

# HIE-SQL: History Information Enhanced Network for Context-Dependent Text-to-SQL Semantic Parsing

Anonymous ACL submission

## Abstract

Recently, context-dependent text-to-SQL semantic parsing which translates natural language into SQL in an interaction process has attracted a lot of attention. Previous works leverage context-dependence information either from interaction history utterances or the previous predicted SQL queries but fail in taking advantage of both since of the mismatch between natural language and logic-form SQL. In this work, we propose a **History Information Enhanced text-to-SQL model (HIE-SQL)** to exploit context-dependence information from both history utterances and the last predicted SQL query. In view of the mismatch, we treat natural language and SQL as two modalities and propose a bimodal pre-trained model to bridge the gap between them. Besides, we design a schema-linking graph to enhance connections from utterances and the SQL query to the database schema. We show our history information enhanced methods improve the performance of HIE-SQL by a significant margin, which achieves new state-of-the-art results on the two context-dependent text-to-SQL benchmarks, the SparC and CoSQL datasets, at the writing time.

## 1 Introduction

Conversation user interfaces to databases have launched a new research hotspot in Text-to-SQL semantic parsing (Zhang et al., 2019; Guo et al., 2019; Wang et al., 2020; Lin et al., 2020; Xu et al., 2021; Cao et al., 2021; Hui et al., 2021; Yu et al., 2021b) and benefited us in industry (Dhamdhere et al., 2017; Weir et al., 2020). Most previous works focus on the context-independent text-to-SQL task and propose many competitive models. Some models (Wang et al., 2020; Scholak et al., 2021) even surprisingly work well on the context-dependent text-to-SQL task by just appending the interaction history utterances to the input. Especially, PICARD (Scholak et al., 2021) achieves state-of-the-art performances both in Spider (Yu et al., 2018b),

$U_1$ : List the name of the teachers and the courses assigned for them to teach.
$S_1$ : <code>SELECT T3.Name, T2.Course FROM course_arrange AS T1 JOIN course AS T2 ON T1.Course_ID = T2.Course_ID JOIN teacher AS T3 ON T1.Teacher_ID = T3.Teacher_ID</code>
$U_2$ : Arrange this list with the teachers name in ascending order.
$S_2$ : <code>SELECT T3.Name, T2.Course FROM course_arrange AS T1 JOIN course AS T2 ON T1.Course_ID = T2.Course_ID JOIN teacher AS T3 ON T1.Teacher_ID = T3.Teacher_ID ORDER BY T3.Name</code>
$U_3$ : Include teachers id in the same list.
$S_3$ : <code>SELECT T3.Name, T2.Course, T1.teacher_ID FROM course_arrange AS T1 JOIN course AS T2 ON T1.Course_ID = T2.Course_ID JOIN teacher AS T3 ON T1.Teacher_ID = T3.Teacher_ID ORDER BY T3.Name</code>

Figure 1: An example of context-dependent text-to-SQL interaction in CoSQL where  $U_i$  is the utterance of turn  $i$  and  $S_i$  is the corresponding SQL query for  $U_i$ . The tokens with red color are the history information that should be considered in later predictions. It is context-independent if we just consider the prediction of  $S_1$ .

a cross-domain context-independent text-to-SQL benchmark, and CoSQL (Yu et al., 2019a), a cross-domain context-dependent text-to-SQL benchmark, before our work. However, every coin has two sides. That implies underachievement of the exploration of context information in context-dependent text-to-SQL semantic parsing.

Compared with context-independent text-to-SQL semantic parsing, context-dependent text-to-SQL semantic parsing are more challenging since of the various types of dependence in utterances which make models vulnerable to parsing errors. As R<sup>2</sup>SQL (Hui et al., 2021) considers, different context dependencies between two adjacent utterances require the model to establish dynamic connections between utterances and database schema carefully. However, context information is not only from the last utterance. Long-range dependence is also the case in CoSQL as the prediction of  $S_3$  depends on "the name of the teachers and the courses" in  $U_1$  in Figure 1. A workable proposition for long-range dependence is to inherit context information

065 from previous predicted SQL queries. But it is  
066 not a piece of cake to take advantage of previously  
067 predicted queries since of the mismatch between  
068 natural language and logic-form SQL. As Liu et al.  
069 (2020) conclude, roughly encoding the last pre-  
070 dicted SQL query and utterances takes the wooden  
071 spoon while easily concatenation of interaction his-  
072 tory utterances and current utterance appears to  
073 be strikingly competitive in their evaluation of 13  
074 existing context modeling methods.

075 In this paper, we propose a history information  
076 enhanced network to make full use of both history  
077 interactive utterances and previous predicted SQL  
078 queries. We first encode the last predicted SQL  
079 query by treating the logic-form query as another  
080 modality with natural language. We present SQL-  
081 BERT, a bimodal pre-trained model for SQL and  
082 natural language which is able to capture the se-  
083 mantic connection and bridge the gap between SQL  
084 and natural language. It produces general-purpose  
085 representations and supports our context-dependent  
086 text-to-SQL semantic parsing. As adopted in a ma-  
087 jority of large pre-trained models, we develop SQL-  
088 BERT with the multi-layer Transformer (Vaswani  
089 et al., 2017). We pre-train it with the objective func-  
090 tion of masked language modeling (MLM) on SQL.

091 Besides, we propose a history information en-  
092 hanced schema-linking graph to represent the rela-  
093 tions among current utterance, interaction history  
094 utterances, the last predicted query, and correspond-  
095 ing database schema. Considering it is weird to  
096 shift a topic back and forth in an interaction, we  
097 assume that the long-range dependence is succes-  
098 sive. For example, that  $S_3$  depends on  $U_1$  implies  
099 that  $S_2$  does too in Figure 1. In that case, we can  
100 leverage the long-range dependence from the last  
101 predicted query. In addition, the corresponding  
102 SQL queries of adjacent utterances tend to over-  
103 lap (Zhang et al., 2019). Therefore, unlike the  
104 previous schema-linking graph just with utterances  
105 and database schema (Hui et al., 2021), the last  
106 predicted query takes part in our graph. Besides,  
107 we distinguish current utterance and interaction his-  
108 tory utterances in the schema-linking graph. We  
109 encode the schema-linking relations with Relative  
110 Self-Attention Mechanism (Shaw et al., 2018).

