# SHARE: An SLM-based Hierarchical Action CorREction Assistant for Text-to-SQL

**Ge Qu** [1], **Jinyang Li** [1], **Bowen Qin** [2*], **Xiaolong Li** [1], **Nan Huo**[1],
**Chenhao Ma** [3], **Reynold Cheng** [1*]
[1]The University of Hong Kong, [2] BAAI
[3]The Chinese University of Hong Kong, Shenzhen
quge@connect.hku.hk, bwqin@baai.ac.cn, ckcheng@cs.hku.hk

## Abstract

Current self-correction approaches in text-to-SQL face two critical limitations: 1) Conventional self-correction methods rely on recursive self-calls of LLMs, resulting in multiplicative computational overhead, and 2) LLMs struggle to implement effective error detection and correction for declarative SQL queries, as they fail to demonstrate the underlying reasoning path. In this work, we propose **SHARE**, an **S**LM-based **H**ierarchical **A**ction cor**RE**ction assistant that enables LLMs to perform more precise error localization and efficient correction. SHARE orchestrates three specialized Small Language Models (SLMs) in a sequential pipeline, where it first transforms declarative SQL queries into stepwise action trajectories that reveal underlying reasoning, followed by a two-phase granular refinement. We further propose a novel hierarchical self-evolution strategy for data-efficient training. Experimental results demonstrate that SHARE effectively enhances self-correction capabilities while proving robust across various LLMs. Furthermore, our comprehensive analysis shows that SHARE maintains strong performance even in low-resource training settings, which is particularly valuable for text-to-SQL applications with data privacy constraints. For reproducibility, we release our code at https://github.com/quge2023/SHARE.

## 1 Introduction

Text-to-SQL, aimed at converting natural language (NL) queries to executable SQL queries (Qin et al., 2022), plays a crucial role in enabling non-technical users to analyze and interact with data in relational databases. In recent years, the application of Large Language Models (LLMs) has improved the performance of text-to-SQL to another level of intelligence (Rajkumar et al., 2022;

---

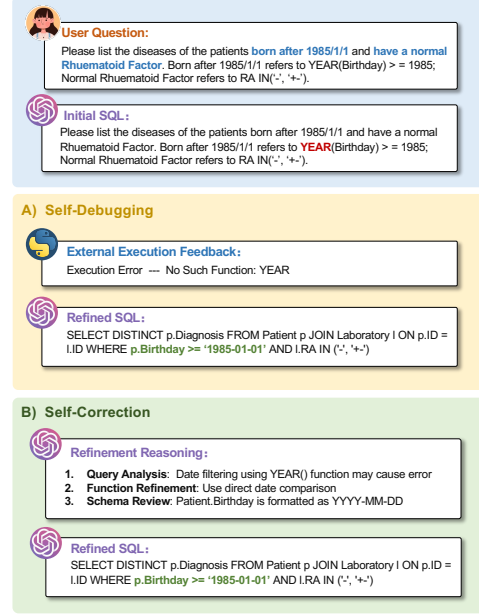* Corresponding authors are Bowen Qin and Reynold Cheng.



Figure 1: Illustrations of self-debugging and self-correction for text-to-SQL.

Talaei et al., 2024; Pourreza et al., 2024; Zhang et al., 2025). A critical component emerging from these architectural developments is the automatic error correction mechanism (Ouyang et al., 2022a; Chen et al., 2023; Pan et al., 2024b), which systematically identifies and rectifies query errors to improve response accuracy.

Contemporary automatic error correction approaches can be categorized along two primary dimensions: self-debugging and self-correction. Self-debugging represents an execution-guided approach in which LLMs iteratively refine their SQL queries based on database execution feedback (Zhong et al., 2023; Li and Xie, 2024; Xie et al., 2024), as illustrated in Figure 1 (A). While demonstrating promising results, this approach faces several fundamental challenges. First, execution feedback from mainstream SQL dialects such as SQLite tends to be concise but insufficiently specific, impeding accurate error localization. This

ambiguity may subsequently induce hallucinations (Dziri et al., 2021; Ji et al., 2023) in the correction process. Furthermore, the fundamental requirement for direct database execution access presents significant operational constraints, particularly in contexts where data privacy and security considerations preclude such direct interaction with the database system (Awan et al., 2023).

These issues have motivated the focus of automatic correction towards self-correcting mechanisms (Liu and Tan, 2024; Askari et al., 2024), where LLMs are prompted to revise their initial outputs through autonomous contextual re-analysis, without relying on external execution feedback. This decoupling from execution environments proves particularly valuable in text-to-SQL applications where database access is typically restricted by privacy constraints. However, achieving effective self-correction often requires multiple inference iterations through proprietary LLM APIs like GPT-4 or Claude-3.5-Sonnet, leading to prohibitive computational costs that scale exponentially. Furthermore, LLMs exhibit a tendency towards self-enhancement bias (Huang et al., 2024), leading them to overestimate the quality of their initial outputs and struggle to effectively identify errors within their self-generated declarative SQL queries. This inherent limitation undermines the overall effectiveness of the correction process.

In this work, we propose an assistant-based framework where generator LLMs create initial outputs and implement self-correction guided by assistants. Our primary contribution, **SHARE** (**S**LM-based **H**ierarchical **A**ction Cor**RE**ction Assistant), orchestrates three specialized Small Language Models (SLMs), each under 8B parameters, in a sequential pipeline. Specifically, the **B**ase **A**ction **M**odel (BAM) transforms raw SQL queries into action trajectories that capture reasoning paths; the **S**chema **A**ugmentation **M**odel (SAM) and the **L**ogic **O**ptimization **M**odel (LOM) further perform orchestrated inference to rectify schema-related and logical errors, respectively, within action trajectories. SHARE improves error detection precision and correction efficacy while reducing computational overhead compared to conventional LLM approaches. Additionally, we also incorporate a novel **hierarchical self-evolution strategy** that enhances data efficiency during training.

Experimental results across 4 diverse text-to-SQL benchmarks demonstrate the effectiveness,

efficiency, and generalizability of SHARE. By enabling SLMs to collaboratively guide LLMs in SQL correction, SHARE achieves substantial execution accuracy improvements over the GPT-4o baseline, with relative gains of 14.80% on BIRD and 11.41% on SPIDER within a single round of correction, while significantly reducing computational costs (Section 4.4). Beyond precision, our framework exhibits strong robustness in varying query complexities, low-resource settings, and various generator models, including both closed-source and open-source LLMs (Section 4.2). Notably, SHARE's learned debugging logic generalizes effectively to previously unseen SQL dialects without dialect-specific supervision (Section 4.5). These findings position our method as a scalable and cost-efficient paradigm for improving the reasoning reliability of LLMs in SQL auto-correction for real-world applications.

## 2 Preliminaries

### 2.1 Task Definition

**Text-to-SQL.** Given a natural language question $q_i \in \mathcal{Q}$, where $\mathcal{Q} = \{q_i\}_{i=1}^n$ with its corresponding database input [1] $d_i \in \mathcal{D}$, where $\mathcal{D} = \{d_i\}_{i=1}^n$, the goal of text-to-SQL is to guide the generator model $\mathcal{G}$ to generate the SQL query $s_i$ by:

$$s_i = f_{\mathcal{G}}(d_i, q_i), \qquad (1)$$

where $f_{\mathcal{G}}(\cdot)$ refers to the mapping function applied by the generator model $\mathcal{G}$.

**Assistant-based Self-Correction.** Given the $d_i$, $q_i$ and the corresponding initial SQL queries $s_i'$ generated by the generator model $\mathcal{G}$, the goal of assistant-based correction aims to utilize feedback $y_i$ generated by the assistant $\mathcal{T}$, which could be a LLM or an agent mechanism, to guide the generator model $\mathcal{G}$ to refine its initial output effectively by:

$$s_i = f_{\mathcal{G}}(d_i, q_i, s_i', y_i), \qquad (2)$$

where $y_i = f_{\mathcal{T}}(d_i, q_i, s_i')$, and $f_{\mathcal{T}}(\cdot)$ is the mapping function applied by the assistant. The generator model $\mathcal{G}$ and assistant $\mathcal{T}$ engage in multiple iterative refinement cycles until a correct SQL query is produced or the predefined maximum number of turns is reached.

---

[1] Database input refers to the schema of a specific database and its corresponding sampled value.
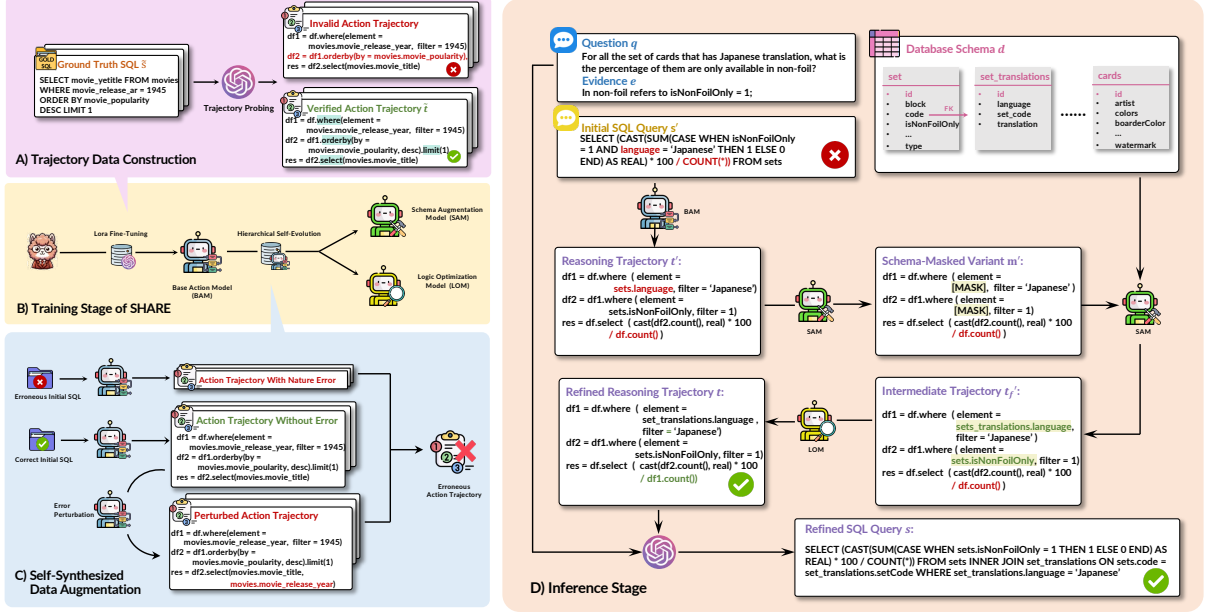
Figure 2: An illustration of SHARE. Figure (A)-(C) illustrate the training architecture of three specialized SLMs in SHARE: the Base Action Model (BAM), Schema Augmentation Model (SAM), and Logic Optimization Model (LOM). Figure D presents SHARE's orchestration of these three models in a sequential pipeline for inference.

