

VET: Verifiable Execution Tracing for Reliable Text-to-SQL Generation

Anonymous ACL submission

Abstract

Large language models (LLMs) have shown remarkable capabilities in text-to-SQL generation, yet existing approaches remain prone to hallucinations and lack verification mechanisms. Current methods such as Chain-of-Thought (CoT) and Program-of-Thought (PoT) typically rely on intermediate reasoning that is either purely textual or executed only as a final step, leaving the reasoning process opaque and prone to grounding and logical hallucinations. In this paper, we introduce Verifiable Execution Tracing (VET), a novel reasoning paradigm that transforms text-to-SQL from unverifiable textual rationales into step-wise executable semantics. VET addresses these limitations by constraining the reasoning process within a candidate schema space and formulating it as a sequence of executable Python steps. Crucially, each step is executed against the real database to produce observable intermediate results, which serve as immediate verification feedback and transform the traditionally opaque generation process into a transparent, debuggable interaction with database reality. Experiments demonstrate superior performance: 70.93% execution accuracy on BIRD benchmark, with exceptional gains on complex queries, validating that executable reasoning fundamentally outperforms textual alternatives.

1 Introduction

Large language models (LLMs) have emerged as a transformative interface for translating natural language into SQL queries, enabling intuitive database interactions across diverse domains. However, despite their remarkable progress, generating reliable and semantically correct SQL queries remains a fundamental challenge, as benchmarks evolve towards rigorous industrial standards (Li et al., 2024b). The inherent ambiguity of natural language, coupled with the structural

complexity of database schemas, creates a gap between user intent and executable queries that current approaches struggle to bridge effectively.

Current text-to-SQL systems face a primary limitation in reliably interpreting user intent within complex database environments. Users frequently provide incomplete or colloquial natural language descriptions, which leads to misinterpretations of field semantics and hallucinated relationships between entities. Recent studies have formalized these failures into a taxonomy of schema-based, logic-based, and content-based hallucinations (Yang et al., 2025; Qu et al., 2024). Even with sophisticated schema linking mechanisms designed to disambiguate retrieval and schema grounding (Wang et al., 2025c), models often struggle to establish accurate correspondences between vague user intents and specific database structures. This misalignment is particularly prevalent in in-context learning settings, where schema errors remain a dominant failure mode (Li et al., 2025b), resulting in queries that are syntactically valid but semantically incorrect.

To address these interpretive challenges, the field has increasingly adopted Chain-of-Thought (CoT) reasoning and its decomposed variants to explicitly structure the generation process. However, relying on these textual reasoning traces introduces a second, fundamental bottleneck: *unverifiability*. As depicted in the left panel of Figure 1, text-based reasoning approaches operate in a purely linguistic space, relying on implicit textual traces where the model essentially "guesses" the mapping between ambiguous intent and schema elements. This detachment from the execution environment makes models susceptible to "faithfully" executing initial errors such as hallucinating data values or violating schema constraints without detection, as the intermediate reasoning cannot be externally validated (Gao et al., 2023a). While some Program-of-Thought (PoT)

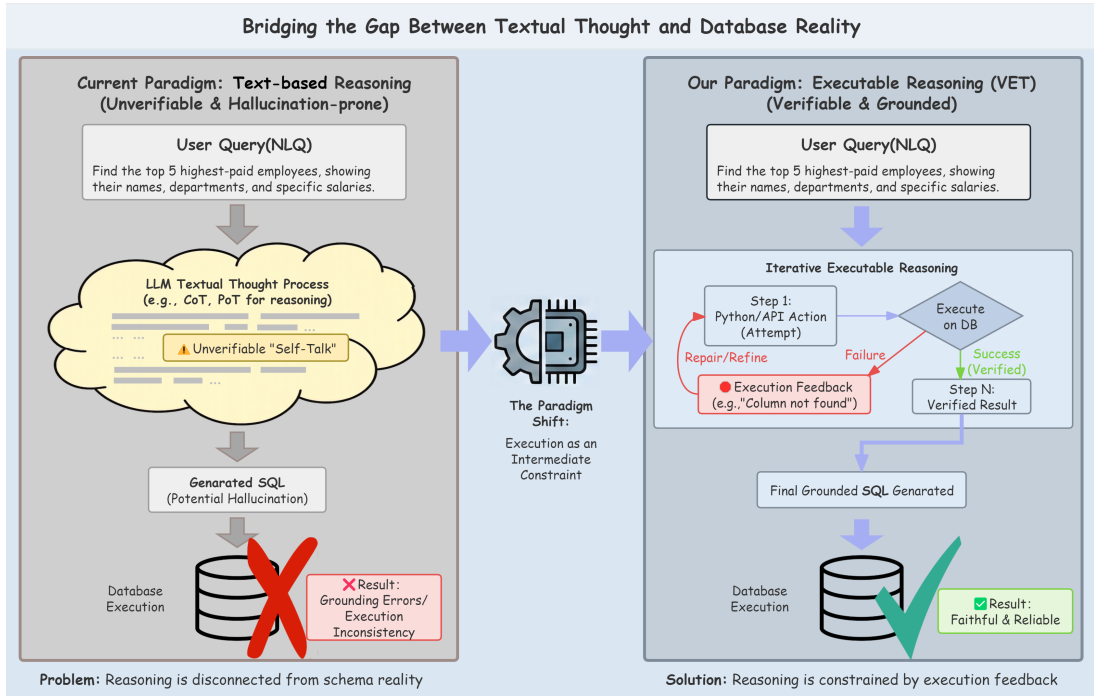


Figure 1: Comparison of reasoning paradigms. **Left:** Traditional text-based approaches operate as a black box, relying on ephemeral and ungrounded textual thoughts that are prone to hallucinations. **Right:** The VET framework utilizes executable code as a tangible carrier of reasoning. By materializing abstract intent into verifiable steps such as probing data existence in temporary DataFrames it transforms the process into a visible, debuggable interaction, ensuring logical consistency prior to SQL synthesis.

084 approaches introduce code execution to mitigate
 085 this (Gao et al., 2023c), they typically employ
 086 a “one-shot” execution paradigm where programs
 087 are validated only at the final stage. This de-
 088 layed verification leaves intermediate logical er-
 089 rors such as redundant filters or incorrect aggrega-
 090 tions undetected until they manifest as incorrect fi-
 091 nal results. Consequently, the reasoning process
 092 remains opaque, lacking the fine-grained, step-by-
 093 step interpretability required to verify and correct
 094 decision-making traces (Nguyen et al., 2025; Liu
 095 et al., 2025).

096 To address the aforementioned limitations and
 097 leverage this execution-based paradigm, we pro-
 098 pose Verifiable Execution Tracing (VET), a novel
 099 framework that transforms text-to-SQL from an
 100 unverifiable generation task into a step-wise ex-
 101 ecutable reasoning process. Unlike emerging multi-
 102 agent frameworks that rely on inter-agent text-
 103 ual critique (Heidari et al., 2025; Wang et al.,
 104 2025a), VET models the generation process as a
 105 function of a validated execution trace: $Y = G(V(E(Q, S)))$.
 106 Contrasting with the black-box failure modes
 107 shown on the left, the right panel of Figure 1
 108 highlights the core innovation of our approach:
 109 using executable code as the tan-

gible carrier of reasoning. Instead of unverifi-
 110 able textual thoughts, code serves as a rigorous
 111 medium that materializes abstract intent into ver-
 112 ifiable logic steps. By executing these intermedi-
 113 ate code steps such as verifying data existence in a
 114 temporary DataFrame VET transforms the reason-
 115 ing process into a visible, debuggable interaction
 116 with the database, ensuring logic consistency be-
 117 fore the final SQL synthesis.

