

# Improving Retrieval-augmented Text-to-SQL with AST-based Ranking and Schema Pruning

Anonymous ACL submission

## Abstract

We focus on Text-to-SQL semantic parsing from the perspective of Large Language Models. Motivated by challenges related to the size of commercial database schemata and the deployability of business intelligence solutions, we propose an approach that dynamically retrieves input database information and uses abstract syntax trees to select few-shot examples for in-context learning.

Furthermore, we investigate the extent to which an in-parallel semantic parser can be leveraged for generating *approximated* versions of the expected SQL queries, to support our retrieval. We take this approach to the extreme—we adapt a model consisting of less than 500M parameters, to act as an extremely efficient approximator, enhancing it with the ability to process schemata in a parallelised manner. We apply our approach to monolingual and cross-lingual benchmarks for semantic parsing, showing improvements over state-of-the-art baselines. Comprehensive experiments highlight the contribution of modules involved in this retrieval-augmented generation setting, revealing interesting directions for future work.

## 1 Introduction

Text-to-SQL semantic parsing aims at translating natural language questions into SQL, to facilitate querying relational databases by non-experts (Zelle and Mooney, 1996). Given their accessibility benefits, Text-to-SQL applications have become popular recently, with many corporations developing Business Intelligence platforms.

The success of Large Language Models (LLMs) in generalising across diverse Natural Language Processing tasks (Ye et al., 2023; OpenAI et al., 2024) has fuelled works that looked at how these multi-billion parameter models can be best employed for Text-to-SQL (Liu et al., 2023; Pourreza and Rafiei, 2023). Recent works in this space have focused on the in-context learning ability of LLMs,

demonstrating that significant improvements can be achieved by selecting suitable (question, SQL) example pairs (Nan et al., 2023; Gao et al., 2023; Guo et al., 2024; Sun et al., 2024). In spite of its underlying benefits, conventional solutions for example selection are usually limited to retrieving examples based solely on the similarity of questions (Nan et al., 2023; An et al., 2023; Guo et al., 2024). Other approaches resort to a preliminary round of parsing which *approximates* expected SQL queries, and directly use these approximations in few-shot prompting (Sun et al., 2024), or to subsequently select (question, SQL) pairs by comparing the approximated query to queries within candidate examples (Gao et al., 2023). The approach proposed by Gao et al. transforms SQL queries into SQL skeletons (Li et al., 2023a) and then filters examples by considering overlap token ratio as the similarity between two skeletons. While incorporating SQL skeleton similarity improves over conventional example selection for Text-to-SQL (Gao et al., 2023), it can result in structural information loss as exemplified in Table 1, where two dissimilar SQL queries are treated as identical. In this paper, we propose a novel approach that selects examples using similarity of normalised SQL Abstract Syntax Trees (ASTs). We argue that considering the similarity of such hierarchical structures can significantly enhance LLMs’ performance for Text-to-SQL parsing.

Apart from example selection, we refine database context input to LLMs by dynamically pruning schemata and selecting values. From the perspective of LLMs, existing studies achieve improvements by including the full database schema in the prompt and additionally *hinting* the importance of particular schema elements or values (Pourreza and Rafiei, 2023; Sun et al., 2024). In this paper, we show that the performance can be boosted with schemata of reduced size.

Inspired by Gao et al. that compute an approxi-

mated query for a given input question, we further explore how combinations of a sparse retriever with such an in-parallel semantic parser (we would refer to it as *approximator*) can be used to retrieve relevant database context input to LLMs. For efficiency, we adapt the semantic parser (a decoder-free model with  $< 500\text{M}$  parameters) proposed by Vougiouklis et al. to process schemata in a parallelised manner. Using this efficient approximator, our schema pruning strategy selects a relevant subschema in order to simplify the task for LLMs and reduce the relevant computational workload. Furthermore, it enables LLM-based Text-to-SQL solutions to handle longer schemata (usually associated with commercial use-cases) exceeding their context window size.

We apply our approach on monolingual (SPIDER, SPIDER-DK, SPIDER-REAL and SPIDER-SYN) and cross-lingual (CSPIDER) benchmarks of different generalisation challenges. We evaluate the applicability of our framework across both closed- and open-source LLMs. Our framework, comprising only a single round of prompting, achieves state-of-the-art performance, outperforming other baselines which may comprise complex prompting and multiple iterations, when LLMs of equal capacity are involved. Through comprehensive experiments, we highlight strengths and limitations. Our contributions can be summarised as follows:

- We propose a novel approach for selecting (question, SQL) examples using similarity of normalised SQL ASTs.
- We take efficient approximation to the extreme, presenting a schema-parallelisable adaptation of the fastest semantic parser to date.
- We introduce a framework for dynamically selecting schema elements and database values, offering substantial execution accuracy improvements over prior works while significantly reducing the computational workload of LLMs.
- We shed light on the benefits of database value selection and its symbiotic relation to schema pruning for Text-to-SQL LLM prompting.

## 2 Preliminaries

Let  $\mathbf{q}$  be the sequence of tokens of a natural language question for database  $\mathbf{D}$  with tables  $\mathbf{t} =$

$\text{SQL}_1$	SELECT T2.name, T2.capacity FROM concert AS T1 JOIN stadium AS T2 ON T1.stadium_id = T2.stadium_id WHERE T1.year >= 2014
	<b>Skeleton:</b> select _ from _ where _
$\text{SQL}_2$	SELECT name FROM highschooler WHERE grade = 10
	<b>Skeleton:</b> select _ from _ where _

Table 1: Two SQL queries with identical SQL skeletons.

$t_1, t_2, \dots, t_T$  and columns  $\mathbf{c} = c_1^1, c_2^1, \dots, c_j^i, \dots, c_{C_T}^T$ , where  $c_j^i$  is the  $j$ -th column of table  $t_i$  and  $C_i \in \mathbb{N}$  is the total number of columns in table  $t_i$ . Furthermore, let  $\mathbf{v}_\mathbf{D} = \{v_{c_1^1}, v_{c_2^1}, \dots, v_{c_{C_T}^T}\}$  be the set of all values associated with the database  $\mathbf{D}$  s.t.  $v_{c_1^1}, \dots, v_{c_{C_T}^T}$  are the DB value sets associated with respective columns  $c_1^1, \dots, c_{C_T}^T \in \mathbf{c}$ . The goal of Text-to-SQL semantic parsing is to predict the SQL query  $\mathbf{s}$  given the  $(\mathbf{q}, \mathbf{D})$  combination, as follows:

$$\mathbf{s} = \arg \max_{\mathbf{s}} p(\mathbf{s} | \mathbf{q}, \mathbf{D}) \quad (1)$$

For in-context learning, we seek to select pertinent input context including few-shot examples, schema, and database values to simplify the task for LLMs.

## 3 Example Selection using Abstract Syntax Trees

Our goal is to identify the most suitable set of  $\mathbb{X}^* = \{(\mathbf{q}_1^*, \mathbf{s}_1^*), \dots, (\mathbf{q}_e^*, \mathbf{s}_e^*)\}$  question-SQL pairs from an index of examples,  $\mathbb{X}$ , s.t.  $\mathbb{X}^* \subseteq \mathbb{X}$ , for maximising the probability of an LLM to predict the correct SQL given  $(\mathbf{q}, \mathbf{D})$ :

$$\mathbb{X}^* = \arg \max_{\mathbb{X}} p(\mathbf{s} | \mathbf{q}, \mathbf{D}, \mathbb{X}) \quad (2)$$

From the perspective of ranking, we consider the relevance score between a candidate example  $(\mathbf{q}_j, \mathbf{s}_j) \in \mathbb{X}$  and the input  $(\mathbf{q}, \mathbf{D})$ . *Vanilla* semantic search is usually based solely on question embeddings, whereas the structure of SQL queries for similar questions is subject to target databases and can thus differ significantly.

To incorporate database context for selecting examples, we propose to re-rank examples retrieved by question embeddings based on normalised SQL ASTs. Inspired by Gao et al., our framework utilises a preliminary model to compute an approximated SQL query  $\mathbf{s}'$ , structurally similar to the

ground truth, given  $(\mathbf{q}, \mathbf{D})$  s.t.  $s' \sim s$ . Examples are then re-ranked by  $\text{score}_{\text{AST}}(s', s_j)$  for each candidate  $s_j$ .

AST represents the hierarchical structure of code in a tree form and can be applied to evaluation metrics for code generation (Tran et al., 2019; Ren et al., 2020). The fact that SQL queries sharing identical abstract meanings may not align with the same syntactic structure poses a challenge for measuring similarity through AST differencing.

**AST Normalisation** Although it is infeasible to exhaustively transform a SQL to another equivalent form, we can normalise ASTs to reduce undesired mismatch. Firstly, nodes of identifiers are lower-cased and unnecessary references are removed (e.g. `<table>.<column>` is substituted with `<column>` if possible). We then delete nodes that create aliases and map each alias to a copy of the subtree to which it references. For cross-domain settings wherein databases at inference time are unseen in the train set, we mask out nodes of values and identifiers after resolving aliases. Otherwise for in-domain settings we further sort nodes associated with JOIN operations(s) to ensure the ordering of tables and keys is consistent.

