R³: "This is My SQL, Are You With Me?" A Consensus-Based Multi-Agent System for Text-to-SQL Tasks

Hanchen Xia^{*}, Feng Jiang^{*}, Naihao Deng⁹, Cunxiang Wang⁹, Guojiang Zhao^{*} Rada Mihalcea⁹, Yue Zhang⁹

School of Mathematical Science, Shanghai Jiao Tong University
 School of Engineering, Westlake University
 University of Michigan
 Carnegie Mellon University

Abstract

Large Language Models (LLMs) have demonstrated exceptional performance across diverse tasks. To harness their capabilities for Textto-SQL, we introduce \mathbf{R}^3 (Review-Rebuttal-Revision), a consensus-based multi-agent system for Text-to-SQL tasks. \mathbf{R}^3 achieves the new state-of-the-art performance of 89.9 on the Spider test set. In the meantime, \mathbf{R}^3 achieves 61.80 on the Bird development set. \mathbf{R}^3 outperforms existing single-LLM and multi-agent Text-to-SQL systems by 1.3% to 8.1% on Spider and Bird, respectively. Surprisingly, we find that for Llama-3-8B, \mathbf{R}^3 outperforms chain-ofthought prompting by over 20%, even outperforming GPT-3.5 on the Spider development set. We open-source our codebase at https: //github.com/1ring2rta/R3.

1 Introduction

Text-to-SQL, the task of converting natural language to SQL queries, enables non-technical users to access databases with natural language (Deng et al., 2022; Katsogiannis-Meimarakis and Koutrika, 2023). Recently, Large Language Models (LLMs) have made significant progress on various tasks (Touvron et al., 2023; OpenAI, 2023).

Although various methods were proposed to enhance the reasoning abilities of LLMs (Wei et al., 2022; Yao et al., 2023; Besta et al., 2024), they are still facing challenges with Text-to-SQL tasks (Li et al., 2023b; Hong et al., 2024). The LLM-based multi-agent system leverages collective intelligence from a group of LLMs and has achieved exceptional performance across various tasks (Park et al., 2023; Hong et al., 2023; Xu et al., 2023), but little work explores using them on Text-to-SQL. The existing multi-agent Text-to-SQL system first decomposes the task into multiple subtasks, which are then accomplished step-by-step in a pipeline by agents (Wang et al., 2023). While achieving remarkable performance, such decomposition-based



Figure 1: \mathbf{R}^3 Architecture. *n* Reviewer agents, each with distinct characteristics, are created to review the generated SQL and its execution result. The process continues until the master node (SQL-Writer agent) and the other nodes reach a consensus, at which point the system outputs the final SQL.

systems necessitate extensive prompt engineering and logic design.

We propose \mathbf{R}^3 , a consensus-based multi-agent system for Text-to-SQL tasks that draws inspiration from the peer-review mechanism. In our designed framework, the LLM does not need to be divided into sub-tasks such as column selection, schema linking, and so on. Instead, it is split into an SQL Writer and multiple Reviewers who provide feedback based on the execution results. Once the generated SQL query is confirmed to be executable, the system enters a review process, where the execution results guide the SQL Writer and reviewers to refine the SQL. Through rounds of "review," "negotiation or rebuttal," and "revision," the SQL Writer and reviewers ultimately reach a consensus and deliver a solution with collective agreement (see Figure 1).

We test \mathbf{R}^3 on the popular Spider and Bird benchmarks. \mathbf{R}^3 outperforms the existing single LLM as well as the multi-agent Tex-to-SQL systems by 1.3% to 8.1% on Spider and Bird, and set new state-of-the-art (SOTA) performance of 89.9 on Spider

dataset. Surprisingly, we find that for Llama-3-8B, R^3 outperforms chain-of-thought prompting by over 20%, even outperforming GPT-3.5 on the Spider-Dev set.

In summary, our contributions are several-fold:

- 1. To the best of our knowledge, \mathbf{R}^3 is the first Text-to-SQL system to use the execution result for SQL refinements, and the first Text-to-SQL system to equip agents with memory sequences to enhance SQL generation.
- R³ offers a consensus-based multi-agent system for Text-to-SQL tasks. Using very succinct prompts, R³ sets the new SOTA performance of 89.9 on the Spider dataset. In the meantime, R³ achieves 61.80 on the Bird-Dev dataset. In addition, R³ effectively helps open-source LLMs such as Llama-3-8B on SQL generation. When using Llama-3-8B as the backbone model, R³ outperforms direct CoT prompting Llama-3-8B by 20%, and outperforms GPT-3.5 on the Spider-Dev set.
- 3. We provide a detailed error analysis of \mathbf{R}^3 on the existing Text-to-SQL benchmarks, shedding light on future research on the Text-to-SQL task.

2 Related Works

Traditional Methods for Text-to-SQL. The Text-to-SQL conversion task has enjoyed a long history dating back to 1970s (Androutsopoulos et al., 1995), and researchers have kept working on this problem for the past few decades (Dahl et al., 1994; Zelle and Mooney, 1996; Popescu et al., 2003; Zhong et al., 2017; Yu et al., 2018). Before the advent of LLMs, systems like RATSQL (Wang et al., 2019) and LGESQL (Cao et al., 2021) adapt BERT (Devlin et al., 2018) architecture to acquire better representations, and carefully design their techniques to link schema in the database system. Later, approaches like PICARD (Scholak et al., 2021), RASAT (Qi et al., 2022), and RESD-SQL (Li et al., 2023a) adapt the T5 model (Raffel et al., 2020) to translate user questions into SQL query in an end-to-end fashion. Additionally, researchers propose a variety of task-specific strategies like relation-aware self-attention (Qi et al., 2022), schema selection (Li et al., 2023a), and constrained decoding (Scholak et al., 2021) to improve the performance of the Text-to-SQL systems.

