SQLong: Enhanced NL2SQL for Longer Contexts with LLMs

Dai Quoc Nguyen, Cong Duy Vu Hoang, Duy Vu, Gioacchino Tangari, Thanh Tien Vu, Don Dharmasiri, Yuan-Fang Li, Long Duong

Oracle Corporation Australia

{dai.nguyen,vu.hoang,duy.vu,gioacchino.tangari,thanh.v.vu,don.dharmasiri,yuanfang.li,long.duong}@oracle.com

Abstract

Open-weight large language models (LLMs) have significantly advanced performance in the Natural Language to SQL (NL2SQL) task. However, their effectiveness diminishes when dealing with large database schemas, as the context length increases. To address this limitation, we present SQLong, a novel and efficient data augmentation framework designed to enhance LLM performance in long-context scenarios for the NL2SQL task. SQLong generates augmented datasets by extending existing database schemas with additional synthetic CREATE TABLE commands and corresponding data rows, sampled from diverse schemas in the training data. This approach effectively simulates long-context scenarios during finetuning and evaluation. Through experiments on the Spider and BIRD datasets, we demonstrate that LLMs finetuned with SQLongaugmented data significantly outperform those trained on standard datasets. These imply SQLong's practical implementation and its impact on improving NL2SQL capabilities in real-world settings with complex database schemas.

CCS Concepts

Computing methodologies → Natural language processing;
 Neural networks;
 Information systems → Social networks.

Keywords

LLMs, NL2SQL, Long Contexts

ACM Reference Format:

1 Introduction

The NL2SQL task focuses on translating natural language questions into SQL queries, enabling non-experts to interact with databases seamlessly [3]. Recent advances leverage LLMs, finetuned on structured input prompts (e.g., task instructions, database schema, and

Permission to make digital or hard copies of all or part of this work for personal or classroom use is granted without fee provided that copies are not made or distributed for profit or commercial advantage and that copies bear this notice and the full citation on the first page. Copyrights for components of this work owned by others than the author(s) must be honored. Abstracting with credit is permitted. To copy otherwise, or republish, to post on servers or to redistribute to lists, requires prior specific permission and/or a fee. Request permissions from permissions@acm.org.

Conference acronym '25, 2025, Woodstock, NY

© 2025 Copyright held by the owner/author(s). Publication rights licensed to ACM. ACM ISBN 978-x-xxxx-xxxx-x/YY/MM https://doi.org/XXXXXXXXXXXXXXXX

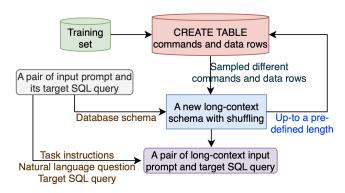


Figure 1: Our proposed SQLong Pipeline.

natural language question), to achieve state-of-the-art performance on benchmarks like Spider [18] and BIRD [9] [10, 17].

Despite significant progress, a critical challenge persists: LLMs finetuned on existing benchmarks still struggle with large database schemas due to limited context handling. Current datasets primarily feature small schemas, failing to represent real-world complexities. Additionally, the absence of publicly available large-schema datasets further hinders progress. Addressing this, we propose SQLong, a data augmentation framework designed to enhance LLM performance in long-context NL2SQL tasks by extending schemas to meet predefined context thresholds.

SQLong constructs augmented data by sampling CREATE TABLE commands and data rows from diverse schemas. These datasets enable LLMs to effectively manage large schemas and maintain robustness in long-context scenarios. Our experiments with *CodeQwen1.5-7B-Chat* [1] and *Llama-3.1-8B-Instruct* [6] show SQLong consistently outperforms baseline finetuning, achieving an average accuracy improvement of over 2.2% on benchmarks like Spider-dev, Spider-test, and BIRD-dev.

Moreover, SQLong enables the creation of 45 long-context test sets, with context lengths up to 128k tokens. Models finetuned with SQLong exhibit significant performance gains, achieving an 11% improvement over base models and a 6% improvement over larger-scale models within the same family. These results highlight SQLong's effectiveness in real-world, large-schema scenarios.

In this paper, we focus on demonstrating that SQLong-augmented models outperform their unaugmented counterparts across varying context lengths. While direct comparisons to retrieval-augmented generation (RAG) schema linking are beyond this paper's scope, our findings suggest combining SQLong with RAG could unlock further gains.

