)`<br>`FROM Patient AS T1`<br>`INNER JOIN Laboratory AS T2`<br>`ON T1.ID = T2.ID`<br>`WHERE T1.SEX = 'M'`<br>`AND T2.WBC BETWEEN 3.5 AND 9.0`<br>`AND (T2.FG < 150 OR T2.FG > 450)` |
| Error Explanation | Incorrect comparison operation. |
| Question | Among all chemical compounds identified in the database, what percent of compounds form a triple-bond. |
| Evidence | triple bond refers to bond_type = '#'; |
| Gold SQL | `SELECT CAST(COUNT(CASE WHEN T.bond_type = '#' THEN T.bond_id ELSE`<br>`NULL END) AS REAL) * 100 / COUNT(T.bond_id)`<br>`FROM bond AS T` |
| Predicted SQL | `SELECT CAST(COUNT(DISTINCT T1.molecule_id) AS REAL) * 100`<br>`/ (SELECT COUNT() FROM molecule) AS percentage`<br>`FROM bond AS T1`<br>`WHERE T1.bond_type = '#'` |
| Error Explanation | Incorrect division operation with the wrong column. |

*F.1.9 SELECT:*

Description: Incorrect formatting in SELECT in the predicted SQL.

| Question | Provide the full name and email address of the Student_Club's Secretary. |
|---|---|
| Evidence | full name refers to first_name, last_name; 'Secretary' is a position of Student Club |
| Gold SQL | `SELECT first_name, last_name, email`<br>`FROM member`<br>`WHERE position = 'Secretary'` |
| Predicted SQL | `SELECT first_name ||' ' || last_name AS full_name, email`<br>`FROM member`<br>`WHERE position = 'Secretary'` |
| Error Explanation | Incorrect column formatting in SELECT. |

## F.2 Vague Question

### F.2.1 Column:

Description: The question and evidence do not specify which columns should be included in the SELECT statement, making it impossible to determine the correct columns from the provided information and database. This leads to incorrect or additional columns in the gold SQL.

| | |
|---|---|
| **Question** | What are the cards for set OGW? State the colour for these cards. |
| **Evidence** | set OGW refers to setCode = 'OGW'; |
| **Gold SQL** | ```SELECT id, colors FROM cards WHERE id IN ( SELECT id FROM set_translations WHERE setCode = 'OGW' )``` |
| **Predicted SQL** | ```SELECT name, colors FROM cards WHERE setCode = 'OGW'``` |
| **Error Explanation** | Unclear if the question is asking for name or id. |
| **Question** | Which of these players performs the best in crossing actions, Alexis, Ariel Borysiuk or Arouna Kone? |
| **Evidence** | player who perform best in crossing actions refers to MAX(crossing); |
| **Gold SQL** | ```SELECT t1.player_name, t2.crossing FROM Player AS t1 INNER JOIN Player_Attributes AS t2 ON t1.player_api_id = t2.player_api_id WHERE t1.player_name IN ('Alexis', 'Ariel Borysiuk', 'Arouna Kone') ORDER BY t2.crossing DESC LIMIT 1``` |
| **Predicted SQL** | ```SELECT T1.player_name FROM Player AS T1 INNER JOIN Player_Attributes AS T2 ON T1.player_fifa_api_id = T2.player_fifa_api_id WHERE T1.player_name IN ('Alexis', 'Ariel Borysiuk', 'Arouna Kone') ORDER BY T2.crossing DESC LIMIT 1``` |
| **Error Explanation** | Unclear from the question if the output should include both the names and the crossing actions. |

## F.3 Incorrect Golden SQL

### F.3.1 Evidence:

Description: The gold SQL does not use the evidence correctly.

| | |
|---|---|
| **Question** | What is the percentage of Story Spotlight cards that also have a text box? List them by their ID. |