The key contributions of our work can be sum-
 118 marized as follows:
 119

- We propose VET (Verifiable Execution Trac-
 121 ing), a novel reasoning paradigm that trans-
 122 forms text-to-SQL from unverifiable genera-
 123 tion to step-wise executable verification. By
 124 ensuring that every intermediate reasoning
 125 step is grounded in database reality, we elimi-
 126 nate the unverifiability inherent in traditional
 127 textual rationales.
 128
- We employ executable code as the primary
 129 vehicle for reasoning, effectively resolving
 130 the semantic ambiguity of natural language
 131 through the precision of programming logic.
 132 Unlike textual CoT which permits vague logi-
 133 cal leaps, code reasoning compels the model
 134

to materialize abstract intent into precise, executable operations, thereby mitigating hallucinations at the source.

- We achieve new state-of-the-art performance on the BIRD benchmark with an execution accuracy of 70.93%. Most notably, VET demonstrates superior robustness on complex queries, outperforming the baseline CoT by a substantial margin (61.11% vs. 27.78%), proving the efficacy of verifiable reasoning in handling intricate logic.

2 Related Work

Evolution of Text-to-SQL. The transformation of natural language into SQL has evolved from pattern matching to semantic parsing. Traditional approaches like RAT-SQL (Wang et al., 2019) and SMBOP (Rubin and Berant, 2020) utilized encoder-decoder models with relation-aware self-attention to align schema elements. With the rise of LLMs, the focus shifted to in-context learning. Methods such as DIN-SQL (Pourreza and Rafiei, 2024) and DAIL-SQL (Gao et al., 2023b) demonstrated that decomposing queries and selecting diverse examples significantly improves performance. However, these methods primarily operate in the textual domain. Recent advancements such as READ-SQL (Gao et al., 2023a) and PARSQL (Dai et al., 2025) have attempted to structure this reasoning process by decomposing queries into sub-components or utilizing parse trees to guide generation. While these decomposition strategies improve logical coherence, they often lack immediate grounding in the actual database state, leaving the “sim-to-real” gap unbridged.

Agentic and Execution-Guided Reasoning. Recent work has shifted towards execution-guided and agentic paradigms. Execution-guided methods such as ExCoT-DPO (Zhai et al., 2025) and Alpha-SQL (Li et al., 2025a) leverage execution feedback or search strategies (e.g., MCTS) to improve SQL generation. Agentic approaches like AgentiQL (Heidari et al., 2025) and MAC-SQL (Wang et al., 2025a) employ multi-agent architectures for planning, coding, and refinement. Unlike these prior methods, VET uniquely performs *inference-time semantic state search*. Python execution is used not only to validate the final output but also to probe intermediate states, ef-

Algorithm 1: Verifiable Execution Tracing

Input: Question Q , schema \mathcal{S} , database \mathcal{D}

Output: Executable SQL query Y

Schema-Constrained Initialization:

Extract candidate column set \mathcal{C}^* from (Q, \mathcal{S}) ;

Initialize data state $I_0 \leftarrow \mathcal{D}$;

Initialize reasoning trace $P \leftarrow []$;

Iterative Executable Reasoning:

while *query intent not yet satisfied* **do**

 Generate executable operation op_t constrained to \mathcal{C}^* ;

 Execute op_t on I_{t-1} to obtain output O_t ;

if *execution error or verification violation* **then**

 Repair op_t using execution feedback;

continue

 Append $p_t = (op_t, I_{t-1}, O_t)$ to P ;

 Update data state $I_t \leftarrow O_t$;

SQL Synthesis and Verification:

Generate SQL query Y conditioned on verified trace P ;

Execute Y on \mathcal{D} to obtain R_{SQL} ;

if $R_{SQL} \neq O_k$ **or** $Sim(Q, Q') < \tau$ **then**

 Reject Y ;

return Y

fectively mitigating content-based hallucinations that prior methods struggle to detect.

Execution-based Verification. Ensuring the correctness of generated SQL is critical. Approaches like R^3 (Xia et al., 2024) assess reliability by executing generated queries and checking for runtime errors. LEVER (Ni et al., 2023) takes this further by training a discriminator to rank SQL candidates based on their execution results. While effective, LEVER acts as a “Discriminative Black-box” relying on statistical likelihoods. VET introduces a “Generative White-box” approach via its Bidirectional Consistency Protocol, verifying validity by semantically reconstructing the user intent from the execution result, offering superior interpretability.

3 Methodology

3.1 Problem Definition

Let Q denote a natural language question, \mathcal{S} a database schema consisting of tables and typed

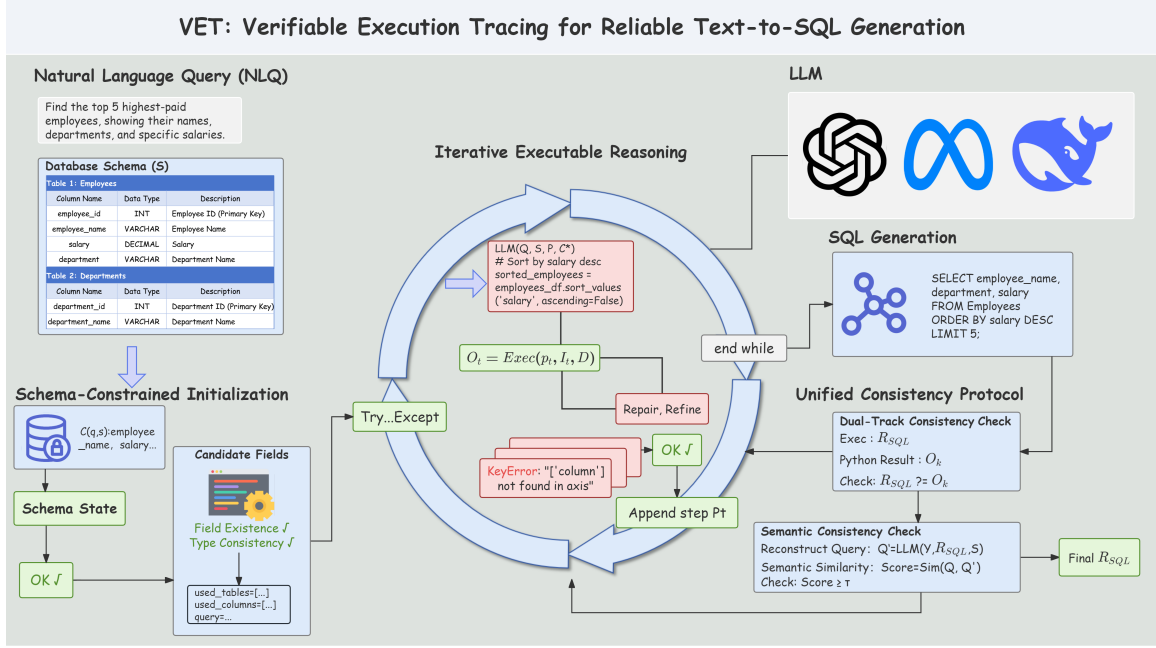


Figure 2: Overview of the proposed Verifiable Execution Tracing (VET) framework. VET first performs schema-constrained initialization to restrict the reasoning space to schema-grounded elements. It then constructs an iterative executable reasoning trace, where each generated operation is immediately executed on the database to produce observable feedback, which is used to verify, refine, or repair subsequent reasoning steps, with data probing operations invoked when ambiguity arises. Finally, VET synthesizes the SQL query from the verified execution trace and applies a unified consistency protocol, enforcing both execution alignment between SQL and reasoning outcomes and semantic consistency with the original question.

columns, and \mathcal{D} the corresponding database instance. The goal of text-to-SQL is to generate an executable SQL query Y such that executing Y on \mathcal{D} yields results that faithfully answer Q .