LLMs for Text-to-SQL. Recent years have witnessed LLMs' breakthroughs in many fields (Ouyang et al., 2022; Touvron et al., 2023; Dubey et al., 2024). Moreover, Brown et al. (2020); Chen et al. (2022); Liu and Liu (2021) have observed that these LLMs can learn in context with a few examples during their inference time. The strong reasoning and in-context learning capabilities of these LLMs have brought a paradigm shift to the Textto-SQL community, which now focuses on leveraging LLMs' ability to handle Text-to-SQL tasks. For instance, Pourreza and Rafiei (2023) propose DIN-SQL to few-shot prompt GPT-4; Dong et al. (2023) introduce C3, which zero-shots GPT-3.5 with hints and checks output consistency; DAIL-SQL (Gao et al., 2023) comprehensively evaluates the efficiency and effectiveness of various prompting techniques.

Output Consistency. Recent works have applied the consistency principle (Wang et al., 2022) to enhance the reasoning ability of LLMs through incontext learning, such as chain-of-thought (CoT) (Wei et al., 2022) or tree-of-thoughts (ToT) (Yao et al., 2023). In addition, Chen et al. (2023) adopt program-of-thoughts (PoT), which uses Python code to assist LLMs in the reasoning process and surpasses CoT on math reasoning.

3 R^3 Architecture

SQL-Writer We task SQL-Writer (SW) agents to: (1) compose the original SQL query based on the user question and database schema; (2) ensure that the SQL query is executable, and correct it when errors occur; (3) respond to reviewer agents' feedback and revise the SQL query accordingly. Specifically, we prompt SW agent through Prompt 1 in Appendix A.6. For task (1), we feed the Prompt 1 to SW agent directly. Given a user question x and the database schema S, task (1) can be formalized as:

$$y = \mathrm{SW}(x, \mathcal{S}),$$

where y is the generated SQL query. For steps (2) and (3), we maintain a truncated dialogue history, denoted as \mathcal{H} , which is initially set to $\mathcal{H} = [(x, S), y]$. Specifically, if an error e occurs during SQL execution, DB(y), we append e to the history, updating $\mathcal{H} \leftarrow \mathcal{H} + e$. Subsequently, we obtain the updated output, y', via the following process:

$$y' = SW(\mathcal{H}).$$

Algorithm 1: R³-Loop

We then concatenate y' to the history, resulting in the update $\mathcal{H} \leftarrow \mathcal{H} + y'$. Furthermore, in consideration of the context window limitations of LLMs, we truncate the dialogue history \mathcal{H} when the length of the prompt exceeds the model's context limit.

Reviewers. We generate the reviewer agent's (REs) professions using an LLM (see Prompt 3 in Appendix A.6) based on the database schema and the content of the SQL query, for instance, "Senior Database Engineer specialized in writing various clauses" and "Data Analyst in the automotive industry", etc. We incorporate these professions in the system prompt for the REs to make them focus on different aspects of the SQL query. These REs are prompted to provide their professional comments based on the database schema, the user's question, the predicted SQL, and its execution result in the table format.

Overall Architecture. After several rounds of "negotiation" between the SQL-writer and REs, we decide whether there is a consensus by checking if the SQL-writer agent generates the same SQL query as in the previous round. When there is a consensus, we terminate the negotiation loop and output the final SQL query. Algorithm 1 depicts the overall process of our system.

Appendix A.6 provides the detailed prompts we use in \mathbf{R}^3 . In addition, we incorporate:

- 1. Program of Thoughts (PoT) (Chen et al., 2023) to prompt the SQL-writer agent to generate Python code before SQL query (see Prompt 2 in Appendix A.6). Therefore, the agents may leverage Python in their reasoning process for better SQL query generation.
- 2. *k*-shots example selection based on similarity of the user question embeddings. Specifically, when our system infers the SQL query in the test

set, we select the k most similar use questions and their corresponding SQL queries from the training set (k-shots) and use them for in-context learning.

4 Experimental Setup

	Spider-Dev (Yu et a	Spider-Test 1., 2018)	Bird-Dev (Li et al., 2023b)
#QA	1,034	2147	1,534
#Domain	138	-	37
#DB	200	206	95
DB Size	879.5 MB	906.5 MB	1.76 GB

Table 1: Statistics of two Text-to-SQL benchmarks we use in our experiments. "#QA", "#Domain" and "#DB" refer to the number of samples, domains and databases, respectively.

Datasets. We conduct experiments on two crossdomain Text-to-SQL benchmarks, Spider and Bird, detailed in Table 1.

Baselines. We conduct our experiments based on LLMs including GPT-3.5-Turbo, GPT-4 (OpenAI, 2023) and Llama-3 (AI@Meta, 2024). As for the compared methods, the raw performance for GPT-3.5 ("-") was evaluated by Li et al. (2023b); C3 employs schema linking filtering (Dong et al., 2023); DAIL selects few-shot demonstrations based on their skeleton similarities (Gao et al., 2023), and "SC" represents Self-Consistency (Wang et al., 2022); PET uses cross-consistency (Li et al., 2024); DIN decomposes the Text-to-SQL task into smaller subtasks (Pourreza and Rafiei, 2023); MAC, as previously mentioned, is the first to apply a Multi-Agent system to Text-to-SQL tasks (Wang et al., 2023).

