

# CREDIT-SQL: Few-shot prompting for context-dependent text-to-SQL with regularized examples from diversity sampling

Anonymous ACL submission

## Abstract

In this paper, we propose a few-shot prompting method called CREDIT-SQL for the context-dependent text-to-SQL problem. CREDIT-SQL converts each question in a multi-turn dialogue into a self-contained question with a fixed few-shot prompt. Once a self-contained question is obtained, CREDIT-SQL converts it into an SQL query using a prompt made of in-context examples selected by diversity sampling and subsequent example voting. CREDIT-SQL with ChatGPT 3.5 achieves 58.6% in terms of the exact set match without values on the dev set of CoSQL, which is the performance comparable to the state-of-the-art models for context-dependent text-to-SQL. We also argue that the example voting we introduced in CREDIT-SQL can serve as an efficient and effective way to mitigate the instability of in-context example selection in general.

## 1 Introduction

Information retrieval from structured knowledge sources is an NLP task widely applicable in many areas. Text-to-SQL is a promising approach to achieve this goal due to the popularity of SQL as an interface between the user and the database. Text-to-SQL systems have shown remarkable improvements (Wang et al., 2020a,b; Lin et al., 2020; Cao et al., 2021; Scholak et al., 2021; Cai and Wan, 2020) along with the rapid advancements of sequence-to-sequence models including the infamous transformer model (Vaswani et al., 2017). The advantages of these advanced sequence-to-sequence models have been often utilized by fine-tuning pre-trained decoder-encoder models. However, these advanced models become so large that they are called large language models (LLMs) which typically have parameter size ranges from a few tens of billion to a few hundreds of billion (Ye et al., 2023; OpenAI, 2023; Touvron et al., 2023a,b; Anil et al., 2023). Because of this large model size,

it takes too much resource to fine-tune these large language models on custom datasets.

To utilize the advantages of advanced sequence-to-sequence models without investing full resources for fine-tuning, in-context learning with zero-shot prompts or few-shot prompts has become popular recently. In the case of few-shot in-context learning, a few in-context examples are listed in the prompt along with a brief instruction, and the LLM outputs the desired sequence as the response to the input prompt. Although the limited context size of available LLMs only allows a handful of in-context examples to be included in each prompt, it has been shown that strategic designs of prompts can perform as well as fine-tuned models in the tasks of text-to-SQL (Pourreza and Rafiei, 2023; Nan et al., 2023; Dong et al., 2023; Gao et al., 2023).

Still, most in-context learning studies on text-to-SQL tasks focus on the context-independent setting where the system needs to answer a single SQL query on the input of a single question. This context-independent setting becomes particularly inconvenient when one needs to develop conversational information retrieval systems where previous questions or answers can implicitly appear in the user’s later questions.

To address this problem, we propose a few-shot prompting method called CREDIT-SQL<sup>1</sup> in this paper (See Figure 1, 2, and 3). CREDIT-SQL does context-dependent question rephrasing to convert a multi-turn text-to-SQL task on each dialogue into a series of question-query pair text-to-SQL tasks. Once all the questions are rephrased, a diversity-sampled prompt is used to address the text-to-SQL tasks. This prompt is composed of examples obtained via multiple trials of diversity sampling and subsequent example voting. Each example in the prompt is represented as a pair of the rephrased question and the regularized SQL query along with

<sup>1</sup>Context-dependent Regularized Examples from DIversity sampling for Text-to-SQL.

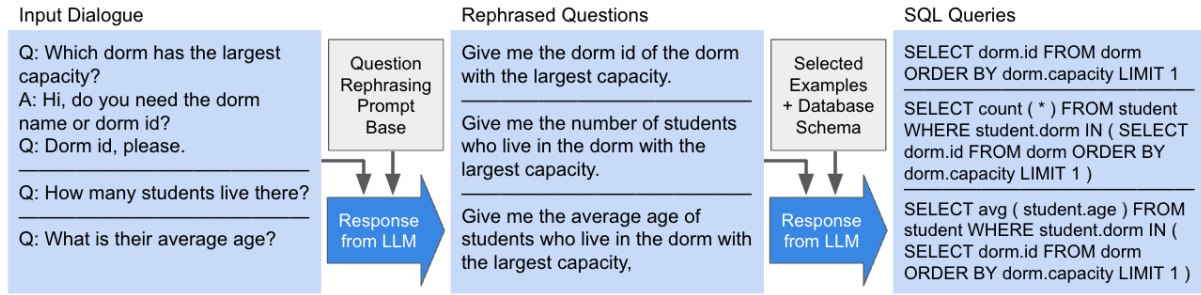


Figure 1: Illustration of overall question to SQL process of CREDIT-SQL. The details of the question rephrasing process are illustrated in Figure 2, and the details of the example selection process are illustrated in Figure 3.

the database schema.

We claim that our approach makes the following contributions. (1) To the best of our knowledge, we propose the first few-shot prompting approach to perform the dialogue state tracking task with systematically selected in-context examples out of the entire training data. (2) We report the performance comparable to the state-of-the-art models in the dialogue state tracking task on the CoSQL dataset, which is ranked 4th on the exact set match without values on the CoSQL dev set among the models reported on the CoSQL leaderboard at the moment of writing. (3) We suggest a new method to mitigate the instability of in-context example selection in the few-shot prompting with LLM by introducing voting on collected examples.

## 2 Related works

### 2.1 Context-dependent text-to-SQL

Most studies on text-to-SQL tasks focused on context-independent settings where a single question is transcribed into a single SQL query. To cope with the complicated scenarios where multiple tables are involved, utilizing graph structures to capture the relations between entities has been the most popular and successful method in text-to-SQL tasks recently (Bogin et al., 2019; Wang et al., 2020a,b; Lin et al., 2020; Cao et al., 2021; Scholak et al., 2021; Cai et al., 2021; Hui et al., 2022).

Unlike its context-independent counterpart, context-dependent text-to-SQL tasks require encoding the context within the dialogue and exploiting this context in the SQL generation. Recently, numerous different approaches have been suggested to tackle this problem. Zhang et al. 2019 used turn attention to edit the SQL query of the previous turn to accommodate the question at the current turn. Cai and Wan 2020 extends the graph structure for the database schema to establish connections be-

tween neighboring turns in the dialogue. Wang et al. 2021 and Hui et al. 2021 suggested using a graph structure state tracker to capture the context of the dialogue at each turn, while Zheng et al. 2022 used BERT to encode the history of the dialogue. Pan et al. 2019, Chen et al. 2021, and Chai et al. 2023 rephrased the question at each turn reflecting the context of the dialogue. Xiao et al. 2022 applied question rephrasing recursively and introduced consistency training to build one of the state-of-the-art models at the time of writing. Other state-of-the-art models used utterance dependency tracking with weighted contrastive learning (Cai et al., 2022) or integrating relational structure through the attention layer into the pre-trained models (Qi et al., 2022).

### 2.2 Prompting with large language models for text-to-SQL

As in its non-prompting counterpart, most of the efforts to perform text-to-SQL tasks focused on context-independent text-to-SQL tasks. Pourreza and Rafiei 2023 used different in-context examples for each difficulty of the question. Dong et al. 2023 achieved one of the best performance among zero-shot prompting efforts. Beyond the arbitrary selection of in-context examples, there have been trials to choose in-context examples in a systematic manner. Liu et al. 2022 introduced question-similarity based example selection using  $k$ -NN algorithm. Nan et al. 2023 noticed that the diversity of the in-context examples in a prompt is indeed important, and suggested methods to balance the similarity and the diversity of in-context examples based on their ground-truth SQL queries. Gao et al. 2023 achieved the best performance on the leaderboard of Spider (Yu et al., 2018) at the time of writing, with systematic in-context example selection using both questions and SQL queries. Beyond

the context-independent setting, [Hu et al. 2022](#) addressed context-dependent text-to-SQL tasks with few-shot prompting and zero-shot prompting, but systematic in-context example selection was still lacking.