Most existing approaches model this task as a direct mapping from Q to Y , optionally augmented with textual intermediate reasoning. However, such formulations implicitly assume that intermediate reasoning steps are correct and grounded, despite lacking any mechanism to verify them against the actual database. As a result, errors such as hallucinated columns, invalid values, or logically inconsistent constraints often remain undetected until final execution. To address this limitation, we reformulate text-to-SQL as a structured reasoning problem with explicit execution-based verification. Instead of treating SQL generation as a one-shot output, we decompose it into reasoning, verification, and synthesis stages, enabling intermediate decisions to be validated before they affect the final query.

3.2 Overview of Verifiable Execution Tracing

We propose **Verifiable Execution Tracing (VET)**, a reasoning framework that transforms text-to-

SQL from an unverifiable generation task into a step-wise executable verification process. VET represents reasoning not as free-form text, but as a sequence of executable operations whose effects can be directly observed on the database.

Formally, VET models SQL generation as a composite function:

$$Y = \mathcal{G}(\mathcal{V}(\mathcal{E}(Q, \mathcal{S}), \mathcal{D})), \quad (1)$$

where \mathcal{E} generates an executable reasoning trace constrained by the schema, \mathcal{V} performs execution-based verification against the database, and \mathcal{G} synthesizes the final SQL query conditioned on the verified trace.

As illustrated in Figure 2, VET consists of three tightly coupled phases: (1) Schema-Constrained Initialization, (2) Iterative Executable Reasoning with Data Probing, and (3) SQL Synthesis with a Unified Consistency Protocol.

3.3 Schema-Constrained Initialization

A fundamental failure mode of text-to-SQL systems is schema hallucination, where models reference tables or columns that do not exist in the database but appear semantically plausible. To

mitigate this issue, VET begins by explicitly constructing a constrained reasoning space over the schema.

Given schema \mathcal{S} and question Q , a schema linking module extracts a candidate column set \mathcal{C}^* containing all columns potentially relevant to the query. This set is treated as a hard constraint: all subsequent reasoning operations must operate exclusively over \mathcal{C}^* .

By externalizing schema constraints as first-class variables, VET enables explicit detection of invalid reasoning steps. If a generated operation references a column outside \mathcal{C}^* , it is immediately identified as a hallucinated schema element and rejected before execution.

3.4 Iterative Executable Reasoning

At the core of VET is the construction of an executable reasoning trace $P = [p_1, p_2, \dots, p_k]$, where each step p_t is defined as a triplet:

$$p_t = (op_t, I_t, O_t), \quad (2)$$

with op_t denoting an executable Python operation (e.g., filtering, joining, aggregation), I_t the input data state, and O_t the observable output obtained by executing op_t on \mathcal{D} . Each operation is required to be type-consistent and executable, ensuring that O_t is deterministically derived from I_t .

Unlike purely textual reasoning, each reasoning step in VET is executed immediately after generation. Execution failures, empty results, or type errors serve as explicit feedback signals that guide refinement of subsequent operations. If an execution error occurs (e.g., due to a non-existent column or invalid value), the system triggers a repair process that leverages the error message to refine the operation, thereby suppressing hallucinations at their source.

Data Probing for Ambiguity Resolution. Natural language questions frequently introduce ambiguity in value expressions and schema mappings. For instance, a categorical filter such as “New York” may correspond to multiple surface forms (e.g., “NY” or “NYC”) in the database, or an entity name such as “Apple” may appear in multiple columns (e.g., Brand vs. Category). To resolve such ambiguity, VET incorporates a data probing mechanism that inspects database contents before committing to filtering constraints.

When ambiguity is detected, the model generates a probing operation (e.g., `df[col].unique()`

or value existence checks) and executes it on \mathcal{D} . The resulting observations are used to refine the original operation, ensuring alignment between user intent and actual database values. This online probing strategy directly mitigates content-based hallucinations, where models generate values that are linguistically plausible but factually absent from the database instance.

Step-wise Verification and Termination. After execution, each step is subjected to lightweight verification conditions that enforce basic logical soundness. For example, filtering operations are expected to satisfy $|O_t| \leq |I_t|$, while aggregation operations should reduce data dimensionality. If a step violates these conditions, it is repaired before being appended to the reasoning trace.

The reasoning process continues iteratively until the accumulated execution state O_k is sufficient to answer the query Q . At this point, the executable trace P serves as a faithful, validated representation of the reasoning process.

3.5 Unified Consistency Protocol

Once a verified reasoning trace is obtained, VET synthesizes the final SQL query Y by conditioning on P . However, generating an executable SQL statement alone does not guarantee semantic correctness. We therefore enforce a unified consistency protocol consisting of two complementary checks.

Dual-Track Execution Consistency. We execute the synthesized SQL Y on \mathcal{D} to obtain result R_{SQL} and compare it against the final execution output of the reasoning trace $R_{SQL} = O_k$. A mismatch indicates that the SQL fails to faithfully represent the verified reasoning process, and the generation is rejected. This dual-track consistency ensures alignment between symbolic SQL execution and executable reasoning semantics.

Semantic Consistency via Question Reconstruction. While execution alignment enforces computational correctness, it does not fully guarantee intent preservation. We therefore reconstruct a natural language question Q' from $(Y, R_{SQL}, \mathcal{S})$ and measure its semantic similarity to the original query Q . The SQL is accepted only if both execution consistency holds and $\text{Sim}(Q, Q') \geq \tau$, ensuring that the final output is both executable and semantically faithful.

Through this unified verification pipeline, VET

ensures that the generated SQL is not only correct by execution, but also grounded in a transparent, verifiable reasoning trajectory aligned with user intent.

4 Experimental Setup

4.1 Datasets and Benchmarks

We evaluate VET on two primary benchmarks that assess distinct reasoning capabilities:

- **BIRD-Dev** (Li et al., 2024b): A large-scale, cross-domain dataset emphasizing complex database values and long-chain reasoning. It is widely considered the most challenging Text-to-SQL benchmark due to its reliance on accurate value grounding. We report Execution Accuracy (EX) on the development set.
- **Spider Variants** (Yu et al., 2019): We employ the standard **Spider** dataset alongside its robustness variants, **Spider-DK** (Domain Knowledge) and **Spider-SYN** (Synonym), to evaluate the model’s resilience against linguistic perturbations and domain-specific terminology.

4.2 Baselines

We compare VET against a comprehensive set of state-of-the-art methods spanning three representative paradigms.

Text-centric prompting methods include Zero-shot and Few-shot Chain-of-Thought prompting, as well as decomposition-based approaches such as DIN-SQL (Pourreza and Rafiei, 2024) and DAIL-SQL (Gao et al., 2023b). These methods rely exclusively on textual reasoning and prompt engineering without external execution feedback.

Program-aided generation methods leverage code as an intermediate reasoning representation. We include TA-SQL (Qu et al., 2024) and Pi-SQL (chi et al., 2025), which employ Python programs to decompose logical structures. Unlike VET, these methods adopt a static generation paradigm without iterative verification or state correction.

