

# RAT-SQL: Relation-Aware Schema Encoding and Linking for Text-to-SQL Parsers ([Wang et al, 2020](#))





# Agenda

- Overview (text-to-SQL)
- Motivation
- Literature
- *RAT-SQL framework*
- *Relation Aware Self-Attention mechanism*
- *Experiments and error analysis*

# What is text-to-SQL?



database: concert singer



Show all *countries* and the number of *singers* in each *country*.



```
SELECT Country , count(*) FROM Singer GROUP BY Country
```

**Task:** translating natural language utterance to SQL queries.

**Application:** give people access to vast amounts of databases

# Why text-to-SQL is a hard problem?

- Generalization to unseen databases and domains.



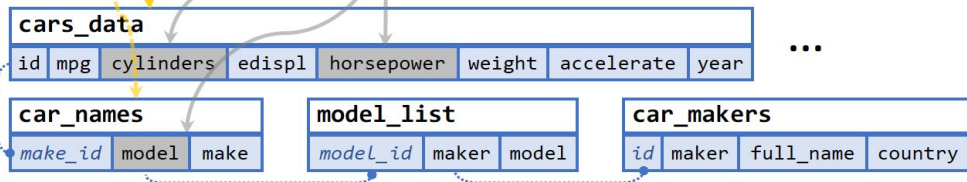
# Why text-to-SQL is a hard problem?

- Schema encoding and linking

## Natural Language Question:

For the cars with 4 cylinders, which model has the largest horsepower?

## Schema:



## Desired SQL:

```
SELECT T1.model
FROM car_names AS T1 JOIN cars_data AS T2
ON T1.make_id = T2.id
WHERE T2.cylinders = 4
ORDER BY T2.horsepower DESC LIMIT 1
```


- Question → Column linking (unknown)
- Question → Table linking (unknown)
- Column → Column foreign keys (known)

Figure 1: A challenging text-to-SQL task from the Spider dataset.



# Motivation

- Address schema encoding and linking problem in text-to-SQL in “RAT-SQL” via *Relation-Aware Self-Attention mechanism*.
- Achieves SOTA performance on ***Spider*** dataset (~8% improvement) for exact match.



Dataset	# Q	# SQL	# DB	# Domain	# Table /DB	ORDER BY	GROUP BY	NESTED	HAVING
ATIS	5,280	947	1	1	32	0	5	315	0
GeoQuery	877	247	1	1	6	20	46	167	9
Scholar	817	193	1	1	7	75	100	7	20
Academic	196	185	1	1	15	23	40	7	18
IMDB	131	89	1	1	16	10	6	1	0
Yelp	128	110	1	1	7	18	21	0	4
Advising	3,898	208	1	1	10	15	9	22	0
Restaurants	378	378	1	1	3	0	0	4	0
WikiSQL	80,654	77,840	26,521	-	1	0	0	0	0
<b>Spider</b>	10,181	5,693	200	138	5.1	1335	1491	844	388

Table 1: Comparisons of text-to-SQL datasets. **Spider** is the *only one* text-to-SQL dataset that contains both databases with multiple tables in different domains and complex SQL queries. It was designed to test the ability of a system to generalize to not only new SQL queries and database schemas but also new domains.