### 3 Methods

#### 3.1 Context-dependent question rephrasing

Inspired by the success of question rephrasing approaches ([Pan et al., 2019](#); [Chen et al., 2021](#); [Chai et al., 2023](#); [Xiao et al., 2022](#)), we summarize the context of each dialogue into a single question with few-shot prompting at every turn. In each multi-turn interaction, rephrased questions from previous turns are appended in the prompt to rephrase questions in further turns. The entire process of question rephrasing is illustrated in Figure 2, and the fixed prompt base we used for question rephrasing is in Appendix A.

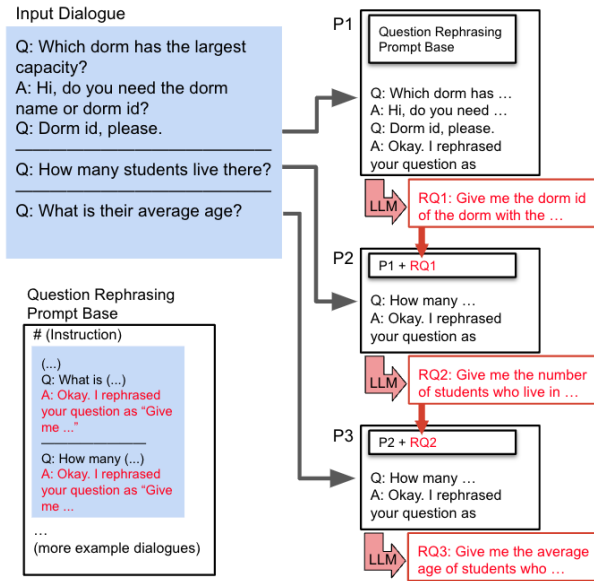


Figure 2: Context-dependent question rephrasing process. A prompt containing examples of multiple dialogues is used to rephrase each question into a rephrased question starting with "Give me ...". Each rephrased question is appended to the prompt for further question rephrasing of later turns.

This question rephrasing process converts the context-dependent text-to-SQL task into the simple text-to-SQL task which is better studied in the literature than the context-dependent counterpart. This process also regularizes the questions in a similar format ("Give me ...") which can help to create consistent in-context examples for the text-to-SQL tasks. This helps LLM to focus more on transcribing

relevant natural language expressions into SQL expressions rather than on deciphering the meaning of the questions written in different styles.

#### 3.2 SQL query regularization

To increase the consistency of the in-context examples for the text-to-SQL tasks, we regularize the SQL queries used in text-to-SQL prompts with rule-based methods. This regularization includes capitalization, spacing, unaliasing, table representation in each column reference, and so on. Examples of affected SQL queries through this regularization are shown in Table 1.

Given SQL	Regularized SQL
select * from tb_1;	SELECT * FROM tb_1
SELECT T1.C1 FROM tb_1 as T1 JOIN tb_2 as T2 on T1.C3=T2.C4	SELECT tb_1.c1 FROM tb_1 JOIN tb_2 ON tb_1.c3 = tb_2.c4
select COUNT(*) from tb_1 where c2=="A"	SELECT count ( * ) FROM tb_1 WHERE tb_1.c2 == 'A'

Table 1: Examples of affected SQL query expressions through the rule-based SQL query regularization.

Similar to the question rephrasing, SQL query regularization helps LLM to focus more on the grammatical structure of SQL queries or links to the database schema rather than on different expression styles of SQL queries.

#### 3.3 Example selection for text-to-SQL

In the few-shot prompting with LLM, the performance of the model is very sensitive to the choice of in-context examples. In particular, strategic sampling of examples out of the training data significantly outperforms the random choices of examples. A natural way to customize in-context examples for each question is to collect the closest examples to the given question, often based on the similarity in the embedding vector space ([Liu et al., 2022](#)). In the meanwhile, [Nan et al. 2023](#) pointed out that keeping the diversity of the example pool can be more important than a mere collection of similar examples. As introduced in [Nan et al. 2023](#), we adopt diversity sampling through the  $k$ -means clustering based on the vectorized SQL queries. Specifically, we perform the clustering with  $k = N$  for the  $N$ -example prompt. In each cluster, we choose the example closest to the centroid in the SQL vector space. In case the number

of resulting clusters  $N_c$  is smaller than  $N$ , we fill the rest of the examples based on the distance with previous examples. To be specific, among training data examples  $\{e_1, e_2, \dots, e_T\}$  with SQL vectors  $\{\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_T\}$ ,  $n$ th selected example  $s_n$  is

$$\begin{aligned} s_n &= e_{i_n}, \\ i_n &= \underset{i}{\operatorname{argmin}} |\mathbf{x}_i - \mathbf{c}_n| \text{ for } n \leq N_c, \\ i_n &= \underset{i \notin \{i_1, \dots, i_{n-1}\}}{\operatorname{argmax}} \left( \min_{j < n} |\mathbf{x}_i - \mathbf{x}_{j}| \right) \\ &\text{for } N_c < n \leq N, \end{aligned}$$

where  $\{\mathbf{c}_1, \mathbf{c}_2, \dots, \mathbf{c}_{N_c}\}$  are the SQL vectors of cluster centroids.

Since the  $k$ -means clustering is a non-deterministic algorithm that depends on the random initial positions of centroids, the selection result of examples varies depending on the choice of random seed. To mitigate this inconsistency of example selection, we adopt example voting. First, we collect  $N$ -example prompt through the  $k$ -means clustering. We repeat this example collection  $M$  times, with a different random seed each time. Then we rank each example by its occurrence among these  $M$  different sets of examples. The entire example selection process is illustrated in Figure 3. With some parameter search, we obtained the best result with  $N = 18$  and  $M = 20$ . For our CREDIT-SQL, we used in-context examples selected by this voting process. We show these selected in-context examples in Appendix B.

### 3.4 Example demonstration

In our few-shot prompt for text-to-SQL, we demonstrate selected in-context examples along with the database ID and the database schema including table names, column names, and foreign keys. A sample text-to-SQL prompt for the CREDIT-SQL approach including the example demonstration is in Appendix C.

## 4 Experiments

### 4.1 Dataset

The most popular benchmark for the text-to-SQL task is Spider dataset (Yu et al., 2018), which is a large-scale, complex, and cross-domain dataset with 10k+ questions with annotated SQL queries and covers 200 different databases across 138 domains. SParC (Yu et al., 2019b) is a multi-turn version of the Spider dataset which covers the

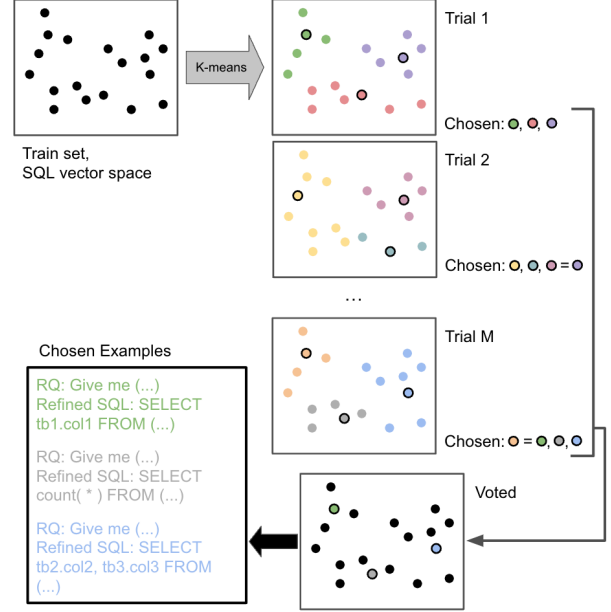


Figure 3: Example selection process using diversity sampling and subsequent example voting. Centroids obtained through the  $k$ -means clustering ( $k = N$ ) are used to pick  $N$  examples. After repeating this selection  $M$  times, aggregated in-context examples are voted by their counts to finally choose  $N$  examples.

same sets of databases as Spider. CoSQL (Yu et al., 2019a) is a dialogue version of Spider and SParC, which includes about 3k dialogues with 10k+ annotated SQL queries over the same sets of databases. CoSQL is different from SParC for it contains turns that does not require immediate SQL query response, such as clarification of the question. This makes CoSQL a more suitable dataset for the development of conversational systems for information retrieval from structured knowledge sources. Since CoSQL is the latest and most complex context-dependent text-to-SQL dataset available at the moment of writing, we benchmark our approach with this dataset.