Metrics. We employ test-suite execution evaluation¹ (Zhong et al., 2020), the standard evaluation protocol for Spider, and the official SQL execution accuracy evaluation for $Bird^2$.

5 Results and Analysis

5.1 General Results

Table 2 compares \mathbf{R}^{3} 's performance with existing baseline methods when we use GPT-3.5-Turbo or GPT-4 as our backbone models. Our best performed system with GPT-4 as the backbone

¹github.com/taoyds/test-suite-sql-eval

²bird-bench.github.io/

Paalthono	Method	Spi	Spider	
Dackoolic		Dev	Test	Dev
	-	72.1	-	37.22
GPT-3.5	C3 (2023)	81.8	82.3	-
Turbo	MAC (2023)	80.6	75.5	50.56
	\boldsymbol{R}^3 (ours)	81.4	81.1	52.15
	DAIL (2023)	83.6	86.6	-
	PET (2024)	82.2	87.6	-
GPT-4	DIN (2023)	82.8	85.3	50.72
	MAC (2023)	86.8	82.8	59.39
	\boldsymbol{R}^3 (ours)	88.1	89.9	61.80

Table 2: Execution accuracy across existing Text-to-SQL systems. We use the GPT-3.5-Turbo in our experiment. The results for plain GPT-3.5-Turbo (first row) are taken from Li et al. (2023b).

Backbone	Method	Spi	Spider	
		Dev	Test	
GPT-3.5	Li et al. (2023b)	72.1	_	
Turbo	R ³	81.4	81.1	
Llama-3-8B	СоТ	52.1	53.5	
Instruct	R^3	72.8	72.6	

Table 3: Execution accuracy comparison when we employ open-source LLMs as the backbone models with \mathbf{R}^3 on Spider-Dev and Spider-Test. We highlight that \mathbf{R}^3 significantly boosts the open-source LLM's capability on SQL generation.

achieves 88.1%, 89.9%, and 61.8% on the Spider-Dev, Spider-Test, and Bird-Dev respectively, surpassing the existing multi-agent Text-to-SQL systems.

5.2 Discussions

Generalizability of \mathbb{R}^3 **framework.** We test our system with open-source Llama-3 models on Spider and report the results in Table 3. To our surprise, with the help of \mathbb{R}^3 , zero-shot Llama-3-8B outperforms GPT-3.5 performance reported by Li et al. (2023b) on Spider-Dev set. This demonstrates the effectiveness of our proposed \mathbb{R}^3 system.

CoT versus PoT. We conduct an ablation study on the impact of CoT, PoT with one or three reviewer agents in the discussion and report the results in Table 4. The results in Table 4 show that the *n*-Reviewer(s) Loop (*n*R-Lp) plays a major role in performance improvement, with the 3R-Lp configuration significantly outperforming the 1R-Lp setup. The proposed \mathbf{R}^3 system achieves a 10.54% improvement over the baseline GPT-4 + CoT. We

	GPT-3.	5-Turbo	GP	Г-4
	Spider	Bird	Spider	Bird
СоТ	78.2	37.22	79.7	53.30
PoT	78.5	36.96	80.0	54.61
1R-Lp + CoT	78.3	44.13	82.3	57.89
1R-Lp + PoT	79.3	46.35	85.4	58.34
<i>R</i> ³ : 3R-Lp + PoT	81.4	52.15	88.1	61.80

Table 4: Ablation Studies on Spider-Dev and Bird-Dev (Execution Accuracy). The 1-Reviewer Loop (1R-Lp) represents that only one reviewer agent participates in the discussion, while the 3-Reviewers Loop (3R-Lp) represents three in the discussion, which is also the default configuration of \mathbf{R}^3 . We conduct all the experiments here under the 5-shot setting.

provide the statistical significant test for these results in Appendix A.1. Appendix A.2 provides a sensitivity analysis of the impacts of the k value in k-shots.

5.3 Error Analysis

In total, GPT-4+ \mathbf{R}^3 fails to generate the gold SQL queries for 123 instances in Spider-Dev. Table 5 shows the error case distribution for our system on Spider-Dev (more in Appendices A.3 and A.4). Note that though we have spotted issues with the gold SQL queries, we still adopt the original set to calculate the performance of our system to ensure a fair comparison.

Gold Error. We notice that though the annotation quality of Spider is good, there are still cases where the gold SQL queries are not correct. Specifically, among the 151 examples, 30.5% are due to incorrect gold SQL queries (4.5% of all the examples in Spider-Dev). To facilitate future research, we catalog the instances with incorrect gold SQL, correct the errors, and share the details.

Ambiguity. We observe that there are a few questions involving ambiguities, a phenomenon spotted on a wide range of NLP tasks (Plank, 2022; Deng et al., 2023). In Table 5.3, both FullName and Maker columns hold the information for the "name of makers", except that FullName holds the full names while Maker holds the name abbreviations. Therefore, both the gold and predicted SQL queries should be considered correct if there is no further clarifications. Such ambiguous requests may be common in real-world applications as the