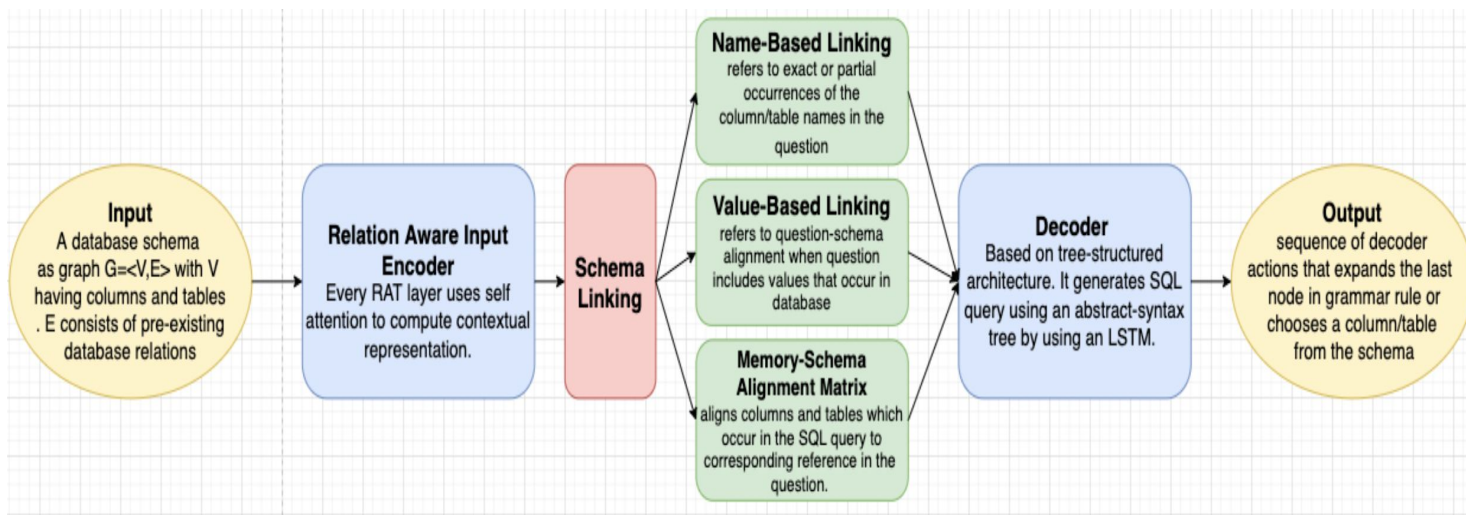


# Literature

- IRNet (Guo et al, 2019)
  - Does not capture binary relations, considers only unary
  - Schema encoder does not exploit schema relations fully.
- GNN (Bogin et al 2019)
  - Does not model context representation of question with schema in encoder.
  - Limits information propagation only to connected nodes defined in predefined graph of foreign keys.



# RAT SQL framework





# Problem Formulation

**Given:** natural language question  $Q$  and schema  $S = \langle C, T \rangle$

**Goal:** Generate SQL program  $P$  represented as abstract syntax tree in the context-free grammar of SQL



# Relation-Aware Self-Attention

Goal is to represent:

- pre-existing relational structure in the input (see later)
- soft relations between sequence elements in the same embedding (self-attention)

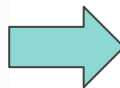
# Relation-Aware Self-Attention

Self-attention in Transformers  
(Vaswani et al)

$$x_i \rightsquigarrow q_i, k_i, v_i$$

$$\alpha_{ij} = \text{softmax}_j \frac{q_i k_j^T}{\sqrt{\text{dim}}}$$

$$y_i = \sum_j \alpha_{ij} v_j$$



Relation-Aware Self-attention in  
RAT-SQL (schema encoding)

$$x_i \rightsquigarrow q_i, k_i, v_i$$

$$\alpha_{ij} = \text{softmax}_j \frac{q_i (k_j + \beta_{ij})^T}{\sqrt{\text{dim}}}$$

$$y_i = \sum_j \alpha_{ij} (v_j + \varepsilon_{ij})$$

~~Relative positional embeddings~~  
Arbitrary edge features



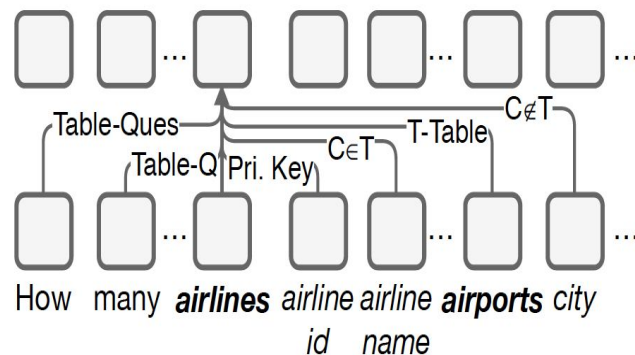
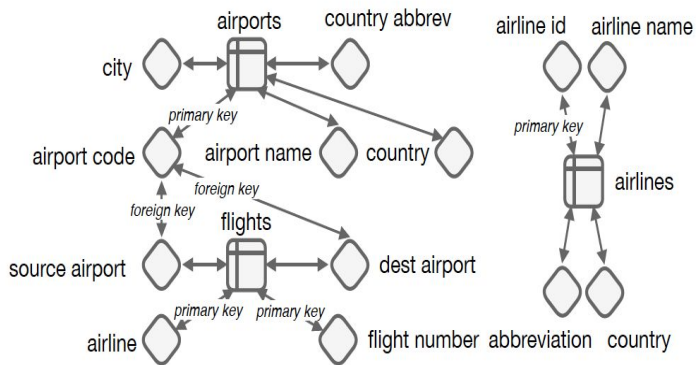
## Pre-existing relations in schema