### 4.2 Models

Throughout the paper, we utilize LLMs of Open AI serviced through Microsoft Azure<sup>2</sup> for the few-shot prompting. We used OpenAI Python API library<sup>3</sup> to access these models. For the baseline approach and the best-performing CREDIT-SQL, we used gpt-3.5-turbo-0301 and we restricted the output size to 600 tokens to accommodate the maximum

<sup>2</sup><https://learn.microsoft.com/en-us/azure/ai-services/openai/concepts/models>

<sup>3</sup><https://github.com/openai/openai-python>. This library is under Apache-2.0 license and we complied to the license.



Model	EM (%)				EX (%)			
	QM		IM		QM		IM	
	Dev	Test	Dev	Test	Dev	Test	Dev	Test
STAR (Cai et al., 2022)	59.7	57.8	30.0	28.2	-	-	-	-
CQR-SQL (Xiao et al., 2022)	58.5	58.3	31.1	27.4	-	-	-	-
RASAT+PICARD (Qi et al., 2022)	58.8	55.7	27.0	26.5	67.0	66.3	39.6	37.4
<i>Few-shot Prompting</i>								
Baseline, randomly sampled dialogues	49.9	-	21.6	-	64.0	-	35.0	-
CREDIT-SQL	58.6	-	25.8	-	66.4	-	38.9	-

Table 2: Results on the CoSQL dataset. Exact set match without values (EM) and execution accuracy with values (EX) are presented for both the question match (QM) and the interaction match (IM). For the few-shot prompting methods, we present the results of the baseline approach (randomly sampled dialogues) as well as the results of CREDIT-SQL, both with the average performance of 5 repeated experiments. Results for other models are as reported in the literature for comparison.

context size of 8192 tokens. For all experiments, we set the temperature to 0 for the consistency of the result.

### 4.3 Evaluation metric

To evaluate the performance of our approach, we use two following metrics as suggested for the SQL-grounded dialogue state tracking task in CoSQL challenge (Yu et al., 2019a):

- **Exact set match without values (EM):** Measures if the predicted SQL query and the ground truth SQL query are equivalent to each other, by comparing the equivalence of each component of the queries. When multiple parallel items are compared, set equivalence is measured so that it does not prefer a particular ordering. It also masks literal/numeral values when comparing each component.
- **Execution accuracy with values (EX):** Measures if the both outputs of the predicted SQL query and the ground truth SQL query are equal to each other. To generate actual outcomes based on the database, it uses the values in each query as they are.

The two accuracy metrics are evaluated at the question level (question match, QM) and at the interaction (dialogue) level (interaction match, IM).

### 4.4 Baseline approach: randomly sampled dialogues

For comparison purposes, we establish a baseline few-shot prompting method. It randomly samples multiple dialogues and presents them along with their database schema. The template of this baseline prompt is in Appendix D.

## 4.5 Experiment results

We report the performances of our baseline approach and CREDIT-SQL on the CoSQL dataset in Table 2, along with the performances of state-of-the-art models. In the dev set, our CREDIT-SQL method outperforms the baseline approach of few-shot prompting using randomly sampled dialogues by 8.7%p in EM for the question match while it outperforms the baseline approach by 2.4%p in EX for the question match. Our approach also shows comparable performance to the state-of-the-art models, by the margin of 0.2%p  $\sim$  1.1%p in EM for the question match on the dev set, and 0.6%p in EX for the question match on the dev set. At the time of writing, our approach ranks to 4th in EM-QM on the CoSQL dev set among the models reported on the CoSQL leaderboard.

## 5 Discussion

### 5.1 Ablation study

We conducted an ablation study to find out the effectiveness of each component of our CREDIT-SQL approach. The study result is reported in Table 3. This study indicates that subtracting SQL regularization from the prompt drops the EM for the question match by 0.6%p. When the prompt with examples collected by diversity sampling is replaced with a prompt with randomly sampled examples, the EM for the question match drops by 5.5%p.

### 5.2 Effectiveness of example voting

To study the effectiveness of example voting in CREDIT-SQL, we investigated the performance of different methods to aggregate multiple sets of

Model	EM (%)	
	QM	IM
CREDIT-SQL	58.6	25.8
w/o SQL Reg.	58.0 (0.6↓)	24.9 (0.9↓)
w/o Div. Prompt	53.1 (5.5↓)	20.1 (5.7↓)

Table 3: Ablation studies for CREDIT-SQL on the CoSQL Dev set. We used gpt-3.5-turbo-0301 with the maximum context size of 8192 tokens. The results without SQL regularization and the results without diversity prompt are presented. Each experiment is repeated 5 times and the average performance is reported. Here we compare the exact set match without values (EM) for both the question match (QM) and the interaction match (IM).

diversity-sampled examples with distinct random seeds. In particular, we plot the EM for the question match on the CoSQL dev set of those different methods versus the number of aggregated sets of diversity-sampled examples in Fig 4. These aggregation methods include: (1) the average of each set’s performance, (2) the maximum of each set’s performance, (3) the performance of EM-based consistency voting from each set’s SQL results, and (4) the performance of the prompt made of voted examples among the all examples of the given sets. The last method is adopted for our CREDIT-SQL. As illustrated in Fig 4, the voted-example prompt outperforms either the average or the maximum of the individual results of distinct diversity-sampled prompts for the number of sets around 12 or more. Moreover, our proposed voted-example prompt performs similar to or better than the popular method of voting on SQL results, for the number of sets around 12 or more. Furthermore, the voted-example prompt is more efficient in the sense that it only uses a single inference of text-to-SQL per each SQL query regardless of the aggregation number  $M$ , while the consistency voting requires  $M$  inferences to aggregate  $M$  sets of examples. This may suggest a new possibility for efficiently mitigating the instability of example selection in the few-shot prompting with LLM in general.

### 5.3 Performance analysis by question difficulty

CoSQL provides the difficulty of each question based on the components of the golden SQL query for that question. Here we analyze the performance of the best-performing CREDIT-SQL on

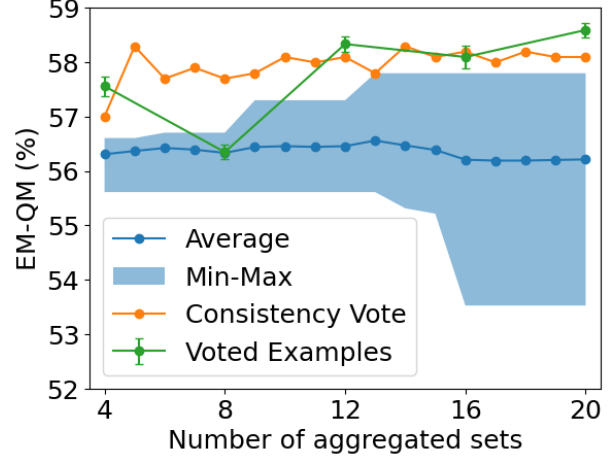


Figure 4: Comparison of different methods of aggregating multiple diversity-sampled sets of examples generated with distinct random seeds. We present an exact set match without values (EM) for question match (QM) on the CoSQL dev set for comparison. For this study, We used gpt-3.5-turbo-0301 with the maximum context size of 8192 tokens. **Blue solid line:** Average EM-QM of the results generated by prompts with distinct random seeds, up to the given number of prompts. **Blue shades:** The minimum to the maximum range for the EM-QM of the results generated by prompts with distinct random seeds, up to the given number of prompts. **Orange:** EM-QM of the consistency voting results on the SQL results of the given number of sets. Voting was done based on the exact set match without values. **Green:** EM-QM of the result of the prompt made of voted examples out of all examples from the given number of prompts. Average performance over 5 repetitions is reported along with the error bar size of the standard deviation.

the CoSQL dev set by the question difficulty (Figure 5). As anticipated, both the execution accuracy and the exact set match decreases as the question difficulty increases.