111 In our experiments, the proposed methods of  
112 SQLBERT and the history information enhanced  
113 schema-linking substantially improve the perfor-  
114 mance of our model. At the time of writing, our  
115 model ranks first on both two large-scale cross-

116 domain context-dependent text-to-SQL leader-  
117 boards, SparC (Yu et al., 2019b) and CoSQL (Yu  
118 et al., 2019a). Specifically, our model achieves  
119 a 64.6% question match and 42.9% interaction  
120 match accuracy on SparC, and a 53.9% question  
121 match and 24.6% interaction match accuracy on  
122 CoSQL.

## 123 2 Related Work

124 Text-to-SQL semantic parsing follows a long  
125 line of research on semantic parsing from natural  
126 language to logical language (Zelle and Mooney,  
127 1996; Zettlemoyer and Collins, 2005; Wong and  
128 Mooney, 2007).

129 Recently, context-independent text-to-SQL se-  
130 mantic parsing has been well studied. Spider (Yu  
131 et al., 2018b) is a famous dataset for the complex  
132 and cross-domain context-independent text-to-SQL  
133 task. Some works (Bogin et al., 2019a,b; Chen  
134 et al., 2021) apply graph neural networks to encode  
135 database schema. Xu et al. (2021) succeed in ap-  
136 plying deep transformers to the context-independent  
137 text-to-SQL task. Yu et al. (2018a) employ a tree-  
138 based decoder to match SQL grammar. Rubin and  
139 Berant (2021) improve the tree-based decoder by  
140 a bottom-up method. Scholak et al. (2021) refine  
141 the sequence-based decoder via carefully designed  
142 restriction rules. Guo et al. (2019) and Gan et al.  
143 (2021) propose SQL intermediate representations  
144 to bridge the gap between natural language and  
145 SQL. Lei et al. (2020) study the role of schema-  
146 linking in text-to-SQL semantic parsing. Wang et al.  
147 (2020) propose a unified framework to capture the  
148 schema-linking. Lin et al. (2020) represent the  
149 schema-linking as a tagged sequence. Cao et al.  
150 (2021) further integrate non-local and local fea-  
151 tures via taking advantage of both schema-linking  
152 graph and its corresponding line graph. Besides,  
153 many previous works (Deng et al., 2021; Yu et al.,  
154 2021a; Shi et al., 2021) focus on pre-train mod-  
155 els for context-independent text-to-SQL semantic  
156 parsing.

157 With more attentions on context-dependent text-  
158 to-SQL semantic parsing, existing works have been  
159 devoted to the context-dependent text-to-SQL task.  
160 SparC (Yu et al., 2019b) and CoSQL (Yu et al.,  
161 2019a) datasets are specially proposed for the task.  
162 EditSQL (Zhang et al., 2019) and IST-SQL (Wang  
163 et al., 2021) focus on taking advantages of the  
164 last predicted query for the prediction of current  
165 query. EditSQL tries to copy the overlap tokens

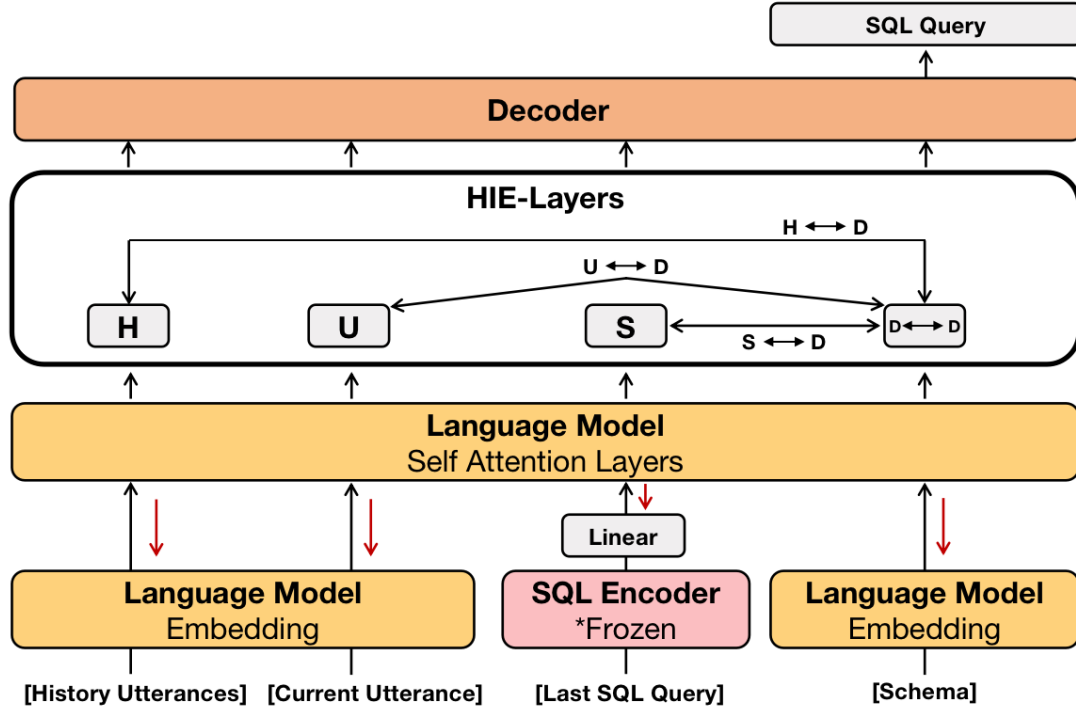


Figure 2: Structure and components of HIE-SQL. During the training stage, the parameters of SQL Encoder will not be updated.

166 from the last predicted query, while IST-SQL pro- 191  
 167 poses an interaction state tracking method to en- 192  
 168 code the information from the last predicted query. 193  
 169 IGSQL (Cai and Wan, 2020) and R<sup>2</sup>SQL (Hui et al., 194  
 170 2021) leverages the contextual information among 195  
 171 the current utterance, interaction history utterances 196  
 172 and database schema via context-aware dynamic 197  
 173 graphs. Notably, R<sup>2</sup>SQL simulates the informa- 198  
 174 tion by connecting the schema graphs with the to- 199  
 175 kens in interactive utterances. Yu et al. (2021b) 200  
 176 creatively propose a context-aware pre-trained lan- 201  
 177 guage model. However, the problem of making 202  
 178 full use of both interaction history utterances and 203  
 179 predicted queries for the context-dependent text-to-  
 180 SQL task remains open.