**AST Similarity** Given two normalised ASTs, we adopt the Change Distilling algorithm (Fluri et al., 2007) that computes a list of tree edit operations to transform the source AST to the target AST. Types of tree edit operations include: insert, delete, alignment, move and update. It is essential to note that move operation relocates a node to a different parent while moving a node within the same parent is an alignment. Therefore, we calculate the similarity between ASTs simply as the ratio of alignments to the total number of operations within the list. Examples of our normalisation and AST similarity are provided in Appendix A.

## 4 Database Context Selection

Apart from relevant question-SQL pairs, prompting for Text-to-SQL parsing requires the context of database schema and values.

### 4.1 Schema Selection

We present a hybrid search strategy that selects a sub-schema given a test question to minimise lengthy and potentially irrelevant schema elements input to LLMs, while maintaining high recall.

Let  $r_j^i$  be a semantic representation of column  $c_j^i$ .

We aggregate the semantic names<sup>1</sup> of  $c_j^i$  and the table it belongs to,  $t_i$ , and its corresponding value set in  $\mathbf{D}$ ,  $v_{c_j^i}$ , as follows:

$$r_j^i = \left\{ t_i \cup c_j^i \cup v_{c_j^i} \mid i \in [1, T] \text{ and } j \in [1, C_i] \right\}.$$

Given question  $\mathbf{q}$ , we retrieve the most relevant columns using  $\text{score}_{\text{BM25}}(\mathbf{q}, r_j^i) \forall i \in [1, T]$  and  $j \in [1, C_i]$  (Robertson et al., 1994). A table is retrieved if any of its columns are retrieved.

#### 4.1.1 Incorporating for Approximated Query

The semantic search for schema selection requires a comprehension of the relevance between heterogeneous database information and natural language questions, in addition to interactions across schema elements. To this end, a trained parser can inherently be a semantic search model for retrieving a sub-schema, where columns and tables are extracted from the approximated query  $s'$ . We argue that a semantic parser which is performing reasonably on the task, can provide us with an  $s'$ , whose structure would assimilate the structure of the expected final query,  $s$ . Consequently, we opt to dynamically determine the number of columns to be retrieved by  $\text{score}_{\text{BM25}}$  as proportional to the number of unique columns in  $s'$ , returned by the approximator. A sub-schema is then obtained by merging schema elements selected by  $\text{score}_{\text{BM25}}$  with elements from the approximated query.

#### 4.1.2 Approximating for Longer Schemata

To further reduce the computational workload, we opt for using a *smaller* model for computing  $s'$ . However, smaller models usually have shorter context windows (i.e.  $< 2\text{k}$  tokens), and, as such, they cannot be easily scaled to the requirements of larger schemata. To this end, we propose an approach that enables transformer-based encoders to process longer schemata, in a parallelised manner.

We start with FastRAT (Vougiouklis et al., 2023), which exploits a decoder-free architecture for efficient text-to-SQL parsing. Given a concatenation of the input natural language question  $\mathbf{q}$  with the column and table names of a database schema, FastRAT computes the SQL operation in which each element of the input schema would participate in the expected SQL query. We refer to these SQL operations as SQL Semantic Prediction (SSP) labels (Yu et al., 2021). SQL queries are then deterministically constructed from the predicted SSP

<sup>1</sup>Semantic name can refer to simply to the name of a particular or to a concatenation of its name and description.

labels. We introduce a schema splitting strategy to scale the model up to the requirements of schemata comprising several columns.

We augment the input embedding matrix of the model, with two special schema-completion tokens, [full\_schema] and [part\_schema], which are used for signalling cases in which a full and a partial schema are provided as input respectively. Our goal is to split a schema consisting of  $\sum_{j=1}^T C_j$  columns into  $r_m$  splits s.t. each split includes a maximum, pre-defined number of columns  $r$ . Each split consists of the question tokens, a single schema-completion token, the table names of the input  $\mathbf{D}$  and up to a maximum  $r$  number of columns allocated to this particular split (see Algorithm 1 for further details).

---

**Algorithm 1:** Algorithm for splitting a schema into smaller splits.

---

```

Input: # col. name tokens concat.
           $\mathbf{c}^{\text{tok}} \leftarrow [c_1^{\text{tok}}, \dots, c_{C_T}^{\text{tok}}]$ 
          # flatten concat. of tab.
          name tokens
           $\mathbf{t}^{\text{tok}} \leftarrow [t_1^{\text{tok}}, t_2^{\text{tok}}, \dots, t_T^{\text{tok}}]$ 
           $\mathbf{q} \leftarrow [q_1, \dots, q_Q]$  # q tokens
          r : int
1 splits  $\leftarrow []$ ;
2 if len( $\mathbf{c}$ ) > r then
3   | prefix  $\leftarrow \mathbf{q} + [“[part\_schema]”]$ ;
4 else
5   | prefix  $\leftarrow \mathbf{q} + [“[full\_schema]”]$ ;
6 end
7 sp  $\leftarrow$  prefix; # one sp per split
8 for  $j \leftarrow 1$  to  $\sum_{j=1}^T C_j$  do
9   | sp  $\leftarrow$  sp +  $\mathbf{c}^{\text{tok}}[j]$ ;
10  | if  $j \bmod r = 0$  or  $j = \sum_{j=1}^T C_j$  then
11    | sp  $\leftarrow$  sp +  $\mathbf{t}^{\text{tok}}$ ;
12    | splits.append(sp);
13    | sp  $\leftarrow$  prefix;
14  | end
15 end
16 Return splits

```

---

The returned schema splits along with the SSP labels corresponding to the schema elements of each split are treated as independent instances during training. At test time, an input schema is split according to Algorithm 1, and the model is input with a batch of the resulting splits. After aggregating the results from all splits, we obtain the SSP label for each column  $\in \mathbf{c}$ . Inconsistencies across

the SSP labels of tables are resolved using majority voting. We refer to this model as FastRAT<sub>ext</sub>.

## 4.2 Value Selection

The inference of LLMs for text-to-SQL parsing can be augmented with column values (Sun et al., 2024). We select values for columns in a schema (or a sub-schema) by simply matching keywords in questions and values. This is based on the assumption that LLMs can generalise to unseen values given a set of representative values; thus, the recall and precision of value selection are less critical. We consider value selection providing additional information for LLMs to discern covert differences among columns. An example of our resulting prompt is shown in Appendix B.

## 5 Experiments

We run experiments using two approximators: FastRAT<sub>ext</sub> and Graphix-T5 (Li et al., 2023b). Graphix-T5 is the approximator used by DAIL-SQL (Gao et al., 2023), and is included to facilitate a fair comparison against the closest work to ours. FastRAT<sub>ext</sub> is trained and tested using  $r = 64$ , unless otherwise stated (cf. Section 5.4). For all experiments, we use 5 (question, SQL) examples.

We test our approach against both closed- and open-source LLMs: (i) gpt-3.5-turbo (gpt-3.5-turbo-0613), (ii) gpt-4 (gpt-4-0613) and (iii) deepseek-coder-33b-instruct. Results using additional models from the DeepSeek family are provided in Appendix C.4.

### 5.1 Datasets

We experiment with several SQL datasets, seeking to explore the effectiveness of our approach on both monolingual and cross-lingual setups. Specifically, we report experiments on CSPIDER (Min et al., 2019) and SPIDER (Yu et al., 2018). Since CSPIDER is a translated version of SPIDER in Chinese, when it comes to the natural language questions, the characteristics of the two with respect to structure and number of examples are identical. We focus our evaluation on the development sets<sup>2</sup>, which are used as test sets in our experiments. These splits consists of 1,034 examples of questions on 20 unique databases that are not met at training time.

We rely on the training splits to maintain an index of (question, SQL) examples, one for each

<sup>2</sup>Appendix D includes results on the test sets.

dataset. Using these splits, we train a monolingual and a cross-lingual version of FastRAT<sub>ext</sub>.

Furthermore, we use popular SPIDER variants: (i) SPIDER-DK (Gan et al., 2021b), (ii) SPIDER-REAL (Deng et al., 2021) and (iii) SPIDER-SYN (Gan et al., 2021a) to evaluate zero-shot domain generalisation in English (leveraging the SPIDER (question, SQL) examples index).

Consistently with the relevant leaderboards<sup>3</sup>, we report results using execution (EX) and exact match (EM) accuracy.<sup>4</sup> Since CSPIDER comes without relevant DB content, we follow previous works, and we focus our evaluation on EM scores (Vougiouklis et al., 2023; Cao et al., 2023).

