Feather-SQL: A Lightweight NL2SQL Framework with Dual-Model Collaboration Paradigm for Small Language Models

Anonymous ACL submission

Abstract

Natural Language to SQL (NL2SQL) has seen significant advancements with large language models (LLMs). However, these models often depend on closed-source systems and high computational resources, posing challenges in data privacy and deployment. In contrast, small language models (SLMs) struggle with NL2SQL tasks, exhibiting poor performance and incompatibility with existing frameworks. To address these issues, we introduce Feather-SQL, a new lightweight framework tailored for SLMs. Feather-SQL improves SQL executability and accuracy through 1) schema pruning and linking, 2) multi-path and multi-candidate generation. Additionally, we introduce the 1+1 Model **Collaboration Paradigm**, which pairs a strong 016 general-purpose chat model with a fine-tuned 017 SQL specialist, combining strong analytical reasoning with high-precision SQL generation. Experimental results on BIRD demonstrate that Feather-SQL improves NL2SQL performance 021 on SLMs, with around 10% boost for models 022 without fine-tuning. The proposed paradigm raises the accuracy ceiling of SLMs to 54.76%, highlighting its effectiveness.

1 Introduction

037

041

Natural Language to SQL (NL2SQL) is the task of converting natural language questions into corresponding SQL queries, allowing users to retrieve structured data from databases without requiring proficiency in SQL language. In recent years, the field has seen significant advancements with the emergence of large language models (LLMs) such as GPT-4 (OpenAI et al., 2024), enabling frameworks like CHASE-SQL (Pourreza et al., 2024) and XiYan-SQL (Gao et al., 2025) to achieve stateof-the-art (SOTA) performance. However, two limitations hinder their practical adoption. First, mainstream methods depend on closed-source models, and their reliance on external APIs introduces data privacy risks in sensitive domains like healthcare



Figure 1: NL2SQL performance on the BIRD DEV dataset. <u>EXE</u> (Executability) measures successful query executions, while <u>ACC</u> (Accuracy) measures correct result matches.

042

043

044

045

051

054

057

059

060

061

062

063

064

065

067

068

and finance (Liu et al., 2024). Second, most opensource research focuses on models with 7B–30B parameters, leaving small language models (SLMs) with 4B or fewer parameters relatively underexplored. Meanwhile, many relational databases are deployed on high-performance systems with limited GPU resources. With efficient inference frameworks (e.g., Intel IPEX-LLM (Intel, 2024)) or quantization techniques, SLMs can help drive the broader adoption of NL2SQL in real-world scenarios while preserving data privacy.

In this paper, we focus on enhancing NL2SQL performance using SLMs. As shown in Figure 1, SLMs face two key challenges: (1) one critical issue is their sharp decline in executability. Unlike LLMs, which can effectively handle longcontext dependencies, SLMs struggle with complex database schema and verbose prompts, often leading to confusion or hallucinated outputs (Nguyen et al., 2024; Qu et al., 2024) (Figure 2); (2) existing frameworks for NL2SQL tasks with LLMs are incompatible with SLMs, as they rely on strong instruction-following capabilities to produce intermediate results, which SLMs lack. As illustrated in Figure 3, SLM outputs frequently violate imposed requirements: they often fail to conform to JSON or array specifications and do not meet predefined

Ω	Schema	Question	Hint
User	 Student (id, name, books, grade) Grade (grade, population) 	Which grade has the most number of books on average?	MAX (AVERAGE (SUM(books grade))
SLM	SELECT g.grade, AVG(g.books) FROM Grade g GROUP BY g.grade ORDER BY AVG(g.books) DESC LIMIT 1;	SELECT g.grade, AVERAGE(s.books) FROM Student s NATURAL JOIN Grade g GROUP BY g.grade;	Find all the students in each grade. Calculate their average number of books, then select the highest grade.
		Execute	x
		+	.
	Syntax Error: no such column: g.books	Syntax Error: no such unction: AVERAGE	Syntax Error: near "Find"

Figure 2: Examples of typical syntax errors produced by small language models (SLMs) in an NL2SQL scenario.

constraints. Directly applying these frameworks to SLMs may further degrade executability.

069

077

085

091

094

097

101

102

103

104 105

106

107

109

To address these challenges, we propose Feather-SQL, a lightweight framework tailored for SLMs to enhance both executability and accuracy in NL2SQL tasks. Feather-SQL consists of six key components: schema pruning, schema linking, multi-path generation, multi-candidate generation, correction, and selection. Designed specifically for SLMs, schema pruning streamlines input processing by discarding irrelevant tables, allowing models to concentrate on essential database elements. Schema linking improves alignment between questions and database schema, ensuring accurate column selection. To mitigate errors from linking and pruning, multi-path generation explores diverse query formulation strategies, enhancing robustness. Multi-candidate generation further improves the model's executability and accuracy by enhancing the variety of generated SQL queries, thereby increasing the likelihood of producing correct candidates. The best candidate is then selected through execution validation and ranking. Complementing these components, we introduce extraction and simplification prompting strategies. Extraction selectively retrieves key information, while simplification removes extraneous prompt details to lower computational overhead. By integrating these techniques, Feather-SQL enables SLMs to generate SQL queries more reliably despite their inherent limitations.

A common approach to enhancing SLMs is finetuning. However, while fine-tuned SLMs for SQL generation tasks (e.g., Prem-SQL (Anindyadeep, 2024), CodeS (Li et al., 2024a)) outperform general-purpose chat models on core NL2SQL tasks, they suffer from catastrophic forgetting (Luo et al., 2025; Kotha et al., 2024) on auxiliary tasks—where task-specific fine-tuning erodes their foundational reasoning abilities. To counter this, we propose **1+1 Model Collaboration Paradigm**,



Figure 3: Experiments conducted on a CHESS-provided BIRD subset for schema linking. Models are required to output a JSON-formatted response containing no more than five relevant columns related to the question, without generating any extraneous content.

in which a general-purpose chat model handles reasoning-intensive auxiliary tasks (e.g., schema linking and candidate selection), while a fine-tuned SQL specialist focuses on query generation. This collaboration leverages both models' strengths: the general model provides broad reasoning ability, while the specialist delivers domain-specific precision. Experiments confirm that the paradigm improves overall performance. Our main contributions are as follows:

110

111

112

113

114

115

116

117

118

119

120

121

122

123

124

125

126

127

128

129

130

131

132

133

134

135

136

- We introduce Feather-SQL, an NL2SQL framework for SLMs to address their unique challenges of low executability and incompatibility with existing LLM-based frameworks.
- We propose a novel 1+1 Model Collaboration paradigm that mitigates catastrophic forgetting in fine-tuned SLMs by delegating reasoning-intensive tasks to a general-purpose chat model.
- Extensive experiments on the Spider and BIRD datasets demonstrate that Feather-SQL consistently achieves strong performance with various SLMs, and when paired with the paradigm, it yields SOTA results on BIRD within the scope of SLMs.