Error Types	Question, Gold & Prediction	Explanation
Gold Error (30.5%)	Q: What are the Asian countries which have a population larger than that of any country in Africa? Gold: X AND population > (SELECT min(population) FROM country WHERE Continent = "Africa") Pred: M AND population > (SELECT max(population) FROM country WHERE Continent = "Africa")	Judged as incorrect because of the incorrect gold SQL query.
Logic (29.8%)	Q: How many owners temporarily do not have any dogs? Gold: SELECT count (*) FROM Owners WHERE owner_id NOT IN (SELECT owner id FROM Dogs) Pred: SELECT (SELECT COUNT (DISTINCT owner_id) FROM Owners) - (SELECT COUNT (DISTINCT owner_id) FROM Dogs WHERE date_departed IS NULL)	The predicted SQL query wrongly assumes that all owners have had dogs.
Ambiguity (13.2%)	Q: What are the names of all makers with more than 3 models? Gold: SELECT T1.FullName HAVING count(*) > 3; Pred: SELECT T1.Maker HAVING count(*) > 3;	Both FullName and Maker columns hold the information for "names".
Inaccuracy (11.3%)	Q: What are the arriving date of the dogs who have gone through a treatment? Gold: SELECT T1.date_arrived, FROM Pred: SELECT T1.date_arrived, T1.Name FROM	The selected Name is not asked by the question.
DB Value (10.6%)	Q: Which city and country is the Alton airport at? Gold: SELECT WHERE AirportName = "Alton" ; Pred: SELECT WHERE AirportName LIKE "%Alton%" ;	Our framework notices there is a space for Alton in the DB, therefore employing a fuzzy match.
Others (4.6%)		

Table 5: Error Analysis of \mathbf{R}^3 on Spider-Dev. We make the part in the question red when it is either annotated incorrectly in the gold SQL query (Gold) or predicted incorrectly in the predicted SQL query (Pred).

lay users may not be familiar with the database schema. This requires future research on interactive Text-to-SQL systems that can understand and deal with such ambiguities in user questions.

Dirty Database Value. We observe that due to the Database (DB) setup for Spider, certain DB values may deviate from what is asked in the question. For instance, in Table 5.5, R^3 notices a space for Alton in DB, therefore employing a fuzzy match. But this deviates the SQL query's execution results from the gold SQL query's results.

Logic. In Table 5.2, we present an example of the logic error made by \mathbf{R}^3 . We notice that LLMs may solve the problems using a more complicated logic, which is prone to mistakes. For instance, in Table 5.2, instead of directly counting the owners who do not own dogs, the LLMs try to subtract the number of dog owners from the total number of owners. This ignores the possibility that some owners may have never had any dogs before. This addresses an issue with the multi-agent system that if the system comes up with a complicated initial SQL query, the following discussion process may try to polish the complicated SQL query instead of switching to an easier solution. In cases like Table 5.2, there is no way to reach a perfect SQL query with the subtraction logic.

Inaccuracy. We observe that the LLMs may incorporate more information than what is asked by the end user. For instance, in Table 5.4, the user does not ask for the name of the dogs, but the LLMs present such information along with the requested arrival date. We hypothesize that since such extra information can potentially be helpful to the end user, LLMs may be biased towards including it.

Our findings indicate that the existing evaluation protocols for Text-to-SQL generation may not authentically capture the capabilities of these sophisticated systems. Therefore, we advocate for a reassessment and enhancement of Text-to-SQL evaluation methods. We provide further error analysis of \mathbf{R}^3 on Bird in Appendix A.4.

6 Conclusion

In this paper, we propose \mathbb{R}^3 , a consensus-based multi-agent system for Text-to-SQL generation. \mathbb{R}^3 sets the new SOTA performance on Spider (89.9) and achieves 61.80 on the Bird Dev set. In addition, we find that \mathbb{R}^3 significantly enhances open-source LLMs such as Llama-3-8B (over 20% improvement on Spider Dev set). Last but not least, we conduct a comprehensive error analysis and identify issues with the current Text-to-SQL evaluation, underscoring the necessity for a more refined evaluation protocol, as the LLMs and LLM-based methods become more powerful than ever.

Limitations

Due to the scope of the study, we only test a limited number of LLMs. In this paper, we study the performance gap between 1R-Lp and 3R-Lp. We leave further studies on the effects of the number of reviewers to future research.

Ethical Statements

In this paper, we propose strategies to improve the SQL generation capabilities of LLMs. To the best of our knowledge, we do not expect our system would have negative impacts on society.

References

AI@Meta. 2024. Llama 3 model card.

- Ion Androutsopoulos, Graeme D Ritchie, and Peter Thanisch. 1995. Natural language interfaces to databases–an introduction. *Natural language engineering*, 1(1):29–81.
- Maciej Besta, Nils Blach, Ales Kubicek, Robert Gerstenberger, Michal Podstawski, Lukas Gianinazzi, Joanna Gajda, Tomasz Lehmann, Hubert Niewiadomski, Piotr Nyczyk, et al. 2024. Graph of thoughts: Solving elaborate problems with large language models. In Proceedings of the AAAI Conference on Artificial Intelligence, volume 38, pages 17682–17690.
- Tom Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared D Kaplan, Prafulla Dhariwal, Arvind Neelakantan, Pranav Shyam, Girish Sastry, Amanda Askell, et al. 2020. Language models are few-shot learners. *Advances in neural information processing systems*, 33:1877–1901.
- Ruisheng Cao, Lu Chen, Zhi Chen, Yanbin Zhao, Su Zhu, and Kai Yu. 2021. Lgesql: line graph enhanced text-to-sql model with mixed local and nonlocal relations. *arXiv preprint arXiv:2106.01093*.
- Mingda Chen, Jingfei Du, Ramakanth Pasunuru, Todor Mihaylov, Srini Iyer, Veselin Stoyanov, and Zornitsa Kozareva. 2022. Improving in-context few-shot learning via self-supervised training. *arXiv preprint arXiv:2205.01703*.
- Wenhu Chen, Xueguang Ma, Xinyi Wang, and William W. Cohen. 2023. Program of thoughts prompting: Disentangling computation from reasoning for numerical reasoning tasks. *Transactions on Machine Learning Research*.
- Deborah A. Dahl, Madeleine Bates, Michael Brown, William Fisher, Kate Hunicke-Smith, David Pallett, Christine Pao, Alexander Rudnicky, and Elizabeth Shriberg. 1994. Expanding the scope of the ATIS task: The ATIS-3 corpus. In *Human Language Technology: Proceedings of a Workshop held at Plainsboro, New Jersey, March 8-11, 1994*.