Agentic and retrieval-augmented frameworks represent the current state-of-the-art in leveraging external knowledge and reranking mechanisms, including MAC-SQL (Wang et al., 2025b), CHESS (Talaei et al., 2024), Super-SQL (Li et al., 2024a), and RSL-SQL (Cao et al., 2024).

Table 1: Execution accuracy on BIRD-Dev for various backbone models and methods, broken down by query difficulty (Simple, Moderate, Challenging).

Backbone	Method	Simple	Moderate	Challenging	Overall
GPT-4	TA-SQL	63.14	48.60	36.11	56.19
	SuperSQL	66.90	46.50	43.80	58.50
	MAC-SQL	65.73	52.69	40.28	59.39
Qwen2.5-32B	Pi-SQL	70.92	56.47	49.66	64.54
	Baseline	56.97	38.49	29.17	48.76
	CoT	57.19	39.78	27.78	49.15
	DIN-SQL	54.38	44.30	31.25	49.15
	DAIL-SQL	59.89	44.52	35.42	52.93
	TA-SQL	61.95	45.38	39.58	54.82
	CHESS	63.68	48.17	42.36	56.98
	RSL-SQL	69.73	54.09	54.48	63.56
	VET (Ours)	72.54	62.37	61.11	68.38
GPT-4o	Baseline	64.65	48.60	40.28	57.50
	RSL-SQL	74.38	57.11	53.79	67.21
	VET (Ours)	76.11	64.09	59.72	70.93

Table 2: Execution accuracy on BIRD-Dev comparing VET (training-free) with fine-tuned models from the leaderboard.

Model	Accuracy (%)			
	Simple	Moderate	Challenging	Overall
Agentar-Scale-SQL	79.35	69.40	64.14	74.90
XiYan-SQL	-	-	-	73.34
CHASE-SQL	-	-	-	73.01
Q-SQL	-	-	-	72.99
OmniSQL-32B	-	-	-	69.23
VET (Ours)	76.11	64.09	59.72	70.93

4.3 Implementation Details

Our primary backbone model is **DeepSeek-Chat (V3)**, chosen for its strong coding capabilities and open-weights availability. To demonstrate model agnosticism, we also conduct experiments using **Qwen2.5-Coder** (7B, 14B, 32B) and **Llama3.1-8B**. For all experiments, we set the temperature to 0 to ensure reproducibility. The maximum number of repair turns in the VET loop is set to 3.

5 Results and Analysis

We conduct a comprehensive evaluation of VET to validate its effectiveness across performance, generalization, robustness, and mechanistic dimensions. All experiments adhere to a reproducible setup with a fixed temperature (0) and a maximum of 3 repair iterations to eliminate stochastic variability in LLM outputs.

5.1 Performance on BIRD-Dev

Table 1 presents the main execution accuracy results on BIRD-Dev across difficulty levels. With DeepSeek-Chat as the backbone, VET achieves an overall EX of 68.38%, outperforming the strongest

open-source baseline RSL-SQL (63.56%) by 4.82 percentage points. When paired with GPT-4o, VET further improves to 70.93%, exceeding GPT-4o + RSL-SQL by 3.72 points.

Notably, the largest gains appear in the *Challenging* subset, where VET attains 61.11% accuracy, surpassing DIN-SQL by 29.86 points and CHESS by 18.75 points. This subset is dominated by queries involving nested aggregation, multi-stage filtering, and value-dependent conditions. Purely text-based decomposition methods tend to accumulate early-stage reasoning errors in such cases, whereas VET detects and corrects erroneous intermediate states through step-wise execution feedback before they propagate to the final SQL.

Beyond direct baselines, Table 2 contextualizes VET against fine-tuned systems reported on the BIRD leaderboard. Despite being entirely training-free, VET (70.93% with GPT-4o) matches or exceeds several fine-tuned models such as OmniSQL-32B (69.23%) and approaches reinforcement-learning-based methods like Q-SQL (72.99%). This comparison highlights that enforcing executable and verifiable reasoning can partially substitute for data-intensive fine-tuning, especially in value-sensitive benchmarks like BIRD.

5.2 Generalization Across Model Scales

Table 3: Performance of VET across different backbone models and model sizes on BIRD-Dev.

Model	Method	Sim.	Mod.	Cha.	Avg.
Qwen2.5-7B	Baseline	45.19	29.25	21.53	38.14
	VET	58.59	42.80	36.11	51.69
Qwen2.5-14B	Baseline	52.97	30.11	34.72	44.33
	VET	68.43	54.19	49.31	62.32
Llama3.1-8B	Baseline	34.70	22.37	22.22	29.79
	VET	59.46	41.94	39.58	52.28

To evaluate generalization across model capacities, we apply VET to backbone models of varying sizes and families, with results shown in Table 3.

Across all backbones, VET consistently improves execution accuracy. Notably, the relative gains are larger for smaller models (e.g., 7B), suggesting that VET functions as an external reasoning scaffold that compensates for limited parametric capacity through structured execution feedback.

5.3 Robustness under Distribution Shifts

Table 4: Execution Accuracy (EX) on Spider, Spider-DK, and Spider-SYN under distribution shifts.

Dataset	Method	Easy	Medium	Hard	Overall
Spider	Baseline	0.927	0.823	0.649	0.747
	+VET	0.948	0.901	0.787	0.849
Spider-DK	Baseline	0.891	0.809	0.514	0.707
	+VET	0.836	0.850	0.730	0.804
Spider-SYN	Baseline	0.847	0.720	0.593	0.658
	+VET	0.923	0.848	0.768	0.807

We evaluate robustness on Spider, Spider-DK, and Spider-SYN, with results summarized in Table 4. These datasets progressively introduce linguistic and domain-level perturbations while preserving the underlying database schemas.

On the standard Spider benchmark, VET improves overall execution accuracy from 74.7% to 84.9%, with the largest gain observed in the Hard subset (64.9% \rightarrow 78.7%). This indicates that even in in-domain settings, executable step-wise verification substantially benefits complex logical queries.

Spider-DK evaluates robustness to domain-specific terminology that deviates from common lexical priors. VET achieves an overall accuracy of 80.4%, outperforming the baseline by 9.7 points. This improvement stems from VET’s data probing mechanism, which validates the existence and semantics of domain-specific values directly against the database, reducing reliance on brittle language priors.

Spider-SYN focuses on synonym substitution at the schema level. VET improves overall accuracy by 14.9 points (65.8% \rightarrow 80.7%). When an incorrect synonym leads to empty execution results, VET automatically explores alternative mappings through execution feedback, eliminating the need for manually curated synonym dictionaries.

5.4 Ablation Study

We conduct an ablation study on BIRD-Dev to analyze the contribution of each component in VET. The full framework consists of Schema-Constrained Initialization (SC), Step-wise Verification (SV), Unified Consistency Checking (UC), and Execution Feedback (EF). Table 5 reports execution accuracy for variants where one component is removed while all others are kept unchanged.

Table 5: Ablation results of VET on BIRD-Dev. Each variant removes one component from the full executable reasoning pipeline while keeping the remaining modules unchanged.