2 Related Work

2.1 Conventional Methods

Extensive research on NL2SQL has been carried137out. Early investigations (Zelle and Mooney, 1996;138Li and Jagadish, 2014; Saha et al., 2016) predominantly employed rule-based or template-based approaches, necessitating considerable manual effort141and thereby limiting both their adaptability and142generalizability.143

To address the shortcomings of earlier methods, sequence-to-sequence models have been proposed. In such models, encoders are responsible for learning the semantic representations of natural language questions and the associated database schema, while decoders generate the corresponding SQL based on these representations. Representative approaches in this category include IRNet (Jha et al., 2019), SQLNet (Xu and et al., 2017), Seq2SQL (Zhong et al., 2017), RyanSQL (Choi et al., 2021), and RESDSQL (Li et al., 2023a), each contributing to advancements in query generation.

144

145

146

147

148

149

150

151

152

153

155

156

157

158

159

161

162

163

166

167

169

170

171

172

173

174

175

176

177

178

179

181

183

184

186

190

191

192

194

Meanwhile, some methods choose to fine-tune pre-trained language models such as BERT (Devlin et al., 2019) and T5 (Raffel et al., 2023), leveraging the broad knowledge captured during pre-training to enhance accuracy and robustness. For instance, Graphix-T5 (Li and et al., 2023) integrates a pretrained transformer with specialized graph-aware layers, improving performance on tasks requiring graph-structured data analysis. Nonetheless, these strategies often demand extensive training data and face considerable challenges when dealing with complex questions and database schema.

2.2 Emerging LLM and SLM Approaches

More recently, the emergence of LLMs has marked a watershed moment. LLM-based NL2SQL methods (Dong et al., 2023; Pourreza and Rafiei, 2023; Gao et al., 2023; Wang et al., 2024; Li et al., 2024b; Qu et al., 2024; Talaei et al., 2024; Ren et al., 2024; Pourreza et al., 2024; Gao et al., 2025) have risen to prominence as leading solutions. For example, DIN-SQL (Pourreza and Rafiei, 2023) decomposes the NL2SQL task into subtasks—such as schema linking, difficulty classification, and SQL generation-thereby streamlining decision-making and enabling more accurate query outputs. CHESS (Talaei et al., 2024) adopts a multi-agent framework in which each agent is assigned a specific model and a few-shot prompting strategy to handle different subtasks. While these approaches offer impressive performance, they introduce security risks and lead to steep increases in computational costs.

Employing SLMs to address NL2SQL tasks has the potential to alleviate the aforementioned challenges. CodeS (Li et al., 2024a) incorporates incremental pre-training and bi-directional data augmentation to fine-tune a series of models (1B, 3B, 7B, and 15B parameters). Models specifically finetuned for NL2SQL tasks, such as premSQL (Anindyadeep, 2024) and SQLCoder (Defog), have likewise demonstrated notable success. Nevertheless, these models remain susceptible to the effects of catastrophic forgetting, which diminishes their capacity to perform general reasoning tasks—such as schema linking—within the broader NL2SQL workflow.

195

196

197

198

199

200

201

202

203

204

205

206

207

208

209

210

211

212

213

214

215

216

217

218

219

220

221

222

223

224

225

226

227

228

229

230

231

232

233

234

235

236

237

238

239

240

241

242

243

244

3 Methodology

3.1 Feather-SQL

As shown in Figure 4, we propose Feather-SQL to enhance the performance of SLMs in NL2SQL. This framework comprises several components, including Schema Pruning, Multi-Path, and Multi-Candidate Generation, which are specifically designed to address the challenges of SLMs. We will elaborate on these components in the following sections.

Schema Pruning. This step dynamically reduces schema complexity by identifying and filtering out irrelevant tables. Given the complete set of Data Definition Language (DDL) statements for all tables, the model analyzes semantic relevance to determine which tables are pertinent to the question. Only the DDLs of these relevant tables are retained in the subsequent processing pipeline. This selective retention mechanism prevents SLMs from processing long inputs, thereby mitigating their limitations in handling long text while preserving essential information.

Schema Linking. This step aligns the question with the database schema by identifying relevant columns through semantic analysis. As a commonly adopted practice, schema linking extracts pertinent columns from the complete schema based on the given question, facilitating downstream processing. By establishing precise mappings between natural language expressions and database elements, this process significantly enhances SQL generation accuracy.

Multi-Path Generation. This step employs four distinct prompt types: (1) with both schema linking and pruning, (2) linking only, (3) pruning only, and (4) without either operation. The multi-path design mitigates the risk of information loss caused by pruning errors and reduces potential misunderstandings arising from linking inaccuracies.

Multi-Candidate Generation. This step generates multiple SQL queries in parallel to increase the likelihood of producing a correct result. To ensure diversity, <u>beam search</u> is employed alongside carefully tuned temperature and top-p parameters.



Figure 4: An overview of the Feather-SQL framework for small language models (SLMs) in NL2SQL tasks. The pipeline comprises six core modules—*schema pruning, schema linking, multi-path generation, multi-candidate generation, correction,* and *selection*—which collaboratively boost query executability and accuracy. Additionally, the *1*+1 *Model Collaboration Paradigm* pairs a general-purpose chat model for auxiliary reasoning with a SQL fine-tuned model for query generation, balancing broad contextual understanding with domain-specific precision.

Each path consistently generates a fixed number of candidate queries, maintaining a balanced exploration of possible solutions. Notably, while LLMs often generate executable answers on the first attempt with minimal accuracy improvement from additional candidates, SLMs benefit significantly from multi-candidate generation, which enhances both executability and accuracy (Appendix B).

256

257

260

264

269

271

Correction. This step executes each generated query and handles it based on the outcome. If a query executes successfully, it is directly added to the array of executable SQL queries. For failed queries, error feedback is used to revise the query through a self-correction approach, generating two new candidate queries. If any of these revised queries are executable, they are also stored in the array of executable SQL queries.

Selection. This step employs a <u>selection ranking</u> method to evaluate all executable queries based on their alignment with the expected answer. If a query yields a limited number of results, the evaluation considers both the query and its execution outcome. In contrast, the evaluation focuses solely on the query itself. The selection process is repeated three times, and the mode of the rankings is returned as the final result.

3.2 Prompting Strategies

Extraction. As mentioned in Section 1, SLMs
struggle to meet structural constraints, thus we propose an extraction strategy to avoid rigid structural outputs by allowing SLMs to freely generate responses. This improves accuracy on reason-

ing tasks by bypassing syntactic constraints. We have two methods to achieve that: (1) Lexical Matching: This method identifies valid schema elements by matching table/column names explicitly mentioned in the natural language response against the database schema. For instance, when the SLM outputs "The required tables include customer and orders", the system verifies and extracts customer/orders only if they exist in the schema. (2) Pattern Matching: This method extracts the final answer by identifying predefined patterns in the model's output, such as "answer is" or "Answer:". For example, if the model generates "The answer is 128", the framework detects the pattern and extracts "128" as the final result.

