Anonymous authorsPaper under double-blind review

000

001

003

006

007

008 009 010

011 012 013

014

015

016

018

019

021

024

025

026

027

028

029

031

034

037 038

040

041

042

043

044

046

047

048

052

ABSTRACT

SQL is the core of data analysis and engineering across industries, powering largescale workflows for data extraction, transformation, and loading. However, in enterprise-level scenarios, it is challenging to generate fully correct SQL code in a single attempt—even for experienced developers or advanced Text-to-SQL LLMs. Multiple iterations of debugging are usually required, yet LLMs often get lost in multi-turn correction. To address this gap, we introduce **Squirrel Benchmark**, the first benchmark designed for enterprise-level SQL reasoning and debugging. Our benchmark is built upon two key innovations: (1) an automated construction workflow that employs reverse engineering to systematically inject realistic bugs into large-scale SQL code, enabling scalable and diverse benchmark generation; and (2) an **execution-free evaluation framework** tailored for enterprise settings, providing fast, accurate, and resource-efficient assessment. Squirrel Benchmark comprises 469 Squirrel-Syntax queries featuring syntax errors with explicit error messages, and 516 Squirrel-Semantic queries targeting semantic errors where SQL fails to meet the requirement. The SQLs are substantially complex, averaging over 140 lines with abstract syntax trees of high complexity (average width > 11, depth > 8.7). We evaluate nearly 30 LLMs on Squirrel Benchmark. Even state-of-the-art reasoning models struggle: Claude-4-Sonnet achieves only 36.46% success on Squirrel-Syntax and 32.17% on Squirrel-Semantic. Most models fail to reach 20% success, underscoring the significant gap between current LLM capabilities and the demands of enterprise SQL debugging. To bridge this gap, we systematically explore four potential solution strategies and conduct extensive experiments to evaluate and compare their effectiveness. Our experiments not only highlight the challenges but also shed light on effective strategies for advancing SQL debugging with LLMs.

1 Introduction

Databases are a cornerstone of modern data infrastructure, powering critical applications across finance, web services, and scientific computing. Structured Query Language (SQL) remains the predominant interface for human—data interaction, enabling large-scale extraction, transformation, and loading (ETL) workflows (Chamberlin & Boyce, 1974; Armbrust et al., 2015). Recent research on Text-to-SQL large language models (LLMs) has sought to help analysts automate routine queries, streamline data workflows, and support advanced business intelligence (Zhong et al., 2017; Yu et al., 2018; Li et al., 2025a).

Enterprise SQL code is often lengthy, complex, and deeply nested, making it extremely challenging for both experienced developers and Text-to-SQL LLMs to generate correct code in a single attempt (Lei et al., 2025). Instead, success typically requires multi-step reasoning and iterative debugging. As shown in Figure 1, debugging generally involves localizing errors, analyzing their causes, consulting schema definitions, applying targeted modifications, and re-executing queries to check whether requirements are satisfied—usually repeating this loop multiple times. Unfortunately, LLMs struggle with this iterative correction process. They frequently fall into anti-patterns such as repeating identical actions without meaningful follow-up, which leads to wasted effort when an initial correction fails (Bouzenia & Pradel, 2025; Laban et al., 2025).

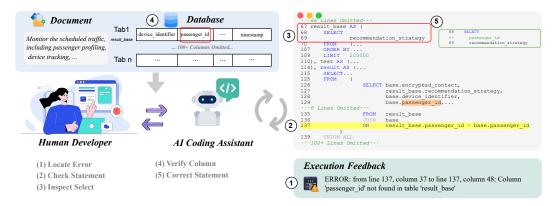


Figure 1: Squirrel Benchmark evaluates LLMs on real-world enterprise-level SQL debugging workflows. It involves multi-step reasoning and actions, including understanding requirements and schemas, diagnosing executor errors, and iteratively refining queries through cycles of reasoning and debugging.

To bridge this gap, we propose moving beyond Text-to-SQL generation and shifting the focus to a model's ability to iteratively debug and self-correct. We introduce Squirrel Benchmark, a benchmark for evaluating LLMs on enterprise-scale SQL debugging. Our construction pipeline utilizes an automated reverse-engineering framework to synthesize realistic and reproducible tasks. This approach minimizes human effort while ensuring high-quality benchmark generation, also providing a foundation for synthetic training data. Furthermore, we design an execution-free evaluation framework tailored to enterprise SQL scenarios. Squirrel Benchmark offers a practical reference point for selecting SQL-focused LLMs in industry. The benchmark comprises 469 Squirrel-Syntax tasks (syntax errors with explicit error messages) and 516 Squirrel-Semantic tasks (semantic errors in which the SQL output does not match the user's requirement). SQL programs in our benchmark are highly complex, averaging over 140 lines (> 420 tokens), with ASTs of width > 11 and depth > 8.7, and incorporating over 15 functions per script on average.

Our evaluation on Squirrel Benchmark indicates significant room for improvement in deploying LLMs within SQL-SWE workflows. Extensive experiments show that even state-of-the-art LLMs struggle: Claude-4-Sonnet achieves only 36.46% success on Squirrel-Syntax and 33.17% on Squirrel-Semantic, while most models fail to reach 20%. These results underscore the difficulty of enterprise SQL debugging and highlight substantial room for improvement. To address this gap, we systematically explore four potential solution strategies and conduct comprehensive experiments to assess their effectiveness. Our results not only illuminate the challenges faced by LLMs in SQL debugging but also provide insights into strategies that can advance performance. Moreover, Squirrel Benchmark exhibits a strong correlation with real-world debugging outcomes, establishing it as a reliable benchmark for aligning models with industrial applications. In summary, this work makes the following contributions:

- We propose an automated reverse-engineering workflow for constructing high-quality SQL debugging benchmarks, which can also be adapted to synthesize realistic training data.
- We present Squirrel Benchmark, a large-scale benchmark comprising 469 syntax and 516 semantic tasks, designed to capture the complexity, diversity, and practicality of enterprise SQL development.
- We conduct a comprehensive evaluation of nearly 30 open-source and proprietary LLMs, showing that even the state-of-the-art LLMs face substantial challenges.
- We introduce three SFT and an agent method as baselines, offering a novel and efficient pathway for further studies.

2 Preliminary

2.1 TASK DEFINITION

SQL debugging is a fundamental but underexplored problem in data-centric machine learning. Existing Text-to-SQL research primarily focuses on translating natural language to SQL queries, but real-world scenarios often involve correcting issues in SQL scripts. The goal of SQL debugging is to automatically repair buggy SQL scripts while preserving the user's intent. This task begins with a

buggy SQL (b), accompanied by auxiliary context \mathcal{C} (e.g., error messages or natural language intent descriptions) and the database schema (σ). The objective is to generate a corrected SQL (\hat{q}):

 $\hat{q} = f_{\theta}(\mathcal{C}, \sigma, b) \tag{1}$

where \hat{q} is executable and faithful to the intent encoded in (C, b, σ) .

We categorize bugs into two primary types: (I) Syntactic errors. b is non-executable. Here, C is the error message E, and the goal is to produce an executable repair while preserving its intended semantics. (II) Semantic errors. b executes successfully but fails to meet the user's requirements. In this case, C is a natural language specification R, and the task is to modify \hat{q} to satisfy R. By covering both types, Squirrel Benchmark unifies execution repair with intent comprehension, offering a challenging and realistic benchmark for SQL debugging.

2.2 CHALLENGES

Despite its practical importance, SQL debugging introduces several unique challenges that are not sufficiently addressed in existing SWE research.

Challenge 1: Lack of Enterprise-level SQL Scripts. Industrial SQL workloads, such as ETL workflows and scheduled analytical jobs, are typically *long*, *complex*, and *schema-heavy*. Scripts can span hundreds of lines, involve deeply nested subqueries and multi-way joins. and reference dozens of tables and columns. This level of intricacy significantly amplifies the challenge for LLMs. In contrast, most existing Text-to-SQL (Li et al., 2024) and SQL-debugging (Li et al., 2025b) benchmarks focus on short, relatively simple queries that are far removed from the scale and complexity of enterprise environments. Unfortunately, such industrial-grade SQL scripts are rarely available in the open-source community, resulting in a pronounced mismatch between academic benchmarks and real-world needs.

Contribution 1: The First Enterprise-level SQL Benchmark

To address this gap, we introduce a large-scale, enterprise-level benchmark that captures the complexity of real-world ETL and analytical workloads (Section 3.1).

Challenge 2: Lack of a Comprehensive Bug Taxonomy. SQL bugs are heterogeneous: some manifest as execution failures (syntax errors), while others silently yield incorrect results (semantic errors). Although recent benchmarks such as BIRD-Critic (Li et al., 2025b) have advanced debugging evaluation, they lack a systematic taxonomy of SQL-specific bug types and their prevalence. Without such categorization, it is difficult to understand where models struggle most and how to target improvements effectively. A comprehensive analysis of SQL bug categories is therefore crucial, not only for benchmarking but also for guiding the design of future bug-fixing models.

Contribution 2: A Hierarchical SQL Bug Taxonomy

We develop a hierarchical taxonomy of SQL bug types derived from an extensive analysis of real-world errors. This provides a structured framework for fine-grained evaluation (Section 3.2).

Challenge 3: Lack of Reliable and Comprehensive SQL debugging Benchmark. High-quality benchmarks for SQL Debugging are scarce. Manually curated datasets are costly to produce and prone to evaluation leakage if models memorize solutions from public templates or repositories. Existing resources often lack diversity, realistic bug patterns, and coverage of enterprise-scale scripts, limiting their usefulness for robust model evaluation. Building a reliable, large-scale benchmark that is both comprehensive and faithful to real-world workflows is therefore a significant challenge.

Contribution 3: An Automated SQL-SWE Synthesis Pipeline

We introduce an automated pipeline for synthesizing and validating SQL bug-fixing examples, ensuring scalability, diversity, and resistance to data leakage (Section 3.3).

3 SQUIRREL BENCHMARK CONSTRUCTION

Figure 2 shows the automated benchmark construction pipeline. It comprises four stages: (1) enterprise-level SQL script generation (Section 3.1), (2) SQL bug taxonomy design (Section 3.2),

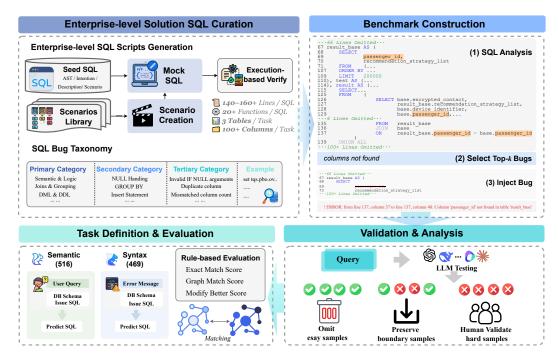


Figure 2: Overview of the Squirrel Benchmark construction and evaluation pipeline. Benchmark construction consists of 4 main stages: (1) Enterprise-level SQL Script Generation, (2) SQL bug taxonomy Design, (3) Issue SQL Construction via reverse engineering, and (4) Validation and Analysis. This pipeline ensures diversity, realism, and rigorous evaluation of SQL Debugging task.

(3) issue SQL construction via reverse engineering (Section 3.3), and (4) validation and analysis (Section 3.4). Section 3.5 further introduces an efficient execution-free evaluation methodology. Section 3.5 presents an execution-free evaluation methodology. All synthetic content is generated with Claude-4-Sonnet (anthropic, 2025) at temperature 0. Examples and prompts are detailed in Appendix F and G.

3.1 Enterprise-level SQL Scripts Generation

Because enterprise SQL scripts are proprietary and rarely accessible, we synthesize realistic, high-quality enterprise SQL.