## 2.2 Action Model

An action model (Zhang et al., 2024) $\mathcal{A}$ is specifically architected to comprehend, plan, and generate stepwise action trajectories $t_i$, given the contextual input $x_i$ by:

$$t_i = \{a_1, a_2, ..., a_n\} = f_{\mathcal{A}}(x_i; \mathbb{A}), \quad (3)$$

where $f_{\mathcal{A}}(\cdot; \mathbb{A})$ refers to the mapping function applied by the action model $\mathcal{A}$ with an action space $\mathbb{A}$ inherently mastered by $\mathcal{A}$, and $a_i \in \mathbb{A}$ denotes as the selected action. Unlike natural language, actions manipulated by action models are typically in the format of API or function calls. For instance, a declarative SQL query `"SELECT .. FROM .. WHERE .."` can be decomposed by an action model into a sequence of functional operations as `[where(column = param_1, value = param_2), select(table = param_1, column = param_2)]`.

## 3 Methodology

### 3.1 Training Pipeline

The general training data used to develop SLMs is often generated through knowledge distillation pipelines, thereby reducing the need for costly human expert annotations (Xu et al., 2024b; Peng and Zhang, 2024). In these pipelines, LLMs serve as **teacher** models, providing annotations and gen-

erating new training instances to guide the learning process of SLMs as **student** models. We employ this approach to train a Base Action Model (BAM), which transforms initially produced declarative SQL queries into structured sequences of actions within a predefined action space, capturing the underlying reasoning steps. Building on this foundation, two action models specialized for two crucial aspects of text-to-SQL are further presented: the Schema-Augmentation Model (SAM) and the Logic Optimization Model (LOM). SAM concentrates on improving schema linking (Dou et al., 2023), while LOM focuses on logical synthesis (Yin and Neubig, 2017). All our models were trained via Lora fine-tuning (Hu et al., 2022).

To train SAM and LOM efficiently, we propose a novel continual learning technique: **hierarchical self-evolve strategy**. Instead of repeatedly querying a teacher LLM (e.g., GPT-4o) for each new training instance, this strategy leverages BAM to synthesize and augment task-specific training data, targeting distinct aspects of text-to-SQL translation. This approach reduces annotation costs while maintaining strong performance, as demonstrated in Section 4.4. Formally, we begin with a seed dataset $\mathcal{C} = \{(d_i, q_i, \tilde{s}_i, s'_i)\}_{i=1}^n$, where each tuple $(d_i, q_i, \tilde{s}_i)$ consists of a database input, a user question, and the associated ground-truth SQL from the BIRD and SPIDER training corpora. A detailed

| Type | Definition | Example |
|------|------------|---------|
| ADD | Inserts an additional action into the original action trajectory. $[a_1, a_2, ..., a_n] \rightarrow$ $[a_1, a_2, \boldsymbol{a_{new}}, ..., a_n]$ | **Before:** df1 = df.orderby(element = movie.likes, desc).limit(1) res = df1.select(element = movie.director) **After:** **df1 = df.groupby(element = movies.id)** df2 = df1.orderby(element = movie.likes, desc).limit(1) res = df2.select(element = movie.director) |
| DELETE | Removes an existing action from the trajectory. $[a_1, \boldsymbol{a_2}, a_3, ..., a_n] \rightarrow$ $[a_1, a_3, ..., a_n]$ | **Before:** df1 = df.where(element = users.country_code,filter = 20) **df2 = df1.where(element = users.gender, filter = 'male')** res = df2.select(users.user_id) **After:** df1 = df.where(element = users.country_code,filter = 20) res = df1.select(users.user_id) |
| SUBSTITUTE | Replaces an existing action with a different action type or modifies the parameters of the existing action. $[a_1, \boldsymbol{a_2}, a_3, ..., a_n] \rightarrow$ $[a_1, \boldsymbol{a_2'}, a_3, ..., a_n] \rightarrow$ | **Before:** df1 = df.where(element = reviews.Date, filter = '2018-09-11') res = df1.select(district.district_id, district.city) **After:** df1 = df.where(element = reviews.Date, filter = '2018-09-11') res = df1.select(district.city, district.district_id ) |

Table 1: Three error perturbation types utilized by Base Action Model (BAM) to implement data augmentation.

distribution of the training data is provided in Appendix A.1. The initial SQL $s_i'$ is generated by GPT-4o using the baseline prompt implemented in BIRD (Li et al., 2024a).

## 3.2 Base Action Model (BAM)

**Training Target.** BAM aims to generate the corresponding action trajectory $t$ given an initial SQL query $s'$. Following prior study (Qu et al., 2024), we design actions in target trajectories as pandas-like APIs (see Refined Reasoning Trajectory in Figure 2 (D)) to present the reasoning process of text-to-SQL transformations. The complete enumeration of actions used in the construction of action trajectories within the SHARE framework is detailed in the Appendix A.2.

**Data Construction.** Given that BAM is the most important model and such reasoning derivations are complex and highly demanding, we employ GPT-4o as a strong teacher LLM to construct the training data for BAM. As shown in Figure 2 (A), GPT-4o is guided to convert each ground-truth (gt) SQL $\tilde{s}$ to a verified action trajectory $\tilde{t}$ by few-shot prompting. The prompt we use is detailed in Appendix D.3. To ensure high-quality training data with fewer hallucinations during conversion, we only contain $(\tilde{s}, \tilde{t})$ pairs in which $\tilde{t}$ can be successfully reverted back to $\tilde{s}$ (Berglund et al., 2024). We ultimately collected 13k query-trajectory pairs as the training data for BAM in this phase.

## 3.3 Schema Augmentation Model (SAM)

**Training Target.** Schema linking is a critical step in identifying the relevant database tables and columns needed to answer user queries (Lei et al., 2020). However, the complexity and heterogeneity of database schemas, spanning from pre-trained language models (PLMs) (Qiu et al., 2020; Li et al., 2023) to large language models (LLMs) (Zhao et al., 2024), often compromise accurate schema linking, thereby introducing substantial burdens on downstream SQL generation. To address this issue, we propose the Schema Augmentation Model (SAM), designed to specifically target and correct schema-linking errors within the input action trajectory. By focusing on schema-related components, SAM aims to isolate and rectify errors before they propagate, ultimately leading to more reliable SQL generation.

**Data Construction.** We begin with a corpus of 13k action trajectories from BAM, each represented as $(d, q, s', \tilde{t})$. For each verified trajectory $\tilde{t}$, we generate a schema-masked variant $\tilde{m}$ by inserting mask symbols [MASK] to isolate schema-specific elements. For each initial SQL $s'$, we extract its referenced tables and columns as an initial schema list $l'$. These paired forms, $(\tilde{t}, \tilde{m})$ and $(d, q, l', \tilde{m}, \tilde{t})$, form the backbone of SAM's two-phase training paradigm. In the first phase, SAM is trained with $(\tilde{t}, \tilde{m})$ to identify and mask schema-related elements accurately. In the second

phase, SAM leverages $(d, q, l', \tilde{m}, \tilde{t})$ to refine the previously masked segments, seamlessly reintegrating the corrected schema links into the schema-masked variant. By employing BAM in a few-shot setting to prepare and orchestrate these training steps, we ensure that SAM efficiently acquires the specialized capabilities needed for robust schema augmentation.

### 3.4 Logic Optimization Model (LOM)

**Training Target.** Given a more precise set of schema-linked database tables or columns, the reasoning logic expressed as an action trajectory should align with both natural language descriptions and valid SQL semantics. Formally, the model takes $d$ and $q$ as input and outputs a refined action trajectory $t$. This trajectory $t$ captures the correct chain of reasoning necessary to accurately resolve the question $q$.

**Data Construction.** As shown in Figure 2 (C), the erroneous action trajectories in the training data for LOM come from two resources. First, we collect the corresponding action trajectory of erroneous initial SQLs. However, the scale of this resource is limited. Therefore, we propose an **action-based perturbation strategy** for data augmentation. We apply three types of perturbation, as illustrated in Table 1, on error-free action trajectories derived from correct initial SQL to reproduce various logic errors in text-to-SQL conversions. Finally, 15k erroneous action trajectories along with their corresponding verified action trajectories are collected for LOM training. The detailed pseudocode of this process is provided in Appendix A.3.

### 3.5 Orchestration Inference

During inference, as shown in Figure 2 (D), SHARE operates in a sequential orchestration that integrates three LoRa fine-tuned models, enabling iterative refinement of action trajectories. Upon receiving an initial SQL query $s'$, SHARE first invokes the BAM to generate a corresponding action trajectory $t'$. The SAM then refines $t'$ by applying schema-based adjustments derived from the given database input $d$, producing the intermediate trajectory $t'_f$. Next, $t'_f$, along with the user query and the refined database content, is forwarded to the LOM for logic-based corrections. The refined trajectory $t$ generated by LOM serves as SHARE's final output and is employed as feedback to guide

the underlying language model in self-correcting $s'$ within a zero-shot setting. Ultimately, with these refined action trajectories serving as self-correction signals, the language model can regenerate more accurate and contextually appropriate SQL queries.

## 4 Experiments

### 4.1 Experiment Settings

**Datasets and Metrics.** The experiments are conducted on four challenging benchmarks for cross-domain text-to-SQLs. 1) BIRD (Li et al., 2024a) is the most challenging large-scale cross-domain text-to-SQL benchmark, which introduces external knowledge as an additional resource in complex scenarios. In this paper, we use its development set for evaluation, which contains 1,534 pairs of text-to-SQL data and 11 complex databases. 2) SPIDER(Yu et al., 2018) is a more standard cross-domain text-to-SQL benchmark. It contains 1,034 examples, covering 20 databases across multiple domains, in the development set. 3) DK (Gan et al., 2021), extended from the SPIDER benchmark, requires text-to-SQL parsers equipped with the ability of domain-knowledge reasoning. 4) RE-ALISTIC removes and switches the obvious mentions of schema items in questions, making it closer to the real scenarios. In this paper, we use widely adopted **Execution Accuracy** (EX) to measure the performance of our framework.

**Compared Methods.** We explore two open-source SLMs, namely Llama-3.1-8B and Phi-3-Mini-3.8B, as backbone models to construct our SHARE. For all baseline and advanced self-correction methods for comparison in Table 2, we employ GPT-4o as the generator model and report the results after a single refinement iteration. Details of these methods are shown in Appendix B.1.

**Implementation Details.** We fine-tune all our models using the LLaMa-Factory library (Zheng et al., 2024) with LoRA (Hu et al., 2022). All our experiments are conducted on 4×A100 GPU with 80GB memory. We detail the hyperparameters for training and inference in Appendix B.2, and claim the reproducibility of this work in Appendix E.