### 181 3 HIE-SQL

182 First, we formally define the conversational text- 205  
 183 to-SQL semantic parsing problem. In the rest of the 206  
 184 section, we detail the architecture of history infor- 207  
 185 mation enhanced text-to-SQL model (HIE-SQL). 208

#### 186 3.1 Preliminaries

187 **Task Definition.** Given the current user utterance 213  
 188  $u_\tau$ , interaction history  $h_\tau = [u_1, u_2, \dots, u_{\tau-1}]$ , the 214  
 189 schema  $D = \langle T, C \rangle$  of the target database such 215  
 190 that the set of tables  $T = \{t_1, \dots, t_{|T|}\}$  and the 216

set of columns  $C = \{c_1, \dots, c_{|C|}\}$ , our goal is to 191  
 generate the corresponding SQL query  $s_\tau$ . 192

193 **Model Architecture.** Figure 2 shows the encoder- 194  
 195 decoder framework of HIE-SQL. We will intro- 196  
 197 duce it in four modules: (i) **Multimodal Encoder**, 198  
 199 which encodes SQL query and natural language 200  
 201 context in a multimodal manner, (ii) **SQLBERT**, 202  
 203 a bimodal pre-trained encoder for SQL and nat- 204  
 205 ural language, (iii) **HIE-Layers**, which encode 206  
 207 pre-defined schema-linking relations between all 208  
 209 elements of the output of Language Model, and 210  
 211 (iv) **Decoder**, which generates SQL query as an 212  
 213 abstract syntax tree. 214  
 215  
 216

#### 217 3.2 Multimodal Encoder

218 Since of the huge syntax structure differences 219  
 220 between SQL and natural language, using a sin- 221  
 222 gle language model to encode both languages at 223  
 224 the same time increases the difficulty and cost of 225  
 226 training the model. Inspired by the efficiency of the 227  
 works (Kiela et al., 2019; Tsimpoukelli et al., 2021) 228  
 to solve the multimodal problems, we build an ad- 229  
 ditional pre-trained Encoder named SQLBERT (we 230  
 will detail it in the following section) to pre-encode 231  
 SQL query. Then we learn weights  $W \in R^{N \times M}$  to 232  
 project the N-dimensional SQL query embeddings 233  
 to M-dimensional token input embedding space of 234  
 235  
 236

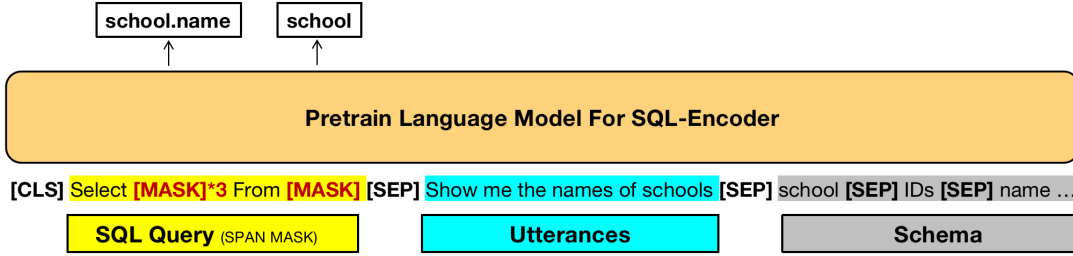


Figure 3: Input format and training objective of SQLBERT.

the language model:

$$S = Wf(s_{\tau-1}), \quad (1)$$

where  $f(\cdot)$  is the last hidden state output of SQLBERT.

We arrange the input format of HIE-SQL as  $x = ([CLS], \mathcal{U}, [CLS], \mathcal{S}, [SEP], \mathcal{T}, [SEP], \mathcal{C})$  in which

$$\begin{aligned} \mathcal{U} &= (u_1, [CLS], u_2, \dots, [CLS], u_\tau), \\ \mathcal{T} &= (t_1, [SEP], t_2, \dots, [SEP], t_{|T|}), \\ \mathcal{C} &= (c_1, [SEP], c_2, \dots, [SEP], c_{|C|}). \end{aligned} \quad (2)$$

All the special separator tokens and language word tokens in  $x$  are converted to the word embedding by embedding layer of the language model. Gathering the embeddings of natural language and SQL, we feed them to self-attention blocks in a language model. In the training stage, we directly take the golden SQL query of the last turn as an input SQL query. As for the inference stage, we apply the SQL query generated by HIE-SQL in the last turn. When predicting the first turn utterance, we just set  $S$  to empty.

### 3.3 SQLBERT

As mentioned above, we treat the SQL query as another modality that can provide information of the SQL query from the previous round as a reference for the model. So we need an encoder to extract the representation of the SQL query.

**Model Architecture.** Considering the success of multi-modal pre-trained models, such as ViLBERT (Lu et al., 2019) for language-image and CodeBERT (Feng et al., 2020) for natural language and programming language, we propose SQLBERT, a bimodal pre-trained model for natural language and SQL. We develop SQLBERT by using the same model architecture as RoBERTa (Liu et al., 2019). The total number of model parameters is 125M.

**Input format.** As the training method showed in Figure 3, we set the same input as CodeBERT (Feng et al., 2020) does. To alleviate the difficulty of training and resolve inconsistencies between natural language and schema, we append the question-relevant database schema to the concatenation of SQL query and question. We represent the whole input sequence into the format as  $x = ([CLS], s_1, s_2, \dots, s_n, [SEP], q_1, q_2, \dots, q_m, [SEP], t_1 : c_{11}, c_{12}, \dots, [SEP], t_2 : c_{21}, \dots, [SEP], \dots)$ , in which  $s$ ,  $q$ ,  $t$ , and  $c$  are the tokens of SQL query, question, tables, and columns respectively.

**Training Objective.** The main training objective of SQLBERT is the masked language modeling (MLM). It’s worth noting that we only mask the tokens of SQL query because we only need SQLBERT to encode SQL query in the downstream task. Specifically, we utilize a special objective referenced span masking (Sun et al., 2019) by sampling 15% independent span in SQL clause except the reserved word (e.g., SELECT, FROM, WHERE), which aims to avoid leaking answers and help SQLBERT learn the information structure of SQL better. In the training stage, we adopt a dynamic masking strategy via randomly shuffling the order of tables and columns in the original schema. We describe the masked span prediction loss as

$$\mathcal{L}(\theta) = \sum_{k=1}^n -\log \mathcal{P}_\theta(s_k^{mask} | s^{\setminus mask}, q, t, c), \quad (3)$$

where  $\theta$  stands for the model parameters,  $s_k^{mask}$  is the masked span of SQL input,  $s^{\setminus mask}$  is the unmasked part.