Model	EX	EM
GraPPa (Yu et al., 2021)	–	73.6
FastRAT (Vougiouklis et al., 2023)	73.2	69.1
FastRAT <sub>ext</sub>	71.5	64.2
Graphix-T5 (Li et al., 2023b)	81.0	77.1
RESDSQL (Li et al., 2023a)	84.1	<b>80.5</b>
<hr/>		
deepseek-coder-33b-instruct		
Ours (w/ FastRAT <sub>ext</sub> )	81.5	62.1
Ours (w/ Graphix-T5)	83.4	64.7
<hr/>		
PaLM2		
Few-shot SQL-PaLM (Sun et al., 2024)	82.7	–
<hr/>		
text-davinci-003		
Zero-shot (Guo et al., 2024)	73.1	–
RAG w/ Rev. Chain (Guo et al., 2024)	<u>85.0</u>	–
<hr/>		
gpt-3.5-turbo		
Zero-shot (Liu et al., 2023)	70.1	–
C3 (Dong et al., 2023)	81.8	–
DAIL-SQL (Gao et al., 2023)	79.0	–
Ours (w/ FastRAT <sub>ext</sub> )	82.0	65.7
Ours (w/ Graphix-T5)	83.0	68.8
<hr/>		
gpt-4		
Zero-shot (Pourreza and Rafiei, 2023)	72.9	40.4
DIN-SQL (Pourreza and Rafiei, 2023)	82.8	60.1
DAIL-SQL (Gao et al., 2023)	83.6	68.7
Ours (w/ FastRAT <sub>ext</sub> )	84.3	73.8
Ours (w/ Graphix-T5)	<b>86.6</b>	<u>77.3</u>

Table 2: EX and EM accuracies on the development split of SPIDER. Fine-tuning-based baselines are listed at the top part of the table. Results of our approach are shown with both FastRAT<sub>ext</sub> and Graphix-T5 as approximators. The **best** model is in bold, the second best is underlined, and the **best prompting-based setup** is in blue.

## 5.2 Baselines

We dichotomize the landscape of baselines in fine-tuning- and prompting-based baselines. Further

<sup>3</sup><https://taolusi.github.io/CSpider-explorer/> and <https://yale-lily.github.io/spider>

<sup>4</sup>EX and EM scores are computed using: <https://github.com/taoyds/test-suite-sql-eval>.

details are provided in Appendix E.

**Fine-tuning-based** (i) **GraPPa**, (ii) **DG-MAML**, (iii) **FastRAT**, (iv) **Graphix-T5**, (v) **RESDSQL** and (vi) **HG2AST**.

**Prompting-based** Zero-shot LLM prompting has been explored by Guo et al.; Liu et al.; Pourreza and Rafiei; (i) **C3** introduces calibration bias for LLM prompting; (ii) **DIN-SQL** uses chain-of-thought prompting with pre-defined prompting templates tailored for the assessed question hardness; (iii) **DAIL-SQL** uses query approximation and SQL skeleton-based similarities for example selection; (iv) **SQL-PaLM** proposes a framework for *soft* column selection and execution-based refinement; (v) **RAG w/ Rev. Chain** augments the input prompt with question skeleton-based example retrieval and an execution-based revision chain.

## 5.3 Text-to-SQL Evaluation

Table 2 and 3 summarise the results of our approach with deepseek-coder-33b-instruct, gpt-3.5-turbo and gpt-4 against the baselines. Our approach, comprising a single-prompting round, surpasses other LLM-based solutions, that incorporate several prompting iterations, for LLMs of the same capacity. We note consistent improvements over DAIL-SQL, the closest work to ours, even when FastRAT<sub>ext</sub> is used as approximator (i.e. a model consisting of < 500M vs the ≥ 3B parameters that DAIL-SQL’s approximator is using). For the same approximator, our framework is able to meet, performance standards of DAIL-SQL (equipped with gpt-4 and an additional self-consistency prompting step) using an open-source model as backbone LLM, by achieving shorter prompts in a single prompting step.

SPIDER results are consistent with the results across the various Spider variants and CSPIDER<sup>5</sup> (Table 3). Our approach leveraging FastRAT<sub>ext</sub> and AST-based re-ranking for example selection outperforms other prompting-based solution, and is in-line with the scores of state-of-the-art fine-tuning-based baselines. While gpt-4 is the most capable model within our framework (with this being more noticeable in the case of SPIDER-SYN), we observe surprising findings with DeepSeek with which in many cases our approach can surpass much more computationally expensive alternatives

<sup>5</sup>In CSPIDER, questions are fully translated in Chinese while the DB content remains in English. Due to this limitation, DB schema and content selection are disabled.

Model	SPIDER-DK		SPIDER-REAL		SPIDER-SYN		CSPIDER
	EX	EM	EX	EM	EX	EM	EM
RAT-SQL + BERT (Wang et al., 2020)	–	40.9	62.1	58.1	–	48.2	–
DG-MAML (Wang et al., 2021)	–	–	–	–	–	–	51.0
FastRAT (Vougiouklis et al., 2023)	–	–	–	–	–	–	<u>61.3</u>
FastRAT <sub>ext</sub>	–	44.1	–	47.8	–	48.5	53.2
HG2AST (Cao et al., 2023)	–	–	–	–	–	–	61.0 <sup>a</sup>
RESDSL (Li et al., 2023a)	66.0	<u>55.3</u>	<b>81.9</b>	<b>77.4</b>	<b>76.9</b>	<b>69.1</b>	–
deepseek-coder-33b-instruct							
Ours (w/ FastRAT <sub>ext</sub> )	<u>70.5</u>	46.4	77.4	59.3	68.7	49.5	55.9
gpt-3.5-turbo							
Zero-shot (Liu et al., 2023)	62.6	–	63.4	–	58.6	–	32.6
DAIL-SQL (Gao et al., 2023)	–	–	67.9	–	–	–	–
Ours (w/ FastRAT <sub>ext</sub> )	68.8	49.3	78.0	60.8	66.9	51.2	54.0
PaLM2							
Few-shot SQL-PaLM (Sun et al., 2024)	66.5	–	77.6	–	<u>74.6</u>	–	–
gpt-4							
DAIL-SQL (Gao et al., 2023)	–	–	76.0	–	–	–	–
Ours (w/ FastRAT <sub>ext</sub> )	<b>72.3</b>	<b>59.1</b>	<u>80.9</u>	<u>66.1</u>	74.4	<u>61.3</u>	<b>64.4</b>

Table 3: Results on SPIDER-DK, SPIDER-REAL, SPIDER-SYN and CSPIDER. Fine-tuning-based baselines are listed at the top. The **best** model is in bold, second best is underlined, and the **best prompt-based setup** is in blue.

<sup>a</sup>Without using question translation; 64.0 EM when question translation is used.

based on larger closed-source LLMs. Our findings remain consistent with (Liu et al., 2023) since the EM scores of prompting-based methods fall behind those of their fine-tuning based counterparts.

### 5.3.1 Schema Selection Evaluation

We evaluate our proposed schema selection strategy in a two-fold manner, given that value selection is applied for selected columns. Firstly, we use recall and schema shortening (rate) to compute averaged metrics across all samples showcasing the extent to which (i) the most relevant schema elements are successfully retrieved, and (ii) the size of the resulting schema, after selection, with respect to its original size. Secondly, we explore how the performance of the end-system changes across different schema pruning settings by reporting EX and EM scores. Recall is the percentage of samples for which all ground-truth schema elements are selected. Schema shortening is the number of schema elements that are excluded divided by the total number of schema elements. Results are summarised in Table 4.

The benefits of schema selection are apparent in the oracle setup, in which only schema elements from the gold query are included in the prompt (cf. Table 8). In this setup, the highest execution accuracy is achieved while filtering out  $> 70\%$  of the original schema on average. We note that

our approach of coupling the schema elements returned in the approximated query with the ones returned by BM25 navigates a healthy trade-off between maximising recall and reducing processing of unnecessary schema elements. We also notice that our strategy of dynamically determining the number of retained schema elements per input (cf. Section 4.1.1) results in improvements compared to static top- $k$  determination. For roughly the same extent of schema shortening (i.e. by comparing scores with dynamic top- $k$  against top-7), the results with the former are higher across all metrics.

### 5.3.2 Ablation Study

Table 5 shows a comprehensive ablation study for the efficacy of our database context selection, and example selection methods including DAIL (Gao et al., 2023) and AST. We consistently notice improvement when selecting examples using AST, for the same approximator. Interestingly, the performance gap is increasing the better the approximator becomes, leading to an improvement  $> 2.4\%$  in the case of an oracle approximator. This finding is in agreement with our hypothesis that AST re-ranking can preserve structural information for more precise example selection when  $s' \sim s$ . The inclusion of combined schema and value selection (SVS) leads to further improvements when coupled with example selection based on AST or DAIL.