### 4.2 Overall Performance

**Overall Results.** Table 2 presents the performance of GPT-4o on the BIRD and SPIDER benchmarks, comparing approaches with baseline and

| Method | BIRD | | | | SPIDER | | | | |
|---|---|---|---|---|---|---|---|---|---|
| | Simple | Moderate | Challenging | Total | Easy | Medium | Hard | Extra Hard | Total |
| GPT-4o | 63.35 | 44.18 | 45.52 | 55.87 | 89.10 | 83.00 | 68.40 | 52.40 | 77.10 |
| *Self Correction w/o Feedback* | | | | | | | | | |
| Self-Correction | 62.70 | 43.75 | 44.83 | 55.28 | 88.30 | 82.70 | 66.10 | 49.40 | 75.90 |
| Self-Consistency | 65.75 | 49.04 | 45.21 | 58.75 | 92.30 | 87.90 | 75.10 | 58.60 | 81.80 |
| Multiple-Prompt | 66.38 | 48.06 | 44.83 | 58.80 | 91.10 | 87.20 | 74.70 | 59.00 | 81.50 |
| *Self Correction w/ Feedback* | | | | | | | | | |
| Self-Debugging ♣ | 65.41 | 47.84 | 46.21 | 58.28 | 91.10 | 85.90 | 74.70 | 60.80 | 81.20 |
| DIN-Correction | 65.62 | 46.98 | 44.83 | 58.02 | 91.90 | 85.20 | 70.10 | 55.40 | 79.50 |
| MAC-Refiner ♣ | 66.27 | 47.41 | <u>46.90</u> | 58.74 | <u>93.10</u> | 85.00 | 73.00 | 56.60 | 80.40 |
| MAGIC | 66.75 | 49.46 | 45.79 | 59.53 | - | - | - | - | <u>85.66</u> |
| **+SHARE-3.8B** | <u>68.00</u> | <u>51.29</u> | 46.21 | <u>60.89</u> | 92.70 | <u>88.30</u> | <u>78.70</u> | <u>65.10</u> | 84.00 |
| **+SHARE-8B** | **70.81** | **56.25** | **46.90** | **64.14** | **94.00** | **90.10** | **78.20** | **70.50** | **85.90** |

Table 2: Self-correction performance of GPT-4o in Execution Accuracy (EX) (%) on BIRD and SPIDER. ♣ means the model uses external execution results as feedback. **Bold** indicates best results, while <u>underlines</u> denote second-best results.

advanced self-correction strategies. Three key observations emerge: 1) When GPT-4o attempts intrinsic Self-Correction, performance actually decreases (55.87% → 55.28% on BIRD). This occurs because the model lacks reliable mechanisms to assess the correctness of its prior reasoning steps, sometimes converting originally correct solutions into incorrect ones (Huang et al., 2024). Moreover, improving prompts for self-correction strategies can inadvertently introduce biases, leading to suboptimal revisions. (Gou et al., 2024). Although advanced methods with carefully crafted designs, such as Multiple-Prompt and Magic, mitigate these biases, they rely heavily on extensive human-engineered prompts and entail significant computational overhead. 2) In contrast, SHARE enables GPT-4o to perform effective self-correction through a single interaction, resulting in a relative improvement of 14.80% in EX on BIRD and 11.41% on SPIDER with SHARE-8B. That is because SHARE introduces a novel mechanism for inferring and analyzing initial hidden reasoning paths, allowing the LLM to identify and rectify errors more precisely. Notably, this process is conducted automatically by SLMs, thereby reducing both the human effort and the computational costs associated with high-level LLMs. 3) Significant performance improvements using SHARE-3.8B and SHARE-8B, which are based on two widely used SLMs (Phi-3-mini and Meta-Llama-8B), demonstrate the generalization and robustness of our training pipeline.

| METHOD | EASY | MED. | HARD | EXTRA | ALL |
|---|---|---|---|---|---|
| *DK* | | | | | |
| GPT-4o | 75.50 | 73.60 | 47.30 | 41.90 | 64.10 |
| +SHARE-3.8B | 84.90 | 73.80 | 48.90 | 56.20 | 69.20 |
| +SHARE-8B | 85.50 | 81.30 | 56.80 | 63.80 | 75.30 |
| *REALISTIC* | | | | | |
| GPT-4o | 81.70 | 82.30 | 69.70 | 49.50 | 73.40 |
| +SHARE-3.8B | 88.10 | 86.70 | 68.70 | 57.70 | 78.00 |
| +SHARE-8B | 87.20 | 88.20 | 74.70 | 68.00 | 81.50 |

Table 3: Execution Accuracy (EX) of GPT-4o + SHARE across queries of varying levels of difficulty on DK and REALISTIC.

**Results on Robust Testing.** Table 3 presents evaluation results in EX on the variant datasets of SPIDER for robustness without any additional training. SHARE-8B effectively enhances the performance of the GPT-4o baseline by 11.20% and 8.10% on DK and REALISTIC, respectively, with improvements across all difficulty levels. SHARE-3.8B also facilitates significant performance gains, achieving a relative increase of 7.96% on DK and 7.18% on REALISTIC.

**Results on Various Generator Models.** Although the teacher LLMs we employ for data generation during training are primarily based on GPT-4o, the performance gains observed with our trained SHARE-8B assistant extend well beyond this single source, as evidenced in Table 4. Notably, our approach significantly improves performance for both proprietary closed-source models, such as Claude-3.5-Sonnet (↑ 28.64% relatively), and for open-source alternatives, such as Llama-3.1-70B (↑ 14.88% relatively). This indicates that our

| MODEL | SIM. | MOD. | CHALL. | TOTAL |
|---|---|---|---|---|
| *Closed-Source LLM* | | | | |
| Claude-3.5-S | 57.08 | 39.87 | 31.03 | 49.41 |
| +SHARE-8B | 68.86 | 57.11 | 50.34 | 63.56 |
| GPT-4o-mini | 55.03 | 41.59 | 35.17 | 49.09 |
| +SHARE-8B | 67.46 | 50.86 | 40.00 | 59.64 |
| GPT-3.5-turbo | 52.51 | 35.99 | 29.66 | 45.35 |
| +SHARE-8B | 57.51 | 38.36 | 31.03 | 49.22 |
| *Open-Source weaker LM* | | | | |
| Llama-3.1-70B | 60.54 | 44.61 | 41.38 | 53.91 |
| +SHARE-8B | 68.86 | 53.88 | 43.45 | 61.93 |
| Qwen-Coder-32B | 65.72 | 46.71 | 45.23 | 58.03 |
| +SHARE-8B | 67.78 | 54.31 | 46.90 | 61.73 |
| Llama-3.1-8B | 41.62 | 27.59 | 20.00 | 35.33 |
| +SHARE-8B | 47.68 | 35.13 | 28.97 | 42.11 |
| DS-Coder-6.7B | 41.30 | 26.94 | 23.45 | 34.57 |
| +SHARE-8B | 56.25 | 45.04 | 36.55 | 51.24 |

Table 4: Self-correction performance of various generator models assisted by SHARE on BIRD. For brevity, we refer to some models using shorthand names and provide their corresponding official model aliases in Appendix B.3. SIM., MOD., CHALL. represent the levels of query difficulty and are the abbreviations of simple, moderate, and challenging, respectively.
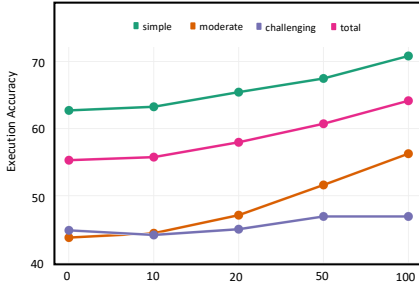


Figure 3: The effect of training data scale on SHARE.

method is not limited to the model bias of error patterns in text-to-SQL, but also general knowledge of SQL correction.

**Results on Low-resource Settings.** To take a closer look at the impact of the amount and quality of training data on self-correction assistance in LLM, we conduct a low-resource training analysis. Specifically, we sampled three subsets of the training data—10%, 20%, and 50%—and repeated each experiment three times to minimize variance. The averaged results are illustrated in Figure 3.

Our findings reveal a strong positive correlation between the amount of training data and model performance, thereby confirming the overall high quality of the training set. In particular, GPT-4o + SHARE-8B with only 50% training data outperforms the state-of-the-art MAGIC baseline, achiev-

| METHOD | SIM. | MOD. | CHALL. | TOTAL |
|---|---|---|---|---|
| SHARE-8B | 70.81 | 56.25 | 46.90 | **64.14** |
| *(a) w/o Schema Aug* | 67.14 | 50.22 | 46.00 | **60.02** $_{(\downarrow 4.08)}$ |
| *(b) w/o Logic Opt* | 64.86 | 46.77 | 39.31 | **56.98** $_{(\downarrow 7.16)}$ |
| *(c) w/o Hierarchy* | 68.34 | 49.89 | 45.02 | **60.55** $_{(\downarrow 3.59)}$ |
| *(d) w/o Error Pert* | 68.08 | 52.79 | 46.20 | **61.38** $_{(\downarrow 2.76)}$ |

Table 5: Ablation study of SHARE. **w/o Schema Aug** and **w/o Logic Opt** denote removing the SAM and LOM, respectively. **w/o Hierarchy** denotes training specialized models sequentially instead of employing the hierarchical evolution strategy. **w/o Error Pert** denotes the removal of action-based error perturbation.

ing 60.71% versus MAGIC's 59.53% in EX in the BIRD data set.

However, this positive trend does not fully hold for more challenging instances. From 0% to 20% of the training data, improvements remain unclear and at 10% there is even a slight decline in performance. Only when the dataset exceeds 20% of the full training set we do observe a clear performance increase. This suggests that performance improvements for harder instances often remain erratic or limited until the dataset size exceeds a certain point, particularly for tasks involving nuanced or rare examples, as demonstrated by (Kaplan et al., 2020).

### 4.3 Ablation Study

Table 5 presents the results of our ablation study. The removal of either the schema or logic refinement module (Table 5 (a–b)) results in substantial performance drops, underscoring the importance of the two-stage refinement architecture for disentangling schema linking from logical reasoning (Lei et al., 2020). Replacing the hierarchical evolution strategy with conventional sequential training (Table 5 (c)) leads to a 3.59% decline, indicating the advantage of the hierarchical approach in preventing error accumulation and bias transfer across stages (Liu et al., 2021). Furthermore, the exclusion of action-based error perturbation (Table 5 (d)) leads to a 2.76% reduction, demonstrating the effectiveness of this lightweight augmentation technique. More detailed analysis is provided in Appendix C.1.

### 4.4 Computational Cost Analysis of SHARE

Table 6 presents the computational cost analysis of SHARE. To the best of our knowledge, our work is the **first** effort in the text-to-SQL domain to implement correction through collaboration between

| Method | LLM InToks | LLM OutToks ↓ | SLM InToks ↓ | SLM OutToks ↓ | Cost / 1K ↓ | EX ↑ |
|---|---|---|---|---|---|---|
| *Inference Stage* | | | | | | |
| MAC-Refiner | 7126.74 | 236.58 | - | - | $20.18 | 58.74 |
| Multiple-Prompt | 21128.65 | 1004.55 | - | - | $62.86 | 58.80 |
| MAGIC | 8245.16 | 1737.98 | - | - | $37.99 | 59.53 |
| GPT-4o + SHARE-8B | **716.30** | **68.32** | 1731.23 | 132.16 | **$2.57** | **64.14** |
| *Training Stage* | | | | | | |
| MAGIC | 4838.63 | 2085.22 | - | - | $32.94 | 59.53 |
| GPT-4o + SHARE-8B | **1623.28** | **66.87** | 2308.75 | 83.04 | **$4.85** | **64.14** |

Table 6: Token usage and computational cost of various self-correction methods on BIRD. In(Out)Toks refers to the average input (output) token length per instance. Cost / 1K refers to the average cost per 1000 instances.



