
The Death of Schema Linking? Text-to-SQL in the Age of Well-Reasoned Language Models

Karime Maamari¹ Fadhil Abubaker¹ Daniel Jaroslawicz¹ Amine Mhedhbi²
¹Distyl AI ²Polytechnique Montreal
{karime,fadhil,daniel}@distyl.ai
amine.mhedhbi@polymtl.ca

Abstract

Schema linking is a crucial step in Text-to-SQL pipelines. Its goal is to retrieve the relevant tables and columns of a target database for a user’s query while disregarding irrelevant ones. However, imperfect schema linking can often exclude required columns needed for accurate query generation. In this work, we revisit schema linking when using the latest generation of large language models (LLMs). We find empirically that newer models are adept at utilizing relevant schema elements during generation even in the presence of large numbers of irrelevant ones. As such, our Text-to-SQL pipeline entirely forgoes schema linking in cases where the schema fits within the model’s context window in order to minimize issues due to filtering required schema elements. Furthermore, instead of filtering contextual information, we highlight techniques such as augmentation, selection, and correction, and adopt them to improve the accuracy of our Text-to-SQL pipeline. Our approach ranks first on the BIRD benchmark achieving an accuracy of 71.83%.

1 Introduction

We address the task of Text-to-SQL: generating a database-executable SQL query given a natural language inquiry (Androustopoulos et al., 1995; Quamar et al., 2022). Text-to-SQL is crucial in democratizing data access as it allows querying databases using natural language. The advent of large language models (LLMs) has significantly advanced Text-to-SQL by simplifying the translation of natural language into SQL.

LLM-based Text-to-SQL approaches typically follow a multi-stage generation pipeline, as shown in Fig. 1 (Hong et al., 2024; Li et al., 2024a; Liu et al., 2024; Zhang et al., 2024). The pipeline begins with a retrieval stage to collect contextual knowledge such as the definition of terms and database schema elements. This is followed by a generation stage, where an LLM produces a candidate SQL query. Finally, the correction stage regenerates the SQL as needed based on encountered errors.

Selecting relevant elements of the database schema (tables and columns) – known as *schema linking* – provides the necessary context for the LLM to produce correct SQL in the downstream generation stage. Effective schema linking implies retrieving *all* the relevant database components associated with the natural language query. Missing even a single required column results in incorrect executable SQL. Thus, it is vital to ensure that all essential columns are retrieved. However, this does not mean that it should be overly inclusive. Research has shown that false positives, *i.e.*, the number of irrelevant columns passed to the LLM, can often degrade text-to-SQL accuracy. For example, Floratou et al. (2024) showed that even when the entire database schema fits into the context of an LLM, it is still advantageous to perform schema linking. At the same time, attempts to prune irrelevant columns may also remove some required ones. Thus, schema linking traditionally contains an inherent trade-off between minimizing false positives while also preserving relevant context.

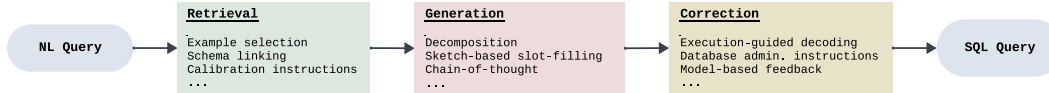


Figure 1: A typical Text-to-SQL pipeline comprised of retrieval, generation and correction stages.

As LLM reasoning capabilities improve, we challenge the conventional wisdom that schema linking is necessary for accurate Text-to-SQL when the schema fits within the model’s context window. We find empirically that as model reasoning improves, the benefits of reducing false positives diminishes, *i.e.*, newer models are more capable of sifting through the schema to identify relevant columns compared to older models (Laban et al., 2024; Li et al., 2024b). This effect might be similar to how the latest LLMs when probed can reliably recall ‘a needle in multiple millions of tokens of haystack’ (Gregory Kamradt, 2023; Reid et al., 2024). For these models, schema linking is unnecessary and can even be detrimental, as it may filter out essential columns. Instead, we present alternatives to schema linking that improve accuracy without schema information loss. Our approach based on these insights currently ranks first in accuracy at 71.83% on the BIRD benchmark (Li et al., 2023).

2 Preliminaries

We first outline the key elements of the Text-to-SQL pipeline, with a focus on schema linking and its implications.

2.1 Text-to-SQL Pipeline

In current state-of-the-art approaches, Text-to-SQL uses multi-stage pipelines comprised of retrieval, generation, and correction stages.

The retrieval stage gathers relevant contextual information, including schema elements, domain knowledge, and example queries (Dong et al., 2023; Gao et al., 2023; Lee et al., 2024; Pourreza & Rafiei, 2023; Wang et al., 2024).

The generation stage often involves more than just producing a candidate SQL query associated given an input context. Rather, approaches frequently augment the generation process through techniques like decomposed generation (Maamari & Mhedhbi, 2024; Pourreza & Rafiei, 2023; Wang et al., 2024) and chain-of-thought prompting (Wei et al., 2023). In addition, most approaches employ methods like self-consistency and multi-choice selection to produce multiple results, selecting the best outcome (Dong et al., 2023; Gao et al., 2023; Lee et al., 2024).