Our main contributions include:

¹The data will be made publicly available upon paper acceptance.

```
Given an input Question, create a syntactically correct
SQLite SQL query to run.
Pay attention to using only the column names that you can
see in the schema description.
Be careful to not query for columns that do not exist. Also,
pay attention to which column is in which table.
Please double check the SQLite SQL query you generate.
DO NOT use alias in the SELECT clauses.
Only use the tables listed below.
CREATE TABLE grades (
  "student_id" INTEGER,
  "student_name" TEXT,
  "subject" TEXT,
  "grade" TEXT,
  PRIMARY KEY ("student_id")
/* 3 rows from grades table:
student_id
              student_name subject
    Alice
                       Α
              math
    Bob
                       В
              math
3
    David
              science
                       В
Question: Show me all the students getting an A in math
SELECT student_name FROM grades WHERE subject =
```

Figure 2: Prompt template for the NL2SQL task.

- Introducing long-context NL2SQL: A challenging new task for evaluating LLM performance on large database schemas.
- SQLong pipeline: A novel, scalable data augmentation approach for generating long-context training and test datasets.
- Empirical insights: Comprehensive experiments validating SQLong's effectiveness in enhancing LLM robustness and accuracy in long-context scenarios.
- Resource sharing: Plans to release SQLong datasets and code to support further research.

The Proposed SQLong Pipeline

'math' AND grade = 'A'

The NL2SQL task aims to translate a natural-language question about a database schema into a corresponding SQL query. Following the standardized prompt template [12], we represent the input prompt to LLMs in the format of (task instructions, database schema, natural language question).² As illustrated in Figure 2, the database schema is represented by CREATE TABLE commands and three sample data rows for each corresponding table.

Using supervised finetuning (SFT) [15], LLMs can be trained on pairs of input prompts and target SQL queries to optimize their performance on the NL2SQL task. Specifically, given a training set T comprising pairs of input prompts x and corresponding target SQL queries s, the supervised finetuning process can be formulated as minimizing the log-likelihood loss [15], as shown below:

$$\mathbb{E}_{(\mathbf{x},\mathbf{s})\sim T}\left[\sum_{i=1}^{|\mathbf{s}|}\log p_{\theta}\left(s_{i}|\mathbf{s}_{< i},\mathbf{x}\right)\right]$$

 $\mathbb{E}_{(\mathbf{x},\mathbf{s})\sim T}\left[\sum_{i=1}^{|\mathbf{s}|}\log p_{\theta}\left(s_{i}|\mathbf{s}_{< i},\mathbf{x}\right)\right]$ wherein $|\mathbf{s}|$ is the length of \mathbf{s} , s_{i} is the i-th token, $\mathbf{s}_{< i}$ is the prefix of s up to the *i*-th position, and θ denotes the given LLM's parameters.

In this work, we introduce SQLong, a novel approach for constructing long-context finetuning and benchmark datasets, as illustrated in Figure 1. SQLong augments database schemas to enable large language models (LLMs) to effectively handle long-context scenarios in natural language to SQL (NL2SQL) tasks.

The SQLong pipeline comprises three main steps:

- (1) Schema Collection. We collect all CREATE TABLE commands and three sample data rows for each table from the training database schemas, compiling them into a comprehensive schema set.
- (2) Schema Augmentation. For each training pair, consisting of an input prompt (task instructions, database schema, natural language question) and its target SQL query, SQLong randomly samples items from the schema set. These sampled items contain table names distinct from those in the given database schema. The sampled items are combined with the original schema, and the resulting schema is randomly shuffled to produce a new, long-context database schema. This shuffling introduces variability in the positions of the original tables and columns.
- (3) Long-Context Prompt Generation. SOLong generates an augmented input prompt in the format of task instructions, the long-context database schema, and the natural language question, while keeping the target SQL query unchanged. It ensures that the combined length of the long-context input prompt and the target SQL query does not exceed a predefined context length (e.g., 32k tokens), maintaining compatibility with the model's tokenizer constraints.

By systematically extending and diversifying the context, SQLong enhances the robustness and effectiveness of LLMs in handling longcontext NL2SQL tasks. We summarise the steps involved in SQLong in Algorithm 1.

Evaluation

We assess the effectiveness of our proposed SQLong model in enhancing NL2SQL performance across both short-context and longcontext scenarios.

Experimental Setup

Datasets. For the short-context evaluation, we utilize widely adopted benchmark datasets, including Spider [18], Spider-realistic [4], Spider-syn [7], and BIRD [9]. ³ It is noted that Spider-Syn is manually created based on Spider training and development sets using synonym substitution in the original questions, while Spiderrealistic is created based on Spider development set by manually removing the explicit mention of column names in the original questions. The BIRD-test set is not publicly available.