Figure 4: The performance of SHARE on BIRD across various SQL dialects, specifically MySQL (left) and PostgreSQL (right).

| Model | Sim. | Mod. | Chall. | Total |
|---|---|---|---|---|
| GPT-4o | 63.35 | 44.18 | 45.52 | 55.87 |
| +SHARE-gpt | 70.81 | 56.25 | 46.90 | 64.14 |
| +SHARE-llama | 71.57 | 57.54 | 48.97 | 65.19 |
| Claude-3.5-S | 57.08 | 39.87 | 31.03 | 49.41 |
| +SHARE-gpt | 68.86 | 57.11 | 50.34 | 63.56 |
| +SHARE-llama | 67.68 | 52.59 | 44.83 | 60.95 |
| Llama-3.1-70B | 60.54 | 44.61 | 41.38 | 53.91 |
| +SHARE-gpt | 68.86 | 53.88 | 43.45 | 61.93 |
| +SHARE-llama | 69.08 | 56.90 | 46.90 | 63.30 |
| Llama-3.1-8B | 41.62 | 27.59 | 20.00 | 35.33 |
| +SHARE-gpt | 47.68 | 35.13 | 28.97 | 42.11 |
| +SHARE-llama | 53.08 | 35.56 | 30.34 | 45.63 |

Table 7: Performance comparison across different teacher models on BIRD.

small and large LMs. This paradigm effectively reduces the inference overhead of SHARE and incurs only **one-tenth** the cost of the most economical baseline. Notably, SHARE remains highly cost-efficient during the training data construction stage. Compared to the In-Context-Learning-based (ICL-based) training method adopted in MAGIC, which is the strongest self-correction method, our self-evolution strategy significantly reduces the reliance on LLMs during the construction of training data, resulting in substantial cost savings in the overall process. Other relevant details of the cost computation are shown in Appendix C.2.

## 4.5 Generalization for Different SQL Dialects

Mainstream text-to-SQL benchmarks predominantly use SQLite as their target SQL dialect, primarily for its accessibility and ease of data collection. However, due to the heightened privacy requirements in data management, MySQL and PostgreSQL, characterized by licensing restrictions and proprietary attributes, are more commonly adopted dialects in real-world implementation scenarios. As demonstrated in Figure 4, SHARE also exhibits effective performance on MySQL and PostgreSQL dialects even without additional training. This can be attributed to SHARE's focus on learning low-level reasoning path corrections, enabling it to generalize across various high-level SQL dialects.

## 4.6 Open-source Teacher Models

During the construction of SHARE, as introduced in Section 3.2, GPT-4o acts as the teacher model for automated data synthesis within the Base Action Model (BAM). To strengthen the flexibility and generalization of our workflow, we further investigate the use of open-source teacher models. Specifically, we replace GPT-4o with Llama-3.1-70B to generate training data for the BAM and retrain SHARE-8B accordingly. To clarify the distinction, we refer to the original version of the SHARE-8B model trained with GPT-4o as **SHARE-gpt**, and the SHARE-8B model trained with Llama-3.1-70B as the teacher model as **SHARE-llama**. Table 7 shows the performance of various generator models assisted by SHARE-gpt and SHARE-llama on the BIRD dev set. It demonstrates that SHARE continues to deliver strong performance even when using an open-source model as the teacher model, and it still works well with a variety of generator models. This suggests that our approach is not restricted by the error patterns of any one model in text-to-SQL, but rather leverages broader knowledge of SQL correction.
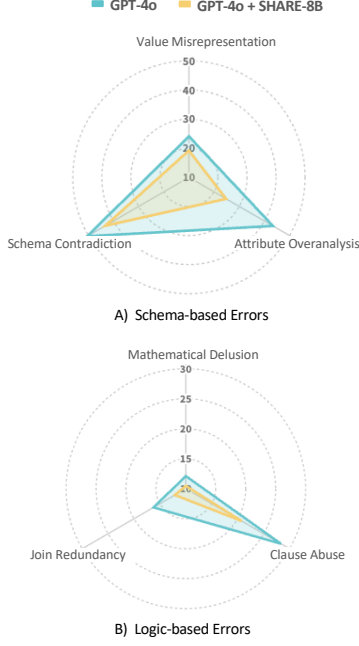
Figure 5: The correction performance across fine-grained error categories on BIRD.

## 4.7 Quantitative Analysis of SQL Error Corrections

To quantify the effectiveness of SHARE in error correction, we analyze the SQL queries generated by the GPT-4o baseline. Following TA-SQL (Qu et al., 2024), we categorize the observed errors into two primary types: **schema-based errors** and **logic-based errors**. Each category is further subdivided into three specific subtypes. For detailed definitions and examples of these error types, we refer the reader to TA-SQL.

Table 5 presents the correction performance across fine-grained error categories on the BIRD dev set. The results indicate that SHARE is effective in mitigating both schema-based and logic-based errors, showing substantial reductions in `Attribute Overanalysis` (↓ 18.61%), `Schema Contradiction` (↓ 7.24%), and `Clause Abuse` (↓ 7.54%). However, SHARE demonstrates limited effectiveness in correcting instances of `Mathematical Delusion` (↓ 1.63%), which may be attributed to the restricted mathematical reasoning capability of the underlying generator model (Mirzadeh et al., 2025). To further elucidate SHARE's correction behavior, we conduct qualitative case studies and present a representative example in Appendix C.6

## 5 Related Work

**LLMs for Text-to-SQL.** In recent years, large language models (LLMs) (Ouyang et al., 2022b; Anthropic, 2024; Team et al., 2024) have attracted considerable attention due to their robust reasoning and domain generalization capabilities. The application of LLMs has improved the performance of text-to-SQL to another level of intelligence. Early research leverages in-context learning capabilities of LLMs to develop text-to-SQL systems through meticulously crafted prompt engineering methodologies (Rajkumar et al., 2022; Pourreza and Rafiei, 2024; Gao et al., 2024; Qu et al., 2024). With the emergence of language agents (Deng et al., 2024; Gu et al., 2024) as a promising paradigm, recent works (Wang et al., 2024; Talaei et al., 2024; Pourreza et al., 2024) leverage multi-agent architectures to construct more reliable and comprehensive frameworks for text-to-SQL conversion, yielding substantial improvements in empirical performance.

**Self-correction in Text-to-SQL.** Self-correction (Pan et al., 2024b), where LLMs evaluate and refine their initial output, has emerged as a crucial technology to enhance the reliability and accuracy of automated code generation tasks (Gu, 2023; Li et al., 2024b). Self-correction has been widely applied in text-to-SQL tasks. Some studies (Lee et al., 2024; Gao et al., 2024) leverage carefully designed prompts to guide LLMs in utilizing their intrinsic reasoning capabilities for more effective self-correction. Additionally, using feedback to effectively guide the self-correction process of LLMs is another promising approach. The sources of feedback are diverse, ranging from human annotations (Pourreza and Rafiei, 2024), external execution environments (Chen et al., 2024a; Wang et al., 2024) to iterative exploration conducted by the LLM itself (Askari et al., 2024; Xia et al., 2024).

## 6 Conclusion

In this research, we propose SHARE, an SLM-based hierarchical action correction assistant designed to enable more precise error localization and effective self-correction for LLMs. We further propose a novel hierarchical evolution strategy for data-efficient training. Experimental results show the effectiveness and robustness of our method, even in low-source settings, unlocking the potential of SLMs in self-correction for text-to-SQL tasks.

## 7 Limitation

In this paper, we demonstrate the effectiveness of SHARE by presenting the performance of SHARE-assisted self-correction in the one-turn setting, where the generator model receives feedback generated by SHARE and performs a single-revision iteration. We leave investigating the performance of our framework in multi-turn interactive self-correction scenarios, where the correction process undergoes multiple refinement cycles, as our future work. Furthermore, our current work exclusively focuses on the text-to-SQL domain. Expanding SHARE to broader code generation tasks represents another key direction for future research.

## 8 Acknowledgement

## 9 Ethical Statement

All datasets employed in this work are publicly accessible. We will also release our models and source code after the review process, ensuring the transparency and reproducibility of our findings. Furthermore, the output generated by our investigations is structured as SQL queries—a programming language format—rather than natural language text, which could potentially involve harmful or biased content. Our team meticulously examines each output to confirm the absence of politically sensitive or biased material. Finally, it is noteworthy that we utilize parameter-efficient LoRA fine-tuning to train our models, which demonstrates superior environmental sustainability compared to full-parameter fine-tuning.

# References

Anthropic. 2024. Claude 3 haiku: our fastest model yet.

Arian Askari, Christian Poelitz, and Xinye Tang. 2024. Magic: Generating self-correction guideline for in-context text-to-sql.

Kamran Ahmad Awan, Ikram Ud Din, Ahmad Almogren, and Joel J. P. C. Rodrigues. 2023. Privacy-preserving big data security for iot with federated learning and cryptography. *IEEE Access*.

Lukas Berglund, Meg Tong, Maximilian Kaufmann, Mikita Balesni, Asa Cooper Stickland, Tomasz Korbak, and Owain Evans. 2024. The reversal curse: Llms trained on "a is b" fail to learn "b is a". In *The Twelfth International Conference on Learning Representations, ICLR 2024, Vienna, Austria, May 7-11, 2024*.

Bei Chen, Fengji Zhang, Anh Nguyen, Daoguang Zan, Zeqi Lin, Jian-Guang Lou, and Weizhu Chen. 2023. Codet: Code generation with generated tests. In *The Eleventh International Conference on Learning Representations, ICLR 2023, Kigali, Rwanda, May 1-5, 2023*.

Xinyun Chen, Maxwell Lin, Nathanael Schärli, and Denny Zhou. 2024a. Teaching large language models to self-debug. In *The Twelfth International Conference on Learning Representations, ICLR 2024, Vienna, Austria, May 7-11, 2024*.

Xinyun Chen, Maxwell Lin, Nathanael Schärli, and Denny Zhou. 2024b. Teaching large language models to self-debug. In *The Twelfth International Conference on Learning Representations, ICLR 2024, Vienna, Austria, May 7-11, 2024*.

Xiang Deng, Yu Gu, Boyuan Zheng, Shijie Chen, Sam Stevens, Boshi Wang, Huan Sun, and Yu Su. 2024. Mind2web: Towards a generalist agent for the web. *Advances in Neural Information Processing Systems*.

Longxu Dou, Yan Gao, Mingyang Pan, Dingzirui Wang, Wanxiang Che, Jian-Guang Lou, and Dechen Zhan. 2023. Unisar: a unified structure-aware autoregressive language model for text-to-sql semantic parsing. *Int. J. Mach. Learn. Cybern.*

Nouha Dziri, Andrea Madotto, Osmar Zaïane, and Avishek Joey Bose. 2021. Neural path hunter: Reducing hallucination in dialogue systems via path grounding. In *Proceedings of the 2021 Conference on Empirical Methods in Natural Language Processing, EMNLP 2021, Virtual Event / Punta Cana, Dominican Republic, 7-11 November, 2021*.

Yujian Gan, Xinyun Chen, and Matthew Purver. 2021. Exploring underexplored limitations of cross-domain text-to-sql generalization. In *Proceedings of the 2021 Conference on Empirical Methods in Natural Language Processing, EMNLP 2021, Virtual Event / Punta Cana, Dominican Republic, 7-11 November, 2021*.