The retrieval and generation stages often contain various techniques chained together. These techniques can be broadly categorized as filtering or augmenting, *i.e.*, techniques can either strip away unnecessary contextual information or attempt to provide additional useful context.

Finally, the correction stage will often employ some combination of execution-based feedback (Wang et al., 2018b; Lin et al., 2020; He et al., 2019; Lyu et al., 2020) or model-based feedback (Talaie et al., 2024; Askari et al., 2024; Wang et al., 2018a) to correct the generated query.

2.2 Schema Linking

Within the retrieval stage, schema linking leverages sophisticated prompting techniques to produce variable-length representations, hierarchically retrieve components, and iteratively process the schema (Talaie et al., 2024; Dong et al., 2023; Pourreza & Rafiei, 2023; Lee et al., 2024; Wang et al., 2024; Gao et al., 2023).

These techniques can vary in i) how they represent the schema and ii) how they perform linking on that representation. For instance, some may represent the schema in natural language, while others utilize code-like structures (Gao et al., 2023). The approach to linking can also differ across techniques, with some directly filtering the schema (Dong et al., 2023; Lee et al., 2024; Talaie et al., 2024), and others using intermediate representations to identify relevant tables and columns (Qu et al., 2024). The choice of schema representation and linking strategy can have a significant influence on accuracy (Gao et al., 2023). This variability underscores the importance of selecting an appropriate

method tailored to the specific requirements of the task, as the degree of filtering can directly impact the loss of schema information incurred during the schema linking process.

Most prior investigations of schema linking regardless of approach have arrived at the same conclusion – schema linking yields meaningful gains in accuracy (Guo et al., 2019; Li et al., 2024a; Talaei et al., 2024). However, these explorations used LLMs that are more sensitive to the presence of irrelevant columns (false positives) as contextual information, where reducing the false positives yielded meaningful gains in performance (Floratos et al., 2024).

3 Experimental Setup

Our experimental analysis guides the design and implementation of our proposed Text-to-SQL pipeline (§4.4). Our experiments aim to answer empirically three research questions (RQs):

- RQ1.* How does the inclusion of irrelevant schema elements impact SQL generation?
- RQ2.* How can the trade-off between precision and recall in schema linking techniques be characterized, and what is its downstream impact on SQL generation?
- RQ3.* How do other techniques and stages within Text-to-SQL pipelines, aside from schema linking, comparatively impact SQL generation?

Next, we cover the details of our setup (datasets and models), and our methodology and metrics.

3.1 Datasets

We conducted our experiments using the BIRD dataset (Li et al., 2023), which is widely considered to be the most challenging Text-to-SQL benchmark. BIRD contains queries from 95 databases spanning a wide breadth of domains, such as education and hockey, and is designed to mimic the complexity of real-world databases. This complexity arises from its “dirty” format, where data, queries, and external knowledge may contain flaws — queries can be incorrect, database columns might be improperly described, and databases can contain null values and unexpected encodings. Our evaluation set consisted of 10% of the entries from each database in the dev set, as done in evaluations in prior work (Talaei et al., 2024). Our training set consisted of 500 of the 9,428 available samples in the BIRD training dataset.

3.2 Models

We used the following language models with context windows sufficiently large to accommodate the entire schema for each query in the BIRD evaluation set:

ft:GPT-4o (fine-tuned)	Llama 3.1-405b	Claude 3.5 Sonnet
GPT-4o	Llama 3.1-70b	Claude 3 Opus
GPT-4o-Mini	Llama 3.1-8b	Mixtral-8x22B
GPT-4-Turbo	Deepseek Coder-V2	Gemini 1.5 Pro

We overview the fine-tuning approach of GPT-4o momentarily under methodology (§3.3).

3.3 Methodology

We design an empirical experiment for each research question:

- Exp 1.* We use a simplified Text-to-SQL pipeline consisting of only schema linking within the retrieval stage followed by a single attempt at generation. For each query, we provide all the required columns and vary the amount of irrelevant columns to look at the impact of irrelevant columns retrieved on the generation.
- Exp 2.* Using the same simplified pipeline, we introduce different implementations of schema linking that vary in precision and recall and analyze their impact on generation.
- Exp 3.* We introduce augmentation, selection, and correction techniques on top of the simplified pipeline without and with schema linking. We run an ablation study to understand the relative impact of each technique on end-to-end accuracy.