Variant	Simple	Mod.	Cha.	Overall
VET (Full)	72.54	62.37	61.11	68.38
w/o Schema Constr.	72.22	58.71	52.78	66.30
w/o Step-wise Verif.	70.92	53.98	47.22	63.56
w/o Unified Consist.	70.59	52.04	52.78	63.30
w/o Execution Feed.	70.81	52.90	45.83	63.04

Schema-Constrained Initialization (SC). Removing SC leads to a consistent but moderate performance drop, with overall accuracy decreasing from 68.38% to 66.30%. This suggests that early schema grounding stabilizes intent interpretation and reduces hallucinated fields, although downstream execution-based mechanisms can partially compensate for its absence.

Step-wise Verification (SV). Disabling SV and reverting to one-shot execution results in a 4.82-point drop in overall accuracy, with the largest degradation on the Challenging subset. This confirms that incremental execution is crucial for early detection of intermediate logical errors that would otherwise remain latent until final SQL execution.

Unified Consistency Checking (UC). Removing UC reduces overall accuracy by 5.08 points, particularly on the Moderate subset. Without enforcing alignment between intermediate execution traces and final SQL semantics, subtle inconsistencies can persist despite successful step-wise execution.

Execution Feedback (EF). EF contributes the most to overall performance. Removing EF causes overall accuracy to drop from 68.38% to 63.04%, with a sharp 15.28-point decline on the Challenging subset, highlighting the importance of grounding reasoning in real database states through execution feedback.

5.5 Qualitative Case Studies Overview

To further illustrate how VET transforms potential failures into successful queries, we provide detailed qualitative analyses in Appendix A. Table 9 (in the appendix) presents three representative scenarios from BIRD-Dev, covering:

1. **Schema Hallucination**, where baseline models make single-table assumptions that VET corrects via KeyError-driven table expansion; 2.

Value Mismatch, where string formatting errors are caught and corrected via execution feedback; 3. **Semantic Ambiguity**, where multi-hop schema paths are discovered through failed execution signals.

These examples highlight the importance of step-wise verification and in-process repair in reducing logic, schema, and value errors.

Additionally, Appendix D presents a full end-to-end case study, demonstrating VET’s reasoning process from initial query to final SQL output in complex multi-step scenarios.

5.6 Error Patterns and Repair Mechanism

To complement quantitative performance, we analyze representative error types and the effect of VET’s repair loop. Detailed qualitative error analysis and repair statistics are provided in Appendix B. Overall, VET substantially reduces logic and schema errors and recovers a meaningful fraction of otherwise failing queries, highlighting the practical benefits of step-wise verification and repair.

5.7 Efficiency Considerations

VET introduces iterative execution overhead due to step-wise verification and repair. However, this overhead is modest relative to the gains in execution accuracy, especially for complex queries. Detailed efficiency statistics, including inference time and token consumption per query difficulty, are provided in Appendix C. Techniques like schema constraints and potential caching can further mitigate runtime in practice.

6 Conclusion

We presented VET, a framework that mitigates ungrounded hallucinations in Text-to-SQL by enforcing an executable and verifiable reasoning process. By constraining the search space, leveraging data probing for ambiguity resolution, and applying unified consistency checks, VET aligns natural language intent with database reality. Experiments on BIRD and Spider demonstrate that VET achieves state-of-the-art performance while generalizing across model sizes and domains. Future work will focus on reducing inference latency and extending VET to settings with restricted database access.

7 Limitations

While VET significantly reduces hallucinations and improves accuracy, it introduces two primary limitations. First, the step-wise execution mechanism increases inference latency compared to direct text-to-SQL generation, as it requires waiting for database feedback at each reasoning step. Although this overhead is justified for high-stakes queries, it may be a bottleneck for real-time applications requiring millisecond-level latency. Second, VET relies on direct interaction with the database during inference. In scenarios involving highly sensitive or private data (e.g., healthcare or finance), strict security protocols might restrict the model’s ability to execute arbitrary Python code against the production database, necessitating sandboxed environments or mirrored schemas.

8 Ethical considerations

The development of execution-augmented text-to-SQL systems raises considerations regarding data privacy and automation. Since VET executes generated code against real databases, there is a risk of unintended data exposure or modification if not properly sandboxed. We strictly employ read-only permissions in our experiments and recommend similar safeguards for deployment. Furthermore, while VET enhances automation in data analysis, we view it as an assistive tool to empower domain experts rather than a replacement for human data analysts. We are committed to transparency and reproducibility; all code and evaluation scripts used in this work will be made publicly available to the research community.

References

Zhenbiao Cao, Yuanlei Zheng, Zhihao Fan, Xiaojin Zhang, Wei Chen, and Xiang Bai. 2024. [Rsl-sql: Robust schema linking in text-to-sql generation](#). *Preprint*, arXiv:2411.00073.

Yongdong chi, Hanqing Wang, Zonghan Yang, Jian Yang, Xiao Yan, Yun Chen, and Guanhua Chen. 2025. [Pi-sql: Enhancing text-to-sql with fine-grained guidance from pivot programming languages](#). *Preprint*, arXiv:2506.00912.

Yaxun Dai, Haiqin Yang, Mou Hao, and Pingfu Chao. 2025. [PARSQL: Enhancing text-to-SQL through SQL parsing and reasoning](#). In *Findings of the Association for Computational Linguistics: ACL 2025*, pages 661–681, Vienna, Austria. Association for Computational Linguistics.

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

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

Luyu Gao, Aman Madaan, Shuyan Zhou, Uri Alon, Pengfei Liu, Yiming Yang, Jamie Callan, and Graham Neubig. 2023c. [Pal: Program-aided language models](#). *Preprint*, arXiv:2211.10435.

Omid Reza Heidari, Siobhan Reid, and Yassine Yaakoubi. 2025. [Agentlql: An agent-inspired multi-expert framework for text-to-sql generation](#). *Preprint*, arXiv:2510.10661.

Boyan Li, Yuyu Luo, Chengliang Chai, Guoliang Li, and Nan Tang. 2024a. [The dawn of natural language to sql: Are we fully ready?](#) *Proceedings of the VLDB Endowment*, 17(11):33183331.

Boyan Li, Jiayi Zhang, Ju Fan, Yanwei Xu, Chong Chen, Nan Tang, and Yuyu Luo. 2025a. [Alpha-sql: Zero-shot text-to-sql using monte carlo tree search](#). *Preprint*, arXiv:2502.17248.

Jinyang Li, Binyuan Hui, Ge Qu, Jiayi Yang, Binhua Li, Bowen Li, Bailin Wang, Bowen Qin, Ruiying Geng, Nan Huo, and 1 others. 2024b. [Can llm already serve as a database interface? a big bench for large-scale database grounded text-to-sqls](#). *Advances in Neural Information Processing Systems*, 36.

Miaoran Li, Jiangning Chen, Minghua Xu, and Xiaolong Wang. 2025b. [Hallucination detection in structured query generation via llm self-debating](#). In *Findings of the Association for Computational Linguistics: EMNLP 2025*, pages 16102–16113.

Hanbing Liu, Haoyang Li, Xiaokang Zhang, Ruo-tong Chen, Haiyong Xu, Tian Tian, Qi Qi, and Jing Zhang. 2025. [Uncovering the impact of chain-of-thought reasoning for direct preference optimization: Lessons from text-to-sql](#). *Preprint*, arXiv:2502.11656.

Giang Nguyen, Ivan Brugere, Shubham Sharma, Sanjay Kariyappa, Anh Totti Nguyen, and Freddy Lecue. 2025. [Interpretable llm-based table question answering](#). *Preprint*, arXiv:2412.12386.