277

278

279

281

282

284

286

287

288

289

290

292

293

294

296

297

298

299

300

301

302

303

304

305

306

307

308

Simplification. The simplification strategy reduces computational overhead by minimizing prompt verbosity. In Feather-SQL, we achieve this by removing superfluous details and using concise instructions with the fewest effective examples (Appendix C). This approach refines the input by eliminating unnecessary complexity, avoiding the need for SLMs to process lengthy inputs while maintaining the clarity of the task.

3.3 1+1 Collaboration Paradigm

Our paradigm categorizes NL2SQL pipeline tasks into two types: reasoning-intensive tasks and SQL generation tasks. Reasoning tasks, such as schema linking and candidate evaluation, require strong contextual understanding and adaptability, while SQL generation demands precision in query synthesis. To optimize performance, we employ two spe-

Mathad	Qwen2	.5-1.5B	Yi-Co	der-1.5B	Phi3-M	ini-3.8B
Methoa	EX (%)	EP (%)	EX (%)	EP (%)	EX (%)	EP (%)
DR	19.36	53.52	15.84	54.82	27.44	71.90
FEQ	<u>21.51</u>	68.25	<u>18.71</u>	73.60	<u>30.12</u>	67.93
MAC-SQL	18.06	52.28	7.63	59.52	29.99	77.64
CHESS	18.71	43.55	2.48	7.82	18.12	39.70
Feather-SQL (Ours)	31.81	88.33	25.23	90.61	36.64	83.70
Mathad	MiniCPM3-4B		Prem-SQL-1.3B		CodeS-3B	
Methoa	EX (%)	EP (%)	EX (%)	EP (%)	EX (%)	EP (%)
DR	27.57	69.30	47.07	88.14	24.19	<u>59.32</u>
FEQ	29.34	63.89	51.63	<u>92.70</u>	25.03	57.50
MAC-SQL	<u>37.35</u>	<u>81.68</u>	8.67 (8.87*)	17.01 (<i>19.23</i> *)	10.10 (<i>13.23</i> *)	40.87 (56.26*)
CHESS	28.42	54.43	24.64	43.22	<u>26.53</u>	56.91
Feather-SQL (Ours)	40.09	87.02	<u>49.28</u>	98.04	33.96	85.31

Table 1: Comparison of EX (Execution Accuracy) and EP (Execution Proportion) across different methods on the BIRD DEV dataset. The best and second-best results are highlighted by **Bold** and <u>underline</u>, respectively. * denotes results with the extraction strategy.

cialized models: the general-purpose chat Model
for reasoning tasks and the SQL fine-tuned model
for SQL generation. By leveraging their complementary strengths, our approach improves overall
NL2SQL accuracy and robustness.

General-purpose Chat Model. This model is designed for reasoning-intensive tasks, leveraging broad linguistic and contextual comprehension without domain-specific fine-tuning, which helps prevent catastrophic forgetting. Compared to the SQL Specialist Model, it is more effective in schema linking and evaluating SQL candidates, ensuring that the SQL generation process is guided by accurate and well-structured contextual information.

> **SQL Fine-tuned Model.** Optimized exclusively for SQL generation, this model is extensively trained on large-scale NL2SQL datasets, allowing it to achieve superior performance on SQL-specific tasks. Its focused training reduces hallucinations and enhances both query executability and accuracy.

4 Experiments

4.1 Settings

324

325

326

328

333

334

335

337

340

4.1.1 Datasets

BIRD (Li et al., 2023b) as a representative and challenging benchmark dataset for NL2SQL, encompasses databases over 37 professional domains. Due to the proprietary nature of the BIRD TEST dataset, we conduct our experiments using the publicly accessible BIRD DEV subset, which contains 1,534 unique question-SQL pairs.

Spider (Yu et al., 2019) is another large-scale

benchmark dataset for cross-domain SQL generation, covering 138 different domains. Compared to BIRD, Spider is relatively simpler, as its SQL structures and schema are generally less complex. Our experiments utilize the SPIDER TEST set, comprising 2,147 question-SQL pairs. 342

343

344

345

346

348

349

350

351

352

353

354

355

356

357

358

359

360

361

362

363

364

365

366

367

368

369

370

371

4.1.2 Evaluation Metrics

Execution Accuracy (EX) (Li et al., 2023b) is a widely adopted metric in NL2SQL evaluations, measuring whether the result of executing the generated query matches the result of the ground truth query. This metric allows for different query formulations that yield the same result. It is calculated as:

$$\mathsf{EX} = \frac{|\{n \in N \mid \mathsf{E}(Q_{gen}) = \mathsf{E}(Q_{gt})\}|}{N} \times 100\%$$

where N denotes the number of questions. Q_{gen} represents the SQL query generated by the model, while Q_{gt} is the ground truth answer. E is the execution function.

Execution Proportion (EP) is an auxiliary metric we proposed, evaluating the proportion of generated SQL queries that can be executed on the corresponding database without syntax errors. This metric reflects the model's upper-bound capability by assuming that any executable query is potentially correct. It is defined as:

$$\mathsf{EP} = \frac{|\{n \in N \mid \mathsf{E}(Q_{\mathsf{gen}}) \neq \mathsf{error}\}|}{N} \times 100\%$$

4.1.3 Baselines

Direct Response (DR) directly generates an SQL query from the natural language question without

Mathad	Qwen2	.5-0.5B	Yi-Code	er-1.5B	DeepSeek-Coder-1.3B	
Wiethoa	EX (%)	EP (%)	EX (%)	EP (%)	EX (%)	EP (%)
DR	28.50	56.45	45.23	87.24	49.28	90.68
FEQ	<u>36.53</u>	67.35	<u>48.30</u>	86.77	45.46	89.89
MAC-SQL	29.06	89.61	13.04	21.70	52.12	<u>93.62</u>
CHESS	15.42	29.16	3.68	10.29	30.18	46.30
Feather-SQL (Ours)	36.98	<u>75.08</u>	49.56	92.04	<u>51.19</u>	94.13
Mathad	MiniCPM3-4B		Prem-SQL-1.3B		CodeS-3B	
Methoa	EX (%)	EP (%)	EX (%)	EP (%)	EX (%)	EP (%)
DR	55.10	<u>93.71</u>	60.92	85.79	47.74	64.23
FEQ	<u>55.75</u>	89.52	64.23	85.75	49.60	64.65
MAC-SQL	25.01	38.47	0.14 (67.91 *)	0.14 (100 *)	0 (74.48 *)	0 (100 *)
CHESS	56.73	89.99	63.86	92.08	<u>66.65</u>	88.54
Feather-SQL (Ours)	58.92	94.18	<u>66.60</u>	<u>92.78</u>	63.25	<u>88.96</u>

Table 2: Comparison of EX (Execution Accuracy) and EP (Execution Proportion) across different methods on the Spider TEST dataset. The best and second-best results for EX are highlighted by **bold** and <u>underline</u>, respectively. * denotes results with the extraction strategy.

applying any refinement techniques. The process follows a single-turn interaction.