| | |
|---|---|
| **Evidence** | Story Spotlight cards that do not have a text box refers to isStorylight = 1 and isTextless = 1; Percentage refer to DIVIDE(SUM(count(id) where isStorylight = 1), SUM(count(id))) * 100 |
| **Gold SQL** | `SELECT CAST(SUM(CASE WHEN isTextless = 0 THEN 1 ELSE 0 END) AS REAL)  100 / COUNT(id)`<br>`FROM cards`<br>`WHERE isStorySpotlight = 1` |
| **Predicted SQL** | `SELECT CAST(SUM(CASE WHEN isStorySpotlight = 1 AND isTextless = 0 THEN 1 ELSE 0 END) AS REAL)  100 / COUNT(id) AS percentage`<br>`FROM cards` |
| **Error Explanation** | The filtering should be done in the numerator as shown in the evidence. |
| **Question** | Who is the illustrator that illustrated the least amount of cards? List the format of play of the cards that he/she illustrated. |
| **Evidence** | format of the cards refers to format; illustrator refers to artist; the least amount of cards refers to MIN(artist) |
| **Gold SQL** | `SELECT T1.artist, T2.format`<br>`FROM cards AS T1`<br>`INNER JOIN legalities AS T2`<br>`ON T2.uuid = T1.uuid`<br>`GROUP BY T1.artist`<br>`ORDER BY COUNT(T1.id) ASC`<br>`LIMIT 1` |
| **Predicted SQL** | `SELECT T1.artist, T2.format`<br>`FROM cards AS T1`<br>`INNER JOIN legalities AS T2`<br>`ON T1.uuid = T2.uuid`<br>`GROUP BY T1.artist, T2.format`<br>`ORDER BY COUNT(T1.artist) ASC`<br>`LIMIT 1` |
| **Error Explanation** | The gold SQL did not follow the evidence instruction to use the right column in MIN(artist). |

*F.3.2 Description:*

Description: The gold SQL does not utilize the column or value descriptions provided in the database, leading to incorrect SQL.

| | |
|---|---|
| **Question** | Which gas station has the highest amount of revenue? |
| **Gold SQL** | `SELECT GasStationID`<br>`FROM transactions_1k`<br>`GROUP BY GasStationID`<br>`ORDER BY SUM(Price) DESC`<br>`LIMIT 1` |

| | |
|---|---|
| **Predicted SQL** | SELECT T2.GasStationID<br>FROM gasstations AS T1<br>INNER JOIN transactions_1k AS T2<br>ON T1.GasStationID = T2.GasStationID<br>WHERE T2.Amount IS NOT NULL<br>AND T2.Price IS NOT NULL<br>GROUP BY T2.GasStationID<br>ORDER BY SUM(T2.Amount * T2.Price) DESC<br>LIMIT 1 |
| **Error Explanation** | The description in the database indicates revenue = amount*price. |
| **Question** | What's the finish time for the driver who ranked second in 2008's Australian Grand Prix? |
| **Gold SQL** | SELECT T1.time<br>FROM results AS T1<br>INNER JOIN races AS T2<br>ON T1.raceId = T2.raceId<br>WHERE T1.rank = 2<br>AND T2.name = 'Australian Grand Prix'<br>AND T2.year = 2008 |
| **Predicted SQL** | SELECT T2.time<br>FROM races AS T1<br>INNER JOIN results AS T2<br>ON T1.raceId = T2.raceId<br>WHERE T1.name = 'Australian Grand Prix'<br>AND T1.year = 2008<br>AND T2.positionOrder = 2 |
| **Error Explanation** | The description in the database indicates the position order is the finishing order of positions while rank is the starting rank positioned by fastest lap speed. |

*F.3.3   Aggregation:*

Description: The gold SQL misses aggregation operations.

| | |
|---|---|
| **Question** | What is the amount of debt that client number 992 has, and how is this client doing with payments? |
| **Gold SQL** | SELECT T3.amount, T3.status<br>FROM client AS T1<br>INNER JOIN account AS T2<br>ON T1.district_id = T2.district_id<br>INNER JOIN loan AS T3<br>ON T2.account_id = T3.account_id<br>WHERE T1.client_id = 992 |

| | |
|---|---|
| **Predicted SQL** | ```SELECT SUM(T3.amount) AS total_debt, T3.status FROM client AS T1 INNER JOIN disp AS T2 ON T1.client_id = T2.client_id INNER JOIN loan AS T3 ON T2.account_id = T3.account_id WHERE T1.client_id = 992 GROUP BY T3.status``` |
| **Error Explanation** | Gold SQL did not use SUM to sum the full amount. |
| **Question** | How many cards have frame effect as extendedart? List out the id of those cards. |
| **Gold SQL** | ```SELECT id FROM cards WHERE frameEffects = 'extendedart' GROUP BY id``` |
| **Predicted SQL** | ```SELECT COUNT(id) FROM cards WHERE frameEffects = 'extendedart' UNION ALL SELECT id FROM cards WHERE frameEffects = 'extendedart'``` |
| **Error Explanation** | The gold SQL did not return a count as asked by the question. |

*F.3.4   Column:*

Description: The gold SQL does not have all the expected columns from the question in SELECT.