Dawei Gao, Haibin Wang, Yaliang Li, Xiuyu Sun, Yichen Qian, Bolin Ding, and Jingren Zhou. 2024. Text-to-sql empowered by large language models: A benchmark evaluation. *Proc. VLDB Endow.*

Zhibin Gou, Zhihong Shao, Yeyun Gong, Yelong Shen, Yujiu Yang, Nan Duan, and Weizhu Chen. 2024. CRITIC: large language models can self-correct with tool-interactive critiquing. In *The Twelfth International Conference on Learning Representations, ICLR 2024, Vienna, Austria, May 7-11, 2024*.

Qiuhan Gu. 2023. Llm-based code generation method for golang compiler testing. In *Proceedings of the 31st ACM Joint European Software Engineering Conference and Symposium on the Foundations of Software Engineering, ESEC/FSE 2023, San Francisco, CA, USA, December 3-9, 2023*.

Yu Gu, Yiheng Shu, Hao Yu, Xiao Liu, Yuxiao Dong, Jie Tang, Jayanth Srinivasa, Hugo Latapie, and Yu Su. 2024. Middleware for llms: Tools are instrumental for language agents in complex environments. *arXiv preprint arXiv:2402.14672*.

Edward J. Hu, Yelong Shen, Phillip Wallis, Zeyuan Allen-Zhu, Yuanzhi Li, Shean Wang, Lu Wang, and Weizhu Chen. 2022. Lora: Low-rank adaptation of large language models. In *The Tenth International Conference on Learning Representations, ICLR 2022, Virtual Event, April 25-29, 2022*.

Jie Huang, Xinyun Chen, Swaroop Mishra, Huaixiu Steven Zheng, Adams Wei Yu, Xinying Song, and Denny Zhou. 2024. Large language models cannot self-correct reasoning yet. In *The Twelfth International Conference on Learning Representations, ICLR 2024, Vienna, Austria, May 7-11, 2024*.

Ziwei Ji, Nayeon Lee, Rita Frieske, Tiezheng Yu, Dan Su, Yan Xu, Etsuko Ishii, Yejin Bang, Andrea Madotto, and Pascale Fung. 2023. Survey of hallucination in natural language generation. *ACM Comput. Surv.*

Jared Kaplan, Sam McCandlish, Tom Henighan, Tom B. Brown, Benjamin Chess, Rewon Child, Scott Gray, Alec Radford, Jeffrey Wu, and Dario Amodei. 2020. Scaling laws for neural language models.

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

Wenqiang Lei, Weixin Wang, Zhixin Ma, Tian Gan, Wei Lu, Min-Yen Kan, and Tat-Seng Chua. 2020. Re-examining the role of schema linking in text-to-SQL. In *Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing (EMNLP)*.

Jinyang Li, Binyuan Hui, Reynold Cheng, Bowen Qin, Chenhao Ma, Nan Huo, Fei Huang, Wenyu Du, Luo

Si, and Yongbin Li. 2023. Graphix-t5: Mixing pre-trained transformers with graph-aware layers for text-to-sql parsing. In *Thirty-Seventh AAAI Conference on Artificial Intelligence, AAAI 2023, Thirty-Fifth Conference on Innovative Applications of Artificial Intelligence, IAAI 2023, Thirteenth Symposium on Educational Advances in Artificial Intelligence, EAAI 2023, Washington, DC, USA, February 7-14, 2023*.

Jinyang Li, Binyuan Hui, Ge Qu, Jiaxi Yang, Binhua Li, Bowen Li, Bailin Wang, Bowen Qin, Ruiying Geng, Nan Huo, et al. 2024a. 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*.

Jinyang Li, Nan Huo, Yan Gao, Jiayi Shi, Yingxiu Zhao, Ge Qu, Yurong Wu, Chenhao Ma, Jian-Guang Lou, and Reynold Cheng. 2024b. Tapilot-crossing: Benchmarking and evolving llms towards interactive data analysis agents.

Zhenwen Li and Tao Xie. 2024. Using llm to select the right sql query from candidates.

Huihui Liu, Yiding Yang, and Xinchao Wang. 2021. Overcoming catastrophic forgetting in graph neural networks. In *Thirty-Fifth AAAI Conference on Artificial Intelligence, AAAI 2021, Thirty-Third Conference on Innovative Applications of Artificial Intelligence, IAAI 2021, The Eleventh Symposium on Educational Advances in Artificial Intelligence, EAAI 2021, Virtual Event, February 2-9, 2021*.

Xiping Liu and Zhao Tan. 2024. Epi-sql: Enhancing text-to-sql translation with error-prevention instructions.

Seyed-Iman Mirzadeh, Keivan Alizadeh, Hooman Shahrokhi, Oncel Tuzel, Samy Bengio, and Mehrdad Farajtabar. 2025. Gsm-symbolic: Understanding the limitations of mathematical reasoning in large language models. In *The Thirteenth International Conference on Learning Representations, ICLR 2025, Singapore, April 24-28, 2025*.

Long Ouyang, Jeffrey Wu, Xu Jiang, Diogo Almeida, Carroll L. Wainwright, Pamela Mishkin, Chong Zhang, Sandhini Agarwal, Katarina Slama, Alex Ray, John Schulman, Jacob Hilton, Fraser Kelton, Luke Miller, Maddie Simens, Amanda Askell, Peter Welinder, Paul F. Christiano, Jan Leike, and Ryan Lowe. 2022a. Training language models to follow instructions with human feedback.

Long Ouyang, Jeffrey Wu, Xu Jiang, Diogo Almeida, Carroll L. Wainwright, Pamela Mishkin, Chong Zhang, Sandhini Agarwal, Katarina Slama, Alex Ray, John Schulman, Jacob Hilton, Fraser Kelton, Luke Miller, Maddie Simens, Amanda Askell, Peter Welinder, Paul F. Christiano, Jan Leike, and Ryan Lowe. 2022b. Training language models to follow instructions with human feedback. In *Advances in Neural Information Processing Systems 35: Annual Conference on Neural Information Processing Systems 2022, NeurIPS 2022, New Orleans, LA, USA, November 28 - December 9, 2022*.

Liangming Pan, Michael Saxon, Wenda Xu, Deepak Nathani, Xinyi Wang, and William Yang Wang. 2024a. Automatically correcting large language models: *Surveying the Landscape of Diverse Automated Correction Strategies*. Trans. Assoc. Comput. Linguistics.

Liangming Pan, Michael Saxon, Wenda Xu, Deepak Nathani, Xinyi Wang, and William Yang Wang. 2024b. Automatically correcting large language models: Surveying the landscape of diverse automated correction strategies. *Transactions of the Association for Computational Linguistics*.

Tianyu Peng and Jiajun Zhang. 2024. Enhancing knowledge distillation of large language models through efficient multi-modal distribution alignment.

Mohammadreza Pourreza, Hailong Li, Ruoxi Sun, Yeounoh Chung, Shayan Talaei, Gaurav Tarlok Kakkar, Yu Gan, Amin Saberi, Fatma Ozcan, and Sercan O. Arik. 2024. Chase-sql: Multi-path reasoning and preference optimized candidate selection in text-to-sql.

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*.

Bowen Qin, Lihan Wang, Binyuan Hui, Bowen Li, Xiangpeng Wei, Binhua Li, Fei Huang, Luo Si, Min Yang, and Yongbin Li. 2022. SUN: exploring intrinsic uncertainties in text-to-sql parsers. In *Proceedings of the 29th International Conference on Computational Linguistics, COLING 2022, Gyeongju, Republic of Korea, October 12-17, 2022*.

XiPeng Qiu, TianXiang Sun, YiGe Xu, YunFan Shao, Ning Dai, and XuanJing Huang. 2020. Pre-trained models for natural language processing: A survey. *Science China Technological Sciences*.

Ge Qu, Jinyang Li, Bowen Li, Bowen Qin, Nan Huo, Chenhao Ma, and Reynold Cheng. 2024. Before generation, align it! A novel and effective strategy for mitigating hallucinations in text-to-sql generation. Association for Computational Linguistics.

Rafael Rafailov, Archit Sharma, Eric Mitchell, Christopher D. Manning, Stefano Ermon, and Chelsea Finn. 2023. Direct preference optimization: Your language model is secretly a reward model. In *Advances in Neural Information Processing Systems 36: Annual Conference on Neural Information Processing Systems 2023, NeurIPS 2023, New Orleans, LA, USA, December 10 - 16, 2023*.

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

Sithursan Sivasubramaniam, Cedric Osei-Akoto, Yi Zhang, Kurt Stockinger, and Jonathan Fuerst. 2024. Sm3-text-to-query: Synthetic multi-model medical text-to-query benchmark.

Shayan Talaei, Mohammadreza Pourreza, Yu-Chen Chang, Azalia Mirhoseini, and Amin Saberi. 2024. Chess: Contextual harnessing for efficient sql synthesis.

Gemini Team, Rohan Anil, Sebastian Borgeaud, Jean-Baptiste Alayrac, Jiahui Yu, Radu Soricut, Johan Schalkwyk, and Andrew M. 2024. Gemini: A family of highly capable multimodal models.

Bing Wang, Changyu Ren, Jian Yang, Xinnian Liang, Ji-aqi Bai, Linzheng Chai, Zhao Yan, Qian-Wen Zhang, Di Yin, Xing Sun, and Zhoujun Li. 2024. Mac-sql: A multi-agent collaborative framework for text-to-sql.

Xuezhi Wang, Jason Wei, Dale Schuurmans, Quoc V. Le, Ed H. Chi, Sharan Narang, Aakanksha Chowdhery, and Denny Zhou. 2023. Self-consistency improves chain of thought reasoning in language models. In *The Eleventh International Conference on Learning Representations, ICLR 2023, Kigali, Rwanda, May 1-5, 2023*.

Hanchen Xia, Feng Jiang, Naihao Deng, Cunxiang Wang, Guojiang Zhao, Rada Mihalcea, and Yue Zhang. 2024. $r^3$: "this is my sql, are you with me?" a consensus-based multi-agent system for text-to-sql tasks.

Yuanzhen Xie, Xinzhou Jin, Tao Xie, Matrixmxlin Matrixmxlin, Liang Chen, Chenyun Yu, Cheng Lei, Chengxiang Zhuo, Bo Hu, and Zang Li. 2024. Decomposition for enhancing attention: Improving llm-based text-to-sql through workflow paradigm. In *Findings of the Association for Computational Linguistics, ACL 2024, Bangkok, Thailand and virtual meeting, August 11-16, 2024*.

Shusheng Xu, Wei Fu, Jiaxuan Gao, Wenjie Ye, Weilin Liu, Zhiyu Mei, Guangju Wang, Chao Yu, and Yi Wu. 2024a. Is dpo superior to ppo for llm alignment? a comprehensive study. In *Proceedings of the 41st International Conference on Machine Learning*.

Xiaohan Xu, Ming Li, Chongyang Tao, Tao Shen, Reynold Cheng, Jinyang Li, Can Xu, Dacheng Tao, and Tianyi Zhou. 2024b. A survey on knowledge distillation of large language models.

Pengcheng Yin and Graham Neubig. 2017. A syntactic neural model for general-purpose code generation. In *Proceedings of the 55th Annual Meeting of the Association for Computational Linguistics, ACL 2017, Vancouver, Canada, July 30 - August 4, Volume 1: Long Papers*.