Seed Enterprise SQL Curation. We curate high-quality SQL scripts q along with corresponding table definitions σ from real-world enterprise applications. To ensure that queries are non-trivial and representative of practical workloads, we filter scripts that fall below a complexity threshold τ . Complexity is quantified via a composite metric:

$$C(q) = \alpha \left(D_{\text{AST}}(q) + W_{\text{AST}}(q) \right) + \beta L(q)$$
 (2)

where D_{AST} , W_{AST} , and L(q) denote AST depth, AST width, and code length, respectively.

For each retained SQL script, we utilize an LLM to abstract its business domain (d), intention (I), and descriptive scenario (S). All scenarios are aggregated into a Scenarios Library, denoted as $\mathcal{D}_{\text{domain}} = \{d\}$. The resulting seed dataset is then defined as:

$$\mathcal{D}_{\text{seed}} = \{ (q_i, \sigma_i, d_i, I_i, S_i, \text{AST}(q_i)) | q_i \in \mathcal{Q}_s, \mathcal{C}(q_i) > \tau \}, \tag{3}$$

where Q_s denotes the candidate SQL pool.

The final seed corpus contains $1{,}000+$ SQL scripts spanning 26 business scenarios, averaging over 120 lines with AST depth > 8 and width > 12. Each script is rigorously validated to be bug-free, resulting in a corpus that accurately captures both the structural complexity and semantic diversity of enterprise SQL.

Solution SQL Synthesis. To expand coverage across domains and code structures, we synthesize new SQL scripts using the seed corpus and the Scenarios Library:

- 1. Seed Sampling. Select $(q_i, \sigma_i, d_i, I_i, S_i, AST(q_i)) \in \mathcal{D}_{seed}$ and a target domain $d_t \in \mathcal{D}_{domain}$.
- 2. Scenario Creation. Conditioned on d_t , the LLM generates a new scenario description S_t together with schema definitions σ_t , following the structure of the seed corpus.
- 3. SQL Synthesis. Given $(I_i, S_i, AST(q_i), S_t, \sigma_t)$, the LLM generates a new SQL script q_t that preserves the complexity of the seed SQL scripts while adapting to the new schema and scenario. This ensures that synthesized queries remain realistic, non-trivial, and representative of enterprise workloads.
- 4. Execution-based Validation. To ensure the correction, each candidate q_t is validated via execution. Specifically, σ_t is instantiated to construct a fake test database, q_t is executed, and only queries that successfully execute are retained:

$$Q_{gt} = \{ (q_t, \sigma_t) \mid exec(q_t, \sigma_t) == passed \}$$
(4)

This synthesis pipeline ensures that the final SQL dataset exhibits (i) enterprise-grade complexity, (ii) broad domain coverage via controlled scenario transfer, and (iii) guaranteed execution correctness.

3.2 SQL Bug Taxonomy

We construct an SQL bug taxonomy by manually annotating 268 erroneous SQL queries collected from real-world applications. Each bug is classified according to a three-level hierarchical error type: (i) *macro categories* (e.g., DML, DDL, semantic, and logic), (ii) *construct-specific subcategories* (e.g., INSERT statements), and (iii) *atomic faults* (e.g., mismatched column counts). This taxonomy organizes common failure patterns and forms a bug library of realistic error templates. The library underpins our controlled bug-injection process (Section 3.3), ensuring that Squirrel Benchmark captures authentic SQL error modes. Table 4 and 3 report the distribution of bug types.

3.3 ISSUE SQL CONSTRUCTION

We construct issue SQL queries through reverse engineering, transforming correct SQL scripts into buggy versions. The process is guided by three principles: structural awareness, taxonomy-guided selection, and minimal-change injection, ensuring that the generated bugs are both realistic and diagnostically useful.

Step 1: Structural Profiling and Taxonomy-Guided Selection For each ground-truth SQL $q_{\rm gt}$, we first analyze its structural and semantic profile, including the AST, function patterns, and clause usage. Based on this profile, we then select the top-k candidate bug types from our hierarchical SQL bug taxonomy. This approach ensures that the injected errors are well-suited to the given SQL while providing broad coverage of real-world error scenarios.

Step 2: Minimal Change-Based Bug Injection. Each injected bug represents the smallest possible modification that induces the targeted error type. This principle preserves maximal similarity between the buggy SQL b and its reference $q_{\rm gt}$, isolating the error signal and reducing confounding factors. As a result, evaluating whether a model can localize and repair the fault becomes both precise and interpretable.

3.4 VALIDATION AND ANALYSIS

We validate Squirrel Benchmark via a model-driven *attack-defense* process. The goal is to filter out trivial cases that most models can easily solve, while retaining challenging but solvable instances that better reflect real-world debugging.

Automated Verification. We first attack the benchmark by evaluating each generated instance with a diverse set of advanced LLMs (including Qwen3-Coder-32B(Yang et al., 2025a), GPT-5(Openai, 2025), DeepSeek-V3.1(DeepSeek, 2025), Claude-4-sonnet(anthropic, 2025), and others). Instances fall into three categories: (i) If the majority of models succeed, the instance is deemed too easy and discarded; (ii) If only a few models succeed, the instance is considered an edge case and retained; (iii) If none of the models succeed, the instance is flagged for manual review. This adversarial filtering ensures that the benchmark emphasizes cases where current models diverge, thereby sharpening its discriminatory power.

Human Verification. Instances flagged as potentially unsolvable are subjected to manual inspection by three expert annotators with extensive SQL experience. Following a cross-validation protocol, annotators assess whether the task is logically inferable from the provided context and whether

multiple valid solutions exist. Instances that fail to meet these criteria are removed. For cases where multiple correct answers are possible, annotators supplement the benchmark with all valid alternative solutions.

Through this *attack-defense* protocol, Squirrel Benchmark removes trivial cases, yielding a challenging yet solvable testbed.

3.5 EVALUATION METRICS

The prevailing metrics for SQL debugging are Exact Match (EM) and Execution Accuracy. However, EM is notoriously strict, failing to credit semantically equivalent queries with divergent syntax. Execution Accuracy, while more forgiving, introduces false positives when test databases lack the necessary content to reveal logical errors (Zhan et al., 2025). Direct execution in production also poses practical barriers, being computationally expensive and raising data privacy concerns. To overcome these challenges, we introduce an execution-free evaluation framework based on three metrics (Detailed definitions and formulas are provided in Appendix C.1.):

- (1) **Exact Match Score (EM)**: This metric assesses strict syntactic correctness by checking for string-level identity between the predicted and reference SQL queries, thereby serving as a baseline for syntactic alignment.
- (2) **Graph Match Score (GM)**: This metric evaluates structural and functional equivalence by comparing the abstract syntax trees (ASTs) of the predicted and reference queries, thereby capturing semantic correctness where EM fails.
- (3) **Modify-Better Score (MB)**: This metric gauges iterative improvement capability by comparing the edit distances from the predicted SQL and the original SQL to the reference, thereby measuring how much closer the refinement is to the target.

4 BENCHMARK STATISTICS

We present a statistical analysis of Squirrel Benchmark, comparing its key features with existing SQL datasets in Table 1 and Figure 3. Our benchmark is designed to emphasize both *complexity* and *realism*, closely mirroring the challenges found in real-world industrial environments—particularly in terms of SQL script structure, error taxonomy, and task diversity.

Table 1: Statistical comparison of Squirrel Benchmark with representative text-to-SQL and SQL debugging benchmarks. The table evaluates benchmarks on scale (# examples), script length (avg. tokens and lines), and structural complexity (avg. function count, AST depth, and width).

	Туре	# Test Examples	Length of SQL		Complex of SQL		
Benchmark			# Tok. /SQL	# Line. /SQL	# Func. /SQL	# AST Depth /SQL	# AST Width /SQL
Spider 1.0 (Yu et al., 2018)	Text-to-SQL	2,147	18.50	_	_	_	_
Spider 2.0-snow (Yu et al., 2018)	Text-to-SQL	121	154.63	56.12	14.90	11.95	9.66
Spider 2.0-lite (Yu et al., 2018)	Text-to-SQL	256	131.79	49.84	13.65	11.97	10.05
BIRD (Li et al., 2024)	Text-to-SQL	1,789	30.90	_	_	_	_
BIRD-Critic-open (Li et al., 2025b)	SQL debugging	600	49.18	9.73	4.30	8.03	6.01
BIRD-Critic-postgresql (Li et al., 2025b)	SQL debugging	530	51.44	6.92	4.78	8.25	6.34
BIRD-Critic-flash (Li et al., 2025b)	SQL debugging	200	34.53	2.84	4.06	7.85	5.20
Squrriel-Syntax	SQL debugging	469	496.90	163.69	21.62	8.93	11.69
Squrriel-Semantic	SQL debugging	516	425.93	141.58	17.34	8.75	11.12

Complexity of SQL Scripts. The SQL scripts in Squirrel Benchmark are not only longer but also structurally more complex, presenting challenges that better mirror real-world enterprise systems. With an average length of 140-160 lines and over 420 tokens, our scripts are an order of magnitude larger than those in BIRD-Critic (which average under 10 lines). This scale directly implies a higher probability of errors and a greater need for models to maintain long-range context and dependency understanding. Additionally, the high number of functions per script (17.34 in Squirrel-Semantic, 21.62 in Squirrel-Syntax) necessitates reasoning across multiple subqueries and nested expressions—a capability that many existing sequence-to-sequence models lack. This scale and functional richness underscore the increased complexity and practical difficulty of the debugging tasks in our benchmark.

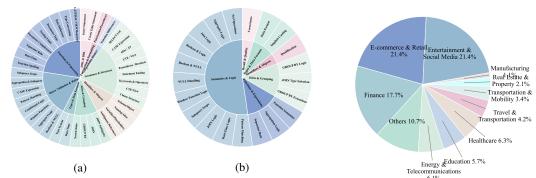


Figure 3: Statistics of errors and domain distribution in Squirrel Benchmark. (a) Two-level error types in Squirrel-Syntax, highlighting the distribution of syntax errors. (b) Two-level error types in Squirrel-Semantic, showing the distribution of semantic errors. (c) Distribution of SQL code across different business domains.

Hierarchical Error Taxonomy. Figures 3(a) and (b) show the two-level error taxonomy for Squirrel-Syntax and Squirrel-Semantic. Detailed error type statistics are in Appendix E. Squirrel Benchmark covers a broad spectrum of common syntax and semantic errors, enabling fine-grained evaluation of model capabilities. Syntax errors include issues related to grammar, structure, and dialect, while semantic errors encompass type mismatches, aggregation errors, and logical inconsistencies. This hierarchical classification allows for detailed insight into model performance across error types, supporting a more rigorous assessment of debugging ability.

Diversity of Task Scenarios. As shown in Figure 3 (c), the domains in Squirrel Benchmark span finance, e-commerce, healthcare, and more than ten additional areas, ensuring that models are evaluated against a broad range of business logic and contextual dependencies. For example, a program from the financial domain may involve complex window functions for time-series analysis, whereas an e-commerce program might require reasoning over multi-table joins across user and product schemas. This diversity tests a model's ability to generalize beyond simplistic syntactic patterns and demands domain-aware reasoning. Consequently, performance on Squirrel Benchmark provides a stronger indicator of a model's practicality and readiness for deployment in heterogeneous real-world environments.

5 EXPERIMENTS

Due to space limitations, we provide detailed experimental settings in Appendix C. This section focuses on the most important results.

5.1 MAIN RESULTS

Existing LLMs are far from being experts on enterprise SQL debugging. As shown in Table 2, we evaluate a diverse set of LLMs on Squirrel, including the Qwen, DeepSeek, Claude, GPT, Gemini, and Doubao families. Claude-4-Sonnet achieves the best performance, with a peak success rate of 36.46% GM score on Squirrel-Syntax and 32.17% GM score on Squirrel-Semantic. Interestingly, although our benchmark is constructed through reverse engineering using Claude-4-Sonnet, it still struggles with forward debugging. Other closed-source LLMs perform even worse, with most failing to exceed 20% GM. Among open-source models, DeepSeek-V3 achieves 30.28% on Squirrel-Syntax, and Qwen-2.5-Coder-32B attains 23.45% on Squirrel-Semantic, demonstrating competitive performance relative to closed-source systems.