For the long-context evaluation, we extend each of the Spider-dev, Spider-test, Spider-realistic, Spider-syn, and BIRD-dev datasets by applying SQLong with a pre-defined context length. Specifically, we

²In datasets with additional complexity, such as BIRD, the question may be supplemented with extra information, such as evidence. For simplicity, this additional information is omitted in Figure 2.

³We use the latest BIRD-dev dataset, updated on June 27, 2024. The BIRD-test set is not publicly available.

Algorithm 1: The algorithm steps involved in the proposed SQLong.

```
1 Input: A training set T of pairs of input prompts and target SQL queries:
     \mathbf{T} = \{((instructions_i, database\_schema_i, question_i), target\_sql_i)\}_{i=1}^N, \text{ wherein each } database\_schema_i \text{ is a set of CREATE TABLE } \}
     commands and three data rows for each corresponding table; a set
     \mathcal{T} = \{((instructions_j, database\_schema_j, question_j), target\_sql_j)\}_{j=1}^{M}; \text{ the base model's tokenizer } tk, \text{ a starting number } s\_n \text{ (default } tk)\}_{j=1}^{M}; \text{ the base model's tokenizer } tk, \text{ a starting number } s\_n \text{ (default } tk)\}_{j=1}^{M}; \text{ the base model's tokenizer } tk, \text{ a starting number } s\_n \text{ (default } tk)\}_{j=1}^{M}; \text{ the base model's tokenizer } tk, \text{ a starting number } s\_n \text{ (default } tk)\}_{j=1}^{M}; \text{ the base model's tokenizer } tk, \text{ a starting number } s\_n \text{ (default } tk)\}_{j=1}^{M}; \text{ the base model's tokenizer } tk)
     4096), an ending number e_n (default 32768), an increasing number i_n (default 512), and a pre-defined number p_n (default 8192).
2 Output: The augmented long-context set \mathcal{T}'.
3 schema\_set \leftarrow collect\_unique\_commands\_and\_data\_rows(\{database\_schema_i\}_{i=1}^{N})
4 table_names ← get_table_names(schema_set)
5 item_lengths \leftarrow \{\}
6 for item \in schema \ set \ do
     item\_lengths \leftarrow item\_lengths \cup \{get\_length(item, tk)\}
s \ \mathcal{T}' \leftarrow \{\}
9 diverse_lengths \leftarrow range(s_n, e_n + 1, i_n)
10 for ((instructions, database_schema, question), target_sql) \in \mathcal{T} do
         original\_length \leftarrow get\_length(instructions + database\_schema + question + target\_sql, tk)
11
         certain\_length \leftarrow randomly\_select\_value(diverse\_lengths)
12
                                                                                                           // This aims to construct long-context
          fine-tuning data with T = \mathcal{T}. Otherwise, certain_length is set to p_n to construct
          long-context benchmark data.
        local\_table\_names \leftarrow get\_table\_names(database\_schema)
13
         augmented\_schema \leftarrow \{\}
14
        for idx \in shuffle\_list(range(0, get\_size(schema\_set))) do
15
             if schema_set[idx] ∉ database_schema and table_names[idx] ∉
16
                local\_table\_names and original\_length + item\_lengths[idx] < certain\_length then
                   original\_length \leftarrow original\_length + item\_lengths[idx]
17
                   augmented\_schema \leftarrow augmented\_schema \cup \{schema\_set[idx]\}
18
         augmented\_long\_context\_schema \leftarrow shuffle\_list(augmented\_schema \cup database\_schema)
19
         \mathcal{T}' \leftarrow \mathcal{T}' \cup \{((instructions, augmented\_long\_context\_schema, question), target\_sql)\}
20
```

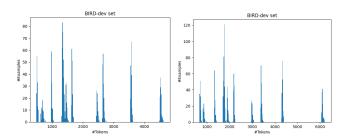


Figure 3: Statistics of input prompt lengths with respect to Llama-3.1-8B-Instruct's tokenizer (left) and CodeQwen1.5-7B-Chat's tokenizer (right) on the original BIRD-dev set. Similarly, the maximum input prompt lengths for the original Spider-related sets are approximately 2,000 tokens for Llama-3.1-8B-Instruct's tokenizer and 2,500 tokens for CodeQwen1.5-7B-Chat's tokenizer.

generate augmented long-context test sets for nine context lengths: 8k, 16k, 24k, 32k, 40k, 48k, 56k, 64k, and 128k. This process results in a total of 45 long-context test sets, constructed in accordance with the tokenizer of the base model.