- Naihao Deng, Yulong Chen, and Yue Zhang. 2022. Recent advances in text-to-SQL: A survey of what we have and what we expect. In *Proceedings of the* 29th International Conference on Computational Linguistics, pages 2166–2187, Gyeongju, Republic of Korea. International Committee on Computational Linguistics.
- Naihao Deng, Xinliang Zhang, Siyang Liu, Winston Wu, Lu Wang, and Rada Mihalcea. 2023. You are what you annotate: Towards better models through annotator representations. In *Findings of the Association* for Computational Linguistics: EMNLP 2023, pages 12475–12498, Singapore. Association for Computational Linguistics.
- Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. 2018. Bert: Pre-training of deep bidirectional transformers for language understanding. *arXiv preprint arXiv:1810.04805*.
- Xuemei Dong, Chao Zhang, Yuhang Ge, Yuren Mao, Yunjun Gao, Jinshu Lin, Dongfang Lou, et al. 2023.
 C3: Zero-shot text-to-sql with chatgpt. *arXiv* preprint arXiv:2307.07306.
- Abhimanyu Dubey, Abhinav Jauhri, Abhinav Pandey, Abhishek Kadian, Ahmad Al-Dahle, Aiesha Letman, Akhil Mathur, Alan Schelten, Amy Yang, Angela Fan, et al. 2024. The Ilama 3 herd of models. *arXiv preprint arXiv:2407.21783*.
- 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. *arXiv preprint arXiv:2308.15363*.
- Sirui Hong, Xiawu Zheng, Jonathan Chen, Yuheng Cheng, Jinlin Wang, Ceyao Zhang, Zili Wang, Steven Ka Shing Yau, Zijuan Lin, Liyang Zhou, et al. 2023. Metagpt: Meta programming for multi-agent collaborative framework. *arXiv preprint arXiv:2308.00352*.
- Zijin Hong, Zheng Yuan, Qinggang Zhang, Hao Chen, Junnan Dong, Feiran Huang, and Xiao Huang. 2024. Next-generation database interfaces: A survey of llmbased text-to-sql. *arXiv preprint arXiv:2406.08426*.
- George Katsogiannis-Meimarakis and Georgia Koutrika. 2023. A survey on deep learning approaches for text-to-sql. *The VLDB Journal*, 32(4):905–936.
- Haoyang Li, Jing Zhang, Cuiping Li, and Hong Chen. 2023a. Resdsql: Decoupling schema linking and skeleton parsing for text-to-sql. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 37, pages 13067–13075.
- Jinyang Li, Binyuan Hui, Ge Qu, Binhua Li, Jiaxi Yang, Bowen Li, Bailin Wang, Bowen Qin, Rongyu Cao, Ruiying Geng, et al. 2023b. Can llm already serve as a database interface. A big bench for large-scale database grounded text-to-sqls. CoRR abs/2305.03111.

- Zhishuai Li, Xiang Wang, Jingjing Zhao, Sun Yang, Guoqing Du, Xiaoru Hu, Bin Zhang, Yuxiao Ye, Ziyue Li, Rui Zhao, et al. 2024. Pet-sql: A promptenhanced two-stage text-to-sql framework with crossconsistency. *arXiv preprint arXiv:2403.09732*.
- Yixin Liu and Pengfei Liu. 2021. Simcls: A simple framework for contrastive learning of abstractive summarization. *arXiv preprint arXiv:2106.01890*.
- Henry B Mann and Donald R Whitney. 1947. On a test of whether one of two random variables is stochastically larger than the other. *The annals of mathematical statistics*, pages 50–60.
- R OpenAI. 2023. Gpt-4 technical report. *arXiv*, pages 2303–08774.
- Long Ouyang, Jeffrey Wu, Xu Jiang, Diogo Almeida, Carroll Wainwright, Pamela Mishkin, Chong Zhang, Sandhini Agarwal, Katarina Slama, Alex Ray, et al. 2022. Training language models to follow instructions with human feedback. *Advances in neural information processing systems*, 35:27730–27744.
- Joon Sung Park, Joseph O'Brien, Carrie Jun Cai, Meredith Ringel Morris, Percy Liang, and Michael S Bernstein. 2023. Generative agents: Interactive simulacra of human behavior. In *Proceedings of the 36th Annual ACM Symposium on User Interface Software and Technology*, pages 1–22.
- Barbara Plank. 2022. The "problem" of human label variation: On ground truth in data, modeling and evaluation. In *Proceedings of the 2022 Conference* on Empirical Methods in Natural Language Processing, pages 10671–10682, Abu Dhabi, United Arab Emirates. Association for Computational Linguistics.
- Ana-Maria Popescu, Oren Etzioni, and Henry Kautz. 2003. Towards a theory of natural language interfaces to databases. In *Proceedings of the 8th international conference on Intelligent user interfaces*, pages 149–157.
- Mohammadreza Pourreza and Davood Rafiei. 2023. Din-sql: Decomposed in-context learning of text-to-sql with self-correction. *arXiv preprint arXiv:2304.11015*.
- Jiexing Qi, Jingyao Tang, Ziwei He, Xiangpeng Wan, Yu Cheng, Chenghu Zhou, Xinbing Wang, Quanshi Zhang, and Zhouhan Lin. 2022. Rasat: Integrating relational structures into pretrained seq2seq model for text-to-sql. *arXiv preprint arXiv:2205.06983*.
- Colin Raffel, Noam Shazeer, Adam Roberts, Katherine Lee, Sharan Narang, Michael Matena, Yanqi Zhou, Wei Li, and Peter J Liu. 2020. Exploring the limits of transfer learning with a unified text-to-text transformer. *The Journal of Machine Learning Research*, 21(1):5485–5551.
- Torsten Scholak, Nathan Schucher, and Dzmitry Bahdanau. 2021. Picard: Parsing incrementally for constrained auto-regressive decoding from language models. *arXiv preprint arXiv:2109.05093*.