Code generation LLMs struggle with SQL debugging. In previous studies, most code LLMs are heavily optimized for code generation, achieving strong performance on benchmarks such as SWE-Bench (Jimenez et al., 2024), BIRD (Li et al., 2024), and Spider (Yu et al., 2018). For example, OmniSQL (Li et al., 2025a), a Text-to-SQL-specialized model, achieves 87.6% on Spider and 64.5% on BIRD. However, its performance on Squirrel-Syntax and Squirrel-Semantic drops sharply to only 6.4% GM, underscoring the substantial gap between SQL generation and SQL debugging.

Reasoning-oriented LLMs (RLMs) exhibit stronger refinement abilities. Comparing RLMs with non-RLMs, we find that RLMs consistently perform better across both open-source and closed-source families. Notably, most RLMs achieve MB scores above 50%, indicating that while their predictions often move closer to the correct solution, they rarely solve the task in a single attempt.

Squirrel-Semantic is more challenging than Squirrel-Syntax. Across all evaluated models, performance on Squirrel-Semantic is consistently lower than on Squirrel-Syntax. This is because Squirrel-

Table 2: Evaluation results of LLMs on Squirrel-Syntax and Squirrel-Semantic. For each section, the best performance is highlighted in **bold**, and the second-best is <u>underlined</u>. EM, GM, and MB denote exact match score, graph match score, and modify-better score, respectively.

				Squirrel-Syntax			Squirrel-Semantic		
Model	Size	Reasoning	MoE	EM	GM	MB	EM	GM	MB
Open Source									
Qwen-2.5-Instruct	7B			2.13	8.53	33.05	1.94	5.62	14.15
Qwen-2.5-Coder	7B			3.20	8.96	37.53	4.84	7.75	18.99
Qwen-2.5-Coder	32B			12.79	20.26	52.88	17.44	23.45	34.69
Qwen-3-Instruct	235B		\checkmark	9.38	20.47	61.19	10.27	15.50	27.57
Qwen-3-Coder-Instruct	30B		\checkmark	5.54	20.90	44.14	6.40	15.12	24.42
Qwen-3-Coder-Instruct	480B		\checkmark	14.93	23.88	61.62	17.05	19.96	31.84
QwQ	32B	✓		8.76	20.51	41.45	10.47	15.31	20.16
Seed-Coder-Instruct	8B			8.53	14.93	42.43	8.72	14.15	24.61
OmniSQL	32B			0.21	6.40	50.75	0.39	6.40	21.17
Deepseek-V3	685B		\checkmark	17.91	30.28	60.34	11.24	21.32	33.27
Deepseek-V3.1	685B		\checkmark	17.91	30.49	63.61	12.02	14.73	32.47
Deepseek-R1	671B	✓		18.34	21.98	58.64	15.89	22.09	30.14
Closed Source									
Claude-4-Sonnet	_	✓		23.88	36.46	68.02	31.78	32.17	43.69
GPT-4o-mini-2024-07-18	_			1.71	4.69	13.01	5.62	6.40	8.74
GPT-4o-2024-11-20	_			2.14	4.69	13.79	2.91	4.84	6.86
GPT-4.1	_			6.40	17.70	61.25	8.52	17.05	30.49
GPT-5	_	✓		13.43	18.55	66.52	16.28	16.47	29.90
Gemini-2.5-Pro	_	✓		15.78	21.54	62.37	14.15	23.06	34.37
Kimi-K2	_	✓	\checkmark	14.07	27.72	61.83	15.70	20.93	31.84
O1-preview	_	✓		8.32	21.11	46.27	8.14	11.43	14.43
O3-mini	_	✓		3.84	19.83	63.54	10.47	28.68	40.78
Doubao-Seed-1.6	230B	✓	\checkmark	19.19	30.92	64.39	16.09	20.93	32.82
Doubao-Seed-1.6-flash	230B	✓	\checkmark	1.50	3.63	9.62	1.55	3.11	6.42
Doubao-Seed-1.6-thinking	230B	✓	\checkmark	15.35	23.24	60.98	16.67	20.93	30.87
Comparison of different SF	Comparison of different SFT method on Qwen-2.5-Coder								
+ SFT				26.44	30.70	48.40	14.34	15.70	18.02
+ diff-SFT	7B			22.17	22.81	34.33	7.95	9.30	12.60
+ DM-SFT				27.27	33.18	55.67	15.12	18.99	24.81

Syntax provides explicit error messages, which help models localize faulty positions, whereas Squirrel-Semantic requires reasoning about deeper semantic inconsistencies without surface-level cues.

5.2 CAN SFT SOLVE THE SQL DEBUGGING?

As detailed in Appendix C.3.1 and Figure 7, we propose 3 representative SFT approaches as baselines: (1) Vanilla SFT, which directly fine-tunes the model on parallel SQL debugging pairs; (2) DM-SFT (Duan et al., 2024), which dynamically masking the loss for unchanged tokens in responses; (3) Diff-SFT, which frames SFT as a search-and-replace task, focusing only on the modified code segments. Results in Table 2 and Figure 4 shows:

(1) Targeted in-domain SFT significantly improves SQL debugging performance. Specifically, Qwen-2.5-Coder-7B + SFT substantially outperforms the base Qwen-2.5-

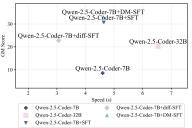


Figure 4: SFT baseline performance on Squirrel-Syntax. The horizontal axis represents the average inference speed, and the vertical axis shows the GM score.

Coder-7B, achieving a 33.17% gain in GM score on Squirrel-Syntax, and even surpasses Qwen-2.5-Coder-32B by 10.44%. (2) DM-SFT improves performance over vanilla SFT by masking the loss on non-diff tokens during training. This design forces the model to focus more on diff segments within pairs, thereby enhancing its effectiveness. (3) Diff-SFT predicts only the diff segments instead of generating the full code, offering a substantial inference speed advantage and reducing generation hallucination. On our benchmark, it requires only half the time of other methods, which is particularly

beneficial for longer code snippets in enterprise applications. However, due to a mismatch between the search-and-replace task and the pretraining/SFT objectives of the base model, its GM score is slightly lower. Overall, these three SFT strategies provide strong baselines for future research on SQL debugging. More analysis is available in Appendix D.1.

5.3 CAN AGENT METHODS SOLVE THE SQL DEBUGGING?

As described in Appendix C.3.2 and Figure 8, we adopt an **agentic framework** for SQL debugging. In this framework, a main agent analyzes error messages and plans SQL modifications, while a code generation sub-agent executes the fixes. Execution results are iteratively fed back to the main agent, forming a ReAct loop (Yao et al., 2022) that continues until the SQL executes successfully.

Figure 5 shows that agent-based systems can significantly boost performance, but results heavily depend on the main agent's capabilities. For example, using Kimi-K2 as the main agent and Qwen3-Coder as the subagent increases EM accuracy by 65% compared to the

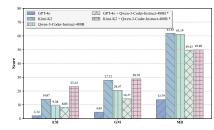


Figure 5: Agent performance on Squirrel-Syntax. '*' denotes agent-based methods, while others are single-model baselines.

Kimi-K2 single-model baseline. In contrast, when GPT-40 serves as the main agent—despite a 300%+ gain over its single-model performance—the combined system still underperforms the single Qwen3-Coder model. We also observe a decline in the MB score of agent-based systems, as multiple rounds of modification can gradually cause the model to deviate from the original SQL. These observations provide initial insights for future exploration of agentic methods in SQL debugging.

6 RELATED WORK

Code Generation and Text-to-SQL Benchmarks. Early text-to-code benchmarks, including HumanEval (Chen et al., 2021b), SQL-Spider (Yu et al., 2018), and BIRD (Li et al., 2024), focus on simple and short code snippets (Zhuo et al., 2025; Jain et al., 2025; Bytedance, 2025). To address the gap with real-world applications, SWE-Bench (Jimenez et al., 2023) evaluates models on complete software issues, which require a comprehensive understanding of codebases. Similarly, Spider2.0 (Lei et al., 2025) extends Text-to-SQL evaluation to enterprise contexts. BIRD-Critic (Li et al., 2024) introduces SQL debugging, but it only handles short, simplified StackOverflow queries that lack enterprise-level complexity. Most of these benchmarks rely on manually curated datasets, which are costly and prone to data leakage (Chou et al., 2025). In this work, we introduce the first enterprise-level SQL debugging benchmark, which is automatically constructed via reverse engineering.

LLMs for Automated Software Engineering. Recent work applies LLMs to automated software engineering through three primary paradigms: (1) **Single-model** approaches, which attempt to produce patches directly from a description and buggy code, often using few-shot prompting or SFT (Huang et al., 2024; Yasunaga & Liang, 2021; Allamanis et al., 2021). These single-model methods are bottlenecked by the need to build large-scale SFT datasets (Pan et al., 2024; Li et al., 2025b; Ma et al., 2024; Yang et al., 2025b; Pham et al., 2025). (2) **Multi-stage Workflows**, which guide models through defect localization, patch generation, and validation (Xia et al., 2024; Zhang et al., 2024). (3) **Agent-based Methods**, which leverage analysis, execution traces, or test feedback for iterative refinement (Yang et al., 2024; Wang et al., 2025; Bouzenia et al., 2024; Chen et al., 2023). In this work, we provide both SFT-based Single-model solutions and Agent-based methods, offering the community a comprehensive understanding of SQL debugging tasks.

7 Conclusion

We introduce Squirrel Benchmark, the first benchmark for enterprise-level SQL debugging. With its automated construction workflow and execution-free evaluation, Squirrel Benchmark enables scalable and reliable assessment of LLMs. Despite recent advances in LLM reasoning, our evaluation of nearly 30 models shows that real-world enterprise SQL debugging remains a significant challenge. To encourage further progress, we highlight four promising directions, including three SFT-based strategies and one agent-driven approach. Importantly, Squirrel Benchmark correlates strongly with practical debugging performance, making it a reliable reference for both academic research and industrial deployment.

CODE OF ETHICS AND ETHICS STATEMENT

Our methodology utilizes publicly accessible resources, including the LLMs and toolkits such as LLaMA-Factory and vLLM. The benchmark datasets used in our evaluation were synthetically generated using these models and are scheduled for public release upon acceptance. While a portion of our SFT data incorporates proprietary enterprise information and is therefore not fully disclosable, we recommend that researchers use our automated benchmark construction pipeline to replicate the training data. This work is centered on the English language and is strictly for research purposes.

REPRODUCIBILITY STATEMENT

To ensure reproducibility, we detail our datasets and annotation process in Section 3 and provide full experimental settings in Appendix C.

REFERENCES

Miltiadis Allamanis, Henry Jackson-Flux, and Marc Brockschmidt. Self-supervised bug detection and repair. *Advances in Neural Information Processing Systems*, 34:27865–27876, 2021.

anthropic. Claude, 2025. URL https://www.anthropic.com/claude/sonnet.

Anthropic. Introducing claude 4, 2025. URL https://www.anthropic.com/news/claude-4.

Michael Armbrust, Reynold S. Xin, Cheng Lian, Yin Huai, Davies Liu, Joseph K. Bradley, Xiangrui Meng, Tomer Kaftan, Michael J. Franklin, Ali Ghodsi, and Matei Zaharia. Spark SQL: relational data processing in spark. In Timos K. Sellis, Susan B. Davidson, and Zachary G. Ives (eds.), Proceedings of the 2015 ACM SIGMOD International Conference on Management of Data, Melbourne, Victoria, Australia, May 31 - June 4, 2015, pp. 1383–1394. ACM, 2015.