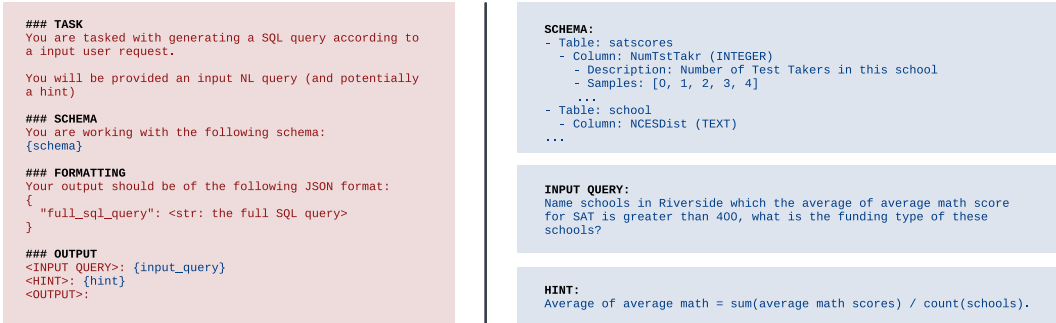


Figure 2: **(Red)** Structure of SQL Generation prompt given input query, hint, and schema; **(Blue)** Examples of a schema, input query, and query hint which act as contextual inputs to a prompt.

Runs and input/output structure. In all runs, the temperature was set to zero and structured output was used whenever possible. Given that not all models are capable of reliable structured output generation, the generated SQL query was fed through an identity call by GPT-4o-Mini in JSON mode to handle any potential issues with output formatting. The relative position of any schema element in an input prompt follows the same ordering as that provided by the schema definition of the benchmark.

Fine-tuning GPT-4o. Fine-tuning is done iteratively. At each iteration, we first fine-tune on a sample of N triples: natural language query, SQL query, and schema elements. For each query, the schema includes all required columns and a random number of irrelevant columns picked uniformly at random. We then evaluate on BIRD’s dev set. For each failed query, we prompt the model to reason about the failure, aggregate the reasoning across queries, and use it to pick the sample for the next iteration. Finally, based on the reasoning, we select a new sample of size N . The iterations stop once a pre-determined accuracy score is reached.

Generation prompts. Fig. 2 shows the structure of the prompt used for SQL generation as well as an example schema, input query, and query hint.

3.4 Metrics

We rely on three different metrics in our experimental analysis:

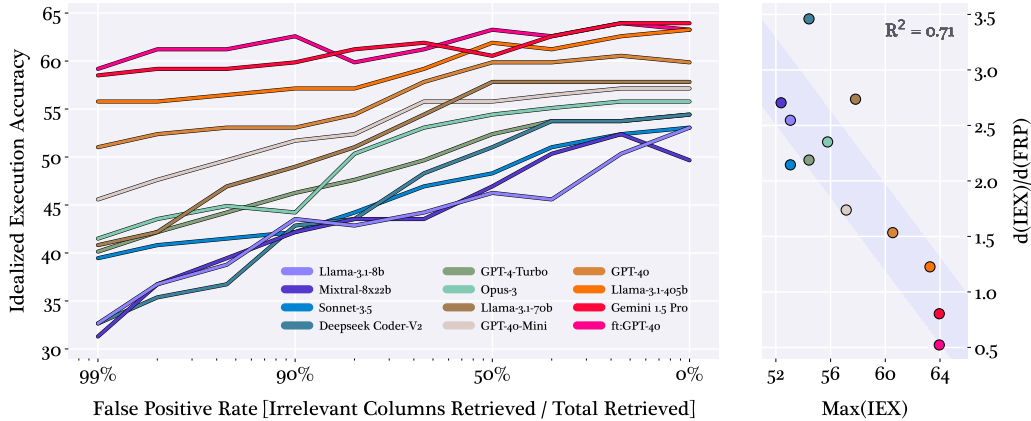
- *Execution Accuracy (EX):* The metric used by the BIRD benchmark to evaluate end-to-end Text-to-SQL pipelines. It is the proportion of queries for which the output of the predicted SQL query is identical to that of the ground truth SQL query. We report EX as a percentage over queries in the evaluation set.
- *False Positive Rate (FPR):* For a given query, the proportion of irrelevant schema columns retrieved over the total number of retrieved columns. We report its average over queries in the evaluation set.
- *Schema Linking Recall (SLR):* The proportion of queries for which all required columns are retrieved over the total number of queries. We use SLR as the downstream generation requires all required columns to be retrieved to be correct.

All queries in evaluation sets are across multiple databases. Note that BIRD also has a second metric: valid efficiency score. It assesses the efficiency of correctly predicted queries by comparing their execution speed to those of the corresponding ground truth queries. In our experiments, we focus on EX, as our research questions are primarily concerned with SQL generation accuracy.

4 Results

4.1 Experiment 1: Impact of False Positives on Accuracy Given Perfect SLR

In this first experiment, we assess the impact of retrieving irrelevant columns on SQL generation accuracy. We create a scenario with perfect schema linking recall such that SQL generation issues are not due to missing required columns.



(a) The idealized execution accuracy (IEX) as the false positive rate (FPR) varies from 99% to 0%. (b) Capability and sensitivity relationship.

Figure 3: Idealized execution accuracy (IEX), *i.e.*, with perfect schema linking recall (SLR), as the false positive rate (FPR) varies for different LMs and the relationship between their capability (maximum IEX) and sensitivity to false positives.