### 5.4 Performances with different OpenAI models

To determine the best OpenAI LLM model to be used for the CREDIT-SQL, we evaluated the CoSQL dev set with different OpenAI models having different maximum context sizes and completion methods (Table 4). While the performance of the CREDIT-SQL varies over different models and context sizes, we observe that the prompting through text completion outperforms the prompting through chat completion significantly.

### 5.5 Error analysis

To understand the cases in which CREDIT-SQL does not perform well, we performed an error anal-

OpenAI model	Maximum context size (tokens)	Completion method	EM (%)	
			QM	IM
gpt-3.5-turbo-0301	8192*	Text	58.6	25.8
gpt-3.5-turbo-0301	4096	Text	57.3	24.6
gpt-3.5-turbo-instruct	4096*	Text	54.9	23.9
gpt-3.5-turbo-16k	8192	Chat	51.7	19.8
gpt-3.5-turbo-16k	16384*	Chat	51.5	18.8
gpt-4-turbo	8192	Chat	54.0	21.5

Table 4: Performances of CREDIT-SQL on the CoSQL dev set, with different OpenAI models. We present the maximum context size we used for the prompt and response as well as the completion method used for each experiment. The asterisk indicates that the presented maximum context size is the maximum capacity allowed by the corresponding model. For the maximum context size of (4096/8192/16384) tokens, we used (12/18/50) examples in the prompt and set the response size to (250/600/600) tokens. For the chat completion method, we input the entire few-shot prompt as a system message.

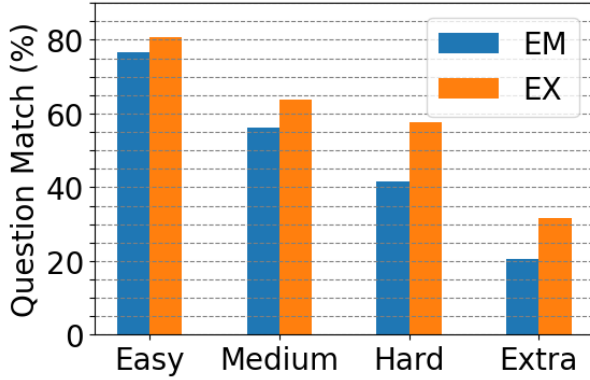


Figure 5: Performance of CREDIT-SQL on the CoSQL dev set by the question difficulty. We present the question match for both the exact set match (EM) and the execution accuracy (EX) along with 4 difficulty categories: easy, medium, hard, and extra.

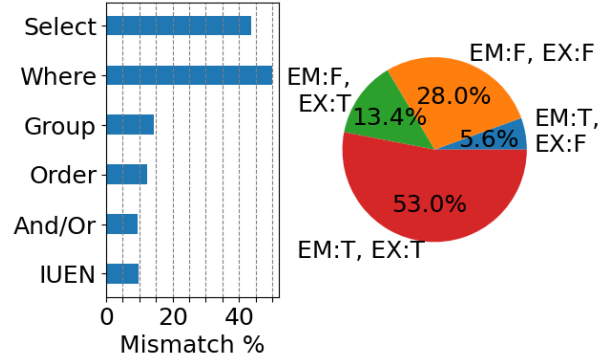


Figure 6: **Left:** Error analysis of CREDIT-SQL on the exact set match errors at the question match level. Here we present the mismatch percentage for each category of SQL query keywords. Here IUEN stands for IN, UNION, EXCEPT, or NOT IN. Since each incorrect SQL query may contain multiple mismatches, mismatch percentages for different categories are not disjoint. **Right:** Correlation of the exact set match (EM) and the execution accuracy (EX), evaluated on the CoSQL dev set at the question match level.

ysis on the results of the best-performing CREDIT-SQL. In particular, we present questions with the incorrect exact set match at the question match level in Figure 6. To categorize the errors, we used the keyword analysis provided by the official evaluation code<sup>4</sup> for the CoSQL dataset (Zhong et al., 2020). To understand how much the exact set match errors and the execution accuracy errors are correlated, we also present the pie chart that describes the correlation of those two metrics in Figure 6.

## 5.6 Limitations and future works

Since our work focused on the CoSQL dataset, the prompts we suggest in this paper might have difficulty in generalizing to the databases and SQL queries outside the CoSQL dataset. Indeed, test-

<sup>4</sup><https://github.com/taoyds/test-suite-sql-eval>

ing the generalizability of CREDIT-SQL to other context-dependent SQL datasets would be an interesting subject for future research. Also, one may test the performance of CRDEIT-SQL on other latest LLMs for future research.

## 6 Conclusion

In this paper, we propose a few-shot prompting method called CREDIT-SQL which is the first few-shot prompting approach to perform the dialogue state tracking task with systematically selected in-context examples out of the entire training data. CREDIT-SQL splits each dialogue state tracking task into multiple question-query pair text-to-SQL

tasks by question rephrasing and utilizes the diversity sampling and subsequent in-context example voting to prepare the few-shot prompts for the text-to-SQL tasks. Experiments demonstrate that CREDIT-SQL achieves a performance comparable to the state-of-the-art models. Also, the technique of example voting used in CREDIT-SQL suggests a new way to mitigate the instability of in-context example selection in the generic few-shot prompting setting.

## References

Rohan Anil, Andrew M. Dai, Orhan Firat, Melvin Johnson, Dmitry Lepikhin, Alexandre Passos, Siamak Shakeri, Emanuel Taropa, Paige Bailey, Zhifeng Chen, Eric Chu, Jonathan H. Clark, Laurent El Shafey, Yanping Huang, Kathy Meier-Hellstern, Gaurav Mishra, Erica Moreira, Mark Omernick, Kevin Robinson, Sebastian Ruder, Yi Tay, Kefan Xiao, Yuanzhong Xu, Yujing Zhang, Gustavo Hernandez Abrego, Junwhan Ahn, Jacob Austin, Paul Barham, Jan Botha, James Bradbury, Siddhartha Brahma, Kevin Brooks, Michele Catasta, Yong Cheng, Colin Cherry, Christopher A. Choquette-Choo, Aakanksha Chowdhery, Clément Crepy, Shachi Dave, Mostafa Dehghani, Sunipa Dev, Jacob Devlin, Mark Díaz, Nan Du, Ethan Dyer, Vlad Feinberg, Fangxiaoyu Feng, Vlad Fienber, Markus Freitag, Xavier Garcia, Sebastian Gehrmann, Lucas Gonzalez, Guy Gur-Ari, Steven Hand, Hadi Hashemi, Le Hou, Joshua Howland, Andrea Hu, Jeffrey Hui, Jeremy Hurwitz, Michael Isard, Abe Ittycheriah, Matthew Jagielski, Wenhao Jia, Kathleen Kenealy, Maxim Krikun, Sneha Kudugunta, Chang Lan, Katherine Lee, Benjamin Lee, Eric Li, Music Li, Wei Li, YaGuang Li, Jian Li, Hyeontaek Lim, Hanzhao Lin, Zhongtao Liu, Frederick Liu, Marcello Maggioni, Aroma Mahendru, Joshua Maynez, Vedant Misra, Maysam Moussalem, Zachary Nado, John Nham, Eric Ni, Andrew Nystrom, Alicia Parrish, Marie Pellat, Martin Polacek, Alex Polozov, Reiner Pope, Siyuan Qiao, Emily Reif, Bryan Richter, Parker Riley, Alex Castro Ros, Aurko Roy, Brennan Saeta, Rajkumar Samuel, Renee Shelby, Ambrose Slone, Daniel Smilkov, David R. So, Daniel Sohn, Simon Tokumine, Dasha Valter, Vijay Vasudevan, Kiran Vodrahalli, Xuezhi Wang, Pidong Wang, Zirui Wang, Tao Wang, John Wieting, Yuhuai Wu, Kelvin Xu, Yunhan Xu, Linting Xue, Pengcheng Yin, Jiahui Yu, Qiao Zhang, Steven Zheng, Ce Zheng, Weikang Zhou, Denny Zhou, Slav Petrov, and Yonghui Wu. 2023. [Palm 2 technical report](#).