Tao Yu, Rui Zhang, Kai Yang, Michihiro Yasunaga, Dongxu Wang, Zifan Li, James Ma, Irene Li, Qingning Yao, Shanelle Roman, Zilin Zhang, and Dragomir R. Radev. 2018. 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, Brussels, Belgium, October 31 - November 4, 2018*.

Jianguo Zhang, Tian Lan, Ming Zhu, Zuxin Liu, Thai Hoang, Shirley Kokane, Weiran Yao, Juntao Tan, Akshara Prabhakar, Haolin Chen, Zhiwei Liu, Yihao Feng, Tulika Awalgaonkar, Rithesh Murthy, Eric Hu, Zeyuan Chen, Ran Xu, Juan Carlos Niebles, Shelby Heinecke, Huan Wang, Silvio Savarese, and Caiming Xiong. 2024. xlam: A family of large action models to empower ai agent systems.

Xuanliang Zhang, Dingzirui Wang, Longxu Dou, Qingfu Zhu, and Wanxiang Che. 2025. A survey of table reasoning with large language models. *Frontiers Comput. Sci.*

Wayne Xin Zhao, Kun Zhou, Junyi Li, Tianyi Tang, Xiaolei Wang, Yupeng Hou, Yingqian Min, Beichen Zhang, Junjie Zhang, Zican Dong, Yifan Du, Chen Yang, Yushuo Chen, Zhipeng Chen, Jinhao Jiang, Ruiyang Ren, Yifan Li, Xinyu Tang, Zikang Liu, Peiyu Liu, Jian-Yun Nie, and Ji-Rong Wen. 2024. A survey of large language models.

Xuanle Zhao, Xianzhen Luo, Qi Shi, Chi Chen, Shuo Wang, Wanxiang Che, Zhiyuan Liu, and Maosong Sun. 2025. Chartcoder: Advancing multimodal large language model for chart-to-code generation.

Yaowei Zheng, Richong Zhang, Junhao Zhang, Yanhan Ye, and Zheyan Luo. 2024. 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)*. Association for Computational Linguistics.

Ruiqi Zhong, Charlie Snell, Dan Klein, and Jason Eisner. 2023. Non-programmers can label programs indirectly via active examples: A case study with text-to-sql. In *Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing, EMNLP 2023, Singapore, December 6-10, 2023*, pages 5126–5152.

| Model | #Input | #Output |
|---|---|---|
| Base Action Model | 196.23 | 25.71 |
| Schema Augmentation Model | 621.46 | 30.24 |
| Logic Optimation Model | 683.69 | 29.80 |

Table 8: The average length of input and output tokens in the training corpus for each model.

| Shorthand Name | Official Model Alias |
|---|---|
| Llama-3.1-70B | `Llama-3.1-70B-Instruct` |
| Qwen-Coder-32B | `Qwen2.5-Coder-32B-Instruct` |
| Llama-3.1-8B | `Llama-3.1-8B-Instruct` |
| DS-Coder-6.7B | `deepseek-coder-6.7b-instruct` |

Table 9: Shorthand model names used throughout the paper and their corresponding official model aliases.

# A SHARE Recipe

## A.1 Training Data Distribution

In this work, we utilize the training set of two mainstream cross-domain text-to-SQL benchmarks as the seed data to construct our own training data. 1) **BIRD**: The training set of BIRD contains 9,428 pairs of text-to-SQL data and 69 big databases across 37 professional domains. 2) **SPIDER**: SPIDER is a more standard text-to-SQL benchmark that contains 8659 training examples across more than 20 domains. Table 8 displays the detailed average length of the input and output tokens for each model in our own training data. Our approach achieved superior performance with notably few training tokens, validating the data-efficient nature of SHARE.

## A.2 Action Space of SHARE

As discussed in Section 3.2, we design actions in target trajectories as pandas-like APIs and employ GPT-4o to convert each ground truth SQL to a verified action trajectory by few-shot prompting for the training data construction. All logically meaningful and validated actions generated by GPT-4o in this process are collected as the action space of SHARE. We further categorize these actions into four types, including clause, dataframe, aggregation, and operator types. Specific actions are presented in Table 11 and 12.

| Approach | Simple | Moderate | Challenging | Total |
|---|---|---|---|---|
| DPO | 65.08 | 48.06 | 42.76 | 57.82 |
| LoRA | **70.81** | **56.25** | **46.90** | **64.14** |

Table 10: The performance of SHARE trained via different fine-tuning approaches on BIRD dev.

## A.3 Erroneous Trajectory Collection

As introduced in Section 3.4, the erroneous action trajectories in the training data for LOM come from two resources: 1) initial erroneous SQLs, and 2) verified action trajectories that are perturbed by our action-based perturbation strategy. To better clarify the data augmentation through error perturbations, we present the pseudocode of this process in Pseudocode 1.

---
**Algorithm 1** Data Augmentation

---
1: **Inputs:**
2: $verified\_corrs$ ▷ A collection of verified correct trajectories $verified\_corr\_1$, $verified\_corr\_2, \dots$
3: $SLM$ ▷ Small Language Model
4: $K$ ▷ Number of error perturbations per ground truth
5: **Output:**
6: $data\_pairs$ ▷ A set of (erroneous trajectory, verified correct trajectory) pairs

7: **Initialize:**
8: Initialize $data\_pairs$ as an empty set.

9: **for** $verified\_corr\_i$ **in** $verified\_corrs$ **do**
10:     **for** $p = 1$ to $K$ **do**
11:         $er\_p \leftarrow$ ErrorPerturb($verified\_corr\_i, SLM$) ▷ Use the smaller model to inject errors into the verified correct trajectory $verified\_corr\_i$.
12:         Add the pair $(er\_p, verified\_corr\_i)$ to $data\_pairs$ ▷ No additional verification is needed, since $verified\_corr\_i$ is correct by definition.
13:     **end for**
14: **end for**
15: **return** $data\_pairs$

---

# B Experiment Setup

## B.1 Baseline Methods

Table 2 presents all the baseline methods we use for comparison. In this work, we consider Self-Correction (Huang et al., 2024), Self-Consistency (Wang et al., 2023), Multiple-Prompt (Lee et al., 2024) as feedback-independent self-correction baselines. These methods leverage the intrinsic capabilities of the LLM through prompt engineering to enable its self-correction. For feedback-

dependent self-correction baselines, we implement Self-Debugging (Chen et al., 2024a), DIN-Correction (Pourreza and Rafiei, 2024), MAC-Refiner (Wang et al., 2024), and MAGIC (Askari et al., 2024), where the Large Language Model (LLM) refines its initial output under the guidance of feedback. Feedback utilized in these approaches is derived from three primary sources: human annotations, external execution environments, and through LLMs iteratively exploring contextual environments.

The details of these methods are as follows: 1) **Self-Correction** is a naive method in which the LLM reconsiders and refines its outputs through vanilla Chain-of-Thought (CoT) prompting. 2) **Self-Consistency** (Wang et al., 2023) is a method that refines the initial output by exploring a broader search space and selecting the most consistent one. We generate five SQL queries using the baseline SQL generation prompt implemented in BIRD (Li et al., 2024a) and consider an instance to be correctly solved if at least one of the generated SQL queries produces the correct result. 3) **Multiple-Prompt** generates diverse queries by systematically reordering candidate tables within the prompt. Following the implementation of (Lee et al., 2024), we generate up to five combinations of prompts for each instance and employ the same evaluation mechanism as self-consistency to determine the results. 4) **Self-Debugging** (Chen et al., 2024b) generates the feedback by investing the execution results and explaining the generated SQL in natural language. The feedback further serves as guidance to instruct the LLM to self-correct. 5) **DIN-Correction** utilizes the human-annotated guideline from DIN-SQL (Pourreza and Rafiei, 2024) for self-correction. 6) **MAC-Refiner**, which is a sub-agent of MAC-SQL (Wang et al., 2024), implements the self-correction process based on multi-dimensional error feedback information by analyzing the execution results, including syntactic correctness, execution feasibility, and retrieval of non-empty results. 7) **Magic** (Askari et al., 2024) collaborates on the failure experiences and automatically distills correction guidelines, employing a crafted LLM-based multi-agent framework for self-correction.

## B.2 Hyper-Parameters

We set the low-rank dimensions as 8, the learning rate as $5e^{-5}$, and the batch size as 8. We train 5 epochs for the Base Action Model (BAM) and train 3 epochs for the Schema Augmentation Model (SAM) and the Logic Optimization Model (LOM). During inference, we set the temperature as 0.1, the top p as 0.95, and the maximum sample length as 1024. We report the experimental results as the average of five repeated trials.

## B.3 Model Reference Mapping

We list in Table 9 the shorthand model names used throughout the paper alongside their corresponding official model aliases.

## C Further Anaylsis

### C.1 Ablation Study

Table 5 presents the results of our ablation study, aimed at isolating and evaluating contributions of each component of SHARE. As shown in Table 5(a)-(b), the substantial performance degradation resulting from the removal of either refinement action model underscores the importance of our **two-stage refinement** design. This two-stage approach, which separately handles schema linking and logical reasoning, proves to be essential for effective text-to-SQL correction (Lei et al., 2020).

To further examine the benefits of the **hierarchical evolution strategy** in SHARE, we compare it against a sequential training pipeline reminiscent of classical continual learning. When trained sequentially, performance declines by 3.59% (Table 5(c)), suggesting that later models may be disproportionately influenced by biases or errors carried over from earlier training stages (Liu et al., 2021). In contrast, SHARE employs hierarchical evolution during training and strategically integrates knowledge at inference time, mitigating these limitations and achieving superior results.

Finally, Table 5(d) highlights the effectiveness of **action-based error perturbation** as a data augmentation strategy. Although simple, it yields a 2.76% improvement in performance, reinforcing its value as a straightforward yet potent enhancement to SHARE's overall text-to-SQL reasoning capability.

### C.2 Computational Cost of SHARE

As illustrated in Section 4.4, we compare the average token usage per instance and average computational cost per 1000 instances during the inference and training stages for GPT-4o assisted by SHARE,

versus other strong LLM-based text-to-SQL correction approaches on the BIRD development set. We present the usage of input and output tokens separately since they have different prices. We take the same price for the calculation as in the previous work (Sivasubramaniam et al., 2024) [2].

Notably, SHARE remains highly cost-efficient during the training data construction stage. As introduced in Section 3.2, SHARE just uses GPT-4o to generate and verify the training data for the Base Action Model (BAM). Afterward, BAM itself creates the training data for all other models, enabling a self-evolution process without further reliance on GPT-4o. To more clearly demonstrate SHARE's cost-efficiency in training, we compare it against the strongest self-correction method, MAGIC. While MAGIC does not undergo a direct fine-tuning step, it adopts an ICL-based training approach, which follows a training-like procedure where the LLM explores the training set, produces a correction guideline for each instance, and memorizes successful corrections as task knowledge for inference. As shown in Table 6, the sharp reduction in token usage directly translates to significant computational savings, underscoring SHARE's superior cost-efficiency in the training stage.

## C.3 Independent Inference of SHARE

While SHARE shows clear improvements among various generator LLMs, **a key question** is how well it performs independently, without any external LLM grounding. This question is especially important for privacy-sensitive scenarios, such as text-to-SQL tasks on confidential relational databases.