Ansong Ni, Srini Iyer, Dragomir Radev, Veselin Stoyanov, Wen-tau Yih, Sida Wang, and Xi Victoria Lin. 2023. [Lever: Learning to verify language-to-code generation with execution](#). In *International Conference on Machine Learning*, pages 26106–26128. PMLR.

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

686 Ge Qu, Jinyang Li, Bowen Li, Bowen Qin, Nan Huo,
687 Chenhao Ma, and Reynold Cheng. 2024. Before
688 generation, align it! a novel and effective strategy
689 for mitigating hallucinations in text-to-sql genera-
690 tion. *arXiv preprint arXiv:2405.15307*.

691 Ohad Rubin and Jonathan Berant. 2020. Smbop: Semi-
692 autoregressive bottom-up semantic parsing. *arXiv
693 preprint arXiv:2010.12412*.

694 Shayan Taleai, Mohammadreza Pourreza, Yu-Chen
695 Chang, Azalia Mirhoseini, and Amin Saberi. 2024.
696 Chess: Contextual harnessing for efficient sql syn-
697 thesis. *arXiv preprint arXiv:2405.16755*.

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

703 Bing Wang, Changyu Ren, Jian Yang, Xinnian Liang,
704 Jiaqi Bai, LinZheng Chai, Zhao Yan, Qian-Wen
705 Zhang, Di Yin, Xing Sun, and Zhoujun Li. 2025a.
706 **MAC-SQL: A multi-agent collaborative framework
707 for text-to-SQL**. In *Proceedings of the 31st Inter-
708 national Conference on Computational Linguistics*,
709 pages 540–557, Abu Dhabi, UAE. Association for
710 Computational Linguistics.

711 Bing Wang, Changyu Ren, Jian Yang, Xinnian Liang,
712 Jiaqi Bai, LinZheng Chai, Zhao Yan, Qian-Wen
713 Zhang, Di Yin, Xing Sun, and Zhoujun Li. 2025b.
714 **Mac-sql: A multi-agent collaborative framework for
715 text-to-sql**. *Preprint*, arXiv:2312.11242.

716 Yihan Wang, Peiyu Liu, and Xin Yang. 2025c.
717 **Linkalign: Scalable schema linking for real-world
718 large-scale multi-database text-to-sql**. *Preprint*,
719 arXiv:2503.18596.

720 Hanchen Xia, Feng Jiang, Naihao Deng, Cunxiang
721 Wang, Guojiang Zhao, Rada Mihalcea, and Yue
722 Zhang. 2024. *r³: "this is my sql, are you with me?"
723 a consensus-based multi-agent system for text-to-sql
724 tasks*. *Preprint*, arXiv:2402.14851.

725 Bo Yang, Yinfen Xia, Weisong Sun, and Yang Liu.
726 2025. **Hallucination detection for llm-based text-
727 to-sql generation via two-stage metamorphic testing**.
728 *Preprint*, arXiv:2512.22250.

729 Tao Yu, Rui Zhang, Kai Yang, Michihiro Ya-
730 sunaga, Dongxu Wang, Zifan Li, James Ma,
731 Irene Li, Qingning Yao, Shanelle Roman, Zilin
732 Zhang, and Dragomir Radev. 2019. **Spider: A
733 large-scale human-labeled dataset for complex and
734 cross-domain semantic parsing and text-to-sql task**.
735 *Preprint*, arXiv:1809.08887.

736 Bohan Zhai, Canwen Xu, Yuxiong He, and Zhewei
737 Yao. 2025. **Excot: Optimizing reasoning for
738 text-to-sql with execution feedback**. *Preprint*,
739 arXiv:2503.19988.

A Qualitative Case Studies 740

741 In this appendix, we provide a detailed qualita-
742 tive analysis of the repair mechanisms in VET to
743 demonstrate how execution feedback transforms
744 potential failures into successful queries. Ta-
745 ble 9 presents three representative scenarios se-
746 lected from the BIRD development set, covering
747 Schema Hallucination, Value Mismatch, and Se-
748 mantic Ambiguity.

749 In the first scenario (**Schema Hallucination**),
750 the baseline model falls victim to ‘single-table
751 bias,’ assuming all necessary columns exist in the
752 primary table. VET detects this via a `KeyError`
753 during execution, which forces the model to ex-
754 pand its search space to related tables (e.g., `frpm`)
755 and generate the correct JOIN condition.

756 The second scenario (**Value Mismatch**) high-
757 lights the fragility of string matching. The base-
758 line blindly assumes a spaced format (‘ = ’), re-
759 sulting in zero matches. VET captures this empty
760 result set as a negative signal (`ValueError` on
761 `argmax`), triggering a re-examination of the lit-
762 eral values and leading to the correct tight format
763 (‘=’).

764 Finally, the third scenario (**Semantic Ambigu-
765 ity**) addresses complex schema linking. The base-
766 line conflates ‘set translations’ with ‘foreign data’,
767 a common semantic drift. The execution failure of
768 the direct merge path compels VET to discover the
769 valid multi-hop path through the `cards` table.

770 These cases collectively illustrate that VET’s in-
771 teractive Python layer serves as a crucial ground-
772 ing mechanism, allowing the model to ‘fail fast’
773 and self-correct before committing to a final SQL
774 query.

B Error Patterns and Repair Efficacy 775

Table 6: Error distribution analysis on 100 sampled error cases. VET significantly reduces errors across all categories, particularly Logic and Schema errors.

Method	Logic Error	Schema Error	Value Error
Baseline (CoT)	66.2%	27.7%	6.1%
VET (Ours)	45.9%	11.4%	2.9%
<i>Reduction (Δ)</i>	↓ 20.3%	↓ 16.3%	↓ 3.2%

776 We complement quantitative results with a qual-
777 itative analysis of 100 sampled error cases (Table
778 6) and repair loop effectiveness (Table 7). VET re-
779 duces logic errors by 20.3 pp (66.2% → 45.9%)—
780 the largest reduction across error types—driven by

Table 7: Effectiveness of the repair loop on BIRD-Dev. “Success Rate” is calculated based on the triggered queries ($N_{trigger} = 156$).

Metric	Count	Global Rate	Success Rate
Total Queries	1534	100%	-
Triggered ($N_{trigger}$)	156	10.17%	-
Code Fixed	135	8.80%	86.54%
Final Correct	82	5.35%	52.56%

the Logical Observability check, which flags non-sensical intermediate states (e.g., empty filter results) and triggers repair before final SQL generation. Schema errors are cut by 16.3 pp (27.7% \rightarrow 11.4%) due to SC and Data Probing, as the model verifies column/table existence before generating operations. Value errors are reduced by 3.2 pp (6.1% \rightarrow 2.9%), with Data Probing resolving format mismatches (e.g., ‘=’ vs. ‘=’) by querying actual data distributions.

The repair loop is triggered for 10.17% of queries (156/1534), with an 86.54% code-fix rate (135/156) and a 52.56% final correction rate (82/156). This means VET recovers 5.35% of queries that would otherwise fail—a critical gain for high-stakes applications like business intelligence. Table 9 highlights three representative failure modes addressed by VET: (1) schema hallucination (resolved via KeyError feedback to switch to the correct table), (2) value format mismatch (fixed via empty sequence error triggering repair), and (3) multi-hop schema path error (corrected via execution failure forcing identification of valid multi-hop paths). These cases illustrate VET’s paradigm shift from post-hoc error correction to in-process validation, transforming unverifiable textual errors into actionable execution feedback.

C Efficiency Statistics

Table 8: Comparison of average inference time (s) and token consumption. VET maintains reasonable efficiency while delivering superior accuracy.