Islem Bouzenia and Michael Pradel. Understanding software engineering agents: A study of thought-action-result trajectories, 2025. URL https://arxiv.org/abs/2506.18824.

Islem Bouzenia, Premkumar Devanbu, and Michael Pradel. Repairagent: An autonomous, llm-based agent for program repair. *arXiv preprint arXiv:2403.17134*, 2024.

Bytedance. Fullstack bench: Evaluating llms as full stack coders, 2025. URL https://arxiv.org/abs/2412.00535.

Donald D. Chamberlin and Raymond F. Boyce. SEQUEL: A structured english query language. In Proceedings of the 1974 ACM SIGMOD Workshop on Data Description, Access and Control, pp. 249–264. ACM, 1974.

Dong Chen, Shaoxin Lin, Muhan Zeng, Daoguang Zan, Jian-Gang Wang, Anton Cheshkov, Jun Sun, Hao Yu, Guoliang Dong, Artem Aliev, et al. Coder: Issue resolving with multi-agent and task graphs. *arXiv preprint arXiv:2406.01304*, 2024.

Mark Chen, Jerry Tworek, Heewoo Jun, Qiming Yuan, Henrique Ponde de Oliveira Pinto, Jared Kaplan, Harri Edwards, Yuri Burda, Nicholas Joseph, Greg Brockman, Alex Ray, Raul Puri, Gretchen Krueger, Michael Petrov, Heidy Khlaaf, Girish Sastry, Pamela Mishkin, Brooke Chan, Scott Gray, Nick Ryder, Mikhail Pavlov, Alethea Power, Lukasz Kaiser, Mohammad Bavarian, Clemens Winter, Philippe Tillet, Felipe Petroski Such, Dave Cummings, Matthias Plappert, Fotios Chantzis, Elizabeth Barnes, Ariel Herbert-Voss, William Hebgen Guss, Alex Nichol, Alex Paino, Nikolas Tezak, Jie Tang, Igor Babuschkin, Suchir Balaji, Shantanu Jain, William Saunders, Christopher Hesse, Andrew N. Carr, Jan Leike, Josh Achiam, Vedant Misra, Evan Morikawa, Alec Radford, Matthew Knight, Miles Brundage, Mira Murati, Katie Mayer, Peter Welinder, Bob McGrew, Dario Amodei, Sam McCandlish, Ilya Sutskever, and Wojciech Zaremba. Evaluating large language models trained on code, 2021a. URL https://arxiv.org/abs/2107.03374.

- Mark Chen, Jerry Tworek, Heewoo Jun, Qiming Yuan, Henrique Ponde De Oliveira Pinto, Jared Kaplan, Harri Edwards, Yuri Burda, Nicholas Joseph, Greg Brockman, et al. Evaluating large language models trained on code. *arXiv preprint arXiv:2107.03374*, 2021b.
- Xinyun Chen, Maxwell Lin, Nathanael Schärli, and Denny Zhou. Teaching large language models to self-debug. *arXiv preprint arXiv:2304.05128*, 2023.
- Jason Chou, Ao Liu, Yuchi Deng, Zhiying Zeng, Tao Zhang, Haotian Zhu, Jianwei Cai, Yue Mao, Chenchen Zhang, Lingyun Tan, Ziyan Xu, Bohui Zhai, Hengyi Liu, Speed Zhu, Wiggin Zhou, and Fengzong Lian. Autocodebench: Large language models are automatic code benchmark generators, 2025. URL https://arxiv.org/abs/2508.09101.
- DeepSeek. Deepseek-v3.1, 2025. URL https://api-docs.deepseek.com/news/news250821.
- DeepSeek-AI. Deepseek-v3 technical report, 2025a. URL https://arxiv.org/abs/2412.19437.
- DeepSeek-AI. Deepseek-r1: Incentivizing reasoning capability in llms via reinforcement learning, 2025b. URL https://arxiv.org/abs/2501.12948.
- Yiwen Duan, Yonghong Yu, Xiaoming Zhao, Yichang Wu, and Wenbo Liu. Pdc dm-sft: A road for llm sql bug-fix enhancing, 2024. URL https://arxiv.org/abs/2411.06767.
- Gemini. Gemini 2.5: Pushing the frontier with advanced reasoning, multimodality, long context, and next generation agentic capabilities, 2025. URL https://arxiv.org/abs/2507.06261.
- Kai Huang, Xiangxin Meng, Jian Zhang, Yang Liu, Wenjie Wang, Shuhao Li, and Yuqing Zhang. An empirical study on fine-tuning large language models of code for automated program repair. In *Proceedings of the 38th IEEE/ACM International Conference on Automated Software Engineering*, ASE '23, pp. 1162–1174. IEEE Press, 2024. ISBN 9798350329964. doi: 10.1109/ASE56229. 2023.00181. URL https://doi.org/10.1109/ASE56229.2023.00181.
- Binyuan Hui, Jian Yang, Zeyu Cui, Jiaxi Yang, Dayiheng Liu, Lei Zhang, Tianyu Liu, Jiajun Zhang, Bowen Yu, Keming Lu, Kai Dang, Yang Fan, Yichang Zhang, An Yang, Rui Men, Fei Huang, Bo Zheng, Yibo Miao, Shanghaoran Quan, Yunlong Feng, Xingzhang Ren, Xuancheng Ren, Jingren Zhou, and Junyang Lin. Qwen2.5-coder technical report, 2024. URL https://arxiv.org/abs/2409.12186.
- Naman Jain, King Han, Alex Gu, Wen-Ding Li, Fanjia Yan, Tianjun Zhang, Sida Wang, Armando Solar-Lezama, Koushik Sen, and Ion Stoica. Livecodebench: Holistic and contamination free evaluation of large language models for code. In *The Thirteenth International Conference on Learning Representations*, 2025. URL https://openreview.net/forum?id=chfJJYC3iL.
- Carlos E Jimenez, John Yang, Alexander Wettig, Shunyu Yao, Kexin Pei, Ofir Press, and Karthik R Narasimhan. Swe-bench: Can language models resolve real-world github issues? In *The Twelfth International Conference on Learning Representations*, 2023.
- Carlos E Jimenez, John Yang, Alexander Wettig, Shunyu Yao, Kexin Pei, Ofir Press, and Karthik R Narasimhan. SWE-bench: Can language models resolve real-world github issues? In *The Twelfth International Conference on Learning Representations*, 2024.
- Kimi-Team. Kimi k2: Open agentic intelligence, 2025. URL https://arxiv.org/abs/2507. 20534.
- Woosuk Kwon, Zhuohan Li, Siyuan Zhuang, Ying Sheng, Lianmin Zheng, Cody Hao Yu, Joseph E. Gonzalez, Hao Zhang, and Ion Stoica. Efficient memory management for large language model serving with pagedattention. In *Proceedings of the ACM SIGOPS 29th Symposium on Operating Systems Principles*, 2023.
- Philippe Laban, Hiroaki Hayashi, Yingbo Zhou, and Jennifer Neville. Llms get lost in multi-turn conversation, 2025. URL https://arxiv.org/abs/2505.06120.

```
Fangyu Lei, Jixuan Chen, Yuxiao Ye, Ruisheng Cao, Dongchan Shin, Hongjin SU, ZHAOQING SUO, Hongcheng Gao, Wenjing Hu, Pengcheng Yin, Victor Zhong, Caiming Xiong, Ruoxi Sun, Qian Liu, Sida Wang, and Tao Yu. Spider 2.0: Evaluating language models on real-world enterprise text-to-SQL workflows. In The Thirteenth International Conference on Learning Representations, 2025.
```

- Haoyang Li, Shang Wu, Xiaokang Zhang, Xinmei Huang, Jing Zhang, Fuxin Jiang, Shuai Wang, Tieying Zhang, Jianjun Chen, Rui Shi, et al. Omnisql: Synthesizing high-quality text-to-sql data at scale. *arXiv preprint arXiv:2503.02240*, 2025a.
- Jinyang Li, Binyuan Hui, Ge Qu, Jiaxi Yang, Binhua Li, Bowen Li, Bailin Wang, Bowen Qin, Ruiying Geng, Nan Huo, et al. Can Ilm already serve as a database interface? a big bench for large-scale database grounded text-to-sqls. *Advances in Neural Information Processing Systems*, 36, 2024.
- Jinyang Li, Xiaolong Li, Ge Qu, Per Jacobsson, Bowen Qin, Binyuan Hui, Shuzheng Si, Nan Huo, Xiaohan Xu, Yue Zhang, Ziwei Tang, Yuanshuai Li, Florensia Widjaja, Xintong Zhu, Feige Zhou, Yongfeng Huang, Yannis Papakonstantinou, Fatma Ozcan, Chenhao Ma, and Reynold Cheng. Swe-sql: Illuminating llm pathways to solve user sql issues in real-world applications, 2025b. URL https://arxiv.org/abs/2506.18951.
- Yingwei Ma, Rongyu Cao, Yongchang Cao, Yue Zhang, Jue Chen, Yibo Liu, Yuchen Liu, Binhua Li, Fei Huang, and Yongbin Li. Lingma swe-gpt: An open development-process-centric language model for automated software improvement, 2024. URL https://arxiv.org/abs/2411.00622.
- OpenAI. Hello GPT-40, 2024. URL https://openai.com/index/hello-gpt-4o/.
- OpenAI. Introducing gpt-4.1 in the api, 2025. URL https://openai.com/index/gpt-4-1/.
- Openai. Introducing gpt-5, 2025. URL https://openai.com/index/introducing-gpt-5/.
- OpenAI. Introducing openai o3 and o4-mini, 2025. URL https://openai.com/index/introducing-o3-and-o4-mini/.
- Jiayi Pan, Xingyao Wang, Graham Neubig, Navdeep Jaitly, Heng Ji, Alane Suhr, and Yizhe Zhang. Training software engineering agents and verifiers with swe-gym, 2024. URL https://arxiv.org/abs/2412.21139.
- Minh V. T. Pham, Huy N. Phan, Hoang N. Phan, Cuong Le Chi, Tien N. Nguyen, and Nghi D. Q. Bui. Swe-synth: Synthesizing verifiable bug-fix data to enable large language models in resolving real-world bugs, 2025. URL https://arxiv.org/abs/2504.14757.
- Qwen. Qwen3-coder: Agentic coding in the world, 2025. URL https://qwenlm.github.io/blog/qwen3-coder/.
- Seed. Introduction to techniques used in seed1.6, 2025. URL https://seed.bytedance.com/en/seed1_6.
- ByteDance Seed, Yuyu Zhang, Jing Su, Yifan Sun, Chenguang Xi, Xia Xiao, Shen Zheng, Anxiang Zhang, Kaibo Liu, Daoguang Zan, Tao Sun, Jinhua Zhu, Shulin Xin, Dong Huang, Yetao Bai, Lixin Dong, Chao Li, Jianchong Chen, Hanzhi Zhou, Yifan Huang, Guanghan Ning, Xierui Song, Jiaze Chen, Siyao Liu, Kai Shen, Liang Xiang, and Yonghui Wu. Seed-coder: Let the code model curate data for itself, 2025. URL https://arxiv.org/abs/2506.03524.
- Xingyao Wang, Boxuan Li, Yufan Song, Frank F. Xu, Xiangru Tang, Mingchen Zhuge, Jiayi Pan, Yueqi Song, Bowen Li, Jaskirat Singh, Hoang H. Tran, Fuqiang Li, Ren Ma, Mingzhang Zheng, Bill Qian, Yanjun Shao, Niklas Muennighoff, Yizhe Zhang, Binyuan Hui, Junyang Lin, Robert Brennan, Hao Peng, Heng Ji, and Graham Neubig. Openhands: An open platform for AI software developers as generalist agents. In *The Thirteenth International Conference on Learning Representations*, 2025.