| | |
|---|---|
| **Question** | List out all post that are related to post ID 61217 and what is the popularity of this post? |
| **Evidence** | post related refers to RelatedPostId; popularity refers to ViewCount |
| **Gold SQL** | ```SELECT T1.ViewCount FROM posts AS T1 INNER JOIN postLinks AS T2 ON T1.Id = T2.PostId WHERE T2.PostId = 61217``` |
| **Predicted SQL** | ```SELECT T2.RelatedPostId, T1.ViewCount FROM posts AS T1 INNER JOIN postLinks AS T2 ON T1.Id = T2.PostId WHERE T2.PostId = 61217``` |
| **Error Explanation** | The gold SQL did not select all the columns asked by the question. |
| **Question** | List down the bond type for molecules from molecule id TR000 to TR050. |
| **Evidence** | double bond refers to bond_type = ' = '; single bond refers to bond_type = '-'; triple bond refers to bond_type = '#'; |

| Gold SQL | SELECT T2.molecule_id, T2.bond_type<br>FROM molecule AS T1<br>INNER JOIN bond AS T2<br>ON T1.molecule_id = T2.molecule_id<br>WHERE T1.molecule_id BETWEEN 'TR000' AND 'TR050' |
|---|---|
| Predicted SQL | SELECT T2.bond_type<br>FROM molecule AS T1<br>INNER JOIN bond AS T2<br>ON T1.molecule_id = T2.molecule_id<br>WHERE T1.molecule_id LIKE 'TR0%'<br>AND T1.molecule_id BETWEEN 'TR000' AND 'TR050' |
| Error Explanation | The gold SQL has more columns selected than what the question asked for. |

*F.3.5   Filtering:*

Description: The gold SQL does not filter using the correct entity or is missing necessary filters.

| Question | How many male customers who are living in North Bohemia have average salary greater than 8000? |
|---|---|
| Evidence | Male means that gender = 'M'; A3 refers to region; A11 pertains to average salary. |
| Gold SQL | SELECT COUNT(T1.client_id)<br>FROM client AS T1<br>INNER JOIN district AS T2<br>ON T1.district_id = T2.district_id<br>WHERE T1.gender = 'M'<br>AND T2.A3 = 'North Bohemia'<br>AND T2.A11 > 8000 |
| Predicted SQL | SELECT COUNT(T1.client_id)<br>FROM client AS T1<br>INNER JOIN district AS T2<br>ON T1.district_id = T2.district_id<br>WHERE T1.gender = 'M'<br>AND T2.A3 = 'north Bohemia'<br>AND T2.A11 > 8000 |
| Error Explanation | The correct entity name has a lower case in north. |
| Question | How many users from New York have a teacher and supporter badge? |
| Evidence | "Supporter" and "Teachers" are both Name of badge; 'New York' is the Location; user refers to UserId |
| Gold SQL | SELECT COUNT(DISTINCT T1.Id)<br>FROM badges AS T1<br>INNER JOIN users AS T2<br>ON T1.UserId = T2.Id<br>WHERE T1.Name IN ('Supporter', 'Teacher')<br>AND T2.Location = 'New York' |

| Predicted SQL | SELECT COUNT(DISTINCT T1.Id)<br>FROM users AS T1<br>INNER JOIN badges AS T2<br>ON T1.Id = T2.UserId<br>WHERE T1.Location = 'New York'<br>AND T2.Name IN ('Supporter', 'Teacher')<br>GROUP BY T1.Id<br>HAVING COUNT(DISTINCT T2.Name) = 2 |
|---|---|
| Error Explanation | The gold SQL missed a filtering that indicates a user is both teacher and a supporter. |

### F.3.6  NULL/DISTINCT/GROUP BY:

Description: The gold SQL missed checking NULL, adding DISTINCT, or adding GROUP BY.

| Question | Calculate the total amount spent on speaker gifts and list the name of the event they were spent on. |
|---|---|
| Evidence | total amount spent = SUM(spent) where category = 'Speaker Gifts' |
| Gold SQL | SELECT SUM(T1.spent), T2.event_name<br>FROM budget AS T1<br>INNER JOIN event AS T2<br>ON T1.link_to_event = T2.event_id<br>WHERE T1.category = 'Speaker Gifts' |
| Predicted SQL | SELECT SUM(T2.spent), T1.event_name<br>FROM event AS T1<br>INNER JOIN budget AS T2<br>ON T1.event_id = T2.link_to_event<br>WHERE T2.category = 'Speaker Gifts'<br>GROUP BY T1.event_name |
| Error Explanation | The gold SQL misses GROUP BY. |