Type of $x$	Type of $y$	Edge label	Description
Column	Column	SAME-TABLE	$x$ and $y$ belong to the same table.
		FOREIGN-KEY-COL-F	$x$ is a foreign key for $y$ .
		FOREIGN-KEY-COL-R	$y$ is a foreign key for $x$ .
Column	Table	PRIMARY-KEY-F	$x$ is the primary key of $y$ .
		BELONGS-TO-F	$x$ is a column of $y$ (but not the primary key).
Table	Column	PRIMARY-KEY-R	$y$ is the primary key of $x$ .
		BELONGS-TO-R	$y$ is a column of $x$ (but not the primary key).
Table	Table	FOREIGN-KEY-TAB-F	Table $x$ has a foreign key column in $y$ .
		FOREIGN-KEY-TAB-R	Same as above, but $x$ and $y$ are reversed.
		FOREIGN-KEY-TAB-B	$x$ and $y$ have foreign keys in both directions.

# Input preprocessing

$$\mathcal{G}_Q = \langle \mathcal{V}_Q, \mathcal{E}_Q \rangle$$

$$\mathcal{V}_Q = \mathcal{V} \cup \mathcal{Q} = \mathcal{C} \cup \mathcal{T} \cup \mathcal{Q}$$





# Schema Encoding

- Representation of every node in G:  $X = (c_1^{\text{init}}, \dots, c_{|C|}^{\text{init}}, t_1^{\text{init}}, \dots, t_{|T|}^{\text{init}}, q_1^{\text{init}}, \dots, q_{|Q|}^{\text{init}})$ .
  - Glove processed through BiLSTM
  - Bert pre-trained embedding
- Initial representations are independent of relational information
- Encoder applies stack of self-attention

- 

For the cars with 4 cylinders, which model has the largest horsepower?

cars_data							
id	mpg	cylinders	edisp1	horsepower	weight	accelerate	year

car_names		
<i>make_id</i>	model	make

model_list		
model_id	maker	model

car_makers			
id	maker	full_name	country

```
SELECT T1.model
FROM car_names AS T1 JOIN cars_data AS T2
ON T1.make_id = T2.id
WHERE T2.cylinders = 4
ORDER BY T2.horsepower DESC LIMIT 1
```

- Question → Column linking (unknown)
- Question → Table linking (unknown)
- Column → Column foreign keys (known)

Figure 1: A challenging text-to-SQL task from the Spider dataset.





## Memory-schema alignment matrix

**Intuition:** Tables and column names that appear in program P will appear in question Q

$$\tilde{L}_{i,j}^{\text{col}} = \frac{y_i W_Q^{\text{col}} (\mathbf{c}_j^{\text{final}} W_K^{\text{col}} + \mathbf{r}_{ij}^K)^\top}{\sqrt{d_x}} \quad (3)$$

$$\tilde{L}_{i,j}^{\text{tab}} = \frac{y_i W_Q^{\text{tab}} (\mathbf{t}_j^{\text{final}} W_K^{\text{tab}} + \mathbf{r}_{ij}^K)^\top}{\sqrt{d_x}}$$

$$L_{i,j}^{\text{col}} = \underset{j}{\text{softmax}} \{ \tilde{L}_{i,j}^{\text{col}} \} \quad L_{i,j}^{\text{tab}} = \underset{j}{\text{softmax}} \{ \tilde{L}_{i,j}^{\text{tab}} \}$$