To implement the experiment, we mock the schema linker as follows. We build an oracle that identifies all columns necessary for a given query using SQLGLOT (SQL Parser and Transpiler). For an input query, the mocked schema linker uses the oracle to retrieve all required columns and injects a pre-defined rate of false positives (irrelevant columns retrieved / total columns retrieved). To inject it, the mocked schema linker samples the irrelevant columns uniformly at random from the target database. If there is not a sufficient number of columns in the target database, *e.g.*, there are 10 required columns and 900 irrelevant columns are needed to produce a 99% false positive rate, yet only 50 columns exist in the target database, it supplements from other databases with columns not conflicting in name. We use a simple pipeline: a retrieval stage containing only the mocked schema linker followed by a zero-shot generation stage.

We run the pipeline using the 12 selected models (§3.2) for generation while applying an equal false positive rate to each query in the evaluation set. We run for 10 different rates that are equally log-spaced from 99% and 0%. We refer to the execution accuracy under perfect schema linking recall as the idealized execution accuracy (IEX). We track the change in IEX as the false positive rate changes. The results in Fig. 3a show the broad trend that IEX improves as the rate of false positives decreases; specifically, the stage of SQL generation improves as less irrelevant columns are included as contextual information. At one extreme with a maximum rate of false positives (99%), there is a $\sim 28\%$ IEX difference between the worst and best performing models. In the other extreme with no false positives (0%), the IEX difference between the worst and best performing models gets reduced to $\sim 14\%$.

We use a model’s maximum IEX in this experiment as a proxy for its SQL generation capability. Models with higher generation capability, such as Gemini 1.5 Pro, demonstrate greater resilience to false positives compared to lower-performing ones like Llama 3.1-8b. We characterize a model’s resilience to false positives by the relative change in its SQL generation capability as the false positive rate changes. We define a metric for sensitivity to false positives as the proportion of change in IEX over the proportion of change in the false positive rate (FPR): $d(IEX)/d(FRP)$. As such, sensitivity to false positives is the slope derived from the model’s (IEX, FPR) data points. Fig. 3b depicts a strong negative correlation between a model’s SQL generation capability (maximum IEX) and its sensitivity to false positives.

Empirical Observation: As the model’s SQL generation capability improves, its sensitivity to the presence of irrelevant columns as contextual information for generation decreases. Perhaps surprisingly, both of these capabilities go hand-in-hand where models that are generally better at SQL generation are also more resilient to large amounts of irrelevant context.

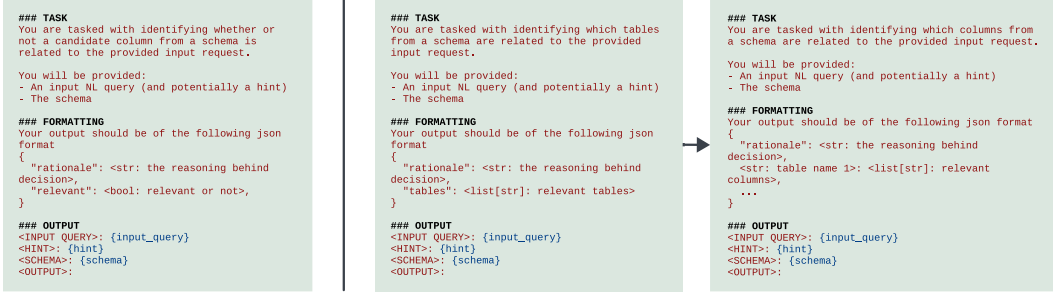


Figure 4: Prompts used for schema linking. **(Left)** Single-Column Schema Linking (SCSL): identifying relevance of a particular column independent of the rest of the schema; **(Middle + Right)** Table-to-Column Schema Linking (TCSL): first identifying relevant tables then relevant columns.

4.2 Experiment 2: Impact of False Positives on Accuracy Given Actual SLR

In Experiment 1, all necessary columns for generation are provided, regardless of the false positive rate. However, in practice, decreasing the amount of false positives requires pruning, which carries the risk of excluding required columns. Here we assess the extent to which schema linking affects recall of required columns and the downstream impact of imperfect recall on generation.

We explore four different schema linking approaches with a broad spectrum of ability to reduce the false positive rate:

- *Single-Column Schema Linking (SCSL)*: Model-determined column-wise relevance. The relevance of each column is assessed without context about other columns and tables. The output is a boolean flag per column indicating the relevance of each column. This is considered a more cautious approach less likely to filter out relevant columns.
- *Hybrid SCSL (HySCSL)*: SCSL with added keyword matching aiming for higher SLR.
- *Table-then-Column Schema Linking (TCSL)*: Model-determined table-then-column filtering approach, as proposed in Talaie et al. (2024) and Pourreza & Rafiei (2023). The model first filters the schema to the relevant tables, then filters the columns within those tables. The output is a set of relevant columns and tables. This is a more aggressive filtering approach.
- *Hybrid TCSL (HyTCSL)*: TCSL with added keyword matching aiming for higher SLR.