Dataset	#DB	#tables	#training	#dev	#test
Spider	200	5 ± 3	6,712	1,034	2,019
Spider-syn	200	5 ± 3	6,712	1,034	_
Spider-realistic	200	5 ± 3	6,712	508	_
BIRD	98	7 ± 3	9,428	1,534	_

Table 1: Statistics of the experimental datasets. #DB denotes the number of databases. #tables denotes the mean and standard deviation of numbers of tables in the databases.

Importantly, the long-context test sets are constructed with distinct database schema alignments. To build Spider-based long-context test sets, we use the database schemas from the BIRD training set, whereas for the BIRD-dev long-context test sets, we use the database schemas from the Spider training set. This ensures a robust evaluation across diverse schema configurations and context lengths. The data statistics of the experimental datasets are presented in Figure 3 and Tables 1 and 3.

Baseline Models and Evaluation Metrics. We evaluate SQLong using two powerful base models: CodeQwen1.5-7B-Chat [1], which supports a context length of up to 64k, and Llama-3.1-8B-Instruct [6], which supports a context length of up to 128k. Following Yu

Model	Spider-dev	Spider-realistic	Spider-syn	Spider-test	BIRD-dev	Average
Qwen2-72B-Instruct	82.7	80.7	73.0	82.9	53.7	74.6
CodeQwen1.5-7B-Chat	76.4	70.1	62.7	75.1	44.3	65.7
Finetuned without SQLong	81.9	76.2	68.7	79.6	51.4	71.6
Finetuned with SQLong	83.4	79.7	71.2	81.3	53.3	73.8
Llama-3.1-70B-Instruct	80.7	78.0	73.0	83.7	61.5	75.4
Llama-3.1-8B-Instruct	71.1	63.8	61.0	65.7	40.9	60.5
Finetuned without SQLong	79.2	76.4	69.6	80.4	51.9	71.5
Finetuned with SQLong	83.2	78.0	73.1	81.8	53.3	73.9

Table 2: Execution-match accuracy results (in %) across different datasets and model configurations. Finetuning with SQLong consistently improves performance, with the best results highlighted in bold.

Length	CodeQwen1.5	5-7B-Chat	Llama-3.1-8B-Instruct		
Length	Spider-related	BIRD-dev	Spider-related	BIRD-dev	
8k	37 ± 4	35 ± 8	48 ± 5	48 ± 8	
16k	72 ± 6	76 ± 8	94 ± 7	102 ± 9	
24k	107 ± 7	118 ± 8	141 ± 8	157 ± 9	
32k	142 ± 8	159 ± 9	186 ± 8	211 ± 9	
40k	177 ± 8	200 ± 9	233 ± 9	269 ± 9	
48k	212 ± 9	242 ± 9	279 ± 9	320 ± 10	
56k	247 ± 9	283 ± 9	326 ± 9	374 ± 9	
64k	283 ± 9	324 ± 9	372 ± 8	429 ± 9	
128k	551 ± 4	639 ± 7	725 ± 9	843 ± 8	

Table 3: Mean and standard deviation statistics of the numbers of tables in input prompts for our augmented long-context test sets with respect to each model's tokenizer.

et al. [18], we report execution-match accuracy on both the original short-context test sets and the augmented long-context test sets.

Training Protocol. For each original training set, we use SQLong to create an augmented *long-context finetuning* dataset with context lengths of up to 32k.⁴ The augmented dataset is combined with the original training set to form the final dataset used for finetuning the base models.⁵

We experiment with two base models: CodeQwen1.5-7B-Chat [1], which supports a 64k context length, and Llama-3.1-8B-Instruct [6], which supports a 128k context length. Finetuning is performed with a batch size of 1, gradient accumulation steps of 8, a learning rate chosen from 1×10^{-6} , 5×10^{-6} , 1×10^{-5} , and up to 5 epochs on $8 \times H100~80 GB~GPUs$.

We use Huggingface's TRL [14] for supervised finetuning, employing 8-bit AdamW [5], Flash Attention v2 [2], and DeepSpeed ZeRO-3 Offload [13]. For a fair comparison, we also finetune the base models on the original training set (i.e., without SQLong) under the same settings.

Inference Protocol. We utilize vLLM [8] for the inference process. For long-context test sets, we employ dynamic NTK RoPE

scaling [11] to extend support up to a 128k context length for CodeQwen1.5-7B-Chat and its finetuned variants.