372

373

374

375

377 378

379

384

386

388

396

397

399

400

401

First Executable Query (FEQ) leverages the model's ability to generate multiple SQL candidates. Among candidates, the first executable query is selected without any refinement. This approach simulates multi-turn query generation.

MAC-SQL (Wang et al., 2024) is an LLM-based multi-stage framework, featuring a core Decomposer agent for SQL generation supported by auxiliary agents for sub-database acquisition and query refinement. It also utilizes few-shot chain-of-thought reasoning to enhance generation processes.
 CHESS (Talaei et al., 2024) comprises four specialized agents: Information Retriever, Schema Selector, Candidate Generator, and Unit Tester. Notably, it employs locality-sensitive hashing and vector databases to efficiently retrieve relevant data from extensive database values and catalogs.

4.1.4 Implementation Details

Backbone Models. Our implementation leverages both general-purpose chat models and SQL fine-tuned models. The chat models include Qwen2.5-0.5B, Qwen2.5-1.5B, Qwen2.5-Coder-1.5B (Hui et al., 2024), Yi-Coder-1.5B (AI et al., 2025), DeepSeek-Coder-1.5B (DeepSeek-AI, 2024), Phi3-mini-3.8B (Abdin et al., 2024), and MiniCPM3-4B (Hu et al., 2024), while the SQL-tuned models consist of Prem-SQL-1.3B (Anindyadeep, 2024) and CodeS-3B (Li et al., 2024a).

402 Candidate Size. In the multi-candidate generation
403 stage, we generate <u>4</u> candidates per path, resulting

in a total candidate pool of <u>16</u>. During the correction stage, the candidate size is reduced to <u>2</u>. **Selection Rounds.** During the selection stage, we perform <u>3</u> rounds for each selection. The final choice is the majority vote across the three rounds, ensuring consistency of the selected candidate. 404

405

406

407

408

409

410

411

412

413

414

415

416

417

418

419

420

421

422

423

424

425

426

427

428

429

430

431

432

433

434

4.2 Main Results

4.2.1 Feather-SQL

To validate the general effectiveness of Feather-SQL for SLMs, we conducted experiments on two datasets across a range of models (all results here were obtained using a unified model without adopting the collaboration paradigm).

BIRD Results. As shown in Table 1, Feather-SQL demonstrates superior performance across all general-purpose chat models, achieving the highest scores in both EX and EP, with EX showing an average increase of approximately 10% and EP exceeding a 20% improvement compared to FEQ. For SQL fine-tuned models, Feather-SQL combined with CodeS achieves substantial gains in both EX and EP, while Prem-SQL shows notable improvements specifically in EP, with an average increase of around 5% compared to FEQ. Besides, we explored the upper bound of Feather-SQL on this dataset (Appendix D).

Moreover, we observe that CHESS and MAC-SQL do not perform effectively on SLMs, with their results on Qwen2.5 and Yi-Coder showing even lower EX and EP scores compared to DR. Their performance also falls behind that of FEQ.

Chat Model	SQL Model	EX (%)	EP (%)
_	Prem-SQL	49.28	98.04
Qwen	Prem-SQL	52.44 ↑	94.08
Qwen Coder	Prem-SQL	52.83 ↑	98.31
Yi Coder	Prem-SQL	54.76 ↑	93.94
_	CodeS	33.96	83.31
Qwen	CodeS	35.79↑	80.05
Qwen Coder	CodeS	37.03 ↑	81.10
Yi Coder	CodeS	39.43 ↑	80.44

Table 3: Paradigm performance under Feather-SQL on the BIRD DEV dataset. When no chat model is specified, the SQL model is also used as the chat model.

Spider Results. Similarly, Table 2 highlights the results on the SPIDER TEST dataset, further confirming that our framework consistently and substantially enhances the NL2SQL performance of SLMs.

Although MAC-SQL and CHESS show inconsistent performance across models, MAC-SQL generally performs well. Notably, SQL fine-tuned models, achieve the best EX when extraction is applied, highlighting the necessity of this step for SLMs. This may be attributed to MAC-SQL's Selector mechanism, which also employs schema pruning. Unlike our table pruning approach, MAC-SQL adopts column pruning, which may be more effective for SPIDER's relatively simple schema structures.

4.2.2 1+1 Collaboration Paradigm

As observed in Table 1, although Feather-SQL improves the EP of Prem-SQL, its EX shows a 2% decrease compared to FEQ. This decline is primarily due to Prem-SQL's inability to handle auxiliary reasoning tasks. To address this limitation, we propose a division of tasks where the general-purpose chat model handles auxiliary reasoning, while the SQL fine-tuned model focuses on SQL generation.

As shown in Table 3, our 1+1 collaboration paradigm under Feather-SQL achieves a 3–6% improvement in EX for both Prem-SQL and CodeS, with Prem-SQL reaching <u>SOTA</u> performance among existing SLMs (Appendix E). However, we observe a decline in EP when paired with a chat model. This is because when the SQL model is also used as the chat model during schema pruning, it returns a query instead of the expected answer. But our extraction strategy sitll retrieves table names from the output, often resulting in an overly pruned schema-containing only one or two tables.

Chat Model	SQL Model	EX (%)	EP (%)
-	Prem-SQL	24.64	43.22
Qwen	Prem-SQL	49.28 ↑	82.07
Qwen Coder	Prem-SQL	49.61 ↑	79.60
Yi Coder	Prem-SQL	47.65 ↑	79.79
-	CodeS	26.53	56.91
Qwen	CodeS	28.55 ↑	56.19
Qwen Coder	CodeS	28.88 ↑	63.04
Yi Coder	CodeS	27.44 ↑	55.22

Table 4: Paradigm performance under CHESS on the BIRD DEV dataset. When no chat model is specified, the SQL model is also used as the chat model.

While a simplified schema can occasionally boost EP, it frequently leads to lower overall EX.

472

473

474

475

476

477

478

479

480

481

482

483

484

485

486

487

488

489

490

491

492

493

494

495

496

497

498

499

500

501

502

503

504

505

506

Additionally, Table 4 shows that our paradigm improves both Prem-SQL and CodeS in CHESS, with EX increasing by ~20% and EP by over ~35% for Prem-SQL, while CodeS sees a smaller but consistent EX gain with no clear trend in EP.

However, the two models benefit differently due to their handling of auxiliary tasks. Prem-SQL attempts to answer linking questions but often does so incorrectly, whereas CodeS, due to severe catastrophic forgetting, fails to provide valid responses. As a result, CHESS defaults to using the original schema with CodeS, reducing linking errors.

Furthermore, since CHESS constructs long prompts without schema pruning, introducing a chat model increases input length and complexity. While this improves reasoning, it does not fully offset CodeS's limitations in processing extended inputs, restricting its EX improvement.