To address this, we evaluate SHARE's standalone capabilities, using the same orchestration prompting strategy as for `Llama-3.1-8B Orchestration`. As shown in Figure 6, SHARE-8B achieves robust performance on its own, surpassing both `Llama-3.1-8B` and `Llama-3.1-8B Orchestration`, and approaching the quality of strong proprietary models like GPT-35-Turbo, Claude-3.5-Sonnet, and GPT-4o-mini.

This improvement results from our parameter-efficient LoRa fine-tuning approach, which augments the model's capabilities without altering its original parameters. By optimizing only a small set of new parameters, we enhance the underlying

---

[2]Pricing of GPT-4o API: https://openai.com/api/pricing/. Price of Llama-3.1 8B usage: https://groq.com/pricing/
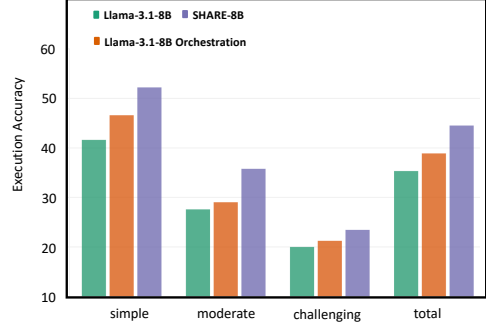


Figure 6: Independent inference performance of SHARE on BIRD.

SLM's text-to-SQL performance without compromising its existing strengths. Moreover, this efficient design supports secure on-device deployment that does not require exposing sensitive database content, relying solely on an 8B-scale SLM and a modest set of additional parameters.

## C.4 Analysis of Overcorrections

Overcorrection (Pan et al., 2024a), which refers to the modification of initially correct SQL queries into incorrect ones, is a notorious challenge in self-correction. To mitigate this issue, we tailored the training data for the Logic Optimization Model (LOM). Specifically, we consider pairs of (erroneous trajectory, verified trajectory) as positive samples and, at a 4:1 ratio, introduce (verified trajectory, verified trajectory) pairs as negative samples.

Analysis of the final correction results shows that, compared to the current leading self-correction method, MAGIC, SHARE reduced the overcorrection rate from **15.52%** to **11.20%**. This demonstrates that our design mitigates overfitting in the correction process in a simple but effective way, preventing the SLMs from treating all trajectories as erroneous and unnecessarily fixing them.

## C.5 Fine-tuning Approach Exploration

Apart from LoRA, we also explored training the SLMs in SHARE using Direct Preference Optimization (DPO) (Rafailov et al., 2023) during our initial attempt. In the data construction stage, verified action trajectories are considered as the chosen responses. We employ the action-based perturbation strategy introduced in Section 3.4 and provide 5-shot examples to guide GPT-4o to generate rejected responses based on chosen responses.

As shown in Table 10, compared to SHARE-

8B trained via LoRA, the performance gains of the DPO-trained version showed a consistent decline across all difficulty levels. We observe that the DPO model outputs many meaningless trajectories, such as invalid pandas output format. We speculate that the suboptimal performance of DPO-based training arises from its limited suitability for low-resource settings. Furthermore, for data science code generation tasks (Zhao et al., 2025) that require complex reasoning, DPO relies on high-quality data and requires crafted rejected responses by human annotation (Xu et al., 2024a). In contrast, LoRA can effectively leverage automatically generated data, making it more efficient under these constraints. Consequently, we adopted LoRA as our primary fine-tuning method, given its lightweight design and effectiveness.

### C.6 Case Study

In order to provide deeper insights into SHARE's effectiveness, we conducted case studies on its correction outputs. Table 13 presents an illustrative case randomly selected from the BIRD development set. By decomposing the declarative SQL query into an action trajectory and leveraging the multi-SLM orchestration during inference, SHARE achieves more granular corrections. Consequently, it produces a correct query that effectively addresses both schema and logic issues. We additionally offer an output log that contains the output of each model during SHARE inference for reference.

## D Implement Prompts

### D.1 Data Construction for BAM

Figure 7 illustrates the prompts employed in our data construction process for BAM. Specifically, we utilize seven examples to guide GPT-4o in generating corresponding action trajectories for SQL queries. The figure presents two representative examples, while the remaining examples can be found in our code.

### D.2 Error Perturbation Prompt

Figure 8 illustrates the prompts we use to implement error perturbations during the data construction process of LOM. Specifically, we utilize four examples to guide BAM to generate perturbed action trajectories. The figure presents a representative example, while the remaining examples can be found in our code.

### D.3 Inference Prompts

In this section, we present prompts that we use in the inference process of SHARE. During inference, we prompt BAM to convert the initial SQL using the prompt in Figure 9. Subsequently, the generated trajectory is forwarded to SAM. The SAM first masks the schema using the prompt in Figure 10 and then generates the schema-refined trajectory using the prompt in Figure 11. The LOM then takes this trajectory as input and refines it as the final output of SHARE using the prompt in Figure 12. Finally, the SHARE-produced trajectory serves as feedback to facilitate the self-correction of LLM, as shown in Figure 13. All the prompts are in the zero-shot setting.

## E Reproducebility

We fine-tune and infer open-source models, including `Llama-3.1-8b-Instruct` and `Phi-3-mini-4K-instruct` (3.8B) via `LlamaFactory`[3]. The `Llama-3.1-70B-Instruct` model is inferred using `vllm`[4]. To accelerate its inference process, we also implemented `deepspeed`[5]. All open-source models are accessed via `huggingface`[6].

We will open-source the source code along with the training data, model checkpoints, and prompts in each stage after the anonymous review phase.

---

[3]https://github.com/hiyouga/LLaMA-Factory
[4]https://github.com/vllm-project/vllm
[5]https://github.com/microsoft/DeepSpeed
[6]https://huggingface.co/models

| Category | Action | Expression | Explanation |
|---|---|---|---|
| Clause | SELECT | select(elements) | Select data from the database. **Parameters:** **elements** - the selected elements. Valid elements include qualified column names in the format of `table.column`, aggregate functions, or any valid SQL-syntax selectable entity. Multiple elements are separated by commas. |
| | WHERE | where(element, filter) | Filter rows by conditions. **Parameters:** **elements** - The qualified column, or expression to be filtered. **filter** - The condition that determines which values from the element are included. |
| | GROUP BY | groupby(elements) | Groups rows that have the same values into summary rows. **Parameters:** **elements** - Qualified columns, or expressions used to group the data. Multiple elements are separated by commas. |
| | HAVEING | having(element, filter) | Filter groups of data with aggregate functions. **Parameters:** **elements** - Qualified columns, or expressions used to filter data. Multiple elements are separated by commas. **filter** - The condition that determines which values from the element are included. |
| | ORDER BY | orderby(by, order) | Sort the result set based on a qualified column or expression. **Parameters:** **by** - The qualified column or expression used to sort the data. **order** - The sorted order. It should be `DESC` or `ASC`. |
| | LIMIT | limit(num) | Restrict the number of rows returned. **Parameters:** **num** - The num is a flexible input that specifies the type of limitation to apply. For instance, `limit(1)` denotes limiting the number of rows to 1. `limit(2,9)` denotes specifying a range of columns to return, which is columns 2 through 9. |
| | DISTINCT | distinct(element) | Remove duplicate rows from the result set. **Parameters:** **element** - The qualified column to be processed. |
| Dataframe | UNION | df1.union(df2) | Union the result set of two dataframes. |
| | INTERSECT | df1.intersect(df2) | Intersect the result set of two dataframes. |
| | EXCEPT | df1.except(df2) | Subtract the result of df2 from the result of df1. |
| Aggregation | SUM | sum(element) | Sum all non-null values of the element or expression. **Parameters:** **element** - The target qualified column or expression. |
| | AVG | average(element) | Calculate the average of all non-null values for the element or expression. **Parameters:** **element** - The target qualified column or expression. |
| | COUNT | count(element) | Returns the total number of the element or expression. **Parameters:** **element** - The target qualified column or expression. |
| | MIN | min(element) | Minimum value in the element or expression except for nULL. **Parameters:** **element** - The target qualified column or expression. |
| | MAX | max(element) | Maximum value in the element or expression except for nULL. **Parameters:** **element** - The target qualified column or expression. |

Table 11: Actions categorized as clause, dataframe, and aggregation types employed in the construction of action trajectories within SHARE. Actions categorized as operator types are shown in Table 12.

| Category | Action | Expression | Explanation |
|---|---|---|---|
| Operator | `CAST` | cast(element, type) | Convert the element or expression to the target data type. **Parameters:** **element** - The qualified column or expression to convert. **type** - The target data type. |
| | `CALCU-LATION` | +, -, *, / | Implement corresponding calculations for two expressions. |
| | `SUBSTR` | substr(element, piv, len) | Extract a substring, starting from a specified position and optionally for a specified length **Parameters:** **element** - The qualified column to extract. **piv** - The start point. **len** - The extracted length. |

Table 12: Actions categorized as operator actions employed in the construction of action trajectories within SHARE. Actions categorized as other types are shown in Table 11.

| | A Detailed Case |
|---|---|
| **Question** | *Which county reported the most number of school closures in the 1980s with school ownership code belonging to Youth Authority Facilities (CEA)?* |
| **Evidence** | *Youth Authority Facilities (CEA) refers to SOC = 11; 1980s = years between 1980 and 1989.* |
| **Gold** | SELECT County FROM schools WHERE strftime('%Y', ClosedDate) BETWEEN '1980' AND '1989' AND SOC = 11 GROUP BY County ORDER BY COUNT(ClosedDate) DESC LIMIT 1 |
| **Initial SQL** | SELECT County, COUNT(*) AS YearCount FROM schools WHERE Year BETWEEN '1980-01-01' AND '1989-12-31' AND SOC = 11 GROUP BY County ORDER BY YearCount DESC LIMIT 1; |
| **MAGIC** | SELECT County, COUNT(*) as ClosureCount FROM schools WHERE ClosedDate BETWEEN '1980-01-01' AND '1989-12-31' AND SOC = 11 GROUP BY County ORDER BY ClosureCount DESC LIMIT 1; |
| **SHARE** | SELECT schools.County FROM schools WHERE schools.SOC = 11 AND schools.ClosedDate BETWEEN '1980-01-01' AND '1989-12-31' GROUP BY schools.County ORDER BY COUNT(schools.ClosedDate) DESC LIMIT 1; |
| | **Output Log** |
| **Initial SQL** | SELECT County, COUNT(*) AS YearCount FROM schools WHERE Year BETWEEN '1980-01-01' AND '1989-12-31' AND SOC = 11 GROUP BY County ORDER BY YearCount DESC LIMIT 1; |
| **BAM:** | `df1 = df.where(element = schools.Year, filter = 'between 1980-01-01 and 1989-12-31'))` `df2 = df1.where(element = schools.SOC, filter = 11)` `df3 = df2.groupby(schools.County).count(schools.Year)` `df4 = df3.orderby(by = count(schools.Year), desc).limit(1)` `res = df4.select(schools.County, count(schools.Year))` |
| **SAM:** | `df1 = df.where(element = schools.ClosedDate, filter = 'between 1980-01-01 and 1989-12-31'))` `df2 = df1.where(element = schools.SOC, filter = 11)` `df3 = df2.groupby(schools.County).count(schools.ClosedDate)` `df4 = df3.orderby(by = count(schools.ClosedDate), desc).limit(1)` `res = df4.select(schools.County, count(schools.ClosedDate))` |
| **LOM:** | `df1 = df.where(element = schools.ClosedDate, filter = 'between 1980-01-01 and 1989-12-31'))` `df2 = df1.where(element = schools.SOC, filter = 11)` `df3 = df2.groupby(schools.County).count(schools.ClosedDate)` `df4 = df3.orderby(by = count(schools.ClosedDate), desc).limit(1)` `res = df4.select(schools.County)` |
| **Final Output** | SELECT schools.County FROM schools WHERE schools.SOC = 11 AND schools.ClosedDate BETWEEN '1980-01-01' AND '1989-12-31' GROUP BY schools.County ORDER BY COUNT(schools.ClosedDate) DESC LIMIT 1; |

Table 13: Case study: an illustrative case from BIRD dev.