Ben Bogin, Matt Gardner, and Jonathan Berant. 2019. [Global reasoning over database structures for text-to-SQL parsing](#). In *Proceedings of the 2019 Conference on Empirical Methods in Natural Language Processing and the 9th International Joint Conference on Natural Language Processing (EMNLP-IJCNLP)*,

pages 3659–3664, Hong Kong, China. Association for Computational Linguistics.

Ruichu Cai, Jinjie Yuan, Boyan Xu, and Zhifeng Hao. 2021. [Sadga: Structure-aware dual graph aggregation network for text-to-sql](#). *Advances in Neural Information Processing Systems*, 34:7664–7676.

Yitao Cai and Xiaojun Wan. 2020. [IGSQL: Database schema interaction graph based neural model for context-dependent text-to-SQL generation](#). In *Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing (EMNLP)*, pages 6903–6912, Online. Association for Computational Linguistics.

Zefeng Cai, Xiangyu Li, Binyuan Hui, Min Yang, Bowen Li, Binhua Li, Zheng Cao, Weijie Li, Fei Huang, Luo Si, and Yongbin Li. 2022. [STAR: SQL guided pre-training for context-dependent text-to-SQL parsing](#). In *Findings of the Association for Computational Linguistics: EMNLP 2022*, pages 1235–1247, Abu Dhabi, United Arab Emirates. Association for Computational Linguistics.

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 non-local relations](#). In *Proceedings of the 59th Annual Meeting of the Association for Computational Linguistics and the 11th International Joint Conference on Natural Language Processing (Volume 1: Long Papers)*, pages 2541–2555, Online. Association for Computational Linguistics.

Linzheng Chai, Dongling Xiao, Jian Yang, Liqun Yang, Qian-Wen Zhang, Yunbo Cao, Zhoujun Li, and Zhao Yan. 2023. [Qurg: Question rewriting guided context-dependent text-to-sql semantic parsing](#). *arXiv preprint arXiv:2305.06655*.

Zhi Chen, Lu Chen, Hanqi Li, Ruisheng Cao, Da Ma, Mengyue Wu, and Kai Yu. 2021. [Decoupled dialogue modeling and semantic parsing for multi-turn text-to-SQL](#). In *Findings of the Association for Computational Linguistics: ACL-IJCNLP 2021*, pages 3063–3074, Online. Association for Computational Linguistics.

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*.

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*.

Yushi Hu, Chia-Hsuan Lee, Tianbao Xie, Tao Yu, Noah A Smith, and Mari Ostendorf. 2022. [In-context learning for few-shot dialogue state tracking](#). *arXiv preprint arXiv:2203.08568*.



546	Binyuan Hui, Ruiying Geng, Qiyu Ren, Binhua Li,	Torsten Scholak, Nathan Schucher, and Dzmitry Bah-	603
547	Yongbin Li, Jian Sun, Fei Huang, Luo Si, Pengfei	danau. 2021. <a href="#">PICARD: Parsing incrementally for</a>	604
548	Zhu, and Xiaodan Zhu. 2021. Dynamic hybrid rela-	<a href="#">constrained auto-regressive decoding from language</a>	605
549	tion exploration network for cross-domain context-	<a href="#">models</a> . In <i>Proceedings of the 2021 Conference on</i>	606
550	dependent semantic parsing. In <i>Proceedings of</i>	<i>Empirical Methods in Natural Language Processing</i> ,	607
551	<i>the AAAI Conference on Artificial Intelligence</i> , vol-	pages 9895–9901, Online and Punta Cana, Domini-	608
552	ume 35, pages 13116–13124.	can Republic. Association for Computational Lin-	609
		guistics.	610
553	Binyuan Hui, Ruiying Geng, Lihan Wang, Bowen Qin,	Hugo Touvron, Thibaut Lavril, Gautier Izacard, Xavier	611
554	Yanyang Li, Bowen Li, Jian Sun, and Yongbin Li.	Martinet, Marie-Anne Lachaux, Timothée Lacroix,	612
555	2022. <a href="#">S<sup>2</sup>SQL: Injecting syntax to question-schema</a>	Baptiste Rozière, Naman Goyal, Eric Hambro, Faisal	613
556	<a href="#">interaction graph encoder for text-to-SQL parsers</a> .	Azhar, Aurelien Rodriguez, Armand Joulin, Edouard	614
557	In <i>Findings of the Association for Computational</i>	Grave, and Guillaume Lample. 2023a. <a href="#">Llama: Open</a>	615
558	<i>Linguistics: ACL 2022</i> , pages 1254–1262, Dublin,	<a href="#">and efficient foundation language models</a> .	616
559	Ireland. Association for Computational Linguistics.		
560	Xi Victoria Lin, Richard Socher, and Caiming Xiong.	Hugo Touvron, Louis Martin, Kevin Stone, Peter Al-	617
561	2020. <a href="#">Bridging textual and tabular data for cross-</a>	bert, Amjad Almahairi, Yasmine Babaei, Nikolay	618
562	<a href="#">domain text-to-SQL semantic parsing</a> . In <i>Findings</i>	Bashlykov, Soumya Batra, Prajjwal Bhargava, Shruti	619
563	<i>of the Association for Computational Linguistics:</i>	Bhosale, Dan Bikel, Lukas Blecher, Cristian Canton	620
564	<i>EMNLP 2020</i> , pages 4870–4888, Online. Association	Ferrer, Moya Chen, Guillem Cucurull, David Esiobu,	621
565	for Computational Linguistics.	Jude Fernandes, Jeremy Fu, Wenyan Fu, Brian Fuller,	622
		Cynthia Gao, Vedanuj Goswami, Naman Goyal, An-	623
566	Jiachang Liu, Dinghan Shen, Yizhe Zhang, Bill Dolan,	thony Hartshorn, Saghar Hosseini, Rui Hou, Hakan	624
567	Lawrence Carin, and Weizhu Chen. 2022. <a href="#">What</a>	Inan, Marcin Kardas, Viktor Kerkez, Madian Khabsa,	625
568	<a href="#">makes good in-context examples for GPT-3?</a> In	Isabel Kloumann, Artem Korenev, Punit Singh Koura,	626
569	<i>Proceedings of Deep Learning Inside Out (DeeLIO</i>	Marie-Anne Lachaux, Thibaut Lavril, Jenya Lee, Di-	627
570	<i>2022): The 3rd Workshop on Knowledge Extrac-</i>	ana Liskovich, Yinghai Lu, Yuning Mao, Xavier Marti-	628
571	<i>tion and Integration for Deep Learning Architectures</i> ,	net, Todor Mihaylov, Pushkar Mishra, Igor Moly-	629
572	pages 100–114, Dublin, Ireland and Online. Associa-	bog, Yixin Nie, Andrew Poulton, Jeremy Reizen-	630
573	tion for Computational Linguistics.	stein, Rashi Rungta, Kalyan Saladi, Alan Schelten,	631
574	Linyong Nan, Yilun Zhao, Weijin Zou, Narutatsu	Ruan Silva, Eric Michael Smith, Ranjan Subrama-	632
575	Ri, Jaesung Tae, Ellen Zhang, Arman Cohan, and	nian, Xiaoqing Ellen Tan, Binh Tang, Ross Tay-	633
576	Dragomir Radev. 2023. Enhancing few-shot text-	lor, Adina Williams, Jian Xiang Kuan, Puxin Xu,	634
577	to-sql capabilities of large language models: A	Zheng Yan, Iliyan Zarov, Yuchen Zhang, Angela Fan,	635
578	study on prompt design strategies. <i>arXiv preprint</i>	Melanie Kambadur, Sharan Narang, Aurelien Ro-	636
579	<i>arXiv:2305.12586</i> .	driguez, Robert Stojnic, Sergey Edunov, and Thomas	637
		Scialom. 2023b. <a href="#">Llama 2: Open foundation and</a>	638
580	OpenAI. 2023. <a href="#">Gpt-4 technical report</a> .	<a href="#">fine-tuned chat models</a> .	639
581	Zhufeng Pan, Kun Bai, Yan Wang, Lianqiang Zhou,	Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob	640
582	and Xiaojiang Liu. 2019. <a href="#">Improving open-domain</a>	Uszkoreit, Llion Jones, Aidan N Gomez, Łukasz	641
583	<a href="#">dialogue systems via multi-turn incomplete utterance</a>	Kaiser, and Illia Polosukhin. 2017. Attention is all	642
584	<a href="#">restoration</a> . In <i>Proceedings of the 2019 Conference</i>	you need. <i>Advances in neural information processing</i>	643
585	<i>on Empirical Methods in Natural Language Pro-</i>	<i>systems</i> , 30.	644
586	<i>cessing and the 9th International Joint Conference</i>		
587	<i>on Natural Language Processing (EMNLP-IJCNLP)</i> ,	Bailin Wang, Richard Shin, Xiaodong Liu, Oleksandr	645
588	pages 1824–1833, Hong Kong, China. Association	Polozov, and Matthew Richardson. 2020a. <a href="#">RAT-</a>	646
589	for Computational Linguistics.	<a href="#">SQL: Relation-aware schema encoding and linking</a>	647
		<a href="#">for text-to-SQL parsers</a> . In <i>Proceedings of the 58th</i>	648
590	Mohammadreza Pourreza and Davood Rafiei. 2023.	<i>Annual Meeting of the Association for Computational</i>	649
591	Din-sql: Decomposed in-context learning of	<i>Linguistics</i> , pages 7567–7578, Online. Association	650
592	text-to-sql with self-correction. <i>arXiv preprint</i>	for Computational Linguistics.	651
593	<i>arXiv:2304.11015</i> .		
594	Jiexing Qi, Jingyao Tang, Ziwei He, Xiangpeng Wan,	Kai Wang, Weizhou Shen, Yunyi Yang, Xiaojun Quan,	652
595	Yu Cheng, Chenghu Zhou, Xinbing Wang, Quanshi	and Rui Wang. 2020b. <a href="#">Relational graph attention</a>	653
596	Zhang, and Zhouhan Lin. 2022. <a href="#">RASAT: Integrating</a>	<a href="#">network for aspect-based sentiment analysis</a> . In <i>Pro-</i>	654
597	<a href="#">relational structures into pretrained Seq2Seq model</a>	<i>ceedings of the 58th Annual Meeting of the Asso-</i>	655
598	<a href="#">for text-to-SQL</a> . In <i>Proceedings of the 2022 Con-</i>	<i>ciation for Computational Linguistics</i> , pages 3229–	656
599	<i>ference on Empirical Methods in Natural Language</i>	3238, Online. Association for Computational Lin-	657
600	<i>Processing</i> , pages 3215–3229, Abu Dhabi, United	guistics.	658
601	Arab Emirates. Association for Computational Lin-		
602	guistics.	Run-Ze Wang, Zhen-Hua Ling, Jingbo Zhou, and Yu Hu.	659
		2021. Tracking interaction states for multi-turn text-	660
		to-sql semantic parsing. In <i>Proceedings of the AAAI</i>	661