3.2 Main Results

Performance on Original Datasets. Table 2 summarizes the results on the original development and test sets, comparing base models with larger LLMs such as Llama-3.1-70B-Instruct [6] and Qwen2-72B-Instruct [16]. Models finetuned using long-context augmentation via SQLong consistently outperform their counterparts finetuned on original contexts. On average, SQLong delivers an absolute improvement of over 2.2% across five benchmark datasets. Additionally, SQLong-finetuned models achieve performance comparable to much larger LLMs on specific datasets, showcasing the scalability and efficiency of the approach.

Performance on Long-Context Datasets. Figure 4 illustrates the experimental results on long-context test sets. Across all datasets, models finetuned with SQLong demonstrate superior performance compared to those trained without SQLong. For instance, on the Spider-test datasets with 8k and 24k context lengths, the Llama-3.1-8B-Instruct model achieves outstanding results of 77.1% and 72.3%, reflecting absolute gains of 7.2% and 13.3%, respectively. Notably, the SQLong-finetuned Llama-8B model outperforms the larger Llama-70B model on 41 out of 45 long-context test sets, with minor exceptions on Spider-realistic 8k and BIRD-dev 8k, 16k, and 24k sets. Similar performance trends are observed with the Qwen models.

On average, SQLong finetuning delivers an 11% absolute improvement over models without SQLong and a 6% advantage over 70B models within the same model family. These results underscore the efficacy of SQLong in handling long-context scenarios and advancing the performance of NL2SQL systems.

Positional robustness. To evaluate the positional robustness of fine-tuned models, we conduct an experiment where each original database schema is placed at varying positions within the input prompt to assess the models' ability to detect them.

We select a set of 124 samples from Spider-dev, Spider-realistic, and Spider-syn, ensuring each sample has a maximum input prompt and target SQL query length of 384 tokens according to CodeQwen1.5-7B-Chat's tokenizer. Using SQLong, we augment this set to a 64k context length. In each augmented set, the original database schemas

⁴Due to computational constraints, we limit finetuning to context lengths of up to 32k. Specifically, for each training example, the context length is randomly sampled from a range starting at 4,096 and increasing by 512 increments up to 32,768.

⁵For Spider, we finetune the base models on the Spider training set and evaluate performance on Spider-dev, Spider-test, Spider-realistic, and Spider-syn.

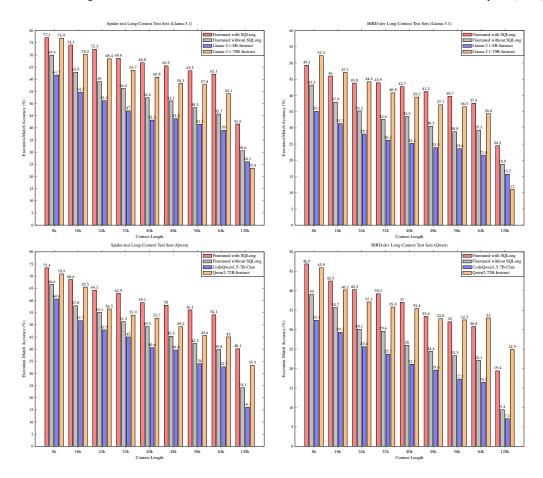


Figure 4: Execution-match accuracy (in %) for Llama-3.1 (top) and Qwen (bottom) families on Spider-test (left) and BIRD-dev (right) long-context test sets.

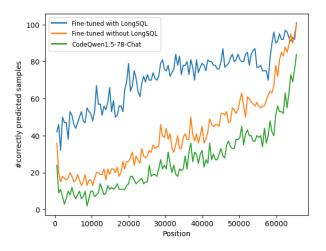


Figure 5: Robust impact of fine-tuned models.

are positioned at specific offsets, starting from 512 and incrementing by 512 up to 64k. This results in 125 new test sets, each containing 124 samples with a 64k context length, corresponding to a distinct schema position.

We compute the number of correctly executed samples for each test set, as shown in Figure 5. The results demonstrate that the long-context fine-tuned model with SQLong is significantly more robust compared to the model without fine-tuning.

4 Conclusion and Future Work

Handling large database schemas poses a significant challenge for NL2SQL models. In this paper, we introduce long-context NL2SQL generation, a novel task that reflects real-world scenarios, and propose SQLong, a simple yet effective augmentation approach for creating long-context finetuning and benchmark datasets. Experiments show that LLMs finetuned with SQLong significantly outperform their counterparts on benchmarks like Spider, BIRD, and our long-context test sets (up to 128k context length).