- Hugo Touvron, Louis Martin, Kevin Stone, Peter Albert, Amjad Almahairi, Yasmine Babaei, Nikolay Bashlykov, Soumya Batra, Prajjwal Bhargava, Shruti Bhosale, et al. 2023. Llama 2: Open foundation and fine-tuned chat models. *arXiv preprint arXiv:2307.09288*.
- Bailin Wang, Richard Shin, Xiaodong Liu, Oleksandr Polozov, and Matthew Richardson. 2019. Rat-sql: Relation-aware schema encoding and linking for textto-sql parsers. arXiv preprint arXiv:1911.04942.
- Bing Wang, Changyu Ren, Jian Yang, Xinnian Liang, Jiaqi Bai, Qian-Wen Zhang, Zhao Yan, and Zhoujun Li. 2023. Mac-sql: Multi-agent collaboration for text-to-sql. *arXiv preprint arXiv:2312.11242*.
- Xuezhi Wang, Jason Wei, Dale Schuurmans, Quoc Le, Ed Chi, Sharan Narang, Aakanksha Chowdhery, and Denny Zhou. 2022. Self-consistency improves chain of thought reasoning in language models. *arXiv preprint arXiv:2203.11171*.
- Jason Wei, Xuezhi Wang, Dale Schuurmans, Maarten Bosma, Fei Xia, Ed Chi, Quoc V Le, Denny Zhou, et al. 2022. Chain-of-thought prompting elicits reasoning in large language models. *Advances in Neural Information Processing Systems*, 35:24824–24837.
- Zelai Xu, Chao Yu, Fei Fang, Yu Wang, and Yi Wu. 2023. Language agents with reinforcement learning for strategic play in the werewolf game. *arXiv* preprint arXiv:2310.18940.
- Shunyu Yao, Dian Yu, Jeffrey Zhao, Izhak Shafran, Thomas L Griffiths, Yuan Cao, and Karthik Narasimhan. 2023. Tree of thoughts: Deliberate problem solving with large language models. *arXiv preprint arXiv:2305.10601*.
- Tao Yu, Rui Zhang, Kai Yang, Michihiro Yasunaga, Dongxu Wang, Zifan Li, James Ma, Irene Li, Qingning Yao, Shanelle Roman, et al. 2018. Spider: A large-scale human-labeled dataset for complex and cross-domain semantic parsing and text-to-sql task. arXiv 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, pages 1050–1055.
- Ruiqi Zhong, Tao Yu, and Dan Klein. 2020. Semantic evaluation for text-to-sql with distilled test suites. *arXiv preprint arXiv:2010.02840*.
- Victor Zhong, Caiming Xiong, and Richard Socher. 2017. Seq2sql: Generating structured queries from natural language using reinforcement learning. *arXiv* preprint arXiv:1709.00103.

A Appendix

A.1 Significance Test

We divided the generated SQL by several strategies in Table 4 into 10 equal parts and calculated the execution accuracy for each. To test whether our strategy can indeed improve execution accuracy, we conduct a significance test between the "CoT" and "3R-Lp+PoT" strategies. The null hypothesis of the test is that the median execution accuracy obtained by the two strategies is the same. The Mann-Whitney U Test (Mann and Whitney, 1947) is a non-parametric statistical method used to compare whether there is a significant difference in the medians of two independent samples. Compared to the Analysis of Variance (ANOVA), it does not require the data to be normally distributed, making it suitable for small samples or data with unknown distribution.

The *p*-value of the test is 0.0024, which is below the commonly accepted significance level of 0.05. Therefore, we have reason to reject the null hypothesis, indicating that the "3R-Lp+PoT" strategy leads to a significant performance improvement.

Effects of the number of "reviewer" agents.

A.2 Effects of k in k-shot.



Figure 2: k-shot Sensitivity Analysis.

We test various k values on 200 random samples

from Spider-Dev. As shown in Figure 2, compared to CoT, the performance of the \mathbf{R}^3 system remains relatively stable regardless of the number of examples, which corroborates our previous findings from the 0-shot experiments with Llama-3.