In our implementation, SCSL and HySCSL use GPT-4o-Mini and TCSL and HyTCSL use GPT-4o. We use GPT-4o-Mini with SCSL and HySCSL as GPT-4o shows a negligible gain in our experiments but is much cheaper. The prompts used in our implementation are shown in Fig. 4.

Table 1 reports the mean (\pm stdev) of the false positive rate (FPR) and the schema linking recall (SLR) across our four schema linking techniques and without schema filtering at all (full schema) from 12 different runs. Table 1 shows that these approaches are robust across runs and that as FPR decreases, *i.e.*, more irrelevant columns are filtered, more required columns can be filtered and SLR decreases. SLR as obtained from schema linking represents an upper bound on possible EX, where $100\% - SLR$ represents the ratio of queries that will fail in the generation stage due to missing required columns.

We run the same simplified pipeline as our first experiment: a retrieval stage containing one of the five schema linking approaches in Table 1 followed by a zero-shot generation stage. We do so across the 12 selected models (§3.2) and track the associated EX. Fig. 5 shows five EX data points given the FPR of the five different schema linking approaches. We interpolate between two adjacent (FPR, EX) data points and show EX as a solid line. We also add the idealized EX from Experiment 1 as a dashed line. The difference between the two lines shows the EX difference due to changes in SLR.

We observe three different classes of models in our results: models where some schema linking method improves performance (*e.g.*, Llama 3.1-8b), models where all schema linking methods degrade performance (*e.g.*, Gemini 1.5 Pro), and models where schema linking has only negligible impact on performance (*e.g.*, GPT-4o-Mini). A model falls into one of these three buckets based on its SQL generation capability. More capable models – *e.g.* Gemini 1.5 Pro, ft:GPT-4o, Llama-3.1-405b – end up with a reduction in EX. Less capable models – *e.g.* Llama 3.1-8b, Mixtral-8x22b, Deepseek Coder-V2 – end up with a gain in EX.

Approach	FPR	SLR
Without schema filtering (Full Schema)	94.62	100.00
Hybrid Single-Column Schema Linking (HySCSL)	82.08 \pm 0.44	90.36 \pm 0.78
Single-Column Schema Linking (SCSL)	67.23 \pm 0.92	88.77 \pm 0.75
Hybrid Table-to-Column Schema Linking (HyTCSL)	19.85 \pm 0.99	83.00 \pm 0.92
Table-to-Column Schema Linking (TCSL)	9.79 \pm 0.73	77.44 \pm 1.34

Table 1: The mean (\pm stddev) false positive rate (FPR) and schema linking recall (SLR), across 12 runs, associated with five schema linking approaches. The approaches are sorted in descending order of their ability to reduce false positives. We report the mean values (\pm stddev) from 12 runs.

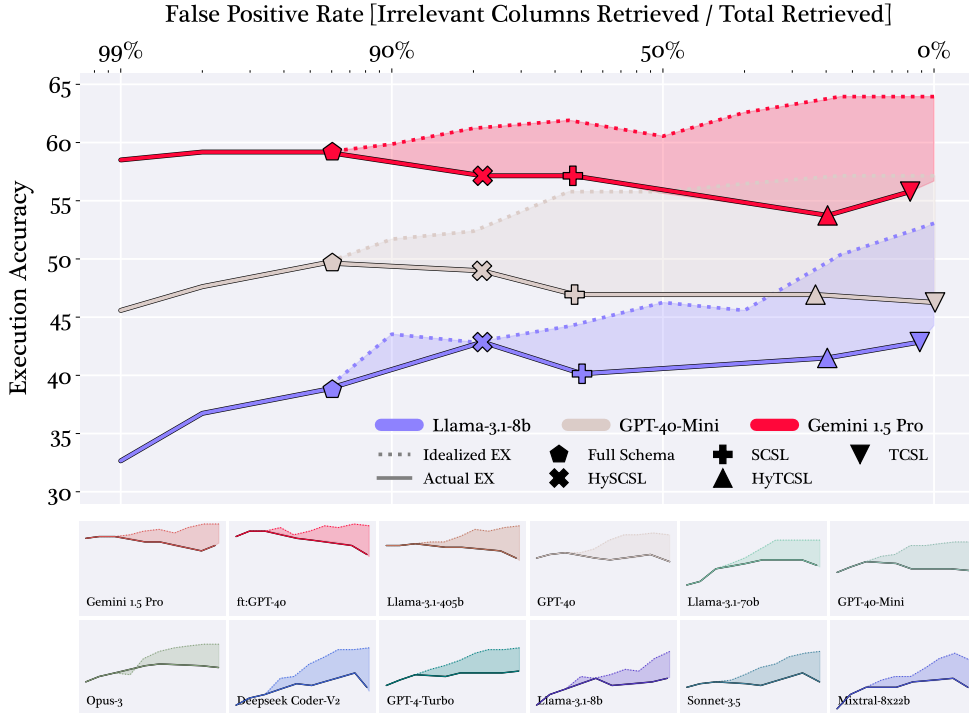


Figure 5: The execution accuracy (EX) given different schema linking and idealized EX, assuming perfect Schema Linking Recall (SLR) as the false positive rate (FPR) varies. EX as a solid line and idealized EX as a dashed line.