- Chunqiu Steven Xia, Yinlin Deng, Soren Dunn, and Lingming Zhang. Agentless: Demystifying llm-based software engineering agents. *arXiv preprint arXiv:2407.01489*, 2024.
- An Yang, Anfeng Li, Baosong Yang, Beichen Zhang, Binyuan Hui, Bo Zheng, Bowen Yu, Chang Gao, Chengen Huang, Chenxu Lv, Chujie Zheng, Dayiheng Liu, Fan Zhou, Fei Huang, Feng Hu, Hao Ge, Haoran Wei, Huan Lin, Jialong Tang, Jian Yang, Jianhong Tu, Jianwei Zhang, Jianxin Yang, Jiaxi Yang, Jing Zhou, Jingren Zhou, Junyang Lin, Kai Dang, Keqin Bao, Kexin Yang, Le Yu, Lianghao Deng, Mei Li, Mingfeng Xue, Mingze Li, Pei Zhang, Peng Wang, Qin Zhu, Rui Men, Ruize Gao, Shixuan Liu, Shuang Luo, Tianhao Li, Tianyi Tang, Wenbiao Yin, Xingzhang Ren, Xinyu Wang, Xinyu Zhang, Xuancheng Ren, Yang Fan, Yang Su, Yichang Zhang, Yinger Zhang, Yu Wan, Yuqiong Liu, Zekun Wang, Zeyu Cui, Zhenru Zhang, Zhipeng Zhou, and Zihan Qiu. Qwen3 technical report, 2025a. URL https://arxiv.org/abs/2505.09388.
- John Yang, Carlos E. Jimenez, Alexander Wettig, Kilian Lieret, Shunyu Yao, Karthik Narasimhan, and Ofir Press. Swe-agent: Agent computer interfaces enable software engineering language models, 2024.
- John Yang, Kilian Lieret, Carlos E. Jimenez, Alexander Wettig, Kabir Khandpur, Yanzhe Zhang, Binyuan Hui, Ofir Press, Ludwig Schmidt, and Diyi Yang. Swe-smith: Scaling data for software engineering agents, 2025b. URL https://arxiv.org/abs/2504.21798.
- Shunyu Yao, Jeffrey Zhao, Dian Yu, Nan Du, Izhak Shafran, Karthik R Narasimhan, and Yuan Cao. React: Synergizing reasoning and acting in language models. In *The Eleventh International Conference on Learning Representations*, 2022.
- Michihiro Yasunaga and Percy Liang. Break-it-fix-it: Unsupervised learning for program repair. In *International conference on machine learning*, pp. 11941–11952. PMLR, 2021.
- 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. In *Proceedings of the 2018 Conference on Empirical Methods in Natural Language Processing*, pp. 3911–3921, 2018.
- Yi Zhan, Longjie Cui, Han Weng, Guifeng Wang, Yu Tian, Boyi Liu, Yingxiang Yang, Xiaoming Yin, Jiajun Xie, and Yang Sun. Towards database-free text-to-SQL evaluation: A graph-based metric for functional correctness. In Owen Rambow, Leo Wanner, Marianna Apidianaki, Hend Al-Khalifa, Barbara Di Eugenio, and Steven Schockaert (eds.), *Proceedings of the 31st International Conference on Computational Linguistics*, pp. 4586–4610, Abu Dhabi, UAE, January 2025. Association for Computational Linguistics. URL https://aclanthology.org/2025.coling-main.308/.
- Yuntong Zhang, Haifeng Ruan, Zhiyu Fan, and Abhik Roychoudhury. Autocoderover: Autonomous program improvement. In *Proceedings of the 33rd ACM SIGSOFT International Symposium on Software Testing and Analysis*, pp. 1592–1604, 2024.
- Yaowei Zheng, Richong Zhang, Junhao Zhang, Yanhan Ye, Zheyan Luo, Zhangchi Feng, and Yongqiang Ma. Llamafactory: Unified efficient fine-tuning of 100+ language models. In *Proceedings of the 62nd Annual Meeting of the Association for Computational Linguistics (Volume 3: System Demonstrations)*, Bangkok, Thailand, 2024. Association for Computational Linguistics. URL http://arxiv.org/abs/2403.13372.
- Victor Zhong, Caiming Xiong, and Richard Socher. Seq2sql: Generating structured queries from natural language using reinforcement learning. *arXiv* preprint arXiv:1709.00103, 2017.
- Terry Yue Zhuo, Vu Minh Chien, Jenny Chim, Han Hu, Wenhao Yu, Ratnadira Widyasari, Imam Nur Bani Yusuf, Haolan Zhan, Junda He, Indraneil Paul, Simon Brunner, Chen GONG, James Hoang, Armel Randy Zebaze, Xiaoheng Hong, Wen-Ding Li, Jean Kaddour, Ming Xu, Zhihan Zhang, Prateek Yadav, Naman Jain, Alex Gu, Zhoujun Cheng, Jiawei Liu, Qian Liu, Zijian Wang, David Lo, Binyuan Hui, Niklas Muennighoff, Daniel Fried, Xiaoning Du, Harm de Vries, and Leandro Von Werra. Bigcodebench: Benchmarking code generation with diverse function calls and complex instructions. In *The Thirteenth International Conference on Learning Representations*, 2025. URL https://openreview.net/forum?id=YrycTjllL0.

A USE OF LLMS

We claim the following regarding the use of LLMs in this work: (1) Claude-4-Sonnet was used for data construction; details are provided in Section 3, with prompts listed in Appendix G. (2) LLMs were employed for evaluation on the benchmark introduced in this paper; the specific models are listed in Appendix C.2. (3) LLMs were used during manuscript preparation solely for text polishing and refinement. (4) We used Cursor for programming assistance; however, all code was manually reviewed. We further claim that all core ideas and intellectual contributions were developed exclusively by the authors, without input from any LLM.

B LIMITATIONS

This work introduces a benchmark for enterprise SQL debugging, supporting future software engineering research. However, the study has several limitations. *First, the benchmark's synthetic nature*. Although the dataset was generated automatically, we manually cross-validated items that all models failed. Nonetheless, the potential for undetected artifacts remains. Developing more sophisticated automated validation techniques is a key future direction. *Second, constraints of the evaluation framework*. Our rule-based, execution-free method uses exact match, graph match, and an edit-direction criterion. This approach is effective for debugging contexts where minimal, correct fixes are expected. However, it is inherently limited by its dependence on reference solutions. For semantic tasks where solution code is diverse, a more flexible evaluation strategy is needed. We acknowledge this as an area for future improvement.

C EXPERIMENTAL SETTINGS

C.1 EVALUATION METRICS

The evaluation of enterprise SQL debugging systems faces distinct challenges: the process is inherently complex, often admitting multiple valid solutions, while practical constraints preclude execution-based assessment in production environments. An effective framework must therefore (i) circumvent the prohibitive costs and privacy concerns of query execution; (ii) evaluate structural and semantic correctness beyond string-level similarity; and (iii) accommodate the legitimate ambiguity of real-world repairs. To address these requirements, we introduce an execution-free evaluation methodology based on three complementary metrics.

Exact Match Score (EM). This metric provides a strict, reproducible measure of syntactic correctness by comparing the predicted SQL string directly against the reference:

$$EM = \frac{1}{N} \sum_{i=1}^{N} \mathbf{1}[\hat{q}_i = q_i]$$
 (5)

where \hat{q}_i is the predicted SQL, q_i is the reference SQL, and $\mathbf{1}[\cdot]$ is the indicator function. While stringent, EM serves as a clear lower bound on model performance.

Graph-Match Score (GM). To overcome the limitations of string-based comparison, we assess semantic equivalence through code structure. Each SQL query is parsed into an abstract syntax tree (AST) and converted into a directed graph. The GM score is then computed based on graph isomorphism:

$$GM = \frac{1}{N} \sum_{i=1}^{N} \mathbf{1}[Graph(\hat{q}_i) \cong Graph(q_i)]$$
 (6)

where \cong denotes graph isomorphism. This approach recognizes semantically equivalent codes that may differ syntactically.

Modify-Better Score (MB). For iterative debugging scenarios, absolute correctness is insufficient; we must measure progressive improvement. The MB metric evaluates whether a prediction moves closer to the correct solution by comparing AST edit distances:

$$MB = \frac{1}{N} \sum_{i=1}^{N} \mathbf{1}[d(\hat{q}_i, q_i) < d(b_i, q_i)]$$
 (7)

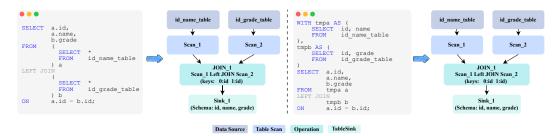


Figure 6: Illustration of Graph Match Score. Although the left and right SQL snippets differ syntactically, their optimized abstract syntax trees are structurally identical. Graph matching evaluates semantic equivalence through tree isomorphism.

where $d(\cdot, \cdot)$ denotes normalized AST edit distance, \hat{q}_i is the predicted repair, q_i is the reference SQL, and b_i is the original buggy query. This metric specifically assesses a model's capacity for incremental repair in debugging workflows.

Together, these metrics provide a comprehensive evaluation framework that balances efficiency, reproducibility, and semantic understanding while avoiding the practical limitations of execution-based assessment.

C.2 LLMs

 This study ensures a robust evaluation by leveraging a diverse set of large language models (LLMs), encompassing both open-source and proprietary architectures to cover a broad range of capabilities. The evaluated models are as follows:

Open-Source Models

- DeepSeek Series: DeepSeek-R1-0528 (DeepSeek-AI, 2025b), DeepSeek-V3-0324 (DeepSeek-AI, 2025a), Deepseek-V3.1
- **Qwen Series:** Qwen-2.5-Instruct, Qwen-2.5-Coder (Hui et al., 2024), Qwen-3-235B-A22B-Instruct-2507 (Yang et al., 2025a), Qwen-3-Coder-480B-A35B-Instruct (Qwen, 2025), QwQ-32B
- Specialized Code Models: Seed-Coder-8B (Seed et al., 2025), OmniSQL-32B (Li et al., 2025a)

Close-Source Models

- Anthropic: Claude-Sonnet-4 (Anthropic, 2025)
- **OpenAI:** GPT-4o-mini-2024-07-18, GPT-4o-2024-11-20 (OpenAI, 2024), GPT-4.1 (OpenAI, 2025), o3-mini (OpenAI, 2025), o1-Preview, GPT-5
- Google: Gemini 2.5 Pro (Gemini, 2025)
- Moonshot AI: Kimi-K2 (Kimi-Team, 2025)
- **ByteDance:** Doubao family (Doubao-Seed-1.6, Doubao-Seed-1.6-flash, Doubao-Seed-1.6-thinking) (Seed, 2025)

C.3 BASELINES

C.3.1 SFT BASELINES

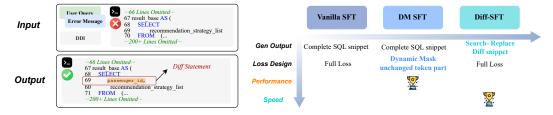


Figure 7: Illustration of three distinct supervised fine-tuning (SFT) methods.

We propose three distinct supervised fine-tuning (SFT) methods as baselines.

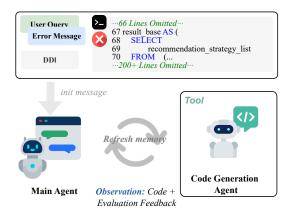


Figure 8: Overview of the agentic method. The *main agent* observes the error message and the issued SQL, analyzes the cause of the failure and the required modification sketch, and outputs an action in the form of code modification instructions. These instructions are executed by a *code generation sub-agent*, and the resulting execution feedback is used to update the main agent's memory. This iterative ReAct loop continues until a termination condition is met.