4.3 Ablation Studies

4.3.1 Component Contribution

We conducted an ablation study to quantify the impact of each framework component by removing them one at a time and measuring changes in EX and EP on the <u>BIRD DEV</u> dataset, using *QWen2.5-1.5B* (Table 5).

We can see from the ablation results that removing any of the components causes a drop in both EX and EP. This underscores that each step in our pipeline contributes to overall performance, and omitting even one module leads to noticeably reduced accuracy or executability.

Among these, schema pruning is shown to be the most critical: when it is removed, EX falls

463

464

465

466

467

468

469

470

471

435

Framework	EX (%)	EP (%)
Full Model	31.81	88.33
-w/o Schema Pruning	-4.63↓	-20.34 ↓
-w/o Schema Linking	-3.45↓	-20.92↓
-w/o Multi-Candidate	-2.47↓	-17.99↓
-w/o Correction	-0.20↓	-12.58↓
-w/o Selection	-2.21↓	-10.36 ↓

Table 5: Ablation Study on Framework Components.

from 31.81% to 27.18%, the single largest drop in our study. This highlights how focusing on only the relevant tables and columns helps the model concentrate on essential schema elements, thereby yielding more accurate SQL generation. In contrast, removing correction only reduces EX by 0.20%, indicating that it has a relatively minor impact on the framework's effectiveness.

4.3.2 Path Contribution

507

508

509

510

511

512

513

514

515

516

517

518

520

522

524

526

528

529

530

531

532

534

535

537

539

541

545

We analyzed the origins of SQL answers from four models to understand how each processing path affects the final output. As shown in Figure 5, our multi-path framework includes four paths: one using both schema linking and pruning, one using only schema linking, one using only schema pruning, and one without either.

For all four models, the path *Full Schema & Linking* is consistently the largest contributor, followed by *Pruned Schema & Linking*. This ranking underscores the critical role of linking in the framework, regardless of whether the schema is pruned or not.

Additionally, we find that schema pruning collectively accounts for over 25% across the models. These observations are consistent with the ablation findings in 4.3.1, further illustrating the essential roles of each component in ensuring executable and accurate query generation.

4.3.3 Candidate Size

We further investigated the impact of different candidate sizes. Figure 6 presents the results based on our four paths. In our experiments, the total candidate size increases from 4 to 24, which corresponds to the number of candidates generated per path increasing from 1 to 6. The figure illustrates how EX changes as the overall candidate size grows from 4 to 24.

We observe a concave trend, consistent with Appendix B: EX steadily increases as the candidate



Figure 5: Distribution of correct SQL answers contributed by each path across four different SLMs.

size rises from 4 to 16 but then plateaus from 16 to 24. Once the model reaches its approximate upper bound, further increases in candidate size result in only a marginal difference in performance. Therefore, we select a candidate size of 16, as it is the earliest point at which EX saturates, thus balancing computational efficiency and model performance.



Figure 6: Effect of candidate size on EX performance.

5 Conclusion

In this work, we introduced Feather-SQL, the first lightweight framework designed to enhance NL2SQL performance for SLMs. We conduct comprehensive evaluations on the challenging BIRD and Spider datasets, where Feather-SQL yields improvements in both executability and accuracy. Additionally, we present the 1+1 Model Collaboration paradigm—a novel approach that pairs a general-purpose chat model with a SQL specialist to combine robust reasoning with precise query generation. Our evaluation results show that this paradigm boosts accuracy across different frameworks, demonstrating its consistent effectiveness. Moreover, the flexibility of our approach provides a robust foundation not only for advancing NL2SQL but also for application to other structured tasks and domains.

551

552

553

554

555

556

557

558

559

560

561

562

563

564

565

566

567

568

6 Limitations

571

587

591

596

597

599

607

608

611

612

613

614

615

616

617

618

620

Despite the promising performance gains achieved by Feather-SQL, our current framework does not 573 yet reach very high absolute accuracy on datasets. 574 For instance, the best cumulative accuracy on BIRD DEV is around 74% (Gao et al., 2025; 576 577 Pourreza et al., 2024). In fact, many LLM-based NL2SQL systems typically report accuracy in the 578 60+% range, while the SOTA results achieved by SLMs remain below 55%. However, our approach is the first to surpass all previous methods at the 581 582 1B-parameter scale. Feather-SQL with the Model Collaboration Paradigm lays a strong foundation for promoting the broader adoption of NL2SQL in real-world applications.

References

- Marah Abdin, Jyoti Aneja, Hany Awadalla, Ahmed Awadallah, Ammar Ahmad Awan, Nguyen Bach, Amit Bahree, Arash Bakhtiari, Jianmin Bao, and Harkirat Behl. 2024. Phi-3 technical report: A highly capable language model locally on your phone. Preprint, arXiv:2404.14219.
- 01. AI, :, Alex Young, Bei Chen, Chao Li, and Chengen Huang. 2025. Yi: Open foundation models by 01.ai. Preprint, arXiv:2403.04652.
- Anindyadeep. 2024. Premsql: End-to-end localfirst text-to-sql pipelines. https://github.com/ premAI-io/premsql. Accessed: 2024-12-10.
- DongHyun Choi, Myeong Cheol Shin, EungGyun Kim, and Dong Ryeol Shin. 2021. Ryansql: Recursively applying sketch-based slot fillings for complex textto-sql in cross-domain databases. <u>Computational</u> Linguistics, 47(2):309–332.
- DeepSeek-AI. 2024. Deepseek llm: Scaling opensource language models with longtermism. <u>Preprint</u>, arXiv:2401.02954.
- Defog. Sqlcoder. https://github.com/defog-ai/ sqlcoder. Accessed: 2024-12-10.
- Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. 2019. Bert: Pre-training of deep bidirectional transformers for language understanding. <u>Preprint</u>, arXiv:1810.04805.
- Xuemei Dong, Chao Zhang, Yuhang Ge, Yuren Mao, Yunjun Gao, lu Chen, Jinshu Lin, and Dongfang Lou. 2023. C3: Zero-shot text-to-sql with chatgpt. <u>Preprint</u>, arXiv:2307.07306.
- Dawei Gao, Haibin Wang, Yaliang Li, Xiuyu Sun, Yichen Qian, Bolin Ding, and Jingren Zhou. 2023.
 Text-to-sql empowered by large language models: A benchmark evaluation. Preprint, arXiv:2308.15363.

Yingqi Gao, Yifu Liu, Xiaoxia Li, Xiaorong Shi, Yin Zhu, Yiming Wang, Shiqi Li, Wei Li, Yuntao Hong, Zhiling Luo, Jinyang Gao, Liyu Mou, and Yu Li. 2025. A preview of xiyan-sql: A multi-generator ensemble framework for text-to-sql. Preprint, arXiv:2411.08599.