662	<i>Conference on Artificial Intelligence</i> , volume 35,	( <i>EMNLP-IJCNLP</i> ), pages 5338–5349, Hong Kong,	720
663	pages 13979–13987.	China. Association for Computational Linguistics.	721
664	Dongling Xiao, LinZheng Chai, Qian-Wen Zhang, Zhao	Yanzhao Zheng, Haibin Wang, Baohua Dong, Xingjun	722
665	Yan, Zhoujun Li, and Yunbo Cao. 2022. <a href="#">CQR-</a>	Wang, and Changshan Li. 2022. <a href="#">HIE-SQL: History</a>	723
666	<a href="#">SQL: Conversational question reformulation en-</a>	<a href="#">information enhanced network for context-dependent</a>	724
667	<a href="#">hanced context-dependent text-to-SQL parsers</a> . In	<a href="#">text-to-SQL semantic parsing</a> . In <i>Findings of the As-</i>	725
668	<i>Findings of the Association for Computational Lin-</i>	<i>sociation for Computational Linguistics: ACL 2022</i> ,	726
669	<i>guistics: EMNLP 2022</i> , pages 2055–2068, Abu	pages 2997–3007, Dublin, Ireland. Association for	727
670	Dhabi, United Arab Emirates. Association for Com-	Computational Linguistics.	728
671	putational Linguistics.		
672	Junjie Ye, Xuanting Chen, Nuo Xu, Can Zu, Zekai	Ruiqi Zhong, Tao Yu, and Dan Klein. 2020. Seman-	729
673	Shao, Shichun Liu, Yuhan Cui, Zeyang Zhou, Chao	tic evaluation for text-to-sql with distilled test suite.	730
674	Gong, Yang Shen, Jie Zhou, Siming Chen, Tao Gui,	In <i>The 2020 Conference on Empirical Methods in</i>	731
675	Qi Zhang, and Xuanjing Huang. 2023. <a href="#">A comprehen-</a>	<i>Natural Language Processing</i> . Association for Com-	732
676	<a href="#">sive capability analysis of gpt-3 and gpt-3.5 series</a>	putational Linguistics.	733
677	<a href="#">models</a> .		
678	Tao Yu, Rui Zhang, Heyang Er, Suyi Li, Eric Xue,		
679	Bo Pang, Xi Victoria Lin, Yi Chern Tan, Tianze		
680	Shi, Zihan Li, Youxuan Jiang, Michihiro Yasunaga,		
681	Sungrok Shim, Tao Chen, Alexander Fabbri, Zifan		
682	Li, Luyao Chen, Yuwen Zhang, Shreya Dixit, Vin-		
683	cent Zhang, Caiming Xiong, Richard Socher, Walter		
684	Lasecki, and Dragomir Radev. 2019a. <a href="#">CoSQL: A</a>		
685	<a href="#">conversational text-to-SQL challenge towards cross-</a>		
686	<a href="#">domain natural language interfaces to databases</a> . In		
687	<i>Proceedings of the 2019 Conference on Empirical</i>		
688	<i>Methods in Natural Language Processing and the</i>		
689	<i>9th International Joint Conference on Natural Lan-</i>		
690	<i>guage Processing (EMNLP-IJCNLP)</i> , pages 1962–		
691	1979, Hong Kong, China. Association for Computa-		
692	tional Linguistics.		
693	Tao Yu, Rui Zhang, Kai Yang, Michihiro Yasunaga,		
694	Dongxu Wang, Zifan Li, James Ma, Irene Li, Qingn-		
695	ing Yao, Shanelle Roman, Zilin Zhang, and Dragomir		
696	Radev. 2018. <a href="#">Spider: A large-scale human-labeled</a>		
697	<a href="#">dataset for complex and cross-domain semantic pars-</a>		
698	<a href="#">ing and text-to-SQL task</a> . In <i>Proceedings of the 2018</i>		
699	<i>Conference on Empirical Methods in Natural Lan-</i>		
700	<i>guage Processing</i> , pages 3911–3921, Brussels, Bel-		
701	gium. Association for Computational Linguistics.		
702	Tao Yu, Rui Zhang, Michihiro Yasunaga, Yi Chern		
703	Tan, Xi Victoria Lin, Suyi Li, Heyang Er, Irene		
704	Li, Bo Pang, Tao Chen, Emily Ji, Shreya Dixit,		
705	David Proctor, Sungrok Shim, Jonathan Kraft, Vin-		
706	cent Zhang, Caiming Xiong, Richard Socher, and		
707	Dragomir Radev. 2019b. <a href="#">SPaC: Cross-domain se-</a>		
708	<a href="#">mantic parsing in context</a> . In <i>Proceedings of the</i>		
709	<i>57th Annual Meeting of the Association for Computa-</i>		
710	<i>tional Linguistics</i> , pages 4511–4523, Florence, Italy.		
711	Association for Computational Linguistics.		
712	Rui Zhang, Tao Yu, Heyang Er, Sungrok Shim, Eric		
713	Xue, Xi Victoria Lin, Tianze Shi, Caiming Xiong,		
714	Richard Socher, and Dragomir Radev. 2019. <a href="#">Editing-</a>		
715	<a href="#">based SQL query generation for cross-domain</a>		
716	<a href="#">context-dependent questions</a> . In <i>Proceedings of the</i>		
717	<i>2019 Conference on Empirical Methods in Natu-</i>		
718	<i>ral Language Processing and the 9th International</i>		
719	<i>Joint Conference on Natural Language Processing</i>		

## A Question Rephrasing Prompt Base

Following is the question rephrasing prompt base we used for CREDIT-SQL:

### Make sure the rephrased questions do not need any implicit information of the previous dialogue.

Database: department\_store

Q1: List the name of all suppliers in alphabetical order of their addresses.