Future work includes leveraging a RAG-based schema linking approach to retrieve relevant schema elements, enabling more concise and efficient inputs for SQLong-tuned models.

References

- [1] Jinze Bai, Shuai Bai, et al. 2023. Qwen Technical Report. arXiv preprint arXiv:2309.16609 (2023).
- [2] Tri Dao. 2023. Flashattention-2: Faster attention with better parallelism and work partitioning. arXiv preprint arXiv:2307.08691 (2023).
- [3] 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 Coling. 2166–2187.
- [4] Xiang Deng, Ahmed Hassan Awadallah, Christopher Meek, Oleksandr Polozov, Huan Sun, and Matthew Richardson. 2020. Structure-Grounded Pretraining for Text-to-SQL. arXiv preprint arXiv:2010.12773 (2020).
- [5] Tim Dettmers, Mike Lewis, Sam Shleifer, and Luke Zettlemoyer. 2021. 8-bit optimizers via block-wise quantization. arXiv preprint arXiv:2110.02861 (2021).
- [6] 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 llama 3 herd of models. arXiv preprint arXiv:2407.21783 (2024).
- [7] Yujian Gan, Xinyun Chen, Qiuping Huang, Matthew Purver, John R Woodward, Jinxia Xie, and Pengsheng Huang. 2021. Towards Robustness of Text-to-SQL Models against Synonym Substitution. In ACL-IJCNLP. 2505–2515.
- [8] Woosuk Kwon, Zhuohan Li, Siyuan Zhuang, Ying Sheng, Lianmin Zheng, Cody Hao Yu, Joseph E. Gonzalez, Hao Zhang, and Ion Stoica. 2023. Efficient Memory Management for Large Language Model Serving with PagedAttention. In SOSP.
- [9] Jinyang Li, Binyuan Hui, Ge Qu, Jiaxi Yang, Binhua Li, Bowen Li, Bailin Wang, Bowen Qin, Ruiying Geng, Nan Huo, et al. 2023. Can LLM Already Serve as A Database Interface? A BIg Bench for Large-Scale Database Grounded Text-to-SQLs. NeurlPS 2023 Track on Datasets and Benchmarks 36 (2023).

- [10] Xinyu Liu, Shuyu Shen, Boyan Li, Peixian Ma, Runzhi Jiang, Yuyu Luo, Yuxin Zhang, Ju Fan, Guoliang Li, and Nan Tang. 2024. A Survey of NL2SQL with Large Language Models: Where are we, and where are we going? arXiv preprint arXiv:2408.05109 (2024).
- [11] Bowen Peng, Jeffrey Quesnelle, Honglu Fan, and Enrico Shippole. 2023. Yarn: Efficient context window extension of large language models. arXiv preprint arXiv:2309.00071 (2023).
- [12] Nitarshan Rajkumar, Raymond Li, and Dzmitry Bahdanau. 2022. Evaluating the text-to-sql capabilities of large language models. arXiv preprint arXiv:2204.00498 (2022).
- [13] Jie Ren, Samyam Rajbhandari, Reza Yazdani Aminabadi, Olatunji Ruwase, Shuangyan Yang, Minjia Zhang, Dong Li, and Yuxiong He. 2021. {Zero-offload}: Democratizing {billion-scale} model training. In USENIX ATC. 551–564.
- [14] Leandro von Werra, Younes Belkada, Lewis Tunstall, Edward Beeching, Tristan Thrush, Nathan Lambert, and Shengyi Huang. 2020. TRL: Transformer Reinforcement Learning. https://github.com/huggingface/trl.
- [15] Jason Wei, Maarten Bosma, Vincent Zhao, Kelvin Guu, Adams Wei Yu, Brian Lester, Nan Du, Andrew M Dai, and Quoc V Le. 2022. Finetuned Language Models are Zero-Shot Learners. In *ICLR*.
- [16] An Yang, Baosong Yang, et al. 2024. Qwen2 Technical Report. arXiv preprint arXiv:2407.10671 (2024).
- [17] Jiaxi Yang, Binyuan Hui, Min Yang, Jian Yang, Junyang Lin, and Chang Zhou. 2024. Synthesizing Text-to-SQL Data from Weak and Strong LLMs. arXiv preprint arXiv:2408.03256 (2024).
- [18] 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. In EMNLP. 3911–3921.