621

622

623

624

625

626

627

628

629

630

631

632

633

634

635

636

637

638

639

640

641

642

643

644

645

646

647

648

649

650

651

652

653

654

655

656

657

658

659

660

661

662

663

664

665

666

667

668

669

670

671

672

673

674

- Shengding Hu, Yuge Tu, Xu Han, Chaoqun He, Ganqu Cui, Xiang Long, Zhi Zheng, Yewei Fang, and Yuxiang Huang. 2024. Minicpm: Unveiling the potential of small language models with scalable training strategies. Preprint, arXiv:2404.06395.
- Binyuan Hui, Jian Yang, Zeyu Cui, Jiaxi Yang, Dayiheng Liu, Lei Zhang, Tianyu Liu, and Jiajun Zhang. 2024. Qwen2.5-coder technical report. <u>Preprint</u>, arXiv:2409.12186.
- Intel. 2024. Intel® llm library for pytorch. https: //github.com/intel/ipex-llm. Accessed: 2024-12-10.
- Dipendra Jha, Logan Ward, Zijiang Yang, Christopher Wolverton, Ian Foster, Wei-keng Liao, Alok Choudhary, and Ankit Agrawal. 2019. Irnet: A general purpose deep residual regression framework for materials discovery. In <u>Proceedings of the 25th ACM</u> <u>SIGKDD International Conference on Knowledge</u> Discovery & Data Mining, pages 2385–2393.
- Suhas Kotha, Jacob Mitchell Springer, and Aditi Raghunathan. 2024. Understanding catastrophic forgetting in language models via implicit inference. In <u>The Twelfth International Conference on Learning</u> <u>Representations</u>.
- Fei Li and Hosagrahar V Jagadish. 2014. Nalir: an interactive natural language interface for querying relational databases. In Proceedings of the 2014 ACM <u>SIGMOD International Conference on Management</u> <u>of Data</u>, SIGMOD '14, page 709–712, New York, <u>NY, USA</u>. Association for Computing Machinery.
- Haoyang Li, Jing Zhang, Cuiping Li, and Hong Chen. 2023a. Resdsql: Decoupling schema linking and skeleton parsing for text-to-sql. In <u>AAAI</u>.
- Haoyang Li, Jing Zhang, Hanbing Liu, Ju Fan, Xiaokang Zhang, Jun Zhu, Renjie Wei, Hongyan Pan, Cuiping Li, and Hong Chen. 2024a. Codes: Towards building open-source language models for text-to-sql. Preprint, arXiv:2402.16347.
- Jinyang Li and et al. 2023. Graphix-t5: Mixing pretrained transformers with graph-aware layers for textto-sql parsing. <u>arXiv:2301.07507</u>.
- Jinyang Li, Binyuan Hui, Ge Qu, Binhua Li, Jiaxi Yang, Bowen Li, Bailin Wang, Bowen Qin, Ruiying Geng, Nan Huo, Xuanhe Zhou, Chenhao Ma, Guoliang Li, Kevin C. C. Chang, Fei Huang, Reynold Cheng, and Yongbin Li. 2023b. Can Ilm already serve as a database interface? a big bench for large-scale database grounded text-to-sqls. <u>Preprint</u>, arXiv:2305.03111.

- 676
- 686 691
- 695

- 701
- 705
- 706
- 710

712

713

715 716

718

- 719
- 720 721

722

723 725

726 727 728

731

- Zhishuai Li, Xiang Wang, Jingjing Zhao, Sun Yang, Guoqing Du, Xiaoru Hu, Bin Zhang, Yuxiao Ye, Ziyue Li, Rui Zhao, et al. 2024b. Pet-sql: A promptenhanced two-stage text-to-sql framework with crossconsistency. arXiv preprint arXiv:2403.09732.
- Xinyu Liu, Shuyu Shen, Boyan Li, Peixian Ma, Runzhi Jiang, Yuyu Luo, Yuxin Zhang, Ju Fan, Guoliang Li, and Nan Tang. 2024. A survey of nl2sql with large language models: Where are we, and where are we going? Preprint, arXiv:2408.05109.
- Yun Luo, Zhen Yang, Fandong Meng, Yafu Li, Jie Zhou, and Yue Zhang. 2025. An empirical study of catastrophic forgetting in large language models during continual fine-tuning. Preprint, arXiv:2308.08747.
- Chien Van Nguyen, Xuan Shen, Ryan Aponte, Yu Xia, Samyadeep Basu, Zhengmian Hu, Jian Chen, Mihir Parmar, Sasidhar Kunapuli, Joe Barrow, Junda Wu, Ashish Singh, Yu Wang, Jiuxiang Gu, Franck Dernoncourt, Nesreen K. Ahmed, Nedim Lipka, Ruiyi Zhang, Xiang Chen, Tong Yu, Sungchul Kim, Hanieh Deilamsalehy, Namyong Park, Mike Rimer, Zhehao Zhang, Huanrui Yang, Ryan A. Rossi, and Thien Huu Nguyen. 2024. A survey of small language models. Preprint, arXiv:2410.20011.
- OpenAI, Josh Achiam, Steven Adler, Sandhini Agarwal, Lama Ahmad, Ilge Akkaya, Florencia Leoni Aleman, Diogo Almeida, Janko Altenschmidt, Sam Altman, Shyamal Anadkat, Red Avila, and Igor Babuschkin. 2024. Gpt-4 technical report. Preprint, arXiv:2303.08774.
- 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. Preprint, arXiv:2410.01943.
- Mohammadreza Pourreza and Davood Rafiei. 2023. Din-sql: Decomposed in-context learning of text-tosql with self-correction. Preprint, arXiv:2304.11015.
- 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. Preprint, arXiv:2405.15307.
- Colin Raffel, Noam Shazeer, Adam Roberts, Katherine Lee, Sharan Narang, Michael Matena, Yanqi Zhou, Wei Li, and Peter J. Liu. 2023. Exploring the limits of transfer learning with a unified text-to-text transformer. Preprint, arXiv:1910.10683.
- Tonghui Ren, Yuankai Fan, Zhenying He, Ren Huang, Jiaqi Dai, Can Huang, Yinan Jing, Kai Zhang, Yifan Yang, and X. Sean Wang. 2024. PURPLE: Making a Large Language Model a Better SQL Writer . In 2024 IEEE 40th International Conference on Data Engineering (ICDE), pages 15-28, Los Alamitos, CA, USA. IEEE Computer Society.

Diptikalyan Saha, Avrilia Floratou, Karthik Sankaranarayanan, Umar Farooq Minhas, Ashish R. Mittal, and Fatma Özcan. 2016. Athena: an ontology-driven system for natural language querying over relational data stores. Proc. VLDB Endow., 9(12):1209-1220.