Approximator	Schema Selection Setup	Recall	Schema Shorten.	EX	EM
<i>Oracle</i>	Gold Query	100.0	71.3	86.3	73.9
FastRAT <sub>ext</sub>	N/A	100.0	0.0	79.3	63.6
FastRAT <sub>ext</sub>	BM25 (top-10)	92.0	36.5	78.9	64.1
FastRAT <sub>ext</sub>	BM25 (top-20)	98.3	14.1	80.7	64.9
FastRAT <sub>ext</sub>	Approx. Query	86.8	71.3	78.4	63.8
FastRAT <sub>ext</sub>	Approx. Query + BM25 (top-7)	93.3	50.4	81.1	65.4
FastRAT <sub>ext</sub>	Approx. Query + BM25 (top-10)	97.0	37.3	81.2	65.6
FastRAT <sub>ext</sub>	Approx. Query + BM25 (dynamic top- <i>k</i> )	97.2	49.0	<b>82.0</b>	<b>65.7</b>
Graphix-T5	N/A	100.0	0.0	79.8	65.6
Graphix-T5	Approx. Query	92.3	71.8	81.8	68.8
Graphix-T5	Approx. Query + BM25 (dynamic top- <i>k</i> )	97.9	49.4	<b>83.0</b>	<b>68.8</b>

Table 4: Recall, Schema Shortening, EX and EM scores (using gpt-3.5-turbo) across different schema selection setups, on the development split of SPIDER. Value selection is enabled for the selected columns across all setups. For the oracle setup, we report performance upper-bounds using *only* the schema elements from the gold query.

Approximator	Selection	EX	EM
N/A	Question Similarity	74.7	52.3
FastRAT <sub>ext</sub>	DAIL	78.6	61.4
	DAIL + SVS	81.3	62.3
	AST	79.3	63.6
	AST + VS	80.4	63.8
	AST + SS	78.9	63.8
	AST + SVS	82.0	65.7
Graphix-T5	DAIL	77.8	61.9
	DAIL + SVS	81.0	63.7
	AST	79.8	65.6
	AST + VS	81.4	66.4
	AST + SS	80.2	66.2
	AST + SVS	83.0	68.8
<i>Oracle</i>	DAIL	79.1	63.2
	DAIL + SVS	82.5	66.0
	AST	81.0	67.6
	AST + VS	82.6	68.1
	AST + SS	82.5	69.6
	AST + SVS	84.6	71.3

Table 5: EX and EM scores on the development set of SPIDER, with gpt-3.5-turbo, across different approximators, and selection setups: example selection (with DAIL or AST), schema selection (SS), value selection (VS), and schema & value selection (SVS). We report results from oracle approximator using gold queries.

## 5.4 Schema Splitting

We evaluate the effect of splitting a schema into  $r_m$  splits, using FastRAT<sub>ext</sub> for schema selection. Figure 1 shows EX scores across different maximum number of columns per schema split ( $r$ ), on the development set of SPIDER. We see that the EX scores of our approach remain consistent across different  $r$ . The performance in the case where particular schemata from the development set are split into  $r_m = 3$  or  $r_m = 4$  splits (i.e. for  $r = 24$  or  $r = 16$  respectively) is identical to the scores

where schemata are split using the default  $r$  with which FastRAT<sub>ext</sub> has been trained.

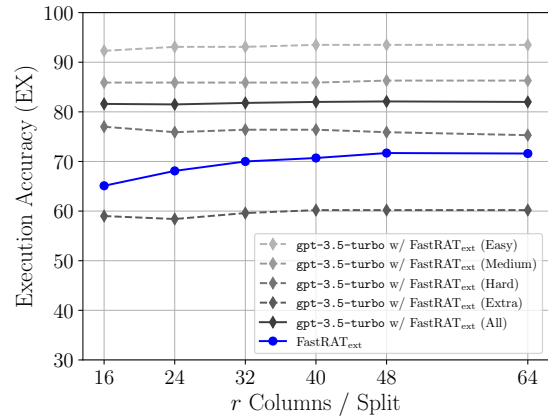


Figure 1: Execution accuracy scores on the development set of SPIDER across different maximum numbers of columns per schema split,  $r$ . The results of our approach, using gpt-3.5-turbo, are presented across different SPIDER-query difficulty levels.

## 6 Discussion

**Are there any theoretical performance upper limits for example selection using AST?** For each data instance in the development set of Spider, we compute the average AST similarity between the approximated query and each SQL query that is included (after example selection) in the corresponding prompt. In Table 6, we measure EX scores on the development set of SPIDER, across different AST-similarity intervals. We see an obvious correlation between execution accuracy and AST scores—execution accuracy is higher for higher AST scores. Besides highlighting an empirical, execution accuracy upper-bound, in the case of test

questions whose SQL structure is well-covered (i.e. with high AST score) in the examples space, our approach can hint data instances that might be challenging for the current configuration, without even requiring to prompt the target LLM or executing the resulting SQL against a database instance. Challenging data instances can be taken into consideration with respect to the existence of insufficient examples in  $\mathbb{X}$  to support the expected SQL structure or a potentially harmful approximator.

AST Interval	Approximator		
	Oracle	Graphix	FastRAT <sub>ext</sub>
[0.0, 1.0]	84.6	83.0	82.0
[0.95, 1.0]	90.8	88.4	88.7
[0.9, 0.95)	80.7	74.1	76.8
[0.85, 0.9)	73.0	68.3	59.4
[0.8, 0.85)	62.4	66.8	63.5
[0.0, 0.8)	50.0	54.1	58.7

Table 6: EX scores on the SPIDER development set using gpt-3.5-turbo, across different average AST-similarity intervals.

**Is the choice of approximator critical?** Although our AST re-ranking and schema selection benefit from more accurate SQL predicted by a stronger approximator, the choice of approximator depends on the desired trade-off between effectiveness and efficiency in practice. FastRAT<sub>ext</sub> is over 600 times faster than Graphix-T5 (Li et al., 2023b) on an A100 80G GPU, while the resulting difference, within our framework, in EX on the SPIDER development set is  $\leq 1\%$  (Table 5).

**Does schema selection improve the performance?** In Table 5, we noted that performing schema selection (i.e. SS) without DB value selection does not necessarily lead to performance improvements. This is in partial agreement with Sun et al. that hard column selection can be harmful for the end-to-end performance, and can be attributed to the drop of recall, when less capable approximators are involved. Nonetheless, as we note in Section 5.3.2, the combination of schema and value selection (SVS) can consistently improve EX and EM, while significantly reducing the LLM token-processing cost due to shortened schema.

## 7 Related Work

Significant number of recent works have looked at how LLMs can be employed in Text-to-SQL scenarios (Rajkumar et al., 2022; Chang and Fosler-

Lussier, 2023; Liu et al., 2023; Pourreza and Rafiei, 2023; Gao et al., 2023; Guo et al., 2024). More recent works have looked at how incorporating examples in the prompt could benefit the performance of LLMs in the end task (Pourreza and Rafiei, 2023; Gao et al., 2023; Guo et al., 2024; Sun et al., 2024). In spite of its underlying benefits, conventional solutions for example selection have focused on retrieving pairs using question similarity (Nan et al., 2023). Other approaches have sought to *approximate* expected SQL queries, and either directly use these approximations in the prompt, in a few-shot setting (Sun et al., 2024) or to filter candidate (question, SQL) pairs by taking into consideration the similarity of their corresponding SQL query skeletons (Li et al., 2023a) against the skeleton of the approximated SQL (Gao et al., 2023). We argue that such example selection strategies can result in information loss, and we propose an approach for re-ranking examples using similarity of normalised SQL ASTs.

The benefits of schema selection for Text-to-SQL have been highlighted across the relevant bibliography (Wang et al., 2020; Li et al., 2023b; Pourreza and Rafiei, 2023). From the LLMs perspective, pruning schema elements from the prompts has been usually leading to performance degradation (Sun et al., 2024). Inspired by Gao et al., we compute a preliminary query for a given  $(q, D)$  by we adapting FastRAT (Vougiouklis et al., 2023), to the requirements of processing longer schemata, in a parallelised manner. We couple the resulting *approximator* with a sparse retriever, and we propose a dynamic strategy for reducing the computational cost of the task while achieving performance improvements.

## 8 Conclusion

In this paper, we augment LLMs for Text-to-SQL semantic parsing by selecting suitable examples and database information. We present a novel AST-based metric to rank examples by similarity of SQL queries. Our hybrid search strategy for schema selection reuses a preliminary query to reduce irrelevant schema elements while maintaining high recall. Extensive experiments demonstrate that our AST-based ranking outperforms previous approaches of example selection and that a symbiotic combination of schema and value selection can further enhance the end-to-end performance of both closed- and open-source LLM solutions.



## 567 Limitations

568 There are limitations with regards to both of our  
569 example selection and schema selection. Our AST-  
570 based ranking can be biased when an approximated  
571 SQL deviates significantly from structurally cor-  
572 rect answers. To address the failure of approxi-  
573 mators, a future direction is to sensibly diversify  
574 selected examples such that LLMs can generalise  
575 compositionally. As for schema selection, our se-  
576 mantic search relies on an approximator which is  
577 essentially a parser with high precision in schema  
578 linking but lack mechanisms to control recall as a  
579 standalone model. Therefore, it is worth extending  
580 cross-encoder architecture such as FastRAT to sup-  
581 port ranking schema elements while being a SQL  
582 approximator in the meantime.