Vanilla SFT. This is the standard sequence-to-sequence fine-tuning approach. The model takes as input the error message, DDL, and the issue SQL (the buggy query), and is trained to generate the complete, corrected reference SQL. While simple, this method establishes a fundamental baseline for performance.

DM-SFT (**Dynamic-Masked SFT**). In enterprise SQL debugging, the differences between an issue SQL and its reference SQL are often minimal within lengthy code snippets. Consequently, Vanilla SFT models can rapidly reduce loss by learning to copy the large, unchanged portions of the input, potentially failing to focus on the critical, erroneous segments. To mitigate this, we adopt Dynamic-Masked SFT (DM-SFT) (Duan et al., 2024), which randomly masks the loss calculation for 50% of the tokens that are identical between the input and output. By increasing the loss contribution of the changed tokens, this method encourages the model to prioritize learning the necessary edits.

Diff-SFT. Generating the complete SQL code increases inference overhead significantly. We propose an alternative method where the model only outputs a "diff" snippet, framing the task as a search-and-replace operation. The model's objective is to identify the erroneous code segment in the input and generate the corresponding corrected snippet.

C.3.2 AGENT BASELINES

As illustrated in Figure 8, we design an **agentic framework** to address SQL debugging by explicitly structuring the process into observation, analysis, and action phases (Yao et al., 2022). The *main agent* receives the error message together with the issued SQL as observations, analyzes the root cause of the failure and identifies the sketch of the required fix, and then outputs an action in the form of modification instructions. These instructions are passed to a *code generation sub-agent*, which produces the concrete SQL patch. The newly generated SQL is executed, and its results—either successful outputs or subsequent error messages—are fed back to update the main agent's memory. Through this iterative loop of error observation, analysis, and guided code generation, the framework incrementally refines the SQL until a correct solution is obtained.

C.4 DATASET

Our training dataset comprises three parts:

- Data synthesized from Text-to-SQL datasets BIRD (Li et al., 2024) and Spider (Yu et al., 2018), containing 1, 054 samples.
- Data constructed via reverse engineering, where bugs were manually injected into correct SQL code, containing 2,015 samples.

• Data mined from online execution logs, containing error messages and the corresponding issue SQL, with the reference SQL manually written and verified by experts, containing 1,971 samples.

C.5 HYPERPARAMETERS

Fine-tuning. For self-supervised fine-tuning, models are trained for 5 epoch with a learning rate of 1e-5 and a per device batch size of 64. We employed the AdamW optimizer and a cosine learning rate scheduler with a warm-up phase corresponding to 3% of the total training steps.

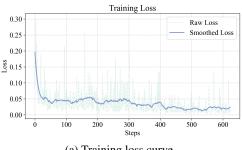
Evaluation. We use the Pass@1 (Chen et al., 2021a) as the default evaluation metric. Following Yang et al. (2024); Chen et al. (2024), we use a temperature of 0.0 for deterministic action decoding and input truncation to manage context length.

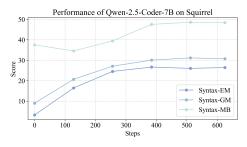
C.6 EXPERIMENTAL ENVIRONMENTS

All experiments are conducted on 32 NVIDIA H20 GPUs. Our code primarily relies on Python 3.12 and PyTorch 2.7.0. Models are self-supervised fine-tuned with LLaMA-Factory (Zheng et al., 2024)¹, and inference is performed with vLLM (Kwon et al., 2023).

D ADDITIONAL EXPERIMENTAL RESULTS

D.1 ADDITIONAL ANALYSIS OF SFT PERFORMANCE ON SQUIRREL BENCHMARK





(a) Training loss curve. (b)

(b) Performance at different training steps.

Figure 9: Analysis of Qwen-2.5-Coder-7B Vanilla SFT on Squirrel Benchmark, showing corresponding training loss and step-wise performance.

Rapid loss decay in SQL debugging fine-tuning. Figure 9a illustrates that the training loss quickly drops below 0.05 within a few steps, approaching zero. This behavior arises because the constructed SQL debugging parallel data contain inputs with error messages and issue SQL statements, and outputs with the corrected SQL. In most cases, only a small portion of tokens differ between the input and output. Consequently, the model primarily copies tokens from the input, leading to extremely low training loss. When the majority of output tokens carry minimal information, the model tends to ignore the truly informative segments that require correction.

Performance improves with increased training steps. Figure 9b shows that as training progresses, model performance steadily improves, particularly during the early steps. Beyond approximately 400 training steps, the gains become marginal, indicating diminishing returns. This trend suggests that while additional in-domain training helps, the benefit of further fine-tuning eventually saturates.

¹https://github.com/hiyouga/LLaMA-Factory.git

D.2 CASE STUDY

org.apache.calcite.sql.parser.SqlParseException: Encountered "AS" at line 14, column 54. Error Message: Was expecting one of: ")" ..."MULTISET" ... "ARRAY" ... Predict SQL Preference SQL 21 SELECT attorney_id, 21 SELECT attorney_id, case_type_id, case_type_id, 23-CAST(consultation_revenue_7d * 100 23+ CAST(consultation_revenue_7d * 100 AS AS BIGINT) AS consultation_revenue), BIGINT) AS consultation_revenue, consultation_bookings_7d AS consultation_bookings_7d AS consultation_bookings, consultation_bookings,

Figure 10: Model Hallucination: After modifying the code according to the error message, the model also inserted an extra ")" in similar fragments, which caused the fix to fail.

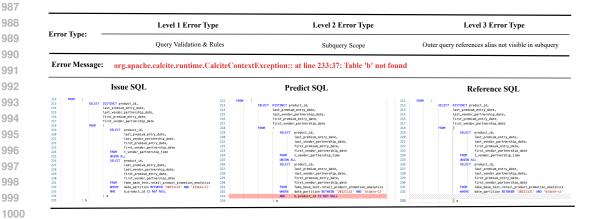


Figure 11: Long Context Reasoning Limitation: The error code uses a non-existent table b (which is usually an alias for a longer table name in SQL), but the model fail to detect this error during the repair process.

E SQL BUG TAXONOMY

E.1 BUG DISTRIBUTION OF SQUIRREL-SEMANTIC

Table 3: Error type distribution in Squirrel-Semantic

Level 1 Error Type	Level 2 Error Type	Level 3 Error Type			
	Aggregate Logic	Using COUNT(column) instead of COUNT(*) and misunderstanding NULL exclusion	43		
Semantics & Logic		Using SUM()/AVG() on a column with NULLs without COALESCE	27		
	Join Logic	JOIN condition placed in WHERE clause (accidental CROSS JOIN)	41		
		Failing to handle NULLs in JOIN keys (causing rows to disappear)	13		
		Missing condition causing Cartesian product	2		
	Boolean & Logic	Three-valued logic error: NOT (a = b) not equivalent to a != b when NULLs present	14		
		Improper Boolean usage (e.g., WHERE col = TRUE)	9		
	NULL Handling	NULL compared with = (should use IS NULL)	33		
		Confusion between IS NULL and =NULL	2		
	Window Function Logic	Using RANK() instead of ROW_NUMBER() or DENSE_RANK() leading to duplicates/skips	31		
		Incorrect partitioning/ordering in window function leading to wrong row assignment	4		
	Subquery Scope	Misplaced LIMIT inside subquery affecting outer results	3		
		Correlated subquery missing correlation condition	2		
	JOIN Logic	Missing condition causing Cartesian product	12		
		Wrong join key used inside nested subquery	2		
	Set Operations	UNION vs. UNION ALL misuse (unintended deduplication)	55		
	Date/Time Logic	Confusion between DATE, TIMESTAMP, and INTERVAL types	23		
	Pattern Matching	Incorrect LIKE usage	2		
E 0.E	Separator Rule	collect_set/concat_ws separator uses semicolon	54		
Functions & Expressions	Function Semantics	Misunderstanding the empty handling of aggregate functions	1		
	an over never	Misuse of ROLLUP / CUBE	14		
	GROUP BY Extensions	Rollup/Cube/Grouping Sets producing unexpected super-aggregate rows	3		
Joins & Grouping	anavin nyvy	Grouping by a functionally dependent column unnecessarily	17		
	GROUP BY Logic	Rollup/Cube/Grouping Sets producing unexpected super-aggregate rows	4		
	JOIN Type Selection	Using INNER JOIN when LEFT JOIN is needed (loss of data)	64		
Result & Quality		Duplicate rows due to many-to-many join not being accounted for	1		
	Correctness	Incorrect output data	1		
	Implicit Casting	Implicit cast changing semantics (e.g., string to number)	15		
Types & Data Formats	Data Format	Misused format placeholder	1		
Identifiers & Objects	Qualification	Qualifying a column with the wrong table alias in a complex join	22		

E.2 Bug Distribution of Squirrel-Syntax

Table 4: Error type distribution in Squirrel-Syntax

Level 1 Error Type	Level 2 Error Type	Level 3 Error Type	Coun	
		Missing parameter for explode	15	
Functions & Expressions	Parameter Completeness	Incorrect explode parameter	9	
		explode(map) requires two aliases	2	
		date_add missing parameter (also typo data_add)	2	
		array_contains wrong argument type	1	
	Parameter Type	get_json_object wrong argument type	4	
		array_contains wrong argument type	3	
		from_json wrong argument type	1	
	LATERAL VIEW Required	Missing LATERAL VIEW	94	
	EATERAE VIEW Required	Missing alias for LATERAL VIEW function output	1	
	Date Difference	datadiff argument/typo error	5	
	Type Conversion	Multiple AS in CAST	15	
	Nesting Limit	Aggregate expressions cannot be nested	2	
	Separator Rule	collect_set/concat_ws separator uses semicolon	11	
	Date/Time	to_unix_timestap typo	1	
	Function Spelling	concat_ws typo	1	
	CASE E	Missing END or THRN in CASE WHEN	72	
Query Validation & Rules	CASE Expression	Multiple END in CASE WHEN	4	
	a	Missing argument in IN	7	
	Conditional Logic	IN subquery returns multiple columns	1	
		Window function misused with GROUP BY	3	
	Window Functions	Window function used inside WHERE/HAVING	3	
		Window function frame clause misuse (e.g., ROWS BETWEEN error)	1	
	Subquery Scope	Outer query references alias not visible in subquery	2	
	Aggregation & Subquery	SELECT list contains non-aggregated column not in GROUP BY	62	
	Pattern Matching	Incorrect LIKE usage	2	
		5		
	Aggregate Usage	Aggregate function in SELECT without GROUP BY	1	
	Boolean & NULL	NULL compared with = (should use IS NULL)	2	
	Clause Structure	Incorrect clause ordering - JOIN after WHERE	7	
		Invalid SELECT clause syntax with subquery	6	
		Missing SELECT before FROM clause	5	
		Multiple WHERE	4	
		Missing partition conditions in WHERE clause	3	
		Missing logical connector in WHERE	26	
		Non-query expression in illegal context	3	
		Missing FROM clause	2	
		Column count mismatch in UNION	1	
Grammar & Structure		WITH AS not first	26	
	CTE/View	Unnecessary WITH AS	13	
		Trailing comma after last view	23	
		Keyword spelling error	3	
	Keywords & Operators	Space in !=	2	
	Keywords & Operators	Missing IN keyword	2	
	Statement Ending	Extra trailing statements	4	
	Parentheses / Brackets		5	
		Missing closing parenthesis		
	Alias / AS	Redundant AS	3	
	SELECT List	Missing column list after SELECT	1	
		Variable error	13	
	Variables/Placeholders Ambiguous References	Missing partition conditions in DELETE statement	2	
		Partition column comparison with numeric type not allowed	2	
Identifiers & Objects		Column exists in multiple tables but alias omitted	8	
ruentiners & Objects		Ambiguous alias in nested subquery with same column name	1	
	Sahama/Ohiast	Field/Table does not exist	11	
	Schema/Object	Missing partition query conditions	2	
	Naming/Alias	Duplicate names (column/alias)	5	
		Missing grouping column	14	
	GROUP BY	Missing HAVING clause for aggregate filtering	1	
Joins & Grouping		Missing condition causing Cartesian product	6	
Joins & Grouping	JOIN Ambiguity	Missing table prefix for duplicate column names in join	35	
	Wissing table prenx for duplicate column names in Joh	Ambiguous column reference due to multiple levels of alias	1	
	Nested Joins	Punctuation error	49	
Dunatuation & Farmatt'	Punatuation/Passathassa			
unctuation & Formatting	Punctuation/Parentheses	Incorrect quote type for column alias with special characters	4	
		Missing semicolon between statements	5	
	Insert Statement	Insert error	37	
DML & DDL		Mismatched column count	5	
	Create Table Statement	Table creation error	10	
Compatibility/Dialact	Function Differences	TRANSFORM with lambda expression not supported in Hive	3	
Compatibility/Dialect	Function Differences	wm_concat function not supported in the current SQL dialect	1	
Types & Data Formats	Type System	Type mismatch	16	
		to_unix_timestap typo	2	