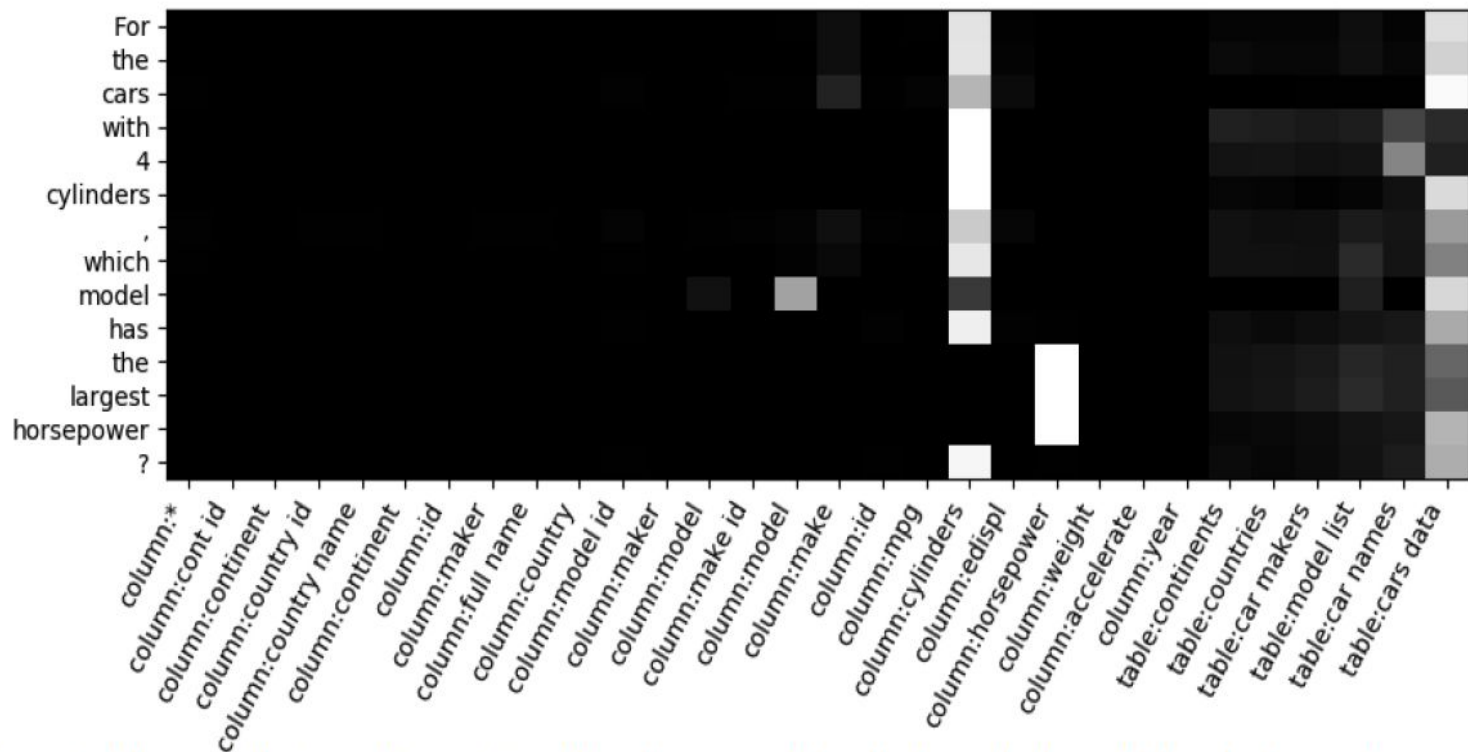


Figure 5: Alignment between the question “For the cars with 4 cylinders, which model has the largest horsepower” and the database `car_1` schema (columns and tables) depicted in Figure 1.



## Decoder/generation of SQL

- Follows the tree structured architecture of [Yin and Neubig \(2017\)](#)
  - expand into a grammar rule : APPLYRULE

$$\text{APPLYRULE}[R] \mid a_{<t}, y) = \text{softmax}_R(g(\mathbf{h}_t))$$

$$f_{\text{LSTM}}([\mathbf{a}_{t-1} \parallel \mathbf{z}_t \parallel \mathbf{h}_{p_t} \parallel \mathbf{a}_{p_t} \parallel \mathbf{n}_{f_t}], \mathbf{m}_{t-1}, \mathbf{h}_{t-1})$$



## Decoder/generation of SQL

- choose a column/table from the schema (terminal):  
SELECTCOLUMN and SELECTTABLE.

$$\tilde{\lambda}_i = \frac{\mathbf{h}_t W_Q^{\text{sc}} (y_i W_K^{\text{sc}})^T}{\sqrt{d_x}} \quad \lambda_i = \underset{i}{\text{softmax}} \{ \tilde{\lambda}_i \}$$

$$\Pr(a_t = \text{SELECTCOLUMN}[i] \mid a_{<t}, y) = \sum_{j=1}^{|y|} \lambda_j L_{j,i}^{\text{col}}$$