583 We demonstrate that schema splitting strategies  
584 within our framework can be applied across various  
585 numbers of splits without noticeable performance  
586 degradation. Nonetheless, given the lack of avail-  
587 able datasets that incorporate longer commercial  
588 schemata, we focus our experiments on the cross-  
589 database setting provided by CSPIDER and SPIDER  
590 variants.

## 591 Ethics Statement

592 We do not make use of any private, proprietary,  
593 or sensitive data. FastRAT<sub>ext</sub> is trained on pub-  
594 licly available Text-to-SQL datasets, using publicly  
595 available encoder-models as base. Our framework  
596 for retrieval-augmented generation builds on-top  
597 of large, pre-trained language models, which may  
598 have been trained using proprietary data (e.g. in the  
599 case of the OpenAI models). Given the nature of  
600 pre-training schemes, it is possible that our system  
601 could carry forward biases present in the datasets  
602 and/or the involved LLMs.

## 603 References

604 Shengnan An, Bo Zhou, Zeqi Lin, Qiang Fu, Bei Chen,  
605 Nanning Zheng, Weizhu Chen, and Jian-Guang Lou.  
606 2023. [Skill-based few-shot selection for in-context  
607 learning](#). In *Proceedings of the 2023 Conference  
608 on Empirical Methods in Natural Language Process-  
609 ing*, pages 13472–13492, Singapore. Association for  
610 Computational Linguistics.

611 Ruisheng Cao, Lu Chen, Jieyu Li, Hanchong Zhang,  
612 Hongshen Xu, Wangyou Zhang, and Kai Yu. 2023.  
613 [A heterogeneous graph to abstract syntax tree frame-  
614 work for text-to-SQL](#). *IEEE Transactions on Pattern  
615 Analysis and Machine Intelligence*, 45(11):13796–  
616 13813.

Shuaichen Chang and Eric Fosler-Lussier. 2023. [How  
617 to prompt llms for text-to-SQL: A study in zero-shot,  
618 single-domain, and cross-domain settings](#). *Preprint*,  
619 arXiv:2305.11853. 620

Xinyun Chen, Maxwell Lin, Nathanael Schärli, and  
621 Denny Zhou. 2024. [Teaching large language models  
622 to self-debug](#). In *The Twelfth International Confer-  
623 ence on Learning Representations*. 624

Xiang Deng, Ahmed Hassan Awadallah, Christopher  
625 Meek, Oleksandr Polozov, Huan Sun, and Matthew  
626 Richardson. 2021. [Structure-grounded pretraining  
627 for text-to-SQL](#). In *Proceedings of the 2021 Con-  
628 ference of the North American Chapter of the Asso-  
629 ciation for Computational Linguistics: Human Lan-  
630 guage Technologies*, pages 1337–1350, Online. As-  
631 sociation for Computational Linguistics. 632

Xuemei Dong, Chao Zhang, Yuhang Ge, Yuren Mao,  
633 Yunjun Gao, lu Chen, Jinshu Lin, and Dongfang Lou.  
634 2023. [C3: Zero-shot text-to-SQL with ChatGPT](#).  
635 *Preprint*, arXiv:2307.07306. 636

Beat Fluri, Michael Wursch, Martin Plnzer, and Harald  
637 Gall. 2007. [Change distilling:tree differencing for  
638 fine-grained source code change extraction](#). *IEEE  
639 Transactions on Software Engineering*, 33(11):725–  
640 743. 641

Yujian Gan, Xinyun Chen, Qiuping Huang, Matthew  
642 Purver, John R. Woodward, Jinxia Xie, and Peng-  
643 sheng Huang. 2021a. [Towards robustness of text-  
644 to-SQL models against synonym substitution](#). In  
645 *Proceedings of the 59th Annual Meeting of the Asso-  
646 ciation for Computational Linguistics and the 11th  
647 International Joint Conference on Natural Language  
648 Processing (Volume 1: Long Papers)*, pages 2505–  
649 2515, Online. Association for Computational Lin-  
650 guistics. 651

Yujian Gan, Xinyun Chen, and Matthew Purver. 2021b.  
652 [Exploring underexplored limitations of cross-domain  
653 text-to-SQL generalization](#). In *Proceedings of the  
654 2021 Conference on Empirical Methods in Natural  
655 Language Processing*, pages 8926–8931, Online and  
656 Punta Cana, Dominican Republic. Association for  
657 Computational Linguistics. 658

Dawei Gao, Haibin Wang, Yaliang Li, Xiuyu Sun,  
659 Yichen Qian, Bolin Ding, and Jingren Zhou. 2023.  
660 [Text-to-SQL empowered by large language models:  
661 A benchmark evaluation](#). *CoRR*, abs/2308.15363. 662

Chunxi Guo, Zhiliang Tian, Jintao Tang, Shasha  
663 Li, Zhihua Wen, Kaixuan Wang, and Ting Wang.  
664 2024. [Retrieval-augmented GPT-3.5-based text-  
665 to-SQL framework with sample-aware prompting  
666 and dynamic revision chain](#). In *Neural Information  
667 Processing*, pages 341–356, Singapore. Springer Na-  
668 ture Singapore. 669

Haoyang Li, Jing Zhang, Cuiping Li, and Hong Chen.  
670 2023a. [RESDSL: Decoupling schema linking  
671 and skeleton parsing for text-to-sql](#). *Proceedings  
672 of the AAAI Conference on Artificial Intelligence*,  
673 37(11):13067–13075. 674