A1: Okay. I rephrased your question as "Give me the names of all suppliers, ordered alphabetically by their addresses."

Q2: And their phone numbers?

A2: Okay. I rephrased your question as "Give me the phone numbers of all suppliers, ordered alphabetically by their addresses."

Q3: What states is Tesco located in?

A3: do you mean the address id or the address details?

Q4: The address details.

A4: Okay. I rephrased your question as "Give me the address details of Tesco."

Q5: How about Lidl?

A5: Okay. I rephrased your question as "Give me the address details of Lidl."

Database: theme\_gallery

Q1: Show the average ticket prices for exhibitions for all years before 2009?

A1: Okay. I rephrased your question as "Give me the average ticket price of all exhibitions whose year before 2009."

Q2: Show minimum ticket prices for exhibitions for all years before 2009?

A2: Okay. I rephrased your question as "Give me the minimum ticket price of all exhibitions whose year before 2009."

Q3: Show maximum ticket prices for exhibitions for all years before 2009?

A3: Okay. I rephrased your question as "Give me the maximum ticket price of all exhibitions whose year before 2009."

Q4: Show average ticket price for exhibitions for year 2005?

A4: Okay. I rephrased your question as "Give me the average ticket price of all exhibitions whose year is 2005."

Q5: Show average ticket price for exhibitions for all years before 2006?

A5: Okay. I rephrased your question as "Give me the average ticket price of all exhibitions whose year before 2006."

Database: apartment\_rentals

Q1: How many bathroom count those building 809 has?

A1: The total number of bathrooms building 809 has is NULL.

Q2: How about building 808?

A2: Okay. I rephrased your question as "Give me the total number of bathroom count of all building 809."

Q3: How about the total number of rooms it has?

A3: Okay. I rephrased your question as "Give me the total number of rooms of all building 809."

Database: wine\_1

Q1: Which wine has the highest price?

A1: Okay. I rephrased your question as "Give me the wine with the highest price."

Q2: Which wine has the highest score?

A2: Okay. I rephrased your question as "Give me the wine with the highest score."

Q3: Which wine has the most cases available?

A3: Okay. I rephrased your question as "Give me the wine with the most cases available."

Q4: How about winery with the least cases?

A4: Okay. I rephrased your question as "Give me the winery with the least cases."

Database: dorm\_1

Q1: How many students live in each dorm?

A1: Okay. I rephrased your question as "Give me the total number of students of each dorm."

Q2: What is the average age of the students living in each dorm?

A2: Okay. I rephrased your question as "Give me the average number of students of each dorm."

Database: roller\_coaster

Q1: what is the highest length of coaster

A1: Okay. I rephrased your question as "Give me the highest length of all costers."

Q2: which roller coaster are top speed

A2: Okay. I rephrased your question as "Give me the roller coaster with highest speed."

Q3: what is the height of Boardwalk Bullet

A3: Okay. I rephrased your question as "Give me the height of Boardwalk Bullet."





936			986
937	Database: activity_1	Database: store_product	987
938	Q1: Hello, what are the names of the faculty mem-	Q1: Please list the headquarters of the store that is	988
939	bers who don't participate in any activity?	called Blackville?	989
940	A1: Okay. I rephrased your question as "Give	A1: Okay. I rephrased your question as "Give	990
941	me the names of all faculty members who do not	me the headquarters of all stores with name	991
942	participate in any activity."	Blackville."	992
943	Q2: Which faculty members participate in activi-	Q2: What is the population of that city?	993
944	ties?	A2: Okay. I rephrased your question as "Give me	994
945	A2: Okay. I rephrased your question as "Give me	the population of all cities where headquarters of	995
946	the names of all faculty members.who participate	all stores with name Blackville are located."	996
947	in any activity."	Q3: What is the name of the store in the city with	997
948	Q3: What activities do Michael Goodrich and Ger-	the greatest population?	998
949	ald Masson participate in?	A3: Okay. I rephrased your question as "Give me	999
950	A3: Okay. I rephrased your question as "Give me	the name of the store in the city with the greatest	1000
951	the activities that both of Michael Goodrich and	population."	1001
952	Gerald Masson participate in."		1002
953	Q4: In total, how many faculty members are there?	Database: school_bus	1003
954	A4: Okay. I rephrased your question as "Give me	Q1: How many drivers are there?	1004
955	the total numbe of faculty members."	A1: Okay. I rephrased your question as "Give me	1005
956		the total number of all drivers."	1006
957	Database: product_catalog	Q2: What is the average age of your drivers?	1007
958	Q1: what are the names of all catalogs	A2: Okay. I rephrased your question as "Give me	1008
959	A1: Okay. I rephrased your question as "Give me	the average age of all drivers."	1009
960	the names of all catalogs."		1010
961	Q2: list the publication dates for all catalogs	Database: customers_and_addresses	1011
962	A2: Okay. I rephrased your question as "Give me	Q1: Can you give me a list of all products?	1012
963	the publication dates of all catalogs."	A1: Okay. I rephrased your question as "Give me	1013
964	Q3: list the catalogs that have level number greater	the all products."	1014
965	than 5	Q2: Could you order them by order quantity?	1015
966	A3: Do you want catalog name or id?	A2: From small to large?	1016
967	Q4: id	Q3: large to small	1017
968	A4: Okay. I rephrased your question as "Give me	A3: Okay. I rephrased your question as "Give	1018
969	the catalog id of catalogs with level number greater	me the all products, ordered by order quantity in	1019
970	than 5."	descending order."	1020
971			1021
972	Database: hr_1	Database: game_1	1022
973	Q1: Show me all the employees in the same depart-	Q1: Please show me the ids of the students who	1023
974	ment as Clara.	play video games and play sports?	1024
975	A1: Here is all the information of all the employees	A1: Okay. I rephrased your question as "Give	1025
976	in the same department as the employee with a first	me the ids of all students who play both of video	1026
977	name Clara.	games and sports."	1027
978	Q2: Can you show me their hire dates?		1028
979	A2: Okay. I rephrased your question as "Give me	Database: protein_institute	1029
980	the hire dates of all employees with department	Q1: show the institute type with lowest number of	1030
981	same as the department of Clara."	institutions	1031
982	Q3: What is their average salary?	A1: Okay. I rephrased your question as "Give	1032
983	A3: Okay. I rephrased your question as "Give	me the institute type with the lowest number of	1033
984	me the average salary of all employees with	institutions."	1034
985	department same as the department of Clara."	Q2: show the institute type with highest number of	1035
		institutions	1036

A2: Okay. I rephrased your question as "Give me the institute type with the highest number of institutions."

Q3: which university is founded first

A3: Okay. I rephrased your question as "Give me the university with the earliest foundation date."

Database: formula\_1

Q1: list the forename of all races

A1: Okay. I rephrased your question as "Give me the forenames of all races."

Q2: what is the forename of the race that occurred most recently?

A2: Okay. I rephrased your question as "Give me the forename of the race with the latest occurrence date."

Q3: what is its code?

A3: Do you mean circuit id?

Q4: yes

A4: Okay. I rephrased your question as "Give me the circuit id of the race with the latest occurrence date."

## B In-context Examples Selected via Voting for CREDIT-SQL Text-to-SQL Prompt

We present the 18 voted examples we used for the text-to-SQL prompt of our CREDIT-SQL approach in Table 5. We performed  $k$ -means clustering with  $k = 18$  and voted examples collected from 20 trials of clustering with distinct random seeds.