Method	Simple		Moderate		Challenging	
	Time	Tok.	Time	Tok.	Time	Tok.
VET	42.4	14k	49.7	16k	56.5	17k
w/o Schema	50.5	33k	47.6	36k	90.5	35k
w/o Step.	41.2	12k	47.7	15k	116.8	15k
w/o Unified	40.3	13k	47.0	15k	52.7	14k
w/o Feed.	40.0	14k	47.8	16k	78.5	18k

While VET introduces iterative execution over-

head (Table 8), the trade-off between performance and efficiency is favorable. Inference time for VET is modestly higher than baselines (e.g., 56.5s for Challenging queries vs. 40.0s for the baseline without EF), but this is offset by a 15.28 pp accuracy gain in the same subset. Token consumption remains competitive (17k tokens for Challenging queries vs. 35k for baselines without SC), as schema constraints reduce the model’s search space and limit redundant token generation. Compared to multi-agent methods like MAC-SQL, VET’s deterministic verification loop avoids consensus-based sampling overhead, resulting in lower token usage despite iterative repair. For real-world applications, the marginal increase in inference time is justified by substantial accuracy gains and the 5.35% recovery of failed queries via repair loops.

810
811
812
813
814
815
816
817
818
819
820
821
822
823
824
825
826
827

781
782
783
784
785
786
787
788
789
790
791
792
793
794
795
796
797
798
799
800
801
802
803
804
805
806
807
808

809

Table 9: **Qualitative Analysis of Error Repair.** We present three representative cases where the Baseline fails due to schema hallucination, value mismatch, and logic errors. In contrast, VET utilizes the Python execution feedback (e.g., `KeyError`, `ValueError`) as a grounded signal to detect inconsistencies and iteratively refine the reasoning process, leading to correct SQL generation.

Query & SQL Comparison	Interactive Repair Process (Python)	Error Analysis
<p>Q: What is the lowest grade for the District Special Education Consortia School (ID 613360)?</p> <p>Baseline (X): SELECT MIN(LowGrade) FROM schools ... <i>→ Hallucinates 'LowGrade' in 'schools'.</i></p> <p>VET (✓): SELECT MIN(T1.'Low Grade') FROM frpm AS T1 JOIN schools AS T2 ...</p>	<p>[Trigger] Execution Failure: # Assume data is in 'schools' table low_grade = schools['Low Grade'].min() Error: <code>KeyError: 'Low Grade'</code> <i>(Probe: Column missing in 'schools')</i></p> <p>↓ <i>Self-Correction via Repair Loop</i></p> <p>[Repair] Cross-Table Retrieval: # Switch to Join 'frpm' table merged = pd.merge(schools, frpm, ...) result = merged['Low Grade'].min() Result: Execution Success ✓</p>	<p>Diagnosis: Single-table Bias. The Baseline fails to locate data in the external frpm table.</p> <p>Effect: Schema Probe. The <code>KeyError</code> invalidates the single-table assumption. This forces the model to switch strategies, retrieving the correct table and generating correct JOIN logic.</p>
<p>Q: Is the molecule with the most double bonds carcinogenic?</p> <p>Baseline (X): ... WHERE bond_type = ' = ' ... <i>→ Value Error: Extra spaces cause zero matches.</i></p> <p>VET (✓): SELECT ... WHERE T1.bond_type = '=' ...</p>	<p>[Trigger] Execution Failure: # Filter by string with spaces bonds = df[df['bond_type'] == ' = '] top_mol = bonds...idxmax() Error: <code>ValueError: attempt to get argmax of an empty sequence</code> <i>(Probe: Query returned 0 rows)</i></p> <p>↓ <i>Self-Correction via Repair Loop</i></p> <p>[Repair] Robust Handling → Correct SQL: # Code robustly handles empty case if not bonds.empty: ... Result: Execution Success ✓ <i>Feedback:</i> The “Empty” signal guided the final SQL to use the correct literal '='.</p>	<p>Diagnosis: String Literal Mismatch. The Baseline incorrectly infers the data format as ' = ' (with spaces), leading to an empty result set.</p> <p>Effect: Execution Feedback. The <code>ValueError</code> (empty sequence) acts as a negative signal, indicating the string literal was incorrect. This feedback steers the model to generate the correct, tight-formatted string '=' in the final SQL.</p>
<p>Q: List top 2 sets in 'Shadowmoor' block that have Italian-language cards.</p> <p>Baseline (X): ... FROM sets JOIN set_translations ON ... WHERE language='Italian' <i>→ Logic Error: Wrong table. Checks set name translation, not card content.</i></p> <p>VET (✓): ... FROM sets T1 JOIN cards T2 ... JOIN foreign_data T3 ... WHERE T3.language='Italian'</p>	<p>[Trigger] Execution Failure: # Attempt direct link between Sets and ForeignData merged = pd.merge(sets, foreign_data, on='id') Error: <code>KeyError: "['name'] not in index"</code> <i>(Direct merge failed; ambiguous columns/no path)</i></p> <p>↓ <i>Self-Correction via Repair Loop</i></p> <p>[Repair] Multi-hop Path Found: # Recognized missing intermediate link # Path: Sets → Cards → ForeignData Result: Execution Success ✓</p>	<p>Diagnosis: Semantic Ambiguity. Baseline confused “translated set names” (‘set_translations’) with “foreign language cards” (‘foreign_data’).</p> <p>Effect: Path Validation. The execution failure invalidates the incorrect or direct join attempts. This constrains the model to search for a valid schema path, correctly identifying ‘cards’ as the necessary intermediate bridge.</p>

D End-to-End Execution Case Study

828

In this section, we present a comprehensive end-to-end case study demonstrating the VET process on a sample query from the BIRD-Dev dataset. This example illustrates how the model initializes the environment, performs step-wise executable reasoning, and validates the final output.

829

830

831

D.1 Phase 1: Initialization & Schema Loading

832

The process begins with the initialization of the database schema and the loading of relevant tables.

833

Query: “What is the highest eligible free rate for K-12 students in the schools in Alameda County?”

834

Evidence: Eligible free rate for K-12 = ‘Free Meal Count (K-12)’ / ‘Enrollment (K-12)’

835

Database: californiaschools

836

Selected Tables: schools, satscores, frpm

837

The system first defines the schema structure and loads the data into Pandas DataFrames.