**Training data.** We train SQLBERT with the open-source Text-to-SQL datasets including Spider, SparC and CoSQL, whose data structures and annotation styles are quite similar. For each sample, we only use its question, SQL query, and the corresponding database schema. As for SparC and CoSQL, which is a context-dependent version, we simply concatenate the current utterance with the



	Current Utterance	Interaction History	SQL Query
Columns	U-C-EM (Exact Match)	H-C-EM	S-C-EC (Equal Columns)
	U-C-PM (Partial Match)	H-C-PM	S-C-UC (Unequal Columns)
	U-C-VM (Value Match)	H-C-VM	
Tables	U-T-EM	H-T-EM	S-T-ET (Equal Tables)
	U-T-PM	H-T-PM	S-T-UT (Unequal Tables)

Table 1: Edge types between current utterance  $U$ , interaction history  $H$ , SQL  $S$ , and database schema  $D$  (Columns  $C$  and Tables  $T$ ). We omit the pre-existing relations in schema such as the foreign-key relation (C-C-FK) in the table. We set a default "no relation" edge type for every node pair.

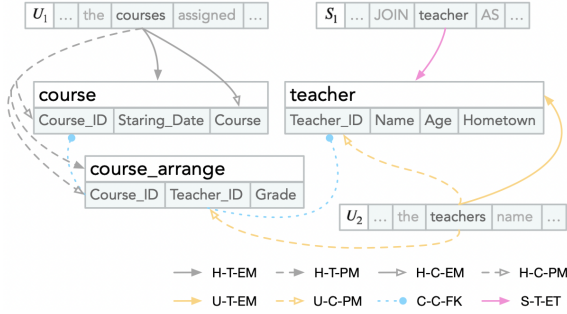


Figure 4: An example of the schema-linking graph for the prediction of  $S_2$  in Figure 1. The graph is a sub-graph of the whole schema-linking graph. We only respectively choose one token in the history utterance ( $U_1$ ), the current utterance ( $U_2$ ), and the last predicted SQL query ( $S_1$ ) in the example. Besides, we omit all unequal relation edges (S-C-UC and S-T-UT) and default "no relation" edges.

history utterances to build the question input. The size of the training dataset is 34,175.

### 3.4 HIE-Layers

**Schema-Linking Graph.** To explicitly encode the complex relational database schema. We convert it to a directed graph  $\mathcal{G} = \langle \mathcal{V}, \mathcal{E} \rangle$ , where  $\mathcal{V} = C \cup T$  and  $\mathcal{E}$  represents the set of pre-existing relations within columns and tables such as the foreign-key relation. In addition, we also consider the unseen linking to the schema in the contexts of current utterance, interaction history utterances, and the last predicted SQL query. Specifically, we define the context-dependent schema-linking graph  $\mathcal{G}_c = \langle \mathcal{V}_c, \mathcal{E}_c \rangle$  where  $\mathcal{V}_c = C \cup T \cup U \cup H \cup S$  and  $\mathcal{E}_c = \mathcal{E} \cup \mathcal{E}_{U \leftrightarrow D} \cup \mathcal{E}_{H \leftrightarrow D} \cup \mathcal{E}_{S \leftrightarrow D}$ . The additional relation edges are listed in Table 1. We set two match types between the language tokens of  $U$ ,  $H$ , and  $D$ : EM for Exact Match, PM for Partial Match. When using database contents, we set VM (Value Match) for exactly matching the value of columns. As for SQL  $S$ , we simply match the words of tables and columns that appear in it to the target database

schema: EC (Equal Columns) for column match, ET (Equal Tables) for table match. In Figure 4, we show an example of the proposed schema-linking graph.

**Graph Encoding.** The work (Wang et al., 2020) shows that Relative Self-Attention Mechanism (Shaw et al., 2018) is an efficient way to encode graphs whose nodes are at the token level. It rebuilds the calculation of the self-attention module in the transformer layers as follows:

$$\begin{aligned}
 e_{ij} &= \frac{x_i W^Q (x_j W^K + r_{ij}^K)^T}{\sqrt{d_z}}, \\
 \alpha_{ij} &= \text{softmax}_j \{e_{ij}\}, \\
 z_i &= \sum_{j=1}^n \alpha_{ij} (x_j W^V + r_{ij}^V).
 \end{aligned} \tag{4}$$

HIE-Layers consist of 8 transformer layers, whose self-attention modules are described above. Specifically, we initialize a learned embedding for each type of edge defined above. For every input sample, we build a relation matrix  $\mathcal{R} \subseteq (L \times L)$  where  $L$  is the length of the input token.  $\mathcal{R}^{(i,j)}$  represents the relation type between  $i$ -th and  $j$ -th input tokens. While computing the relative attention, we set the  $r_{ij}^K = r_{ij}^V = \mathcal{R}_e^{(i,j)}$  where  $\mathcal{R}_e^{(i,j)}$  is the corresponding embedding of  $\mathcal{R}^{(i,j)}$ .

### 3.5 Decoder

To build the decoder of HIE-SQL, we apply the same work (Yin and Neubig, 2017) as Wang et al. (2020) propose, which generates SQL as an abstract syntax tree in depth-first traversal order by using LSTM (Hochreiter and Schmidhuber, 1997) to output sequences of decoder actions. We recommend the reader to refer to the work (Yin and Neubig, 2017) for details.

Dataset	System Response	Interaction	Train	Dev	Test	User Questions	Vocab	Avg Turn
CoSQL	✓	3007	2164	293	551	15598	9585	5.2
SparC	✗	4298	3034	422	842	12726	3794	3.0

Table 2: Details of SparC and CoSQL datasets.

Model	SparC				CoSQL			
	Dev		Test		Dev		Test	
	QM	IM	QM	IM	QM	IM	QM	IM
EditSQL + BERT (Zhang et al., 2019)	47.2	29.5	47.9	25.3	39.9	12.3	40.8	13.7
IGSQL + BERT (Cai and Wan, 2020)	50.7	32.5	51.2	29.5	44.1	15.8	42.5	15.0
IST-SQL + BERT (Wang et al., 2021)	-	-	-	-	44.4	14.7	41.8	15.2
R <sup>2</sup> SQL + BERT (Hui et al., 2021)	54.1	35.2	55.8	30.8	45.7	19.5	46.8	17.0
RAT-SQL <sup>†</sup> + SCoRe (Yu et al., 2021b)	62.2	42.5	62.4	38.1	52.1	22.0	51.6	21.2
T5-3B + PICARD <sup>†</sup> (Scholak et al., 2021)	-	-	-	-	<b>56.9</b>	24.2	<b>54.6</b>	23.7
HIE-SQL + GraPPa (ours)	<b>64.7</b>	<b>45.0</b>	<b>64.6</b>	<b>42.9</b>	56.4	<b>28.7</b>	53.9	<b>24.6</b>

Table 3: Performances of various models in SparC and CoSQL. QM and IM stand for question match and interaction match respectively. The models with † are proposed for the context-independent text-to-SQL task and applied to the context-dependent text-to-SQL task by just appending interaction history utterances to the input.