## C Sample Text-to-SQL Prompt for CREDIT-SQL

Following is a sample text-to-SQL prompt for CREDIT-SQL in the case of  $N = 2$ :

---

# Generate the SQL query for each question about the given database schema.

Database: party\_host

Table host, columns = [\*, host\_id, name, nationality, age]

Table party, columns = [\*, party\_id, party\_theme, location, first\_year, last\_year, number\_of\_hosts]

Table party\_host, columns = [\*, party\_id, host\_id, is\_main\_in\_charge] Foreign\_keys = [party\_host.party\_id = party.party\_id, party\_host.host\_id = host.host\_id]

Q: Give me the names of all nations with hosts younger than 35.

SQL: SELECT DISTINCT host.nationality FROM host WHERE host.age < 35

Database: company\_1

Table department, columns = [\*, name, number, mgr\_ssn, mgr\_start\_date]

Table dependent, columns = [\*, essn, dependent\_name, sex, bdate, relationship]

Table dept\_locations, columns = [\*, number, location]

Table employee, columns = [\*, fname, minit, lname, ssn, bdate, address, sex, salary, super\_ssn, dno]

Table project, columns = [\*, pname, number, location, dnum]

Table works\_on, columns = [\*, essn, pno, hours]

Foreign\_keys = []

Q: Give me the total number of employees with salary lower than 50000.

SQL: SELECT count ( \* ) FROM employee WHERE employee.salary < 50000

## D Sample Prompt for the Baseline Approach

Following is a sample prompt for the baseline approach with 2 dialogues:

---

# Generate the SQL query for each question about the given database schema.

Database: inn\_1

Table reservations, columns = [\*, code, room, check-in, checkout, rate, lastname, firstname, adults, kids]

Table rooms, columns = [\*, roomid, roomname, beds, bedtype, maxoccupancy, baseprice, decor]

Foreign\_keys = [reservations.room = rooms.roomid]

Q: which rooms have a king bed?

A: SQL query you requested is [SELECT \* FROM Rooms WHERE bedtype = 'King']

Q: what are the decor of those rooms?

A: SQL query you requested is [SELECT distinct decor FROM Rooms WHERE bedtype = 'King']

Q: show me the names and the decor of the rooms that have a king bed

A: SQL query you requested is [SELECT room-Name, decor FROM Rooms WHERE bedtype = 'King']

Q: show me the price of the rooms also

A: SQL query you requested is [SELECT room-Name, decor, basePrice FROM Rooms WHERE

bedtype = 'King']  
 Q: now sort the list by price  
 A: SQL query you requested is [SELECT room-  
 Name, decor, basePrice FROM Rooms WHERE  
 bedtype = 'King' order by basePrice]  
 Database: restaurant\_1  
 Table restaurant, columns = [\*, resid, resname, ad-  
 dress, rating]  
 Table restaurant\_type, columns = [\*, restypeid,  
 restypename, restypedescription]  
 Table student, columns = [\*, stupid, lname, fname,  
 age, sex, major, advisor, city\_code]  
 Table type\_of\_restaurant, columns = [\*, resid,  
 restypeid]  
 Table visits\_restaurant, columns = [\*, stupid, resid,  
 time, spent]  
 Foreign\_keys = [type\_of\_restaurant.resid =  
 restaurant.resid, type\_of\_restaurant.restypeid =  
 restaurant\_type.restypeid, visits\_restaurant.stuid  
 = student.stuid, visits\_restaurant.resid = restau-  
 rant.resid]  
 Q: Which city code of Smith?  
 Q: Is Smith first name or last name?  
 Q: first name  
 A: SQL query you requested is [SELECT city\_code  
 from Student where fname = "Smith"]  
 Q: What age of last name Kim?  
 A: SQL query you requested is [SELECT age from  
 Student where lname = "Kim"]  
 Q: What are the student id of city code WAS?  
 A: SQL query you requested is [SELECT stuid  
 from Student where city\_code = "WAS"]  
 Q: Which major has least number of students? List  
 the major and the number of students.  
 A: SQL query you requested is [SELECT Major  
 , count ( \* ) FROM Student GROUP BY Major  
 ORDER BY count ( Major ) ASC LIMIT 1]

1	Give me the invoice dates of all customers. [chinook_1] SELECT customer.firstname , customer.lastname , invoice.invoicedate FROM customer JOIN invoice ON customer.customerid = invoice.customerid
2	Give me the event id of all events that have a participant with the detail Kenyatta Kuhn. [local_govt_in_alabama] SELECT events.event_id FROM events JOIN participants_in_events ON participants_in_events.event_id = events.event_id JOIN participants ON participants_in_events.participant_id = participants.participant_id WHERE participants.participant_details = 'Kenyatta Kuhn'
3	Give me the patient id of the appointment with the most recent start date. [hospital_1] SELECT appointment.patient FROM appointment ORDER BY appointment.start DESC LIMIT 1
4	Give me the minimum and maximum number of bathrooms and bedrooms of all the apartments. [apartment_rentals] SELECT min ( apartments.bathroom_count ) , max ( apartments.bathroom_count ) , min ( apartments.bedroom_count ) , max ( apartments.bedroom_count ) FROM apartments
5	Give me the total number of all professors. [college_1] SELECT count ( * ) FROM professor
6	Give me the names of 5 products that are not in any event. [solvency_ii] SELECT products.product_name FROM products WHERE products.product_id NOT IN ( SELECT products_in_events.product_id FROM products_in_events )
7	Give me the total number of students who play football. [game_1] SELECT count ( * ) FROM sportsinfo JOIN student ON sportsinfo.stuid = student.stuid WHERE sportsinfo.sportname = 'Football'
8	Give me the names of all instructors who are advising more than one student. [college_2] SELECT instructor.name FROM instructor JOIN advisor ON instructor.id = advisor.i_id GROUP BY advisor.i_id HAVING count ( * ) > 1
9	Give me the account types of all customers whose credit score is above 100. [loan_1] SELECT customer.acc_type FROM customer
10	Give me the total number of students who have behavior incident reports with recommendations. [behavior_monitoring] SELECT count ( * ) FROM ( SELECT * FROM behavior_incident JOIN students ON behavior_incident.student_id = students.student_id GROUP BY behavior_incident.student_id )
11	Give me the total number of residents of each property, and the property id. [local_govt_and_lot] SELECT properties.property_id , count ( * ) FROM properties JOIN residents ON properties.property_id = residents.property_id GROUP BY properties.property_id
12	Give me the claim id of the claim that incurred the most number of settlements. [insurance_policies] SELECT claims.claim_id FROM claims JOIN settlements ON claims.claim_id = settlements.claim_id GROUP BY claims.claim_id ORDER BY count ( * ) DESC LIMIT 1
13	Give me the date of ceremony of all music festivals with category 'best song' and 'awarded'. [music_4] SELECT music_festival.date_of_ceremony FROM music_festival WHERE music_festival.category = 'Best Song' AND music_festival.result = 'Awarded'
14	Give me the ids of all employees with role Role_Code. [cre_Doc_Tracking_DB] SELECT employees.employee_id , employees.role_code FROM employees
15	Give me the first names of all students who have a dorm id of 160. [dorm_1] SELECT student.fname FROM student JOIN lives_in ON student.stuid = lives_in.stuid WHERE lives_in.dormid = 160
16	Give me the college name of the employee with name Reggie Lewis. [company_employee] SELECT people.graduation_college FROM people WHERE people.name = 'Reggie Lewis'
17	Give me the names of all nations with hosts younger than 35. [party_host] SELECT DISTINCT host.nationality FROM host WHERE host.age < 35
18	Give me the total number of customers who pay by Credit card. [customers_campaigns_ecommerce] SELECT count ( * ) FROM customers WHERE customers.payment_method = 'Credit Card'

Table 5: Voted examples used for CREDIT-SQL. The first row of each example: rephrased question [database ID], the second row of each example: regulated SQL query. 16