732

733

734

735

736

737

738

739

740

741

742

743

744

745

746

747

748

749

750

751

752

753

754

755

756

757

758

759

762

- Shayan Talaei, Mohammadreza Pourreza, Yu-Chen Chang, Azalia Mirhoseini, and Amin Saberi. 2024. Chess: Contextual harnessing for efficient sql synthesis. Preprint, arXiv:2405.16755.
- Bing Wang, Changyu Ren, Jian Yang, Xinnian Liang, Jiaqi 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. Preprint, arXiv:2312.11242.
- Xiaojun Xu and et al. 2017. Sqlnet: Generating structured queries from natural language without reinforcement learning. arXiv:1711.04436.
- Tao Yu, Rui Zhang, Kai Yang, Michihiro Yasunaga, Dongxu Wang, Zifan Li, James Ma, Irene Li, Qingning Yao, Shanelle Roman, Zilin Zhang, and Dragomir Radev. 2019. Spider: A large-scale human-labeled dataset for complex and cross-domain semantic parsing and text-to-sql task. Preprint, arXiv:1809.08887.
- John M. Zelle and Raymond J. Mooney. 1996. Learning to parse database queries using inductive logic programming. In Proceedings of the Thirteenth National Conference on Artificial Intelligence -Volume 2, AAAI'96, page 1050–1055. AAAI Press.
- Victor Zhong, Caiming Xiong, and Richard Socher. 2017. Seq2sql: Generating structured queries from natural language using reinforcement learning. Preprint, arXiv:1709.00103.

A Experimental Settings

All experiments were conducted on 4 NVIDIA A6000 GPUs using the <u>vLLM</u> inference acceleration framework to improve model efficiency. For stages that produce multiple answers, such as candidate generation and selection, we primarily used a temperature of 0.2 and a top_p of 0.8 to balance diversity and accuracy. In contrast, for tasks requiring a single answer, such as schema pruning and schema linking, we employed greedy search to ensure deterministic outputs.

B Multi-Candidate Motivation

Top N	Yi-Coder-1.5B		MiniCPM3-4B		Prem-SQL-1.3B	
Tob-W	ACC (%)	EXE (%)	ACC (%)	EXE (%)	ACC (%)	EXE (%)
1	15.65	46.26	26.53	65.31	55.78	92.52
3	24.49	70.75	35.37	76.87	59.86	97.28
5	30.61	78.91	36.05	82.31	62.59	97.96
7	33.33	82.31	37.41	84.35	65.31	97.96
Top N	CodeS-3B		GPT-40		Claude-3.5-Sonnet	
Top-IN	ACC (%)	EXE (%)	ACC (%)	EXE (%)	ACC (%)	EXE (%)
1	24.49	61.90	51.70	93.20	40.82	86.39
3	27.21	68.71	53.74	94.56	41.50	87.76
5	29.93	72.11	56.46	94.56	42.18	88.44
7	29.93	73.47	56.46	94.56	42.18	88.44

Table 6: Comparison of Accuracy (ACC) and Execution (EXE) on the BIRD DEV Subset from CHESS using multi-candidate generation strategy.



Figure 7: Improvement in Accuracy (ΔACC) and Executable Rate (ΔEXE) compared to Top-1 candidates

The results demonstrate that SLMs exhibit a performance gap between TOP-1 and TOP-7 results. This indicates that employing a multi-candidate generation strategy can effectively improve the accuracy and execution rates by selecting the best result. In contrast, larger models already perform robustly with TOP-1 outputs, and therefore, the additional benefit from multi-candidate generation is limited. Additionally, the fine-tuned SQL model CodeS-3B shows some improvement, but the gains are not as pronounced as those observed in the other SLMs.

768

764

765

766

769

770

771

772

774

775

777 C Prompt Length Comparison

On average, *CHESS* uses notably longer prompts due to detailed instructions and complex examples,
 while *MAC-SQL* has fewer words overall. *Feather-SQL* demonstrates the smallest average prompt length,
 indicating that concise design can effectively balance context and complexity.

Method	Stage	Word Count
	Information Retriever	423
CHESS	Schema Selector	2522
CHE55	Generate Candidate	4888
	Revise	1835
	Selector	552
MAC-SQL	Decomposer	836
	Reviser	174
	Schema Pruning	267
	Schema Linking	287
Feather-SQL	Generation	190
	Correction	106
	Selection	271

Table 7: Stages and corresponding word counts for each baseline.

781 D Framework Upper Bound

787

788

789

To explore the upper bound of the Feather-SQL framework, we also evaluated its performance using cumulative accuracy, which measures whether the correct SQL query is present within the Top-n generated results. Specifically, we retained the top 4 candidates after the selection ranking in this experiment, rather than solely selecting the top 1 candidate in default.

As indicated in Table 8, Top-3 is approximately 10% higher than Top-1 (EX). This suggests that there is room for further improvement in the selection mechanism. If the selection can be refined to accurately identify the optimal SQL query, the performance gap between Top-N and Top-1 could be considerably reduced.

Model	Top-1 (%)	Top-2 (%)	Top-3 (%)
Qwen	31.8	39.0	40.5
Yi Coder	25.2	32.6	34.5
Prem-SQL	49.2	60.2	62.6

Table 8: Cumulative Accuracy on BIRD DEV.

E SOTA Result Illustration



Figure 8: Accuracy (%) versus model size (in billions of parameters) for various small language models. Fine-tuned models are shown in red, general-purpose chat models in blue, and ours (Feather-SQL + Model Collaboration Paradigm) is marked with a purple star.

F Prompts

Schema Pruning Prompt

F.1

```
prompt_pruning_system = """
You are an agent designed to find all related tables to generate SQL query
for question based on the database schema and hint.
## Requirements
1. You don't need to answer the question, your task is only finding all related tables .
2. Consider all constraints of each table, including primary keys, foreign keys, and data
    types.
3. You can generate chain of thoughts, but ensure all tables mentioned truly exist.
4. Successfully answer related columns could help you win $100000 dollars.
......
prompt_pruning = """
## Instructions
1. Prioritize the table that most directly contains the information needed to answer the
    question, considering:
    - Table relationships such as foreign keys.
   - Whether the table has columns directly related to the entities or actions in the
    question.
2. Reasoning like two shown examples.
  -----Example-----
## Database Schema
CREATE TABLE Employees (
   employee_id INT PRIMARY KEY,
   name VARCHAR(100),
   department VARCHAR(100),
```

```
salary DECIMAL(10, 2)
   );
   CREATE TABLE Departments (
       department_id INT PRIMARY KEY,
       department_name VARCHAR(100),
       location VARCHAR(100)
   );
   ## Question
   What is the salary of the employee named 'Alice'?
   ## Relevant Tables
   This table directly contains the columns name and salary, which are the only necessary fields
        to answer the question.
   The name column is used to locate the specific employee named 'Alice', and the salary column
       provides the required
   salary information. The Departments table is irrelevant because it does not store employee-
       level data like salaries
   or names, and its information is unrelated to this specific query.
   The relevant table is Employees.
     -----Task------
   ## Database Schema
   You are provided with the structure of the database "{database_name}":
   {database_schema}
   ## Question
   {question}
   ## Hint
   {hint}
   Among the following tables: {tables}, which tables are relevant for addressing the question?
   ## Relevant Tables
   ,, ,, ,,
F.2
   Schema Linking Prompt
   prompt_linking_system="""
   You are an agent designed to find all related columns to generate SQL query for question based
        on the database schema and the hint.
   ## Requirements
   1. You don't need to answer the question, your task is only finding all related columns.
   2. Hint could help you to find the correct related columns.
   3. Consider all constraints of each table, including primary keys, foreign keys, and data
       types.
   4. You can generate chain of thoughts, but ensure all columns mentioned truly exist.
   7. Successfully answer related columns could help you win $100000 dollars.
   prompt_linking="""
   ## Instructions
   1. Select columns that relates to information requested by the question, considering:
       - Whether the column is key to filtering results (used in WHERE clauses).
```