### 3.6 Regularization Strategy

We introduce R-Drop (Liang et al., 2021), a simple regularization strategy, to prevent the overfitting of the model. Concretely, we feed every input data  $x_i$  to go through our model twice and the loss function is as follows:

$$\begin{aligned}
\mathcal{L}_{NLL}^i &= -\log\mathcal{P}_1(y_i|x_i) - \log\mathcal{P}_2(y_i|x_i), \\
\mathcal{L}_{KL}^i &= \frac{1}{2}(D_{KL}(\mathcal{P}_1(y_i|x_i)||\mathcal{P}_2(y_i|x_i)) \\
&\quad + D_{KL}(\mathcal{P}_2(y_i|x_i)||\mathcal{P}_1(y_i|x_i))), \\
\mathcal{L}^i &= \mathcal{L}_{NLL}^i + \mathcal{L}_{KL}^i,
\end{aligned} \tag{5}$$

where  $-\log\mathcal{P}_1(y_i|x_i)$  and  $-\log\mathcal{P}_2(y_i|x_i)$  are two output distributions for input  $x_i$  at all decoder steps,  $\mathcal{L}_{NLL}^i$  is the negative log-likelihood learning objective of decoder actions, and  $\mathcal{L}_{KL}^i$  is the bidirectional Kullback-Leibler (KL) divergence between these two output distributions.

## 4 Experiment

### 4.1 Setup

**Setting.** We initialize the weights of Language Model with GraPPa (Yu et al., 2021a), an effective pre-training model for table semantic parsing that performs well on the context-independent text-to-SQL datasets (e.g. Spider). We stack 8 HIE-layers,

which are introduced in section 3.4, on top of the Language Model. When training the model with R-Drop, we set the Dropout rate of 0.1 for the Language Model and HIE-Layers, 0.3 for the decoder. We use Adam optimizer to conduct the parameter learning and set the learning rate of  $1e^{-5}$  for fine-tuning GraPPa and  $1e^{-4}$  for HIE-Layers and Decoder. The learning rate linearly increases to the setting point at first  $max\_steps/8$  steps, then decreases to 0 at  $max\_steps$ , where  $max\_steps = 50000$  with 24 training batch-size. As for SQLBERT, we fine-tune CodeBERT<sub>BASE</sub> (Feng et al., 2020) on the dataset we described in Section 3.3. We set the learning rate as  $1e^{-5}$ , a batch size of 64, and train SQLBERT for 10 epochs. The shape of learned weights of the linear layer applied to the output of SQLBERT is  $768 \times 1024$ . While inferring, we set the beam size to 3.

**Datasets.** We conduct experiments on two cross-domain context-dependent text-to-SQL datasets, SparC (Yu et al., 2019b) and CoSQL (Yu et al., 2019a). Table 2 depicts the statistic information of them.

**Evaluation Metrics.** The main metric we used to measure model performance in SparC and CoSQL is interaction match (IM), which requires all output SQL queries in interaction to be correct. We also

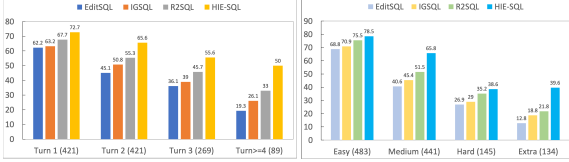


Figure 5: Performances of previous works and HIE-SQL in different turns (left) and different difficulty levels (right) on SparC.

use question match (QM) to evaluate the accuracy of every single question.

## 4.2 Experiment Result.

Results of our proposed HIE-SQL model are shown in Table 3. In terms of interaction match, our model achieves state-of-the-art performances on both development set and test set of SparC and CoSQL. For the test set of SparC, HIE-SQL outperforms the prior state-of-the-art (Yu et al., 2021b) by 4.8% in IM and 2.2% in QM. For CoSQL, compared with the previous state-of-the-art (Scholak et al., 2021), a rule-based auto-regressive method based on the large pre-trained model-T5-3B (Raffel et al., 2020) which is optimized for a GPU with 40GB of memory, HIE-SQL improves IM of development set by 4.5% and IM of the test set by 0.9%. Besides, HIE-SQL surpasses RAT-SQL + SCoRe in all metrics of SparC and CoSQL. This demonstrates that properly integrating interaction utterances and predicted SQL queries is an effective way to enhance the model’s ability for Context-Dependent Text-to-SQL Semantic Parsing.

To further explore the advantages of HIE-SQL, we test the performance on different turns and at different difficulty levels of utterances. As shown in Figure 5, with the increase of turns, the lead of our model gets greater and greater. When the indexes of turns are greater than or equal to 4, the accuracy of HIE-SQL is 17% higher than that of R<sup>2</sup>SQL. It demonstrates that the main contribution of introducing SQL query is to improve the robustness of the model to long interaction. HIE-SQL is also robust to the varying difficulty levels of utterances. Our model performs equally in hard and extra hard levels, and achieves 39.6% accuracy on the extra hard level, which is 17.8% higher than that of R<sup>2</sup>SQL.

## 4.3 Ablation Study

We provide ablation studies to examine the contribution of each component of HIE-SQL. We want

Model	SparC		CoSQL	
	QM	IM	QM	IM
HIE-SQL	64.7	<b>45.0</b>	56.4	<b>28.7</b>
w/o SQL query	<b>65.8</b>	44.3	<b>56.5</b>	23.9
w/o SQLBERT	63.9	44.7	54.8	26.3
w/o $\mathcal{E}_{H \leftrightarrow D}$	64.0	44.3	56.0	26.3

Table 4: Ablation study of HIE-SQL in development sets of SparC and CoSQL. As for ablation on SQL query, we drop the SQL query and only feed utterances and database schema to the model. As for ablation on SQLBERT, we directly concatenate the tokens of SQL query and other context tokens for the input of the language model. And w/o  $\mathcal{E}_{H \leftrightarrow D}$  means we treat historical utterances like the current utterance in our schema-linking.

Dataset	Model	T-F	F-T	T-T
SparC	HIE-SQL	125	88	<b>383</b>
	w/o SQL query	132	104	379
CoSQL	HIE-SQL	140	106	<b>278</b>
	w/o SQL query	161	128	254

Table 5: The counts of different switches in the pairs of adjacent predicted SQL queries. T-F stands for the match of the former predicted query and unmatched of the later predicted query with golden queries. F-T stands for the reverse case. T-T is the case of both matching.

to identify whether introducing the last SQL query has a significant impact on performance. Also, we would like to investigate whether the pre-trained SQL encoder, SQLBERT, can improve the model’s ability to understand SQL queries. What’s more, we conduct another ablation study regarding additional graph edges between historical utterances and database schema  $\mathcal{E}_{H \leftrightarrow D}$  to check the necessity of the join of historical utterance information in schema-linking.