675	Jinyang Li, Binyuan Hui, Reynold Cheng, Bowen Qin,	Jin, Denny Jin, Shino Jomoto, Billie Jonn, Hee-	736
676	Chenhao Ma, Nan Huo, Fei Huang, Wenyu Du, Luo	woo Jun, Tomer Kaftan, Łukasz Kaiser, Ali Ka-	737
677	Si, and Yongbin Li. 2023b. <a href="#">Graphix-T5: mixing</a>	mali, Ingmar Kanitscheider, Nitish Shirish Keskar,	738
678	<a href="#">pre-trained transformers with graph-aware layers for</a>	Tabarak Khan, Logan Kilpatrick, Jong Wook Kim,	739
679	<a href="#">text-to-SQL parsing</a> . In <i>Proceedings of the Thirty-</i>	Christina Kim, Yongjik Kim, Jan Hendrik Kirchner,	740
680	<i>Seventh AAAI Conference on Artificial Intelligence</i>	Jamie Kiros, Matt Knight, Daniel Kokotajlo,	741
681	<i>and Thirty-Fifth Conference on Innovative Applica-</i>	Łukasz Kondraciuk, Andrew Kondrich, Aris Kon-	742
682	<i>tions of Artificial Intelligence and Thirteenth Sympo-</i>	stantinidis, Kyle Kopic, Gretchen Krueger, Vishal	743
683	<i>sium on Educational Advances in Artificial Intelli-</i>	Kuo, Michael Lampe, Ikai Lan, Teddy Lee, Jan	744
684	<i>gence</i> , AAAI’23/IAAI’23/EAAI’23. AAAI Press.	Leike, Jade Leung, Daniel Levy, Chak Ming Li,	745
685	Aiwei Liu, Xuming Hu, Lijie Wen, and Philip S.	Rachel Lim, Molly Lin, Stephanie Lin, Mateusz	746
686	Yu. 2023. <a href="#">A comprehensive evaluation of Chat-</a>	Litwin, Theresa Lopez, Ryan Lowe, Patricia Lue,	747
687	<a href="#">GPT’s zero-shot text-to-SQL capability</a> . <i>Preprint</i> ,	Anna Makanju, Kim Malfacini, Sam Manning, Todor	748
688	arXiv:2303.13547.	Markov, Yaniv Markovski, Bianca Martin, Katie	749
689	Qingkai Min, Yuefeng Shi, and Yue Zhang. 2019. <a href="#">A</a>	Mayer, Andrew Mayne, Bob McGrew, Scott Mayer	750
690	<a href="#">pilot study for Chinese SQL semantic parsing</a> . In	McKinney, Christine McLeavey, Paul McMillan,	751
691	<i>Proceedings of the 2019 Conference on Empirical</i>	Jake McNeil, David Medina, Aalok Mehta, Jacob	752
692	<i>Methods in Natural Language Processing and the</i>	Menick, Luke Metz, Andrey Mishchenko, Pamela	753
693	<i>9th International Joint Conference on Natural Lan-</i>	Mishkin, Vinnie Monaco, Evan Morikawa, Daniel	754
694	<i>guage Processing (EMNLP-IJCNLP)</i> , pages 3652–	Mossing, Tong Mu, Mira Murati, Oleg Murk, David	755
695	3658, Hong Kong, China. Association for Computa-	Mély, Ashvin Nair, Reiichiro Nakano, Rajeev Nayak,	756
696	tional Linguistics.	Arvind Neelakantan, Richard Ngo, Hyeonwoo Noh,	757
697	Linyong Nan, Yilun Zhao, Weijin Zou, Narutatsu	Long Ouyang, Cullen O’Keefe, Jakub Pachocki, Alex	758
698	Ri, Jaesung Tae, Ellen Zhang, Arman Cohan, and	Paino, Joe Palermo, Ashley Pantuliano, Giambat-	759
699	Dragomir Radev. 2023. <a href="#">Enhancing text-to-SQL capa-</a>	tista Parascandolo, Joel Parish, Emy Parparita, Alex	760
700	<a href="#">bilities of large language models: A study on prompt</a>	Passos, Mikhail Pavlov, Andrew Peng, Adam Perel-	761
701	<a href="#">design strategies</a> . In <i>Findings of the Association</i>	man, Filipe de Avila Belbute Peres, Michael Petrov,	762
702	<i>for Computational Linguistics: EMNLP 2023</i> , pages	Henrique Ponde de Oliveira Pinto, Michael, Poko-	763
703	14935–14956, Singapore. Association for Computa-	rny, Michelle Pokrass, Vitchy H. Pong, Tolly Pow-	764
704	tional Linguistics.	ell, Alethea Power, Boris Power, Elizabeth Proehl,	765
705	OpenAI, Josh Achiam, Steven Adler, Sandhini Agarwal,	Raul Puri, Alec Radford, Jack Rae, Aditya Ramesh,	766
706	Lama Ahmad, Ilge Akkaya, Florencia Leoni Ale-	Cameron Raymond, Francis Real, Kendra Rimbach,	767
707	man, Diogo Almeida, Janko Alvenschmidt, Sam Alt-	Carl Ross, Bob Rotsted, Henri Roussez, Nick Ry-	768
708	man, Shyamal Anadkat, Red Avila, Igor Babuschkin,	der, Mario Saltarelli, Ted Sanders, Shibani Santurkar,	769
709	Suchir Balaji, Valerie Balcom, Paul Baltescu, Haim-	Girish Sastry, Heather Schmidt, David Schnurr, John	770
710	ing Bao, Mohammad Bavarian, Jeff Belgum, Ir-	Schulman, Daniel Selsam, Kyla Sheppard, Toki	771
711	wan Bello, Jake Berdine, Gabriel Bernadett-Shapiro,	Sherbakov, Jessica Shieh, Sarah Shoker, Pranav	772
712	Christopher Berner, Lenny Bogdonoff, Oleg Boiko,	Shyam, Szymon Sidor, Eric Sigler, Maddie Simens,	773
713	Madelaine Boyd, Anna-Luisa Brakman, Greg Brock-	Jordan Sitkin, Katarina Slama, Ian Sohl, Benjamin	774
714	man, Tim Brooks, Miles Brundage, Kevin Button,	Sokolowsky, Yang Song, Natalie Staudacher, Fe-	775
715	Trevor Cai, Rosie Campbell, Andrew Cann, Brittany	lipe Petroski Such, Natalie Summers, Ilya Sutskever,	776
716	Carey, Chelsea Carlson, Rory Carmichael, Brooke	Jie Tang, Nikolas Tezak, Madeleine B. Thompson,	777
717	Chan, Che Chang, Fotis Chantzis, Derek Chen, Sully	Phil Tillet, Amin Tootoonchian, Elizabeth Tseng,	778
718	Chen, Ruby Chen, Jason Chen, Mark Chen, Ben	Preston Tuggle, Nick Turley, Jerry Tworek, Juan Fe-	779
719	Chess, Chester Cho, Casey Chu, Hyung Won Chung,	lipe Cerón Uribe, Andrea Vallone, Arun Vijayvergiya,	780
720	Dave Cummings, Jeremiah Currier, Yunxing Dai,	Chelsea Voss, Carroll Wainwright, Justin Jay Wang,	781
721	Cory Decareaux, Thomas Degry, Noah Deutsch,	Alvin Wang, Ben Wang, Jonathan Ward, Jason Wei,	782
722	Damien Deville, Arka Dhar, David Dohan, Steve	CJ Weinmann, Akila Welihinda, Peter Welinder, Ji-	783
723	Dowling, Sheila Dunning, Adrien Ecoffet, Atty Eleti,	ayi Weng, Lilian Weng, Matt Wiethoff, Dave Willner,	784
724	Tyna Eloundou, David Farhi, Liam Fedus, Niko Felix,	Clemens Winter, Samuel Wolrich, Hannah Wong,	785
725	Simón Posada Fishman, Juston Forte, Isabella Ful-	Lauren Workman, Sherwin Wu, Jeff Wu, Michael	786
726	ford, Leo Gao, Elie Georges, Christian Gibson, Vik	Wu, Kai Xiao, Tao Xu, Sarah Yoo, Kevin Yu,	787
727	Goel, Tarun Gogineni, Gabriel Goh, Rapha Gontijo-	Qiming Yuan, Wojciech Zaremba, Rowan Zellers,	788
728	Lopes, Jonathan Gordon, Morgan Grafstein, Scott	Chong Zhang, Marvin Zhang, Shengjia Zhao, Tian-	789
729	Gray, Ryan Greene, Joshua Gross, Shixiang Shane	hao Zheng, Juntang Zhuang, William Zhuk, and Bar-	790
730	Gu, Yufei Guo, Chris Hallacy, Jesse Han, Jeff Harris,	ret Zoph. 2024. <a href="#">GPT-4 technical report</a> . <i>Preprint</i> ,	791
731	Yuchen He, Mike Heaton, Johannes Heidecke, Chris	arXiv:2303.08774.	792
732	Hesse, Alan Hickey, Wade Hickey, Peter Hoeschele,	Mohammadreza Pourreza and Davood Rafiei. 2023.	793
733	Brandon Houghton, Kenny Hsu, Shengli Hu, Xin	<a href="#">DIN-SQL: Decomposed in-context learning</a>	794
734	Hu, Joost Huizinga, Shantanu Jain, Shawn Jain,	<a href="#">of text-to-SQL with self-correction</a> . <i>Preprint</i> ,	795
735	Joanne Jang, Angela Jiang, Roger Jiang, Haozhun	arXiv:2304.11015.	796
		Colin Raffel, Noam Shazeer, Adam Roberts, Katherine	797