F EXAMPLES

1134

11351136

11371138

1139

1140

1141

1148

1149

1150

1152

1153

1154

115511561157

1158

11591160

1161 1162 1163

1169

1170

11711172

1173

1174

1175

1176

1177

11781179

F.1 SQUIRREL-SYNTAX EXAMPLE

```
org.apache.calcite.runtime.CalciteContextException:: from line 139, column 37 to line 139, column 48: Column 'passenger' id' not found in table
  'result_base'
                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                             [omit 50 Lines]
retention pipeline AS (
SELECT encrypted contact,
SELECT encrypted contact,
SELECT encrypted contact,
passenger id, teller,
passenger id, teller,
pourney score, le id,
passenger id collection,
device ld collection,
device ld collection,
reflection reason,
recommendation strategy list,
retention pipeline' AS data source
NHEEE is_retained = '1'
}
                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                   correct SQL
                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                               Issue SQL
                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                              ),
result base AS (
SELECT
                      ),
result_base AS (
                                                                                                                                    SELECT passenger id,
SAX(journey_score) AS journey_score,
COLLECT_SET(data_source) AS recommendation_strategy_list
                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                              SELECT passenger_id,
MAX(fourney score) AS journey score,
COLLECT_SET(data_source) AS recommendation_strategy_list
                                                                                                                                                                                                                        LLECT SET(data source) As recommendation

SELECT encrypted contact,
device identifier,
recommendation timestamp,
journey.score,
traveler_profile_id,
passenger id collection,
device_id_rollection,
rejection_reson,
recommendation_strategy_list,
data source
                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                   SELECT encrypted contact, device identifier, processes, contact, device identifier, processes, contact, device identifier, processes, contact, device identifier, processes, contact, device in the contact in the co
                                                                                                                                                                                                                                                                                                 data_source
recommendation_model
                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                           data_source
recommendation_model
                                                                                                                                                                                                                                                                                                   L
encrypted_contact,
device identifier,
                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                encrypted_contact,
device identifier,
                                                                                                                                                                                                                                                                                              device_identifier,
passenger id,
journey_score,
traveler profile id,
passenger id collection,
device_id_collection,
is_selected,
                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                           device identifier, passenger_id, recommendation_timestamp, journey_score, traveler_profile_id, passenger_id_collection, passenger_id_collection, is_selected_research_research_research_research_research_research_research_research_research_research_research_research_research_research_research_research_research_research_research_research_research_research_research_research_research_research_research_research_research_research_research_research_research_research_research_research_research_research_research_research_research_research_research_research_research_research_research_research_research_research_research_research_research_research_research_research_research_research_research_research_research_research_research_research_research_research_research_research_research_research_research_research_research_research_research_research_research_research_research_research_research_research_research_research_research_research_research_research_research_research_research_research_research_research_research_research_research_research_research_research_research_research_research_research_research_research_research_research_research_research_research_research_research_research_research_research_research_research_research_research_research_research_research_research_research_research_research_research_research_research_research_research_research_research_research_research_research_research_research_research_research_research_research_research_research_research_research_research_research_research_research_research_research_research_research_research_research_research_research_research_research_research_research_research_research_research_research_research_research_research_research_research_research_research_research_research_research_research_research_research_research_research_research_research_research_research_research_research_research_research_research_research_research_research_research_research_research_research_research_research_research_research_research_research_research_research_research_research_research_research
                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                 FROM retention_pipeline
)
GROUP BY
passenger_id
                                                                                                                                    GROUP BY passenger_id
                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                  ORDER BY

ORDER BY

LIMIT JOURNEY, SCORE DESC

INTERIOR SELECT

SELECT encrypted contact,
device identifier,
percommendation timestamp,
journey, score,
passenger id collection,
device id collection,
device id collection,
contact user,
recommendation strategy_list
                   passen,

JOHN TO PASSEN,

JOURNEY SCORE DESC

LIMIT 200000 

), result AS (

SELECT eners
                                                                                                      20000 S (
encrypted contact, device identifier, device identifier, device interest of the contact of the contac
                                                                                                                                    Commendation_strategy_list

SELECT base.encrypted contact,
result base.reCommendation strategy_list,
base.device_identifier,
base.passenger_id,
base.passenger_id,
base.passenger_id,
base.passenger_id collection,
base.passenger_id collection,
base.device_id_collection

FROM
DATE to base.passenger_id = base.passenger_id

ON result_base_passenger_id = base.passenger_id
                        experiment mapping AS (
SELECT passenger id,
                                                                                                                                      CAST (
conv(
                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                 ubstr(
	md5(CAST(passenger_id AS STRING)),
	-15
                                                                                                        -15
),
16,
10
) AS BIGINT
) % 2 AS BIGINT
) AS experiment_group_label
(
                                                                                                                                         SELECT passenger_id
FROM result
GROUP BY
passenger_id
                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                   GROUP BY
passenger_id
                        )
INSERT OVERWRITE TABLE fake base test.passenger_journey_recommendations
PARTITION(processing_date = '${date}', model_version = 'v2')
SELECT result.',
                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                              | NESERT OVERWRITE TABLE fake base test-passenger journey recommendations PARTITION(processing date - "$(date)', model version - "v2') SELECT EXPERIENCE TABLE (SELECT EXPERIENT EMPORTMENT (ADDITION OF TABLE ) AS active_period result ("$[MTE-2]", '-', '5(DXTE-8)") As active_period ("$[MTE-2]", '-', '5(DXTE-8)") As 
                                                                               OWENWITE IMEA TAKE DASE test_passenger_journey_recomes

ON(processing date = "s(date)", model version = "v2")

result.

experiment_mapping_experiment_group_label,

concast("s[DATE+2]", '--', 'S[DATE+8]") AS active_period

result
                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                FROM
LEFT JO
                      FROM
LEFT JO
                                                                                    experiment_mapping
result.passenger_id = experiment_mapping.passenger_id
                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                            experiment_mapping result.passenger_id = experiment_mapping.passenger_id
                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                ON
```

Figure 12: The example of Squirrel-Syntax, where an explicit error message exists.

F.2 SQUIRREL-SEMANTIC EXAMPLE

1188

1189

```
1190
                           In the original table, passenger_id is not a unique key. Could you help me check why the output contains a large number of duplicate rows? Please fix the bug.
1191
1192
                                 -- 60+ Lines Omited
result_base AS (
                                                                                                                                                                      -- 60+ Lines Omited. result_base AS (
                                                                                                                            correct SQL
                                                                                                                                                                                                                                                             Issue SQL
1193
                                       SELECT
                                                    recommendation_strategy_list
                                                                                                                                                                                         recommendation_strategy_list
1194
                                                           SELECT passenger id,
                                                                                                                                                                                                 SELECT passenger id,
                                                                         passenger_id,
MAX(journey_score) AS journey_score,
COLLECT_SET(data_source)
    AS recommendation_strategy_list
1195
                                                                                                                                                                                                               passenger_id,
MAX(journey_score) AS journey_score,
COLLECT_SET(data_source)
AS recommendation_strategy_list
1196
                                                                               SELECT encrypted_contact,
device_identifier,
passenger_id,
recommendation_timestamp,
                                                                                                                                                                                                                     SELECT encrypted_contact,
device_identifier,
passenger_id,
recommendation_timestamp,
1197
1198
                                                                                               journey_score,
traveler_profile_id,
                                                                                                                                                                                                                                    journey_score,
traveler_profile_id,
1199
                                                                                              passenger_id_collection,
device_id_collection,
is_selected,
                                                                                                                                                                                                                                   passenger_id_collection,
device_id_collection,
1200
                                                                                                                                                                                                                                   is_selected,
rejection_reason,
recommendation_strategy_list,
                                                                                              rejection_reason,
recommendation_strategy_list,
1201
                                                                                              data_source recommendation_model
                                                                                                                                                                                                                                   data_source recommendation_model
1202
1203
                                                                                                                                                                                                                      SELECT encrypted_contact,
                                                                                             encrypted_contact,
                                                                                              device_identifier,
passenger_id,
recommendation_timestamp,
                                                                                                                                                                                                                                   device_identifier,
passenger_id,
recommendation_timestamp,
1204
                                                                                               journey_score,
traveler_profile_id,
                                                                                                                                                                                                                                    journey_score,
traveler_profile_id,
                                                                                              passenger id collection,
device_id_collection,
is_selected,
rejection_reason,
recommendation_strategy_list,
                                                                                                                                                                                                                                   passenger id collection,
device_id_collection,
is_selected,
rejection_reason,
recommendation_strategy_list,
1206
1207
1208
                                                                                              data_source
retention_pipeline
                                                                                                                                                                                                                                   data_source
retention_pipeline
1209
1210
                                       journey_score DESC
1211
                                                                                                                                                                                         journey_score DESC
200000
                                ), result AS (
SELECT en
                                                                                                                                                                     ), result AS (
SELECT end
1212
                                                    encrypted_contact,
device_identifier,
                                                                                                                                                                                         encrypted_contact, device_identifier,
1213
                                                     passenger_id,
recommendation_timestamp,
                                                                                                                                                                                          passenger_id,
recommendation_timestamp,
1214
                                                                                                                                                                                         recommendation_timestamp,
journey_score,
traveler_profile_id,
passenger_id_collection,
device_id_collection,
device_id_tollection,
0 AS is_test_user,
recommendation_strategy_list
                                                     journey_score,
traveler_profile_id,
1215
                                                    passenger_id_collection,
device_id_collection,
0 AS is_test_user,
recommendation_strategy_list
1216
1217
                                                           SELECT base.encrypted_contact,
    result_base.recommendation_strategy_list,
    base.device_identifier,
    base.passenger_id,
    base.recommendation_timestamp,
                                                                                                                                                                                                SELECT base.encrypted_contact,
    result_base.recommendation_strategy_list,
    base.device_identifier,
    base.pasenger_id,
    base.recommendation_timestamp,
1218
1219
                                                                         base.recommendation timestamp,
base.journey_score,
base.traveler_profile_id,
base.passenger_id_collection,
base.device_id_collection
result_base
base
result_base_passenger_id = base.passenger_id
                                                                                                                                                                                                              base.recommendation timestamp,
base.journey_score,
base.traveler_profile_id,
base.passenger_id_collection,
base.device_id_collection
result_base
base
result_base.passenger_id = base.passenger_id
1221
1222
1223
                                 experiment_mapping AS (
SELECT passenger_id,
                                                                                                                                                                      experiment_mapping AS (
SELECT passenger_id,
                                                                                                                                                                                          pasa
CAST (
CAST (
                                                     CAST (
1225
                                                                                                                                                                                                                     md5(CAST(passenger_id AS STRING)),
-15
                                                                                md5(CAST(passenger_id AS STRING)),
-15
1226
1227
                                                                                                                                                                                                              ),
16,
1228
                                                     ) AS BIGINT
) % 2 AS BIGINT
) AS experiment_group_label
                                                                                                                                                                                          ) % 2 AS BIGINT
) AS experiment_group_label
1229
                                                           SELECT passenger_id
                                                                                                                                                                                                 SELECT passenger_id
1230
1231
1232
                                            OVERWRITE TABLE fake_base_test.passenger_journey_recomme
PARTITION(processing_date = '${date}', model_version = result.*,
                                                                                                                                                                      1233
1234
                                                                                                                                                                                          eriment_mapping.experiment_group_label,
eat('${DATE+2}', '~', '${DATE+8}') AS active_period
                                              experiment_mapping.experiment_group_label,
concat('${DATE+2}', '~', '${DATE+8}') AS active_period
                                                                                                                                                                                   experime
                                                                                                                                                                              result
JOIN
                                              result
1235
                                              experiment_mapping
result.passenger_id = experiment_mapping.passenger_id
                                                                                                                                                                                   experiment_mapping
result.passenger_id = experiment_mapping.passenger_id
1236
1237
```

Figure 13: The example of Squirrel-Semantic.