Empirical Observation: As the model’s SQL generation capability improves, the benefit of schema linking diminishes. In some cases, this can result in a net reduction in accuracy due to missing required columns for generation.

4.3 Experiment 3: Impact of Non-Filtering Stages and Techniques

Instead of filtering contextual information through schema linking, we focus on techniques that preserve information. We assess the gains of using *augmentation* and *selection* techniques as well as adding a correction stage, which are detailed as follows:

- *Augmentation:* We add contextual information by: (i) expanding column descriptions and add query hints and (ii) adding structural expectations of the output, *e.g.*, expected orderings and aggregations, using CoT planning.
- *Correction:* After generating a candidate SQL query, we iteratively apply corrections through re-generation based on database execution errors (Wang et al., 2018b), revision through database administrator instructions (Talaie et al., 2024), and model-based feedback similar to Reflex-

Method	Execution Accuracy (EX)		
	ft:GPT-4o	Gemini 1.5 Pro	Llama 3.1-405b
Full pipeline	67.35	60.54	59.18
w/o Augmentation	64.63 (↓ 2.72)	60.54	59.86 (↑ 0.68)
w/o Selection	65.31 (↓ 2.04)	57.82 (↓ 2.72)	58.50 (↓ 0.68)
w/o Correction	65.99 (↓ 1.36)	57.14 (↓ 3.40)	55.78 (↓ 3.40)
w/ TCSL	62.58 (↓ 4.77)	55.78 (↓ 4.76)	56.46 (↓ 2.72)
w/ SCSL	55.78 (↓ 11.57)	55.10 (↓ 5.44)	54.42 (↓ 4.76)
Base model	59.18 (↓ 8.17)	57.82 (↓ 2.72)	53.74 (↓ 5.44)

Table 2: Ablation of different methods reporting the execution accuracy (EX) for fine-tuned GPT-4o, Gemini 1.5 Pro, and Llama 3.1-405b. The table compares the full pipeline to variations without augmentation, selection, and correction techniques; with Table-to-Column Schema Linking (TCSL) and Single-Column Schema Linking (SCSL); and as base model performance. Reductions (↓) or increases (↑) in accuracy compared to the full pipeline are indicated.

ion (Shinn et al., 2023). We use these corrections to generate instructions to also augment contextual information.

- *Selection*: We use self-consistency (Wang et al., 2023) to generate multiple responses and select the *most consistent* result. We use selection across the whole pipeline for augmentation, SQL generation, and SQL correction.

We implemented a full pipeline using zero-shot generation. The pipeline contained augmentation, correction and selection as described above and no schema linking, *i.e.*, always providing the full schema. We ran an ablation to understand the relative impact of each technique or stage using the three top performing models from Experiment 1: Llama 3.1-405b, Gemini 1.5 Pro, and ft:GPT-4o. Table 2 shows the execution accuracy (EX) of the full pipeline and 6 other variations: without augmentation, correction, or selection, with schema linking (TCSL and SCSL from Experiment 2), and without any of these techniques *i.e.*, providing the full schema alone to a base model. We find that all techniques improve accuracy to varying degrees except schema linking with a noticeable difference when evaluating the full pipeline.

Empirical Observation: Each of augmentation, selection, and correction have a noticeable positive impact on generation accuracy. Note that even though base models such as Gemini 1.5 Pro may show comparable performance to a fine-tuned GPT-4o in SQL generation, they differ when evaluated within end-to-end Text-to-SQL pipelines. For instance, augmentation leads to major benefits with GPT-4o when compared with Gemini 1.5 Pro. An interesting research direction is understanding whether the relative benefits from augmentation between models hold across other tasks as well.

4.4 Discussion and Proposed Approach

Our empirical observations reveal several key insights that guide our proposed approach.

First, as a model’s SQL generation capability improves, its ability to retrieve relevant schema elements from a full schema in its input context also improves. As such, schema linking for state-of-the-art LLMs is less important if the schema fully fits within the context window. However, it is still helpful for models with lower SQL generation accuracy as they struggle with false positives. In our approach, *we maximize the use of the LLM context window to minimize filtering required columns.*

Second, we find that combining augmentation, selection, and correction techniques heavily impacts the accuracy of a Text-to-SQL pipeline. However, the impact is not the same when evaluated in an end-to-end fashion even when comparing models with similar generation capability. In our approach, *we adopt augmentation, selection, and correction as detailed from the implementation in Experiment 3.* We further use ft:GPT-4o as the model of choice for generation since it provides the best end-to-end accuracy and makes the best use of augmentation.