838

[Python]

839

```
import pandas as pd
import numpy as np
from numpy import nan
pd.options.mode.chained_assignment = None
#`schools` data
schools_data = [
    ('37679916037576', 'San Diego', 'Bostonia Language Academy', 'K-3'),
    ('15635291530369', 'Kern', 'Centennial High', '9-12'),
    ('41689990134197', 'San Mateo', 'Aspire East Palo Alto Charter', None)
]

#`schools` primary key: `CDSCode`
schools_columns = ['CDSCode', 'County', 'School', 'GSserved']

schools = pd.DataFrame(schools_data, columns=schools_columns)
#`satscores` data
satscores_data = [
    ('10623800000000', 'Fresno', 291),
    ('1612000133397', 'Alameda', 454),
    ('39686270130864', 'San Joaquin', 116)
]

#`satscores` primary key: `cds`
#`satscores` foreign key: ['cds']
satscores_columns = ['cds', 'cname', 'enroll12']

satscores = pd.DataFrame(satscores_data, columns=satscores_columns)
#`frpm` data
frpm_data = [
    ('19647330129858', 'Los Angeles', 170.0, 149.0, 0.876470588235294),
    ('09737836110498', 'El Dorado', 19.0, 19.0, 1.0),
    ('33670903330412', 'Riverside', 1608.0, 958.0, 0.595771144278607)
]

#`frpm` primary key: `CDSCode`
#`frpm` foreign key: ['CDSCode']
frpm_columns = ['CDSCode', 'County Name', 'Enrollment (K-12)', 'Free Meal Count (K-12)', 'Percent (%) Eligible Free (K-12)']

frpm = pd.DataFrame(frpm_data, columns=frpm_columns)
```

840

841

842

843

844

845

846

847

848

849

850

851

852

853

854

855

856

857

858

859

860

861

862

863

864

865

866

867

868

869

870

871

872

873

874

875

876

877

878

879

880

[Python]

881

```
import sqlite3

conn = sqlite3.connect("dev/dev_databases/california_schools/california_schools.sqlite")

query = "SELECT * FROM `satscores`;"
satscores = pd.read_sql_query(query, conn)
satscores = satscores[['cds', 'cname', 'enroll12']]

query = "SELECT * FROM `schools`;"
```

882

883

884

885

886

887

888

889

890

891

892

```

893 schools = pd.read_sql_query(query, conn)
894 schools = schools[['CDSCode', 'County', 'School', 'GSserved']]
895
896 query = "SELECT * FROM `frpm`;"
897 frpm = pd.read_sql_query(query, conn)
898 frpm = frpm[['CDSCode', 'County Name', 'Enrollment (K-12)', 'Free Meal Count (K-12)', 'Percent (%)
900 Eligible Free (K-12)']]

```

[Python]

```

901 query = """ What is the highest eligible free rate for K-12 students in the schools in Alameda County
902 ? """
903
904 evidence = """ Eligible free rate for K-12 = `Free Meal Count (K-12)` / `Enrollment (K-12)` """
905

```

907 D.2 Phase 2: Iterative Reasoning

908 The model then proceeds with a step-by-step reasoning process, executing Python code at each stage to
 909 verify intermediate results.

910 D.2.1 Step 1: Filter Data by County

[Markdown]

```

911 ### Step 1: Filter target data for Alameda County
912
913
914

```

915 To answer the question "What is the highest eligible free rate for K-12 students in the schools in
 916 Alameda County", we first need to focus on relevant data for ****Alameda County****:

- 917 - Core Dataset: The `frpm` table contains key metrics such as K-12 enrollment and free meal counts
 918 for each county.
- 919 - Filtering Logic: Filter rows where the `County Name` field equals `Alameda` to obtain a dataset
 920 `alameda_frpm` containing only Alameda County, narrowing the scope for subsequent ratio
 921 calculations.
 922

[Python]

```

923
924 # Step 1: Filter data for Alameda County
925 alameda_frpm = frpm[frpm['County Name'] == 'Alameda']
926
927
928 # Print filtered result
929 print(alameda_frpm)
930

```

[Execution Output]

```

931
932
933
934   CDSCode County Name Enrollment (K-12) Free Meal Count (K-12) \
935 0 01100170109835 Alameda 1087.0 565.0
936 1 01100170112607 Alameda 395.0 186.0
937 2 01100170118489 Alameda 244.0 134.0
938 3 01100170123968 Alameda 191.0 113.0
939 4 01100170124172 Alameda 257.0 14.0
940 .. ... ..
941 374 01751016108971 Alameda 1258.0 98.0
942 375 01751016114318 Alameda 662.0 25.0
943 376 01751016117592 Alameda 1164.0 45.0
944 377 01751016117600 Alameda 723.0 33.0
945 378 01751196002711 Alameda 278.0 14.0
946

```

```

947   Percent (%) Eligible Free (K-12)
948 0 0.519779
949 1 0.470886
950 2 0.549180
951 3 0.591623
952 4 0.054475
953 .. ...
954 374 0.077901
955 375 0.037764
956 376 0.038660
957 377 0.045643
958 378 0.050360
959

```

```

960 [379 rows x 5 columns]

```

D.2.2 Step 2: Calculate Eligible Free Rate

[Markdown]

Step 2: Calculate eligible free rate per school

Based on the filtered `alameda_frpm` dataset, the next step is to calculate the eligible free meal rate for K-12 students in each school:

- Calculation Basis (Core Formula):

> Eligible free rate for K-12 = `Free Meal Count (K-12)` / `Enrollment (K-12)`

- Operational Logic: Add a new column `Eligible Free Rate` to `alameda_frpm` to store the free meal eligibility rate for each row (individual school). This is the core calculation step for finding the "highest rate".

[Python]

Step 2: Calculate the eligible free rate for each school

```
alameda_frpm['Eligible Free Rate'] = alameda_frpm['Free Meal Count (K-12)'] / alameda_frpm['Enrollment (K-12)']
```

Print calculation result

```
print(alameda_frpm)
```

[Execution Output]

	CDSCode	County Name	Enrollment (K-12)	Free Meal Count (K-12)	\
0	01100170109835	Alameda	1087.0	565.0	
1	01100170112607	Alameda	395.0	186.0	
...					
378	01751196002711	Alameda	278.0	14.0	

	Percent (%) Eligible Free (K-12)	Eligible Free Rate
0	0.519779	0.519779
1	0.470886	0.470886
...		
378	0.050360	0.050360

[379 rows x 6 columns]

D.2.3 Step 3: Extract Maximum Rate

[Markdown]

Step 3: Extract highest eligible free rate

After calculating the rate for individual schools, we need to extract the maximum value from the `Eligible Free Rate` column in `alameda_frpm`:

- Core Goal: Find the **highest value** of eligible free meal rates for K-12 students among all schools in Alameda County.

- Operational Logic: Use the Pandas `max()` method to directly obtain the maximum value of this column and assign the result to the variable `highest_eligible_free_rate`, precisely addressing the core requirement of the original question.

[Python]

Step 3: Find the maximum eligible free rate

```
highest_eligible_free_rate = alameda_frpm['Eligible Free Rate'].max()
```

Print highest rate

```
print(highest_eligible_free_rate)
```

[Execution Output]

```
1.0
```

D.3 Phase 3: Unified Consistency Protocol

Based on the verified reasoning trace, the system generates the final SQL query and performs consistency checks.

1031
1032
1033
1034
1035
1036
1037
1038
1039

```
[SQL]
SELECT
  MAX(
    CAST(`Free Meal Count (K-12)` AS REAL) / `Enrollment (K-12)`
  ) AS highest_eligible_free_rate
FROM frpm
WHERE
  `County Name` = 'Alameda';
```

1041
1042
1043
1044
1045
1046

```
[Execution Output]
highest_eligible_free_rate
0                            1.0
[1 rows x 1 columns]
```

1048
1049
1050
1051

Dual-Track Consistency Check: The execution result of the generated SQL ('1.0') is compared with the output of the final Python reasoning step ('highest_eligible_free_rate = 1.0'). Since the two values match, the SQL execution is consistent with the Python-based intermediate computation, confirming result reliability.

1052
1053

Semantic Consistency Check: We reconstruct the natural language query from the generated SQL and schema:

1054
1055
1056

```
[Query]
What is the highest eligible free rate for K-12 students in the schools in Alameda County?
```

1058
1059
1060

A cosine similarity calculation between the reconstructed query and the original question yields: Cosine Score = $0.8926 > \tau$, confirming semantic consistency.
Since both consistency checks pass, the final SQL is validated.