798	Lee, Sharan Narang, Michael Matena, Yanqi Zhou, Wei Li, and Peter J. Liu. 2020. Exploring the limits of transfer learning with a unified text-to-text transformer. <i>J. Mach. Learn. Res.</i> , 21(1).	
799		
800		
801		
802	Nitarshan Rajkumar, Raymond Li, and Dzmitry Bahdanau. 2022. <a href="#">Evaluating the text-to-SQL capabilities of large language models</a> . <i>Preprint</i> , arXiv:2204.00498.	
803		
804		
805		
806	Nils Reimers and Iryna Gurevych. 2019. <a href="#">SentenceBERT: Sentence embeddings using Siamese BERT-networks</a> . In <i>Proceedings of the 2019 Conference on Empirical Methods in Natural Language Processing and the 9th International Joint Conference on Natural Language Processing (EMNLP-IJCNLP)</i> , pages 3982–3992, Hong Kong, China. Association for Computational Linguistics.	
807		
808		
809		
810		
811		
812		
813		
814	Shuo Ren, Daya Guo, Shuai Lu, Long Zhou, Shujie Liu, Duyu Tang, Neel Sundaresan, Ming Zhou, Ambrosio Blanco, and Shuai Ma. 2020. <a href="#">CodeBLEU: a method for automatic evaluation of code synthesis</a> . <i>Preprint</i> , arXiv:2009.10297.	
815		
816		
817		
818		
819	Stephen E. Robertson, Steve Walker, Susan Jones, Micheline Hancock-Beaulieu, and Mike Gattford. 1994. <a href="#">Okapi at TREC-3</a> . In <i>Proceedings of The Third Text REtrieval Conference, TREC 1994, Gaithersburg, Maryland, USA, November 2-4, 1994</i> , volume 500-225 of <i>NIST Special Publication</i> , pages 109–126. National Institute of Standards and Technology (NIST).	
820		
821		
822		
823		
824		
825		
826		
827	Kaitao Song, Xu Tan, Tao Qin, Jianfeng Lu, and Tie-Yan Liu. 2020. <a href="#">MPNet: Masked and permuted pre-training for language understanding</a> . In <i>Advances in Neural Information Processing Systems</i> , volume 33, pages 16857–16867. Curran Associates, Inc.	
828		
829		
830		
831		
832	Ruoxi Sun, Sercan Ö. Arik, Alex Muzio, Lesly Miculicich, Satya Gundabathula, Pengcheng Yin, Hanjun Dai, Hootan Nakhost, Rajarishi Sinha, Zifeng Wang, and Tomas Pfister. 2024. <a href="#">SQL-PaLM: Improved large language model adaptation for text-to-sql (extended)</a> . <i>Preprint</i> , arXiv:2306.00739.	
833		
834		
835		
836		
837		
838	Ngoc Tran, Hieu Tran, Son Nguyen, Hoan Nguyen, and Tien Nguyen. 2019. <a href="#">Does bleu score work for code migration?</a> In <i>2019 IEEE/ACM 27th International Conference on Program Comprehension (ICPC)</i> , pages 165–176.	
839		
840		
841		
842		
843	Pavlos Vougiouklis, Nikos Papisarantopoulos, Danna Zheng, David Tuckey, Chenxin Diao, Zhili Shen, and Jeff Pan. 2023. <a href="#">FastRAT: Fast and efficient cross-lingual text-to-SQL semantic parsing</a> . In <i>Proceedings of the 13th International Joint Conference on Natural Language Processing and the 3rd Conference of the Asia-Pacific Chapter of the Association for Computational Linguistics (Volume 1: Long Papers)</i> , pages 564–576, Nusa Dua, Bali. Association for Computational Linguistics.	
844		
845		
846		
847		
848		
849		
850		
851		
852		
	Bailin Wang, Mirella Lapata, and Ivan Titov. 2021. <a href="#">Meta-learning for domain generalization in semantic parsing</a> . In <i>Proceedings of the 2021 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies</i> , pages 366–379, Online. Association for Computational Linguistics.	853
		854
		855
		856
		857
		858
		859
	Bailin Wang, Richard Shin, Xiaodong Liu, Oleksandr Polozov, and Matthew Richardson. 2020. <a href="#">RAT-SQL: Relation-aware schema encoding and linking for text-to-SQL parsers</a> . In <i>Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics</i> , pages 7567–7578, Online. Association for Computational Linguistics.	860
		861
		862
		863
		864
		865
		866
	Junjie Ye, Xuanting Chen, Nuo Xu, Can Zu, Zekai Shao, Shichun Liu, Yuhan Cui, Zeyang Zhou, Chao Gong, Yang Shen, Jie Zhou, Siming Chen, Tao Gui, Qi Zhang, and Xuanjing Huang. 2023. <a href="#">A comprehensive capability analysis of GPT-3 and GPT-3.5 series models</a> . <i>Preprint</i> , arXiv:2303.10420.	867
		868
		869
		870
		871
		872
	Tao Yu, Chien-Sheng Wu, Xi Victoria Lin, Bailin Wang, Yi Chern Tan, Xinyi Yang, Dragomir Radev, Richard Socher, and Caiming Xiong. 2021. <a href="#">GraPPa: Grammar-augmented pre-training for table semantic parsing</a> . In <i>International Conference on Learning Representations</i> .	873
		874
		875
		876
		877
		878
	Tao Yu, Rui Zhang, Kai Yang, Michihiro Yasunaga, Dongxu Wang, Zifan Li, James Ma, Irene Li, Qingning Yao, Shanelle Roman, Zilin Zhang, and Dragomir Radev. 2018. <a href="#">Spider: A large-scale human-labeled dataset for complex and cross-domain semantic parsing and text-to-SQL task</a> . In <i>Proceedings of the 2018 Conference on Empirical Methods in Natural Language Processing</i> , pages 3911–3921, Brussels, Belgium. Association for Computational Linguistics.	879
		880
		881
		882
		883
		884
		885
		886
		887
	John M. Zelle and Raymond J. Mooney. 1996. Learning to parse database queries using inductive logic programming. In <i>Proceedings of the Thirteenth National Conference on Artificial Intelligence - Volume 2, AAAI’96</i> , page 1050–1055, Portland, Oregon. AAAI Press.	888
		889
		890
		891
		892
		893
	<b>A SQL Similarity using Normalised Abstract Syntax Trees</b>	894
		895
	Table 7 shows the corresponding SQL queries after each step of our AST normalisation as explained in Section 3. An example of the similarity between normalised ASTs is provided in Figure 2, where tables, columns and values are masked out for cross-domain settings.	896
		897
		898
		899
		900
		901
	<b>B Prompt Formulation</b>	902
	Table 8 shows an example of our prompt, after example and DB context selection (i.e. schema and value selection). This prompt is provided as input	903
		904
		905

<b>SQL</b>	<code>SELECT T1.Category, COUNT(*) AS Num FROM Products AS T1 JOIN Orders AS T2 ON T1.id = T2.pid GROUP BY T1.Category ORDER BY Num ASC</code>
<b>1. Unify Identifiers</b>	<code>SELECT t1.category, COUNT(*) AS num FROM products AS t1 JOIN orders AS t2 ON t1.id = t2.pid GROUP BY t1.category ORDER BY num ASC</code>
<b>2. Resolve Aliases</b>	<code>SELECT products.category, COUNT(*) FROM products JOIN orders ON products.id = orders.pid GROUP BY products.category ORDER BY COUNT(*) ASC</code>
<b>3. Reorder JOIN</b>	<code>SELECT products.category, COUNT(*) FROM orders JOIN products ON orders.id = products.pid GROUP BY products.category ORDER BY COUNT(*) ASC</code>
<b>4. Mask Identifiers &amp; Values (cross-domain only)</b>	<code>SELECT _, COUNT(*) FROM _ JOIN _ ON _ = _ GROUP BY _ ORDER BY COUNT(*) ASC</code>

Table 7: An example of the effects to corresponding SQL after each step of our AST normalisation.

to LLMs. Following the latest OpenAI example<sup>6</sup> for Text-to-SQL parsing, we represent a schema with CREATE TABLE statements in SQL. Semantic names or descriptions of tables and columns are included as COMMENT along with the corresponding columns or tables. Note that we filter out comments that can be obtained by simply lowercasing original names and/or removing underscores. To maintain a compact representation of database information, we append selected values of columns into their COMMENT rather than introducing additional lines as in the work by Chen et al. (2024). Example (question, SQL) pairs are provided in a similar manner to DAIL-SQL (Gao et al., 2023), followed by an instruction to prompt LLMs to generate SQL for the test question.

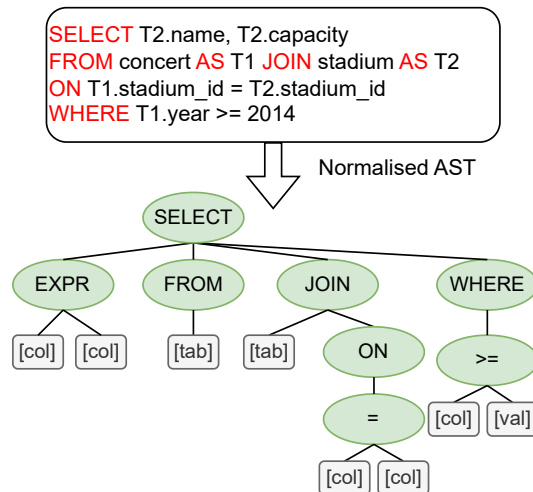
## C Implementation Details

We use this section to provide further details about the implementation of our approach.

### C.1 Example Selection

Following (Gao et al., 2023), we employ the pre-trained all-mpnet-base-v2 model (Song et al., 2020) from Sentence Transformer (Reimers

<sup>6</sup><https://platform.openai.com/examples/default-sql-translate>



Similarity: 0.52

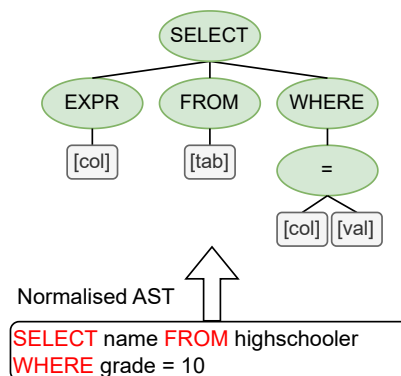


Figure 2: Example of how the similarity between two different SQL queries is computed using normalised ASTs.

and Gurevych, 2019) to compute dense question embeddings for English datasets including SPIDER, SPIDER-DK, SPIDER-REAL, and SPIDER-SYN. For CSPIDER, paraphrase-multilingual-MiniLM-L12-v2 is used instead. SQL queries are parsed into AST by using SQLGlot<sup>7</sup> and are then normalised as explained in Section 3. SQLGlot provides an implementation of the Change Distilling algorithm for AST differencing. We refer readers to SQLGlot’s documentation<sup>8</sup> for more details. For selecting example question-SQL pairs, we first retrieve top 500 examples by question similarity and rerank them in terms of the similarity of normalised SQL ASTs. For relevant experiments in Table 5, we reproduced the implementation of

<sup>7</sup><https://github.com/tobymao/sqlglot>

<sup>8</sup>[https://github.com/tobymao/sqlglot/blob/main/posts/sql\\_diff.md](https://github.com/tobymao/sqlglot/blob/main/posts/sql_diff.md)

	# Given SQLite database schema student_transcripts: CREATE TABLE Departments( department_id number, department_name text COMMENT 'department name (e.g. engineer, statistics, medical) ', ...);
Selected	CREATE TABLE Degree_Programs( degree_program_id number, degree_summary_name text COMMENT 'degree summary name (e.g. PHD, Master, Bachelor) ', ... PRIMARY KEY (degree_program_id), FOREIGN KEY (department_id) REFERENCES Departments(department_id));
Schema	# Your task is to translate Question into SQL. # Some examples are provided based on similar problems:
w/	Question: How many courses does the department of Computer Information Systems offer? SQL: SELECT count(*) FROM department AS T1 JOIN course AS T2 ON T1.dept_code = T2.dept_code WHERE dept_name = "Computer Info.Systems"
Selected	Question: ...
Examples	SQL: ...
	# Complete the following SQL for schema stu- dent_transcripts: Question: How many degrees does the engineering de- partment have? SQL:
Test	
Question	

Table 8: An example of the resulting prompt, after example and schema and DB content selection.