These design choices yield an approach that ranks first on execution accuracy and second on valid efficiency score on the BIRD benchmark.

5 Conclusion

Is it the death of schema linking? For state-of-the-art models, if the schema fits within the context length – yes. However, for smaller or prior generation models, the accuracy gains from schema linking often justify the potential loss. Additionally, in real-world data-warehousing scenarios, the entire schema often exceeds the context window, requiring a multi-stage information retrieval pipeline. In such cases, we argue for maximally using the LLM context window in picking the top-K relevant columns to retain the necessary schema elements. We conclude that while the need for schema linking is highly use-case dependent, its importance is diminishing as costs decrease, context windows widen, and generation capabilities improve.

Acknowledgements

We thank Jinyang Li, Yongbin Li and the rest of the BIRD team for evaluating our approach on the BIRD test set. We also thank John Allard for valuable insights regarding fine-tuning.

References

- I. Androutsopoulos, G. D. Ritchie, and P. Thanisch. Natural language interfaces to databases - an introduction. *CoRR*, abs/cmp-lg/9503016, 1995.
- Arian Askari, Christian Poelitz, and Xinye Tang. Magic: Generating self-correction guideline for in-context text-to-sql. *CoRR*, abs/2406.12692, 2024.
- Xuemei Dong, Chao Zhang, Yuhang Ge, Yuren Mao, Yunjun Gao, lu Chen, Jinshu Lin, and Dongfang Lou. C3: Zero-shot text-to-sql with chatgpt. *CoRR*, abs/2307.07306, 2023.
- Avrilia Floratou, Fotis Psallidas, Fuheng Zhao, Shaleen Deep, Gunther Hagleither, Wangda Tan, Joyce Cahoon, Rana Alotaibi, Jordan Henkel, Abhik Singla, Alex Van Grootel, Brandon Chow, Kai Deng, Katherine Lin, Marcos Campos, K. Venkatesh Emani, Vivek Pandit, Victor Shnayder, Wenjing Wang, and Carlo Curino. Nl2sql is a solved problem... not! In *Conference on Innovative Data Systems Research*, 2024.
- Dawei Gao, Haibin Wang, Yaliang Li, Xiuyu Sun, Yichen Qian, Bolin Ding, and Jingren Zhou. Text-to-sql empowered by large language models: A benchmark evaluation. *CoRR*, abs/2308.15363, 2023.
- Gregory Kamradt, 2023, 2023. https://github.com/gkamradt/LLMTest_NeedleInAHaystack/blob/main/README.md.
- Jiaqi Guo, Zecheng Zhan, Yan Gao, Yan Xiao, Jian-Guang Lou, Ting Liu, and Dongmei Zhang. Towards complex text-to-sql in cross-domain database with intermediate representation. *CoRR*, abs/1905.08205, 2019.
- Pengcheng He, Yi Mao, Kaushik Chakrabarti, and Weizhu Chen. X-sql: reinforce schema representation with context. *CoRR*, abs/1908.08113, 2019.
- Zijin Hong, Zheng Yuan, Qinggang Zhang, Hao Chen, Junnan Dong, Feiran Huang, and Xiao Huang. Next-generation database interfaces: A survey of llm-based text-to-sql. *CoRR*, abs/2406.08426, 2024.
- Philippe Laban, Alexander R. Fabbri, Caiming Xiong, and Chien-Sheng Wu. Summary of a haystack: A challenge to long-context llms and rag systems. *CoRR*, abs/2407.01370, 2024.
- Dongjun Lee, Choongwon Park, Jaehyuk Kim, and Heesoo Park. Mcs-sql: Leveraging multiple prompts and multiple-choice selection for text-to-sql generation. *CoRR*, abs/2405.07467, 2024.
- Boyan Li, Yuyu Luo, Chengliang Chai, Guoliang Li, and Nan Tang. The dawn of natural language to sql: Are we fully ready? *CoRR*, abs/2406.01265, 2024a.
- Jinyang Li, Binyuan Hui, Ge Qu, Jiayi Yang, Binhua Li, Bowen Li, Bailin Wang, Bowen Qin, Rongyu Cao, Ruiying Geng, Nan Huo, Xuanhe Zhou, Chenhao Ma, Guoliang Li, Kevin C. C. Chang, Fei Huang, Reynold Cheng, and Yongbin Li. Can llm already serve as a database interface? a big bench for large-scale database grounded text-to-sqls. *CoRR*, abs/2305.03111, 2023.
- Mo Li, Songyang Zhang, Yunxin Liu, and Kai Chen. Needlebench: Can llms do retrieval and reasoning in 1 million context window? *CoRR*, abs/2407.11963, 2024b.
- Xi Victoria Lin, Richard Socher, and Caiming Xiong. Bridging textual and tabular data for cross-domain text-to-sql semantic parsing. *CoRR*, abs/2012.12627, 2020.
- Xinyu Liu, Shuyu Shen, Boyan Li, Peixian Ma, Runzhi Jiang, Yuyu Luo, Yuxin Zhang, Ju Fan, Guoliang Li, and Nan Tang. A survey of nl2sql with large language models: Where are we, and where are we going? *CoRR*, abs/2408.05109, 2024.
- Qin Lyu, Kaushik Chakrabarti, Shobhit Hathi, Souvik Kundu, Jianwen Zhang, and Zheng Chen. Hybrid ranking network for text-to-sql. *CoRR*, abs/2008.04759, 2020.
- Karime Maamari and Amine Mhedhbi. End-to-end text-to-sql generation within an analytics insight engine. *CoRR*, abs/2406.12104, 2024.