G PROMPTS TEMPLATE

All data synthesis and evaluation using the LLM-as-a-Judge methodology are performed with Claude-4-Sonnet (Anthropic, 2025), with the temperature parameter set to 0.0. The detailed prompts are described below.

G.1 Enterprise-level SQL Scripts Generation Prompts

Prompt for Scenario Creation

Instruction

1242

1243 1244

1245

1246

1247

1248 1249

1250 1251

1252 1253

1254

1255

1256

1257

1259

1260

1261

1262

1263

1264

1265

1266

1267

1268

1269

1270

1271

1272

1273 1274

1276

1277 1278

1279 1280

1281

1284 1285

1286

1287

1290 1291

1293

1294

1295

You are a professional SQL ETL and schema generation expert. Your task is to transfer a database schema from a source domain to a target domain, preserving structural complexity and table relationships, but fully adapting table names, field names, and semantics to the target domain.

Steps

- 1. Analyze Source DDL:
- Examine the number of tables, fields, data types, relationships, and naming patterns.
- Treat this as a structural seed for generating an equivalent schema.
- 2. Generate Target Schema:
- Create a logically equivalent schema under the target domain.
- Rules:
- Use the database fake_base_test.
- Format: CREATE TABLE IF NOT EXISTS fake_base_test.table_name (...);
- Avoid SQL reserved keywords as column names.
- Reflect business meaning in the target domain.
- Optionally add auxiliary fields to maintain equivalent complexity.
- All names, comments, and logic must be consistent with the target domain and unrelated to the source domain.
- 3. Validation:
- Ensure DDL syntax is correct.
- Ensure schema and scenario are fully adapted to the target domain, with no remnants from the source.

Notes

- Do not reuse proprietary identifiers or field names from the source domain.
- Only use the user-provided target domain.
- Preserve the structural pattern, complexity, and relationships of the source schema.

```
## Input Data
```

```
Source DDL: DDL
```

Target Domain: SCENARIO

1282 1283 ## Output Format(JSON)

```
{
    "mock scenario": "Scenario description",
    "mock ddl": "Corresponding CREATE TABLE statements"
}
```

Prompt for Generating Enterprise-level SQL

Instruction

You are a professional SQL ETL code generation expert. Using the provided source SQL as a reference, and given the target domain scenario and its corresponding DDL, generate an SQL ETL script for the target domain that preserves the logical structure and complexity of the source code while adapting it

```
1296
             fully to the target domain.
1297
1298
             ## Requirement
1299
1300
             1. Logical structure equivalence:
1301
             - Analyze the ETL workflow, table relationships, and processing steps in the source SQL code.
             - Preserve the overall structure, complexity, and transformation logic, but replace all table names, field
1302
             names, and data types to match the target domain.
1303
1304
             2. Strictly match the target DDL:
1305
             - All SQL must be fully based on the provided target DDL.
1306
             - Table names and field names must match the target DDL exactly.
             - Do not retain any original business terms, identifiers, or domain concepts from the source code.
1307
             3. Output requirements:
1309

    The code must be executable, and SQL syntax must be correct.

1310
             - Maintain a clear hierarchy and readability (include appropriate comments).
1311

    Naming should reflect the target business domain, ensuring a one-to-one correspondence between SQL

             and the target DDL.
1312
1313
             ## Input Data
1314
1315
             Source SQL: SQL
1316
             Target Domain Scenario: SCENARIO
             Target DDL: DDL
1317
1318
             ## Output Format
1319
1320
1321
                    'mock code': 'Generated target domain SQL ETL code'
1322
1323
```

G.2 ISSUE SQL CONSTRUCTION PROMPTS

1324

```
1326
             Prompt for Error Type Selection
1327
1328
             ## Role:
             You are an expert SQL engineer specializing in designing realistic SQL bugs for testing and debugging
1330
             scenarios.
1331
             ## Task:
1332
             Given a correct SQL query, your job is to:
1333
             Select the top {TOP_K} appropriate error type from the provided taxonomy.
1334
1335
             ##Key Guidelines:
             - Minimal Change: Only introduce the chosen bug. Do not alter the original query's structure or intent
1336
             more than necessary.
1337
             - Realism: The bug should reflect mistakes that real developers are likely to make.
1338
1339
             ##Input:
1340
             1. Correct SQL: \{SQL\}
             2. DDL (optional): {DDL}
1341
             3. Original Intent: {CODE INTENTION}
1342
             4. Error Type Taxonomy: {SEMANTIC ERROR TYPES}
1343
1344
             ##Output Requirements:
             Your output must include:
             - The selected error type(s) at Level 1-3 granularity.
1346
1347
             ##Output Format:
1348
1349
```

```
1350
                 candidate_errors:
1351
1352
                            "level1_error_type": Level 1 error type,
1353
                           "level2_error_type": Level 2 error type,
1354
                            "level3_error_type": Level 3 error type
1355
                      },
1356
                           "level1_error_type": Level 1 error type,
1357
                            "level2_error_type": Level 2 error type,
1358
                            "level3_error_type": Level 3 error type
1359
                      },
1360
                 ]
            }
1361
1362
1363
            Prompt for Squirrel-Syntax Issue SQL Construction
1364
            ## Role:
1365
            You are an expert SQL engineer specializing in designing realistic SQL bugs for testing and debugging
1366
            scenarios.
1367
1368
            ## Task:
1369
            Given a correct SQL query, your task is to introduce an error into the correct query with the smallest
1370
            possible change.
1371
            ## Key Guidelines:
1372
            - Minimal Change: Only introduce the chosen bug. Do not alter the original query's structure or intent
1373
            more than necessary.
1374
            - Realism: The bug should reflect mistakes that real developers are likely to make.
1375
            ## Input:
1376
            1. Correct SQL: {SQL}
1377
1378
            2. DDL (optional): {DDL}
1379
1380
```

- 3. Original Intent: {CODE INTENTION}
- 4. Error Type Taxonomy: {SEMANTIC ERROR TYPES}

Output Requirements:

Your output must include:

1381

1382 1383

1384

1385

1386 1387

1388 1389 1390

1391

1392

1393 1394 1395

1396 1397

1398

1399

1400 1401

1402

1403

- The selected error type(s) at Level 1–3 granularity.
- The modified SQL query with the injected bug.

```
## Output Format:
```

```
"level1_error_type": Level 1 error type,
"level2_error_type": Level 2 error type,
"level3_error_type": Level 3 error type,
"issue_sql": SQL query with the injected bug
```

Prompt for Squirrel-Semantic Issue SQL Construction

Role:

You are an expert SQL engineer specializing in designing realistic SQL bugs for testing and debugging scenarios.

Task:

Given a correct SQL query, your job is to:

1. Introduce the error into the SQL query with the smallest possible change.

```
1404
            2. Write a realistic user-style issue report describing how the bug causes the query to behave incorrectly,
1405
            and the user's real intention.
1406
1407
            ## Key Guidelines:
1408
            - Minimal Change: Only introduce the chosen bug. Do not alter the original query's structure or intent
1409
            more than necessary.
            - Realism: The bug should reflect mistakes that real developers are likely to make.
1410
1411
            ## Input:
1412
            1. Correct SQL: {SQL}
1413
1414
            2. DDL (optional): {DDL}
1415
            3. Original Intent: {CODE INTENTION}
1416
1417
            4. Error Type Taxonomy: {SEMANTIC ERROR TYPES}
1418
1419
            ## Output Requirements:
            Your output must include:
1420
            - The selected error type(s) at Level 1–3 granularity.
1421
           - The modified SQL query with the injected bug.
1422
            - A natural-language user bug report describing the mismatch between expected and actual results
1423
            (without exposing SQL code, since the user does not know the root cause).
1424
            ## Output Format:
1425
1426
1427
                 "level1_error_type": Level 1 error type,
1428
                 "level2_error_type": Level 2 error type,
1429
                 "level3_error_type": Level 3 error type,
1430
                 "user_query": Bug report written in natural language.
1431
                 Describe the expected vs. actual outcome clearly.
                 "issue_sql": SQL query with the injected bug
1432
1433
1434
1435
1436
```

G.3 BENCHMARK EVALUATION PROMPT

```
1460
            Prompt for Squirrel-Syntax Generation
1461
1462
            You are an SQL assistant.
1463
            ## Task
1464
1465
            Based on the error messages and table schema, your task is to fix the issue in the SQL and write the
1466
            correct SQL.
1467
            Remember that you can not change any existing comments and SQL code without errors.
1468
            ## Input Data
1469
            The issue SQL: BUG SQL
1470
            Related tables schema: DDL
1471
            Error Messages: ERROR MESSAGE
1472
1473
            ## Output (JSON):
1474
1475
                  'predict_sql': The fixed SQL.
1476
1477
1478
```

Prompt for Squirrel-Semantic Generation

You are an SQL assistant.

Task

Based on the user query and input table schema, please fix the bugs in the Issue SQL and write the corresponding correct SQL code.

Remember that you can not change any existing comments and SQL code without errors.

```
## Input Data
User Query:USER QUERY
Related tables schema: DDL
Error Messages: ERROR MESSAGE
## Output (JSON):
{
    'predict_sql': The fixed SQL.}
```

Prompt for diff Generation

```
1512
1513
'``last_edit
<><<<<< SEARCH
LAST_EDIT_BEFORE_PLACEHOLDER
=====

1516
LAST_EDIT_AFTER_PLACEHOLDER
>>>>>> REPLACE

1517
>>>>>> REPLACE

1518
```

G.4 AGENT PROMPT

Prompt for Main Agent

You are a SQL expert. Please review the SQL code (with the table DDL) and the error message reported. Your task is to analyze the error and provide fixing edit instructions.

Input:

- Tables DDL

1528 DDL_PLACEHOLDER

- Hive SQL Code:
- " sql SQL_CODE_PLACEHOLDER "
 - Error Message:

ERROR_MESSAGE_PLACEHOLDER

Output Requirements:

You must strictly follow this XML format in your response:

1536 <analysis>

Examine the error message and identify the root cause. Explain what is wrong with the current code and why the error occurred.

</analysis>

<instructions>

Provide clear, step-by-step instructions on how to fix the code. Explain what changes need to be made and where they should be applied.

</instructions>

<sketch_sql>

Provide the edit sketch using the special comment `...` to represent unchanged code between edited lines. Specify each edit in sequence, minimizing unchanged SQL code while making it clear what the edit is and where it should be applied.

</sketch_sql>

Ensure your instructions(in Chinese) and sketch are clear enough that another model can apply them correctly without accidentally deleting or modifying unintended parts of the code.