DAIL selection from the original paper (Gao et al., 2023). The number of few-shot examples is set to 5 across all experiments.

## C.2 Schema & Value Selection

Each database schema is treated as an independent collection of columns that are analogous to documents to be retrieved by using BM25. As mentioned in Section 4.1, we represent a column by concatenating semantic names of both the column and its table, and the column values in the database. Semantic names and values are tokenized using spaCy<sup>9</sup> and preprocessed by lowercasing and stemming<sup>10</sup>. At inference time, the same processing is applied to questions. We adopt the implementation of Okapi BM25 (Robertson et al., 1994) from Rank-BM25<sup>11</sup>. The number of columns to retrieve is dynamically set to  $\lfloor 1.5 \times \gamma \rfloor$  where  $\gamma$  is the number of unique columns in an approximated query. We limit the resulting number to a range between 6 and 20. By retrieving at column level, a table is selected if any of its columns are selected. We merge retrieved schema elements with schema elements from the approximated query to construct a sub-schema. To further increase the recall, we add additional primary keys and foreign keys that

<sup>9</sup><https://github.com/explosion/spaCy>

<sup>10</sup><https://www.nltk.org/api/nltk.stem.porter.html>

<sup>11</sup>[https://github.com/dorianbrown/rank\\_bm25](https://github.com/dorianbrown/rank_bm25)

are not selected but valid based on selected tables, except for experiments where only approximated queries are used (see Table 4). In such cases, however, if the SQL query involves only tables (e.g. SELECT \* FROM books), primary keys of selected tables are still included to ensure that corresponding CREATE TABLE statements (see Table 8) are meaningful and consistent.

For selecting values, similarly, we match the input question and the set of values for each (selected) column that has a non-numeric type. The top 3 results are added to the prompt as exemplified in Table 8. The same setting of schema and value selection is used for all datasets we experimented with except CSPIDER. Due to the cross-lingual nature of CSPIDER, schema selection and value selection are simply disabled.

**Training FastRAT<sub>ext</sub>** We follow the original hyper-parameters provided by (Vougiouklis et al., 2023) for training FastRAT<sub>ext</sub>. The monolingual version of FastRAT is based on BERT<sub>LARGE</sub> while its cross-lingual variant on XLM-RoBERTa-large.

## C.3 OpenAI Models

We use gpt-4 (gpt-4-0613) and gpt-3.5-turbo (gpt-3.5-turbo-0613) for our experiments. For decoding, sampling is disabled and the maximum number of tokens to generate is set to 256. A single experiment, on the SPIDER development set using our approach with FastRAT<sub>ext</sub> as the approximator and dynamic database context selection costs around \$0.8 and \$16.5 in the case of gpt-3.5-turbo-0613 and gpt-4-0613 respectively.

## C.4 Experiments with Open-Source LLMs

We further conduct experiments with open-source models from the DeepSeek family<sup>12</sup>, that specialise in code generation. Prompting and decoding setups remain consistent across all LLMs (cf. Section C of the Appendix). Table 9 summarises the results. We see that our approach can generalise even in the case of open-source LLM alternatives. Interestingly, our scores using deepseek-coder-33b-instruct are comparable to the scores when using gpt-3.5-turbo-0613 across all approximators. Inference experiments are conducted on a machine using 8×NVIDIA-V100 32G GPUs.

<sup>12</sup>We use the implementations provided by <https://huggingface.co/deepseek-ai>.

Model	SPIDER		SPIDER-DK		SPIDER-REAL		SPIDER-SYN		CSPIDER
	EX	EM	EX	EM	EX	EM	EX	EM	EM
deepseek-coder-6.7b-instruct									
Ours (w/ FastRAT <sub>ext</sub> )	78.6	64.8	66.7	46.4	73.2	55.7	66.0	49.5	53.0
Ours (w/ Graphix-T5)	79.5	64.8	—	—	—	—	—	—	—
deepseek-coder-33b-instruct									
Ours (w/ FastRAT <sub>ext</sub> )	81.5	62.1	70.5	46.4	77.4	59.3	68.7	49.5	55.9
Ours (w/ Graphix-T5)	83.4	64.7	—	—	—	—	—	—	—

Table 9: Execution (EX) and exact match (EM) accuracy scores of our approach using DeepSeek family models, on the development splits of SPIDER and CSPIDER, and the SPIDER-DK, SPIDER-REAL and SPIDER-SYN test splits. CSPIDER results are using only FastRAT<sub>ext</sub> as approximator.

## D SPIDER and CSPIDER Experiments

We report experiments on CSPIDER (Min et al., 2019) and SPIDER (Yu et al., 2018), which contain database schema information and examples in Chinese and English respectively. Since CSPIDER is a translated version of the SPIDER dataset, the characteristics of the two with respect to structure and number of examples are identical. Both datasets contain 8,659 examples of questions and SQL queries along with their relevant SQL schemata (i.e. 146 unique databases). The development and test<sup>13</sup> sets consist of 1,034, on 20 unique databases and 2,147, on 40 unique databases, respectively, and none of the relevant databases are seen in the training set. Due to the scarcity of works reporting test scores on these benchmarks, we chose not to include our results in the main body of our manuscript. Table 10 shows the performance of our framework with respect to execution and exact match accuracy scores on the test splits of SPIDER and CSPIDER.

## E Baselines

We compare the performance of our approach against several baselines. We dichotomize the landscape of baselines in fine-tuning- and prompting-based baselines.

**Fine-tuning-based** **GraPPa** uses synthetic data constructed via induced synchronous context-free grammar for pre-training an MLM on the SSP-label classification; **DG-MAML** applies meta-learning targeting zero-shot domain generalization. **FastRAT** incorporates a decoder-free framework, by directly predicting SQL queries from SSP labels; **Graphix-T5** inserts a graph-aware layer into T5

(Raffel et al., 2020) to introduce structural inductive bias; **RESDSQL** decouples schema linking and SQL skeleton parsing using a framework based on a ranking-enhanced encoder and skeleton-aware decoder; **HG2AST** proposes a framework to integrate dedicated structure knowledge by transforming heterogeneous graphs to abstract syntax trees.

**Prompting-based** Zero-shot prompting with LLMs has been explored by Guo et al.; Liu et al.; Pourreza and Rafiei; **C3** introduces calibration bias prompting to alleviate LLMs’ biases; **DIN-SQL** uses chain-of-thought prompting with pre-defined prompting templates tailored for the assessed question hardness; **DAIL-SQL** uses query approximation and SQL skeleton-based similarities for example selection; **SQL-PaLM** proposes a framework for *soft* column selection and execution-based refinement; **RAG w/ Rev. Chain** augments the input prompt with question skeleton-based example retrieval and an execution-based revision chain.

<sup>13</sup>Since the 1st of March 2024, the test sets of both Spider and CSpider have become publicly available.

Model	Easy		Medium		Hard		Extra		All	
	EX	EM	EX	EM	EX	EM	EX	EM	EX	EM
<b>SPIDER</b>										
FastRAT <sub>ext</sub>	86.2	81.3	72.2	66.0	60.0	51.6	48.5	33.9	68.7	60.9
deepseek-coder-33b-instruct										
Ours (w/ Graphix-T5)	89.6	85.5	88.6	66.4	73.9	50.1	58.8	28.0	80.7	60.7
gpt-4										
Ours (w/ Graphix-T5)	91.9	87.4	90.3	80.4	81.2	66.1	74.2	47.6	86.0	73.4
<b>CSPIDER</b>										
FastRAT <sub>ext</sub>	–	67.2	–	49.9	–	41.5	–	12.9	–	45.5
deepseek-coder-33b-instruct										
Ours (w/ FastRAT <sub>ext</sub> )	–	79.7	–	58.2	–	40.2		22.4	–	52.9
gpt-4										
Ours (w/ FastRAT <sub>ext</sub> )	–	81.2	–	67.7	–	53.2	–	29.9	–	61.1

Table 10: Execution (EX) and exact match (EM) accuracy scores of our framework, on the test splits of SPIDER and CSPIDER.