- Mohammadreza Pourreza and Davood Rafiei. Din-sql: Decomposed in-context learning of text-to-sql with self-correction. *CoRR*, abs/2304.11015, 2023.
- Ge Qu, Jinyang Li, Bowen Li, Bowen Qin, Nan Huo, Chenhao Ma, and Reynold Cheng. Before generation, align it! a novel and effective strategy for mitigating hallucinations in text-to-sql generation. *CoRR*, abs/2405.15307, 2024.
- Abdul Quamar, Vasilis Efthymiou, Chuan Lei, and Fatma Özcan. Natural language interfaces to data. *Found. Trends Databases*, 11(4):319–414, 2022.
- Machel Reid, Nikolay Savinov, Denis Teplyashin, Dmitry Lepikhin, Timothy P. Lillicrap, Jean-Baptiste Alayrac, Radu Soricut, Angeliki Lazaridou, Orhan Firat, Julian Schrittwieser, Ioannis Antonoglou, Rohan Anil, Sebastian Borgeaud, Andrew M. Dai, Katie Millican, Ethan Dyer, Mia Glaese, Thibault Sottiaux, Benjamin Lee, Fabio Viola, Malcolm Reynolds, Yuanzhong Xu, James Molloy, Jilin Chen, Michael Isard, Paul Barham, Tom Hennigan, Ross McIlroy, Melvin Johnson, Johan Schalkwyk, Eli Collins, Eliza Rutherford, Erica Moreira, Kareem Ayoub, Megha Goel, Clemens Meyer, Gregory Thornton, Zhen Yang, Henryk Michalewski, Zaheer Abbas, Nathan Schucher, Ankesh Anand, Richard Ives, James Keeling, Karel Lenc, Salem Haykal, Siamak Shakeri, Pranav Shyam, Aakanksha Chowdhery, Roman Ring, Stephen Spencer, Eren Sezener, and et al. Gemini 1.5: Unlocking multimodal understanding across millions of tokens of context. *CoRR*, abs/2403.05530, 2024.
- Noah Shinn, Federico Cassano, Edward Berman, Ashwin Gopinath, Karthik Narasimhan, and Shunyu Yao. Reflexion: Language agents with verbal reinforcement learning. *CoRR*, abs/2303.11366, 2023.
- Shayan Talaei, Mohammadreza Pourreza, Yu-Chen Chang, Azalia Mirhoseini, and Amin Saberi. Chess: Contextual harnessing for efficient sql synthesis. *CoRR*, abs/2405.16755, 2024.
- Bing Wang, Changyu Ren, Jian Yang, Xinnian Liang, Jiaqi Bai, Linzheng Chai, Zhao Yan, Qian-Wen Zhang, Di Yin, Xing Sun, and Zhoujun Li. Mac-sql: A multi-agent collaborative framework for text-to-sql. *CoRR*, abs/2312.11242, 2024.
- Chenglong Wang, Kedar Tatwawadi, Marc Brockschmidt, Po-Sen Huang, Yi Mao, Oleksandr Polozov, and Rishabh Singh. Robust text-to-sql generation with execution-guided decoding. *CoRR*, abs/1807.03100, 2018a.
- Chenglong Wang, Kedar Tatwawadi, Marc Brockschmidt, Po-Sen Huang, Yi Mao, Oleksandr Polozov, and Rishabh Singh. Robust text-to-sql generation with execution-guided decoding. *CoRR*, abs/1807.03100, 2018b.
- Xuezhi Wang, Jason Wei, Dale Schuurmans, Quoc Le, Ed Chi, Sharan Narang, Aakanksha Chowdhery, and Denny Zhou. Self-consistency improves chain of thought reasoning in language models. *CoRR*, abs/2203.11171, 2023.
- Jason Wei, Xuezhi Wang, Dale Schuurmans, Maarten Bosma, Brian Ichter, Fei Xia, Ed Chi, Quoc Le, and Denny Zhou. Chain-of-thought prompting elicits reasoning in large language models. *CoRR*, abs/2201.11903, 2023.
- Weixu Zhang, Yifei Wang, Yuanfeng Song, Victor Junqiu Wei, Yuxing Tian, Yiyang Qi, Jonathan H. Chan, Raymond Chi-Wing Wong, and Haiqin Yang. Natural language interfaces for tabular data querying and visualization: A survey. *CoRR*, abs/2310.17894, 2024.