As shown in Table 4, Our full model achieves about 5 points and 1 point improvement of IM in CoSQL and SparC respectively compared with the model without the last SQL query input. The pre-encoding SQL query by SQLBERT can further improve the performance. It confirms SQLBERT’s ability to efficiently represent SQL features. In addition,  $\mathcal{E}_{H \leftrightarrow D}$  also plays a positive role.

Table 5 shows the continuity of performance of our model compared with that of the model without the last SQL query input. Our model has a higher rate of continuous match, but a lower rate

$U_1$	Which cartoon aired <b>first</b> ?
HIE-SQL	SELECT title FROM cartoon ORDER BY original_air_date asc LIMIT 1
RAT-SQL	SELECT title FROM cartoon ORDER BY original_air_date asc LIMIT 1
$U_2$	What was the <b>last</b> cartoon to air?
HIE-SQL	SELECT title FROM cartoon ORDER BY original_air_date desc LIMIT 1
RAT-SQL	SELECT title FROM cartoon ORDER BY original_air_date desc LIMIT 1
$U_3$	What channel was it on?
HIE-SQL	SELECT channel FROM cartoon ORDER BY original_air_date desc LIMIT 1
RAT-SQL	SELECT channel FROM cartoon ORDER BY original_air_date desc LIMIT 1
$U_4$	What is the production code?
HIE-SQL	SELECT production_code FROM cartoon ORDER BY original_air_date <b>desc</b> LIMIT 1
RAT-SQL	SELECT production_code FROM cartoon ORDER BY original_air_date <b>asc</b> LIMIT 1

Table 6: An example in CoSQL.  $U_i$  is the input utterance of turn  $i$  with corresponding predictions of HIE-SQL and RAT-SQL following. All predictions of HIE-SQL are the ground truth queries in the case.

of switching from mismatch to match. It illustrates that our model does use the SQL information and is sensitive to the accuracy of the last predicted SQL query which explains the higher question match without SQL query input.

We regard R-Drop as a simple means of data augmentation which can improve the generalization of the model. As shown in Figure 6, the model with R-drop outperforms the model without R-Drop in both QM and IM. Additionally, the standard deviations of the IM in the last 20k steps are 0.014 and 0.015 of HIE-SQL and the one without R-Drop respectively even the curve of HIE-SQL has a more obvious upward trend. It shows that R-Drop improves the robustness of our model and stabilizes its performance in IM.

#### 4.4 Case Study

In Table 6, we show the predictions of HIE-SQL and RAT-SQL in an example of CoSQL. Here, HIE-SQL and RAT-SQL both fine-tune GraPPa on CoSQL. As the example shows, RAT-SQL fails to distinguish the right one from two long-range dependences in  $U_1$  and  $U_2$  in Table 6. By contrast, HIE-SQL inherits the right context-dependence from the last predicted query to avoid confusion between  $U_1$  and  $U_2$ .

## 5 Conclusion

We present HIE-SQL, a history information enhanced context-dependent text-to-SQL model, which targets at explicitly capturing the context-

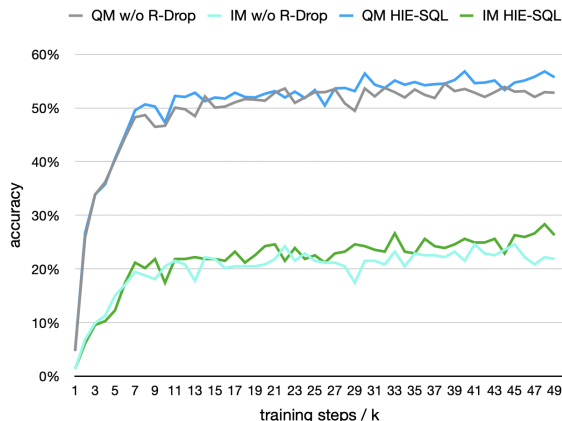


Figure 6: Ablation study result of regarding R-Drop in development set of CoSQL. We show the performances in QM and IM of two models at different training steps. We set the beam size = 1 in the inference stage.

dependence from both interaction history utterances and the last predicted SQL query. With the help of the proposed bimodal pre-trained model, SQLBERT, HIE-SQL bridge the gap between the utterances and predicted SQL despite the mismatch of natural language and logic-form SQL. Moreover, we also introduce a method of schema-linking to enhance the connections among utterances, SQL query, and database schema.

Taken together, HIE-SQL achieves consistent improvements on the context-dependent text-to-SQL task, especially in the interaction match metric. HIE-SQL achieves new state-of-the-art results on two famous context-dependent text-to-SQL datasets, SparC and CoSQL.