**Data Construction prompt for BAM**

You are a text-to-SQL expert. Action Trajectory (AT) is a piece of stepwise action-based code to show the underlying reasoning of text-to-SQL conversion. Actions utilized in the trajectory are pandas-like functions.
Given the database schema, question, evidence and SQL, your task is to convert the SQL to AT which reflect the accurate logic in the SQL. I'll provide several example to you to help you understand the syntax of the AT and the conversion logic. AT ignore 'join' action. Do not generate 'join' action.

```Examples
{example_shots}
```
Now convert the following SQL to valid AT based on the database schema, question and evidence.

question = "{question}"
schema = "{schema}"
evidence = "{evidence}"
SQL = "{sql}"

``` AT
{at}
``

**Examples: A Partial Showcase**

**#Example 1**
question = "How many movies directed by Francis Ford Coppola have a popularity of more than 1,000? Please also show the critic of these movies. "
schema = [movies.movie_title, ratings.critic, movies.director_name, movies.movie_popularity, ratings.movie_id, movies.movie_id']
evidence = "Francis Ford Coppola refers to director_name; popularity of more than 1,000 refers to movie_popularity >1000"
SQL = "SELECT T2.movie_title, T1.critic FROM ratings AS T1 INNER JOIN movies AS T2 ON T1.movie_id = T2.movie_id WHERE T2.director_name = 'Francis Ford Coppola' AND T2.movie_popularity > 1000 "
```AT
df1 = df.where(element = movies.director_name, filter = 'Francis Ford Coppola')
df2 = df1.where(element = movies.movie_popularity, filter = '> 1000')
res = df2.select(movies.movie_title, ratings.critic)"
```

**#Example 2**
question = "Among the professors who have more than 3 research assistants, how many of them are male? "
schema = [prof.gender, RA.student_id, RA.prof_id, prof.prof_id]
evidence = "research assistant refers to the student who serves for research where the abbreviation is RA; more than 3 research assistant refers to COUNT(student_id) > 3;"
SQL = "SELECT COUNT(*) FROM ( SELECT T2.prof_id FROM RA AS T1 INNER JOIN prof AS T2 ON T1.prof_id = T2.prof_id WHERE T2.gender = 'Male' GROUP BY T1.prof_id HAVING COUNT(T1.student_id) > 3 )"

```AT
df1 = df.groupby(prof.prof_id).having(element = count(RA.student_id), filter = '> 3')
df2 = df1.where(element = 'prof.gender', filter = 'Male')
res = df2.count()
```

⋮

Figure 7: The prompt for the data construction process of BAM.

**Error perturbation** prompt for BAM

You are a text-to-SQL expert. Action Trajectory (AT) is a piece of stepwise action-based code to show the underlying reasoning of text-to-SQL conversion. Actions utilized in the trajectory are pandas-like functions.
I will provide you:
1. Schema: A python list and each element is a `table_name`.`column_name` string. It indicates the table and column you could use in the AT.
2. Column description: For each column in the schema, a column description is given to describe the column meaning, column type and example values in this column.
3. Question: the natural language answer you need to answer in the text-to-SQL process
4. Evidence: the oracle knowledge to help you understand the AT
5. AT: The AT that show the correct logic of the text-to-SQL process in the context of the schema, question and evidence.

Your task is to modify the AT and imitate reasonable errors that might occur in text-to-SQL, generating an erroneous AT. It can be implemented by one or several types of error perturbation.
The type of error perturbation as references:
1. Add: Inserts an additional action into the original action trajectory.
2. Delete: Removes an existing action from the trajectory.
3. Substitute: Replaces an existing action with a different action type or modifies the parameters of the existing action.

I'll also provide you some examples:
{example_shots}

Now generating the AT that contains error for the following text-to-SQL instances:
schema = {schema}
```column description
{column_description}
```
question = "{question}"
evidence = "{evidence}"
```input AT
{at}
```

Fill in the following template using your answer.
```erroneous AT
[Your Answer]
```

**Examples: A Partial Showcase**

**#Example 1**
schema = ['movie_crew.job', 'movie.movie_id', 'movie.revenue', 'person.person_name']
```column description
# movie_crew.job: The 'job' column in the 'movie_crew' table (db id: movies_4) stores text descriptions of various crew members' roles within a movie production, acknowledging that multiple individuals can share the same job title. Example roles include 'Stand In', 'Consulting Producer', and 'Simulation & Effects Production Assistant'.
# movie.movie_id: The unique integer identifier for each movie in the 'movie' table.
# movie.revenue: The 'revenue' column in the 'movie' table (db id: movies_4) records the movie's earnings as an integer, reflecting its financial success.
# person.person_name: The 'person name' column in the 'person' table of the 'movies_4' database stores text entries representing individual names, with examples like 'Matthew Ferguson', 'Joe Guzman', and 'Reilly Dolman'.
```

question = "What is the average revenue of the movie in which Dariusz Wolski works as the director of photography?"
evidence = "director of photography refers to job = 'Director of Photography'; average revenue = divide(sum(revenue), count(movie_id))"
```input AT
df1 = df.where(element = person.person_name, filter = 'Dariusz Wolski')
df2 = df1.where(element = movie_crew.job, filter = 'Director of Photography')
res = df2.select(cast(df2.sum(movie.revenue), real) / df2.count(movie.movie_id))
```

Output1:
```erroneous AT
df1 = df.where(element = person.person_name, filter = 'Dariusz Wolski')
res = df1.select(df1.avg(movie.revenue) / df1.count(movie.movie_id))
```

⋮

Figure 8: The prompt for the error perturbation strategy implemented by BAM.

Figure 9: The prompt for BAM to convert SQL to action trajectory.

Figure 10: The prompt for SAM to generate the schema-masked variant given the input trajectory.

**Fillin prompt for SAM**

You are a text-to-SQL expert. Action Trajectory (AT) is a piece of stepwise action-based code to show the underlying reasoning of text-to-SQL conversion. Actions utilized in the trajectory are pandas-like functions.

I will provide you:

1. Schema: For each table, we will have a python list and each element is a `table_name`.`column_name` string to show all the schema in the database. It indicates that the table and column you could use in the AT.

2. Highlighted Schema: a subset of Schema. You can consider it as a guess about the schema that used in the ground-truth SQL in the context of this text-to-SQL process. However, it is not always correct. It may contain irrelevant schema which could lead to errors in the subsequent SQL generation or miss truly related schema.

3. Question: the natural language answer you need to answer in the text-to-SQL process

4. Evidence: the oracle knowledge to help you generate the AT

5. Masked AT: An AT with the schema masked, leaving only the reasoning steps in text-to-SQL.

Your task is to refer to all the provided information and fill in the correct schema at the [MASK] positions in the masked AT. The complete AT should accurately reflect the reasoning process that generates the SQL capable of correctly answering the question. DO NOT modify the logical template in the masked AT; you are only allowed to fill in the schema.

```Schema
{full_schema}
```

highlighted_schema = {partial_schema}
question = "{question}"
evidence = "{evidence}"
```Masked AT
{masked_at}
```

Now, fill in the masked AT and give me the final AT:
```AT
[Your Answer]
```

Figure 11: The prompt for SAM to reinsert the correct schema in the schema-based variant.

**Inference** prompt for LOM

You are a text-to-SQL expert. Action Trajectory (AT) is a piece of stepwise action-based code to show the underlying reasoning of text-to-SQL conversion and help the subsequent generation of the SQL that can answer the question accurately. Actions utilized in the trajectory are pandas-like functions.
I will provide you:
1. Schema: A python list and each element is a `table_name`.`column_name` string. It indicates that the table and column you could use in the AT.
2. Column description: For each column in the schema, a column description is given to describe the column meaning, column type and example values in this column.
3. Question: the natural language answer you need to answer in the text-to-SQL process
4. Evidence: the oracle knowledge to help you generate the AT
5. AT: AT that show the logic reasoning of the text-to-SQL process in the context of the schema, question and evidence. It may contain errors which could lead to errors in the subsequent SQL generation.

Your task is to check the given AT and modify it when needed. The final goal is to generate valid AT which reflect the accurate logic reasoning in the text-to-SQL based on the schema, column description, question and evidence. Later, the modified AT will be converted to SQL.
Please pay attention that:
1. AT ignore 'join' action. Do not generate 'join' action.
2. In the generated AT, only select the thing that request in the question. Do not select any non-requested stuff.
3. The filter condition in the 'where' function doesn't directly match the text in the question. To find the correct value for the 'where' function, you need to reference the example values or all possible values in column description.

schema = {schema}
```column description
{column_description}
```
question = "{question}"
evidence = "{evidence}"
```AT
{at}
```

Now generate the valid AT that display the reasoning process of generating SQL that can accurately answer the question:
```refined AT
[Your Answer]
```

Figure 12: The prompt for LOM to rectify logic-related errors in the input trajectory.

**SQL Generation** prompt for LLM

You are a text-to-SQL expert. I will provide you:
1. Database schema: Includes schema and forerign_keys. Schema is a python list and each element is a `table_name`.`column_name` string. It indicates that the table and column you could use in the AT. foreign_keys is a dictionary that shows foreign key relationships among tables.
2. Column description: For each column in the schema, a column description is given to describe the column meaning, column type and example values in this column.
3. Question: the natural language answer you need to answer in the text-to-SQL process.
4. Evidence: the oracle knowledge to help you generate the AT.
5. AT: a piece of stepwise action-based code to show the underlying reasoning of text-to-SQL conversion and help the subsequent generation of the SQL that can answer the question accurately.

Your task is to convert the given AT to the valid SQL accordding to the database schema, column description, question and evidence. The SQL should be valid in syntax and can answer the question accurately.
To help you better understand AT, please pay attention that:
1. AT ignore 'join' action. Do not generate 'join' action. You need to generate correct join clause by yourself. Use 'INNER JOIN' in your generated SQL.
2. The filter condition in the 'where' function doesn't directly match the text in the question. To find the correct value for the 'where' function, you need to reference the example values or all possible values in column description.
3. In the generated AT, only select the thing that request in the question. Do not select any non-requested stuff.

schema = {schema}
```column description
{column_description}
```
question = "{question}"
evidence = "{evidence}"
```AT
{at}
```

Now generate the valid SQL:
```sqlite
[Your Answer]
```

Figure 13: The prompt for LLM to generate refined SQL given the trajectory outputted by SHARE.