- Whether the column should be part of the SELECT statement to fulfill the user query.
- The relationship of the column to other parts of the question, such as groupings,
- aggregations, or direct match to entities mentioned.
- 2. Reasoning like two shown examples.

------Example------## Database Schema CREATE TABLE Employees (

employee_id INT PRIMARY KEY, name VARCHAR(100), department VARCHAR(100), salary DECIMAL(10, 2)); CREATE TABLE Departments (department_id INT PRIMARY KEY, department_name VARCHAR(100), location VARCHAR(100)); ## Ouestion What is the salary of the employee named 'Alice'? ## Relevant Columns The name column is essential to filter the employee named 'Alice' in the WHERE clause, ensuring we identify the correct individual. The salary column is needed to extract the requested information, which is the employee's salary. Since the question does not involve departments, the Departments table and its columns are irrelevant. The related columns are Employees.name and Employees.salary. -----Task------## Database Schema You are provided with the structure of the database "{database_name}": {schema} ## Question {question} ## Hint {hint} Among the columns, which are relevant for addressing the question? ## Relevant Columns

F.3

Multi-path Generation Prompt

system_prompt_sql_generation = """ You are an expert SQL assistant tasked with generating precise SQL queries based on given database schemas, questions, and hint. ## Responsibilities 1. Analyze the **database schema** and **hint** to determine relationships, including ** primary keys, foreign keys, data types, and constraints**. 2. Generate a single, valid **SQLite SQL query** to answer the question, using provided schema linking information for table and column selection. 3. Your response should contain only the **SQL query**, using standard SQL syntax with correct use of table/column names and SQL clauses. ## Requirements - Respond with only one SQL query, formatted as ```SQL```. - Use clauses like **SELECT**, **FROM**, **WHERE**, **JOIN**, **GROUP BY**, **ORDER BY**, etc. - Ensure SQL is efficient and respects **Important Columns**, table relationships, and relevant constraints. ,,,,, prompt_generation_with_linking = """ You are given a database schema, question, important columns and hint. Generate a valid SQLite query that answers the question. ## Instructions 1. Your response should only contain one SQL query, in standard SQL syntax.

799

798

```
2. Consider all **table relationships**, **primary/foreign keys**, **data types**, and **
    Important Columns** while generating the query.
## Database Schema
Database "{database_name}":
{database_schema}
## Important Columns
{schema_linking}
## Question
{question}
## Hint
{hint}
## Output Requirement
Format the response as:
··`sql
[SQL query]
.....
prompt_generation_without_linking = """
You are given a database schema, question, and hint. Generate a valid SQLite query that
    answers the question.
## Instructions
1. Your response should only contain one SQL query, in standard SQL syntax.
2. Consider all **table relationships**, **primary/foreign keys**, **data types** while
    generating the query.
## Database Schema
Database "{database_name}":
{database_schema}
## Question
{question}
## Hint
{hint}
## Output Requirement
Format the response as:
···sql
[SQL query]
.....
```

800

801

F.4

Correction Prompt prompt_answer_correction_system =""" Suppose you are an expert in SQLite and database management. ## Instructions 1. Based on the database structure provided, previous answer and its error messages, generate one SQL query that answers the question. 2. You should try to fix the error of the previous answer and avoid it from happening again. ## Requirements

1. Your response should consist of only one SQL query, don't generate anything else.

3. Consider all constraints of each table, including primary keys, foreign keys, and data types.

4. Provide your query in standard SQL format with appropriate use of SQL functions, joins, and

conditions.

prompt_answer_correction = """
Database Schema
Given the structure of database:
{schema}

Question {question}

Hint {hint}

Previous answer
{prev_ans}

Error
{errorMsg}

New Answer

F.5

```
Selection Prompt
system_prompt_query_selection = """
You are an expert in analyzing SQL queries and determining their relevance to a given question.
     Your task is to evaluate multiple SQL queries and select the one that best answers the
    question based on the provided database schema and context.
## Responsibilities
1. Analyze the given question: Understand the intent of the question and its expected output.
2. Evaluate each SQL query: Consider the correctness, relevance, and completeness of each
    query in relation to the question.
3. Select the best query: Choose the query that most accurately answers the question, while
    considering database structure, table relationships, and query efficiency.
## Requirements
- Respond with the most relevant SQL query, and nothing else.
- Ensure the selected query is valid for the given database schema and directly addresses the
    question.
,,,,,,
query_selection_prompt = """
You are given a question, a database schema, and multiple SQL queries. Your task is to select
    the SQL query that is most relevant and best answers the question.
## Instructions
1. Analyze the Question: Understand what the user is asking and identify the information that
    needs to be extracted from the database.
2. Evaluate SQL Queries: For each provided SQL query, determine its relevance based on:
   - Accuracy: Does the query correctly match the question's intent?
   - Completeness: Does the query retrieve all the necessary information without omitting
    important details?
    - Efficiency: Is the query optimized for the task, avoiding unnecessary joins or
    conditions?
3. Select the Most Relevant Query: Choose the query that is the best match for the question.
## Database Schema
Database "{database_name}":
{database_schema}
## Ouestion
```

```
{question}
## Hint
{hint}
## SQL Queries
{queries}
## Output Requirement
Reply the query Index in the format of "Index: ".
## Output
.....
query_with_response_selection_prompt = """
You are given a question, a database schema, multiple SQL queries, and their execution results.
     Your task is to select the SQL query that best answers the question based on the query
    and its result.
## Instructions
1. Understand the Question: Determine what the user is asking and identify the specific
    information that needs to be retrieved.
2. Evaluate Each Query and Response Pair: For each provided SQL query and its result,
    determine:
    - Query Accuracy: Does the query correctly represent the user's intent?
    - Result Relevance: Does the result contain the data needed to answer the question
    completely and correctly?
    - Efficiency: Is the query optimized, avoiding unnecessary complexity?
## Database Schema
Database "{database_name}":
{database_schema}
## Ouestion
{question}
## Hint
{hint}
## SQL Queries and Execution Results
{queries}
## Output Requirement
Only reply the query Index in the format of "Index: ".
```