## References

- 498 Ben Bogin, Jonathan Berant, and Matt Gardner. 2019a. **Rep-**  
499 **resenting schema structure with graph neural networks**  
500 **for text-to-SQL parsing.** In *Proceedings of the 57th Con-*  
501 *ference of the Association for Computational Linguistics*  
502 *(ACL)*, volume 1, pages 4560–4565.
- 503 Ben Bogin, Matt Gardner, and Jonathan Berant. 2019b. **Global**  
504 **reasoning over database structures for text-to-SQL**  
505 **reasoning.** In *Proceedings of the 2019 Conference on Empirical*  
506 *Methods in Natural Language Processing and the 9th Inter-*  
507 *national Joint Conference on Natural Language Processing*  
508 *(EMNLP/IJCNLP)*, pages 3657–3662.
- 509 Yitao Cai and Xiaojun Wan. 2020. **IGSQL: database schema**  
510 **interaction graph based neural model for context-dependent**  
511 **text-to-SQL generation.** In *Proceedings of the 2020 Con-*  
512 *ference on Empirical Methods in Natural Language Pro-*  
513 *cessing (EMNLP)*, pages 6903–6912.
- 514 Ruisheng Cao, Lu Chen, Zhi Chen, Yanbin Zhao, Su Zhu,  
515 and Kai Yu. 2021. **LGESQL: line graph enhanced text-to-**  
516 **SQL model with mixed local and non-local relations.** In  
517 *Proceedings of the 59th Annual Meeting of the Associa-*  
518 *tion for Computational Linguistics and the 11th Interna-*  
519 *tional Joint Conference on Natural Language Processing*  
520 *(ACL/IJCNLP)*, volume 1, pages 2541–2555.
- 521 Zhi Chen, Lu Chen, Yanbin Zhao, Ruisheng Cao, Zihan Xu,  
522 Su Zhu, and Kai Yu. 2021. **ShadowGNN: Graph projection**  
523 **neural network for text-to-SQL parser.** In *Proceedings*  
524 *of the 2021 Conference of the North American Chapter*  
525 *of the Association for Computational Linguistics: Human*  
526 *Language Technologies (NAACL-HLT)*, volume 1, pages  
527 5567–5577.
- 528 Xiang Deng, Ahmed Hassan Awadallah, Christopher Meek,  
529 Oleksandr Polozov, Huan Sun, and Matthew Richardson.  
530 2021. **Structure-grounded pretraining for text-to-SQL.** In  
531 *Proceedings of the 2021 Conference of the North American*  
532 *Chapter of the Association for Computational Linguistics:*  
533 *Human Language Technologies (NAACL-HLT)*, volume 1,  
534 pages 1337–1350.
- 535 Kedar Dhamdhere, Kevin S McCurley, Ralfi Nahmias,  
536 Mukund Sundararajan, and Qiqi Yan. 2017. **Analyza: Ex-**  
537 **ploring data with conversation.** In *Proceedings of the 22nd*  
538 *International Conference on Intelligent User Interfaces*  
539 *(ACM IUI)*, pages 493–504.
- 540 Zhangyin Feng, Daya Guo, Duyu Tang, Nan Duan, Xiaocheng  
541 Feng, Ming Gong, Linjun Shou, Bing Qin, Ting Liu, Daxin  
542 Jiang, and Ming Zhou. 2020. **CodeBERT: A pre-trained**  
543 **model for programming and natural languages.** In *Findings*  
544 *of the Association for Computational Linguistics: EMNLP*  
545 *2020 (EMNLP-Findings)*, pages 1536–1547.
- 546 Yujian Gan, Xinyun Chen, Jinxia Xie, Matthew Purver, John R  
547 Woodward, John Drake, and Qiaofu Zhang. 2021. **Natural**  
548 **SQL: Making SQL easier to infer from natural language**  
549 **specifications.** In *Findings of the Association for Computa-*  
550 *tional Linguistics: EMNLP 2021 (EMNLP-Findings)*.
- 551 Jiaqi Guo, Zecheng Zhan, Yan Gao, Yan Xiao, Jian-Guang  
552 Lou, Ting Liu, and Dongmei Zhang. 2019. **Towards com-**  
553 **plex text-to-SQL in cross-domain database with intermedi-**  
554 **ate representation.** In *Proceedings of the 57th Conference*  
555 *of the Association for Computational Linguistics (ACL)*,  
556 volume 1, pages 4524–4535.
- Sepp Hochreiter and Jürgen Schmidhuber. 1997. **Long short-**  
term memory. *Neural Comput.*, 9(8):1735–1780. 557 558
- Binyuan Hui, Ruiying Geng, Qiyu Ren, Binhua Li, Yongbin  
Li, Jian Sun, Fei Huang, Luo Si, Pengfei Zhu, and Xiaodan  
Zhu. 2021. **Dynamic hybrid relation exploration network**  
for cross-domain context-dependent semantic parsing. In  
*Proceedings of the Thirty-Fifth AAAI Conference on Artificial*  
*Intelligence (AAAI)*, volume 35, pages 13116–13124. 559 560 561 562 563 564
- Douwe Kiela, Suvrat Bhooshan, Hamed Firooz, and Davide  
Testuggine. 2019. **Supervised multimodal bitransformers**  
for classifying images and text. In *Visually Grounded In-*  
*teraction and Language (ViGIL), NeurIPS 2019 Workshop.* 565 566 567 568
- Wenqiang Lei, Weixin Wang, Zhixin Ma, Tian Gan, Wei Lu,  
Min-Yen Kan, and Tat-Seng Chua. 2020. **Re-examining**  
the role of schema linking in text-to-sql. In *Proceedings*  
*of the 2020 Conference on Empirical Methods in Natural*  
*Language Processing (EMNLP)*, pages 6943–6954. 569 570 571 572 573
- Xiaobo Liang, Lijun Wu, Juntao Li, Yue Wang, Qi Meng,  
Tao Qin, Wei Chen, Min Zhang, and Tie-Yan Liu. 2021.  
R-Drop: Regularized dropout for neural networks. *arXiv*  
*preprint arXiv:2106.14448.* 574 575 576 577
- Xi Victoria Lin, Richard Socher, and Caiming Xiong. 2020.  
**Bridging textual and tabular data for cross-domain text-**  
**to-SQL semantic parsing.** In *Findings of the Association*  
*for Computational Linguistics: EMNLP 2020 (EMNLP-*  
*Findings)*, pages 4870–4888. 578 579 580 581 582
- Qian Liu, Bei Chen, Jiaqi Guo, Jian-Guang Lou, Bin Zhou,  
and Dongmei Zhang. 2020. **How far are we from effec-**  
**tive context modeling? an exploratory study on semantic**  
**parsing in context.** In *Proceedings of the Twenty-Ninth*  
*International Joint Conference on Artificial Intelligence*  
*(IJCAI)*, pages 3580–3586. 583 584 585 586 587 588
- Yinhan Liu, Myle Ott, Naman Goyal, Jingfei Du, Man-  
dar Joshi, Danqi Chen, Omer Levy, Mike Lewis, Luke  
Zettlemoyer, and Veselin Stoyanov. 2019. **Roberta: A ro-**  
**busly optimized bert pretraining approach.** *arXiv preprint*  
*arXiv:1907.11692.* 589 590 591 592 593
- Jiasen Lu, Dhruv Batra, Devi Parikh, and Stefan Lee. 2019.  
**ViLBERT: Pretraining task-agnostic visiolinguistic rep-**  
**resentations for vision-and-language tasks.** In *Advances in*  
*Neural Information Processing Systems 32: Annual Confer-*  
*ence on Neural Information Processing Systems (NeurIPS)*,  
pages 13–23. 594 595 596 597 598 599
- Colin Raffel, Noam Shazeer, Adam Roberts, Katherine Lee,  
Sharan Narang, Michael Matena, Yanqi Zhou, Wei Li, and  
Peter J. Liu. 2020. **Exploring the limits of transfer learning**  
**with a unified text-to-text transformer.** *J. Mach. Learn.*  
*Res.*, 21:140:1–140:67. 600 601 602 603 604
- Ohad Rubin and Jonathan Berant. 2021. **SmBoP: Semi-**  
**autoregressive bottom-up semantic parsing.** In *Proceedings*  
*of the 2021 Conference of the North American Chapter of*  
*the Association for Computational Linguistics: Human*  
*Language Technologies (NAACL-HLT)*, volume 1, pages  
311–324. 605 606 607 608 609 610
- Torsten Scholak, Nathan Schucher, and Dzmitry Bahdanau.  
2021. **PICARD: Parsing incrementally for constrained**  
**auto-regressive decoding from language models.** In *Pro-*  
*ceedings of the 2021 Conference on Empirical Methods in*  
*Natural Language Processing (EMNLP)*. 611 612 613 614 615