A.3 Spider Error Cases

Error Type	Question, Gold & Prediction	Reason
DB Value	Q: Find the last name of the students who currently live in the state of North Carolina but have not registered in any degree program. Gold: SELECT WHERE T2.state_province_county = 'NorthCarolina' EXCEPT Pred: SELECT WHERE T2.state_province_county = 'North Carolina' EXCEPT	The filtering condition in the question does not match the database value, string "NorthCalifornia" in database do not have a space in between.
Gold Error	<pre>Q: What are the first names of all players, and their average rankings? Gold: SELECT avg(ranking), T1.first_name FROM players AS T1 JOIN rankings AS T2 ON T1.player_id = T2.player_id GROUP BY T1.first_name Pred: SELECT avg(ranking), T1.first_name FROM players AS T1 JOIN rankings AS T2 ON T1.player_id = T2.player_id GROUP BY T1.player_id</pre>	The individuals in the table can be uniquely determined by column player_id not first_name, when GROUP BY.
Gold Error	<pre>Q: Find the id and cell phone of the professionals who operate two or more types of treatments. Gold: SELECT T1.professional_id, T1.cell_number FROM Professionals AS T1 JOIN Treatments AS T2 ON T1.professional_id = T2.professional_id GROUP BY T1.professional_id HAVING count(*) >= 2 Pred: SELECT T1.professional_id, T1.cell_number FROM Professionals AS T1 JOIN Treatments AS T2 ON T1.professional_id = T2.professional_id GROUP BY T1.professional_id HAVING COUNT (DISTINCT T2.treatment_type_code) >= 2</pre>	The gold only finds professionals who have two or more records in the treatment table does not ensure that the records are for different types of treatments
Ambiguity	<pre>Q: What are the names and ids of all makers with more than 3 models? Gold: SELECT T1.FullName, T1.Id FROM CAR_MAKERS AS T1 JOIN MODEL_LIST AS T2 ON T1.Id = T2.Maker GROUP BY T1.Id HAVING count(*) > 3; Pred: SELECT T1.Maker, T1.Id FROM CAR_MAKERS AS T1 JOIN MODEL_LIST AS T2 ON T1.Id = T2.Maker GROUP BY T1.Id HAVING count(*) > 3;</pre>	Both column "Maker" and column "FullName" can answer the question about the "names of makers" in the query.
Imprecise	Q: What are the arriving date and the departing date of the dogs who have gone through a treatment? Gold: SELECT DISTINCT T1.date_arrived, T1.date_departed FROM Dogs AS T1 JOIN Treatments AS T2 ON T1.dog_id = T2.dog_id Pred: SELECT DISTINCT T1.date_arrived, T1.date_departed, T1.Name FROM Dogs AS T1 JOIN Treatments AS T2 ON T1.dog_id = T2.dog id	The question do not require listing the specific names of the students, but only ask to list the students' arrival and departure dates. This falls under information redundancy.

Table 6: Spider error cases.

A.4 BIRD Error Cases

Error Type	Question, Gold & Prediction	Reason
DB Value	Q: How many cards with unknown power that can't be found in foil is in duel deck A? Gold: SELECT SUM(CASE WHEN power LIKE '%*%' OR power IS NULL THEN 1 ELSE 0 END) FROM cards WHERE hasFoil = 0 AND duelDeck = 'a' Pred: SELECT COUNT(*) FROM cards WHERE (power IS NULL OR power = '*') AND hasFoil = 0 AND duelDeck = 'a'	Values in database cannot exact match with the question and evidence.
Gold Error	<pre>Q: How many artists have designed a card with a black border color and is available in both "arena" and "mtgo" printing type? Gold: SELECT COUNT (CASE WHEN availability LIKE '%arena, mtgo%' THEN 1 ELSE NULL END) FROM cards Pred: SELECT COUNT (DISTINCT artist) FROM cards WHERE borderColor = 'black' AND availability LIKE '%arena, mtgo%'</pre>	The gold missed one filtering condition.
Ambiguity	Q: Among black card borders, which card has full artwork? Gold: SELECT id FROM cards WHERE borderColor = 'black' AND isFullArt = 1 Pred: SELECT name FROM cards WHERE borderColor = 'black' AND isFullArt = 1;	Both column "name" and column "id" can answer the question about the "the card has full artwork" in the query.
Logic	Q: Which cities have the top 5 lowest enrollment number for students in grades 1 through 12? Gold: SELECT T2.City FROM frpm AS T1 INNER JOIN schools AS T2 ON T1.CDSCode = T2.CDSCode GROUP BY T2.City ORDER BY SUM(T1.'Enrollment (K-12)') ASC LIMIT 5 Pred: SELECT T2.City, T1."Enrollment (K-12)" FROM frpm AS T1 JOIN schools AS T2 ON T1.CDSCode = T2.CDSCode ORDER BY T1."Enrollment (K-12)" ASC LIMIT 5;	The Pred query is incorrect because it lacks grouping and aggregation. It directly selects cities and their corresponding "Enrollment (K-12)" without summing up the enrollment numbers across all schools within each city. As a result, it returns the cities of the five schools with the lowest enrollment numbers, which does not accurately answer the question seeking the top 5 cities with the lowest total enrollment for students in grades 1 through 12.

Table 7: Bird error cases.

A.5 Interactive Correction Cases

Before	After	Reason		
Q : What is the phone number of the school that has the highest number of test takers with an SAT score of over 1500?				
Before: SELECT schools.Phone FROM satscores JOIN schools ON satscores.cds = schools.CDSCode WHERE satscores.NumGE1500 > 1500 GROUP BY schools.School ORDER BY COUNT(*) DESC LIMIT 1;	After: SELECT schools.Phone FROM satscores JOIN schools ON satscores.cds = schools.CDSCode ORDER BY satscores.NumGE1500 DESC LIMIT 1;	The WHERE clause is checking if the number of test takers with an SAT score over 1500 is greater than 1500, which is not what the question is asking. The question is asking for the school with the highest number of test takers who scored over 1500. Additionally, the GROUP BY clause is grouping by the school name, which may not be unique. It would be better to group by the school code (CDSCode) which is		

Q: List the top five schools, by descending order, from the highest to the lowest, the most number of Enrollment (Ages 5-17). Please give their NCES school identification number.