## Results on Spider dataset

Model	Dev	Test
IRNet (Guo et al., 2019)	53.2	46.7
Global-GNN (Bogin et al., 2019b)	52.7	47.4
IRNet V2 (Guo et al., 2019)	55.4	48.5
<b>RAT-SQL (ours)</b>	<b>62.7</b>	<b>57.2</b>
<i>With BERT:</i>		
EditSQL + BERT (Zhang et al., 2019)	57.6	53.4
GNN + Bertrand-DR (Kelkar et al., 2020)	57.9	54.6
IRNet V2 + BERT (Guo et al., 2019)	63.9	55.0
RYANSQL V2 + BERT (Choi et al., 2020)	<b>70.6</b>	60.6
<b>RAT-SQL + BERT (ours)</b>	69.7	<b>65.6</b>

Split	Easy	Medium	Hard	Extra Hard	All
<i>RAT-SQL</i>					
<b>Dev</b>	80.4	63.9	55.7	40.6	62.7
<b>Test</b>	74.8	60.7	53.6	31.5	57.2
<i>RAT-SQL + BERT</i>					
<b>Dev</b>	86.4	73.6	62.1	42.9	69.7
<b>Test</b>	83.0	71.3	58.3	38.4	65.6



### Easy

What is the number of cars with more than 4 cylinders?

```
SELECT COUNT(*)  
FROM cars_data  
WHERE cylinders > 4
```

### Meidum

For each stadium, how many concerts are there?

```
SELECT T2.name, COUNT(*)  
FROM concert AS T1 JOIN stadium AS T2  
ON T1.stadium_id = T2.stadium_id  
GROUP BY T1.stadium_id
```

### Hard

Which countries in Europe have at least 3 car manufacturers?

```
SELECT T1.country_name  
FROM countries AS T1 JOIN continents  
AS T2 ON T1.continent = T2.cont_id  
JOIN car_makers AS T3 ON  
T1.country_id = T3.country  
WHERE T2.continent = 'Europe'  
GROUP BY T1.country_name  
HAVING COUNT(*) >= 3
```

### Extra Hard

What is the average life expectancy in the countries where English is not the official language?

```
SELECT AVG(life_expectancy)  
FROM country  
WHERE name NOT IN  
(SELECT T1.name  
FROM country AS T1 JOIN  
country_language AS T2  
ON T1.code = T2.country_code  
WHERE T2.language = "English"  
AND T2.is_official = "T")
```

Figure 3: SQL query examples in 4 hardness levels.



## Results on WikiSQL

Model	Dev		Test	
	LF Acc%	Ex. Acc%	LF Acc%	Ex. Acc%
IncSQL (Shi et al., 2018)	49.9	84.0	49.9	83.7
MQAN (McCann et al., 2018)	76.1	82.0	75.4	81.4
RAT-SQL (ours)	73.6	79.5	73.3	78.8
Coarse2Fine (Dong and Lapata, 2018)	72.5	79.0	71.7	78.5
PT-MAML (Huang et al., 2018)	63.1	68.3	62.8	68.0



## Ablation Study

Model	Accuracy (%)
<b>RAT-SQL + value-based linking</b>	<b>60.54 <math>\pm</math> 0.80</b>
RAT-SQL	55.13 $\pm$ 0.84
w/o schema linking relations	40.37 $\pm$ 2.32
w/o schema graph relations	35.59 $\pm$ 0.85





# Error Analysis

- 39% of errors -> a limitation of schema linking
- 29% of errors -> Need of in-domain fine-tuning
  - ‘Older than 21’ -> Age > 21 or age < 21
- 18% of errors -> equivalent implementations of NL but a different SQL syntax

Model	Exact Match	Correctness
<b>RAT-SQL</b>	0.59	0.81
<b>RAT-SQL + BERT</b>	0.67	0.86

Table 7: Consistency of the two RAT-SQL models.



## Key takeaways

- RAT-SQL presented a unified framework to address schema representation and schema linking challenges.
- Contextual representation of question with schema in encoder helps.
- Combining predefined hard schema relations with soft alignment on sequence elements (different from GNN) in encoder added value!