616	Peter Shaw, Jakob Uszkoreit, and Ashish Vaswani. 2018. <a href="#">Self-attention with relative position representations</a> . In <i>Proceedings of the 2018 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies (NAACL-HLT)</i> , volume 2, pages 464–468.	676
617		677
618		678
619		679
620		680
621		681
622	Peng Shi, Patrick Ng, Zhiguo Wang, Henghui Zhu, Alexander Hanbo Li, Jun Wang, Cícero Nogueira dos Santos, and Bing Xiang. 2021. <a href="#">Learning contextual representations for semantic parsing with generation-augmented pre-training</a> . In <i>Proceedings of the Thirty-Fifth AAAI Conference on Artificial Intelligence (AAAI)</i> , pages 13806–13814.	682
623		683
624		684
625		685
626		686
627		687
628	Yu Sun, Shuohuan Wang, Yukun Li, Shikun Feng, Xuyi Chen, Han Zhang, Xin Tian, Danxiang Zhu, Hao Tian, and Hua Wu. 2019. Ernie: Enhanced representation through knowledge integration. <i>arXiv preprint arXiv:1904.09223</i> .	688
629		689
630		690
631		691
632	Maria Tsimpoukelli, Jacob Menick, Serkan Cabi, SM Eslami, Oriol Vinyals, and Felix Hill. 2021. Multimodal few-shot learning with frozen language models. <i>arXiv preprint arXiv:2106.13884</i> .	692
633		693
634		694
635		695
636	Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N. Gomez, Lukasz Kaiser, and Illia Polosukhin. 2017. <a href="#">Attention is all you need</a> . In <i>Advances in Neural Information Processing Systems 30: Annual Conference on Neural Information Processing Systems (NeurIPS)</i> , pages 5998–6008.	696
637		697
638		698
639		699
640		700
641		701
642	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 (ACL)</i> , pages 7567–7578.	702
643		703
644		704
645		705
646		706
647	Runze Wang, Zhen-Hua Ling, Jingbo Zhou, and Yu Hu. 2021. <a href="#">Tracking interaction states for multi-turn text-to-SQL semantic parsing</a> . In <i>Proceedings of the Thirty-Fifth AAAI Conference on Artificial Intelligence (AAAI)</i> , pages 13979–13987.	707
648		708
649		709
650		710
651		711
652	Nathaniel Weir, Prasetya Utama, Alex Galakatos, Andrew Crotty, Amir Ilkhechi, Shekar Ramaswamy, Rohin Bhusan, Nadja Geisler, Benjamin Hättasch, Steffen Eger, et al. 2020. DBPal: A fully pluggable nl2sql training pipeline. In <i>Proceedings of the 2020 ACM SIGMOD International Conference on Management of Data</i> , pages 2347–2361.	712
653		713
654		714
655		715
656		716
657		717
658		718
659	Yuk Wah Wong and Raymond J. Mooney. 2007. <a href="#">Learning synchronous grammars for semantic parsing with lambda calculus</a> . In <i>Proceedings of the 45th Annual Meeting of the Association for Computational Linguistics (ACL)</i> , pages 960–967.	719
660		720
661		721
662		722
663		723
664	Peng Xu, Dhruv Kumar, Wei Yang, Wenjie Zi, Keyi Tang, Chenyang Huang, Jackie Chi Kit Cheung, Simon J. D. Prince, and Yanshuai Cao. 2021. <a href="#">Optimizing deeper transformers on small datasets</a> . In <i>Proceedings of the 59th Annual Meeting of the Association for Computational Linguistics and the 11th International Joint Conference on Natural Language Processing (ACL/IJCNLP)</i> , volume 1, pages 2089–2102.	724
665		725
666		726
667		727
668		728
669		729
670		730
671		731
672	Pengcheng Yin and Graham Neubig. 2017. <a href="#">A syntactic neural model for general-purpose code generation</a> . In <i>Proceedings of the 55th Annual Meeting of the Association for Computational Linguistics (ACL)</i> , pages 440–450.	732
673		733
674		734
675		735
		736
		737
		738
		739
		740
		741
		742
		743
		744
		745
		746
		747
		748
		749
		750
		751
		752
		753
		754
		755
		756
		757
		758
		759
		760
		761
		762
		763
		764
		765
		766
		767
		768
		769
		770
		771
		772
		773
		774
		775
		776
		777
		778
		779
		780
		781
		782
		783
		784
		785
		786
		787
		788
		789
		790
		791
		792
		793
		794
		795
		796
		797
		798
		799
		800
		801
		802
		803
		804
		805
		806
		807
		808
		809
		810
		811
		812
		813
		814
		815
		816
		817
		818
		819
		820
		821
		822
		823
		824
		825
		826
		827
		828
		829
		830
		831
		832
		833
		834
		835
		836
		837
		838
		839
		840
		841
		842
		843
		844
		845
		846
		847
		848
		849
		850
		851
		852
		853
		854
		855
		856
		857
		858
		859
		860
		861
		862
		863
		864
		865
		866
		867
		868
		869
		870
		871
		872
		873
		874
		875
		876
		877
		878
		879
		880
		881
		882
		883
		884
		885
		886
		887
		888
		889
		890
		891
		892
		893
		894
		895
		896
		897
		898
		899
		900
		901
		902
		903
		904
		905
		906
		907
		908
		909
		910
		911
		912
		913
		914
		915
		916
		917
		918
		919
		920
		921
		922
		923
		924
		925
		926
		927
		928
		929
		930
		931
		932
		933
		934
		935
		936
		937
		938
		939
		940
		941
		942
		943
		944
		945
		946
		947
		948
		949
		950
		951
		952
		953
		954
		955
		956
		957
		958
		959
		960
		961
		962
		963
		964
		965
		966
		967
		968
		969
		970
		971
		972
		973
		974
		975
		976
		977
		978
		979
		980
		981
		982
		983
		984
		985
		986
		987
		988
		989
		990
		991
		992
		993
		994
		995
		996
		997
		998
		999
		1000