Before:	After:	The question asks for the top five
SELECT T1."School Name",	SELECT T2.NCESSchool	schools with the most number of
T2.NCESSchool	FROM frpm AS T1	Enrollment (Ages 5-17), and their
FROM frpm AS T1	JOIN schools AS T2 ON	NCES school identification number.
JOIN schools AS T2 ON	T1.CDSCode = T2.CDSCode	However, the SQL query provided
T1.CDSCode = T2.CDSCode	ORDER BY T1."Enrollment	returns the school name, which is not
ORDER BY T1."Enrollment	(Ages 5-17)" DESC	required by the question.
(Ages 5-17)" DESC	LIMIT 5;	
LIMIT 5;		

Q: find the package option of the tv channel that do not have any cartoon directed by Ben Jones.

Before: After: The current query is selecting all SELECT SELECT channels where the director is not Ben Jones, but it should exclude channels TV_Channel.Package_Option TV_Channel.Package_Option FROM TV_Channel that have any cartoons directed by Ben FROM TV_Channel LEFT JOIN Cartoon ON WHERE TV_Channel.id NOT IN Jones. This can be achieved by using a subquery to find all channels that have TV_Channel.id = cartoons directed by Ben Jones, and Cartoon.Channel SELECT Cartoon.Channel WHERE Cartoon.Directed_by FROM Cartoon then excluding these channels in the IS NOT 'Ben Jones'; main query. Here is the corrected SQL: WHERE Cartoon.Directed_by = 'Ben Jones');

Q: List the earnings of poker players in descending order.

 Before:
 After:
 The question only asks for the earnings

 SELECT * FROM poker_player
 SELECT Earnings FROM
 of poker players. Therefore, we should

 ORDER BY Earnings DESC;
 Description
 only select the 'Earnings' column to

 make the result more concise and
 relevant to the question.

Q: How many countries speak both English and Dutch?

Before:	After:	The SQL query provided uses the 'IN'
SELECT COUNT (DISTINCT	SELECT COUNT(*)	clause, which selects countries that
CountryCode)	FROM (speak either English or Dutch. However,
FROM countrylanguage	SELECT CountryCode	the question asks for countries speaking
WHERE Language = 'English'	FROM countrylanguage	both languages. Using IN results in a
AND CountryCode IN (WHERE Language = 'English'	union of sets (either English or Dutch),
SELECT CountryCode	INTERSECT	not an intersection (both English and
FROM countrylanguage	SELECT CountryCode	Dutch). This could be done with two
WHERE Language = 'Dutch')	FROM countrylanguage	separate queries joined by an
	WHERE Language = 'Dutch')	'INTERSECT' clause.

Table 8: Interactive correction cases.

A.6 Prompt Examples

Describe how you understand the question based on the evidence, and help me write an SQL to answer the question. ### EVIDENCE: { evidence } ### USER_QUESTION: { question } ### RELATED SQL: {related_sql} ### DATABASE STRUCTURE: {schema} Prompt 1: CoT-SQL-Writer Write an to answer the question. Program of Thoughts (PoT) is a variant of Chain of Thought (CoT), pre-generating Python code to assist in the creation of SQL. Please apply PoT (and PoT only) before generating an SQL. In your python code, `Table %s` is stored in `db_dict['%s']`, ` db_dict` is of type dict[pandas.DataFrame]. ### RELATED SQL: {related_sqls} ### DATABASE STRUCTURE: {schema} ### EXAMPLES: QUESTION: What is %s in the earliest year and what year was it? SQL: earliest_year = db_dict[%s]['Year'].min() year_filtered_data = step1_result[step1_result['Year'] == earliest_year] result = year_filtered_data[[%s, 'Year']] ```sql SELECT T1.%s, T2.Year FROM %s AS T1 JOIN %s AS T2 ON T1.Id = T2.Id WHERE T2.Year = (SELECT min(YEAR) FROM %s); . . . QUESTION: Show names for all %s except for %s having a %s in year 2023. SQL: %s_2023 = db_dict['%s'][db_dict['%s']['year'] == '2023'] result = db_dict[%s][~db_dict[%s][%s].isin(%ss_2023[%s])] ```sql SELECT name FROM %s EXCEPT SELECT T2.name FROM %s AS T1 WHERE T1. year = 2023

QUESTION: Find the %s that %s is A and B? SQL:

. . .

```
condition_a_data = db_dict[%s][db_dict['Cartoon'][%s] == 'A']
condition_b_data = db_dict[%s][db_dict['Cartoon'][%s] == 'B']
result = pd.merge(condition_a_data, condition_b_data, how='inner')
```sql
SELECT T1.%s FROM %s AS T1 WHERE %s = 'A'
INTERSECT
SELECT T1.%s FROM %s AS T1 WHERE %s = 'B'
. . .
EVIDENCE: { evidence }
USER_QUESTION: { question }
SQL:
 Prompt 2: PoT-SQL-Writer
You are the manager of a Database project. You are going to invite
{n} experts to review an SQL query.
Who would you invite?
considering:
(1) the domain of this database;
(2) the structure of this SQL.
Please write your invitation as a JSON format dictionary, Enclose
the JSON within ```json...```.
DATABASE STRUCTURE:
{schema}
QUESTION: { question }
SQL:
{pred_sql}
EXAMPLES:
```json
{
  "Reviewer PVsg": "Data Analyst in automotive industry",
  "Reviewer 2KtR": "Senior Database Engineer specialized in writing
  various clauses",
  "Reviewer LmN3": "Senior Database Engineer specialized in writing
  filtering conditions"
}
...
### INVITATION:
```

Prompt 3: Invitation