

# Grounding Natural Language to SQL Translation with Data-Based Self-Explanations



Yuankai Fan, Tonghui Ren, Can Huang, Zhenying He, X. Sean Wang  
Fudan University, TeleAI

## Background

Natural Language Interfaces for Databases empower non-technical users to interact with data using natural language (NL). Advanced approaches, utilizing either sequence-to-sequence or more recent sophisticated large-scale language models (LLM), typically implement NL to SQL (NL2SQL) translation in an **end-to-end fashion**. However, like humans, these end-to-end translation models may not always generate the best SQL output on their first try.

**Unsatisfied Results:** The plateauing accuracy (below 80%) observed when the beam size or the number of chat completions is set to 1, they may **fail** to generate best-quality translations in their **initial attempts**.

**However**, expanding the search space through wider beams or more chat completions steadily improves accuracy, without necessitating modifications to the underlying model architectures, which could be the potential.

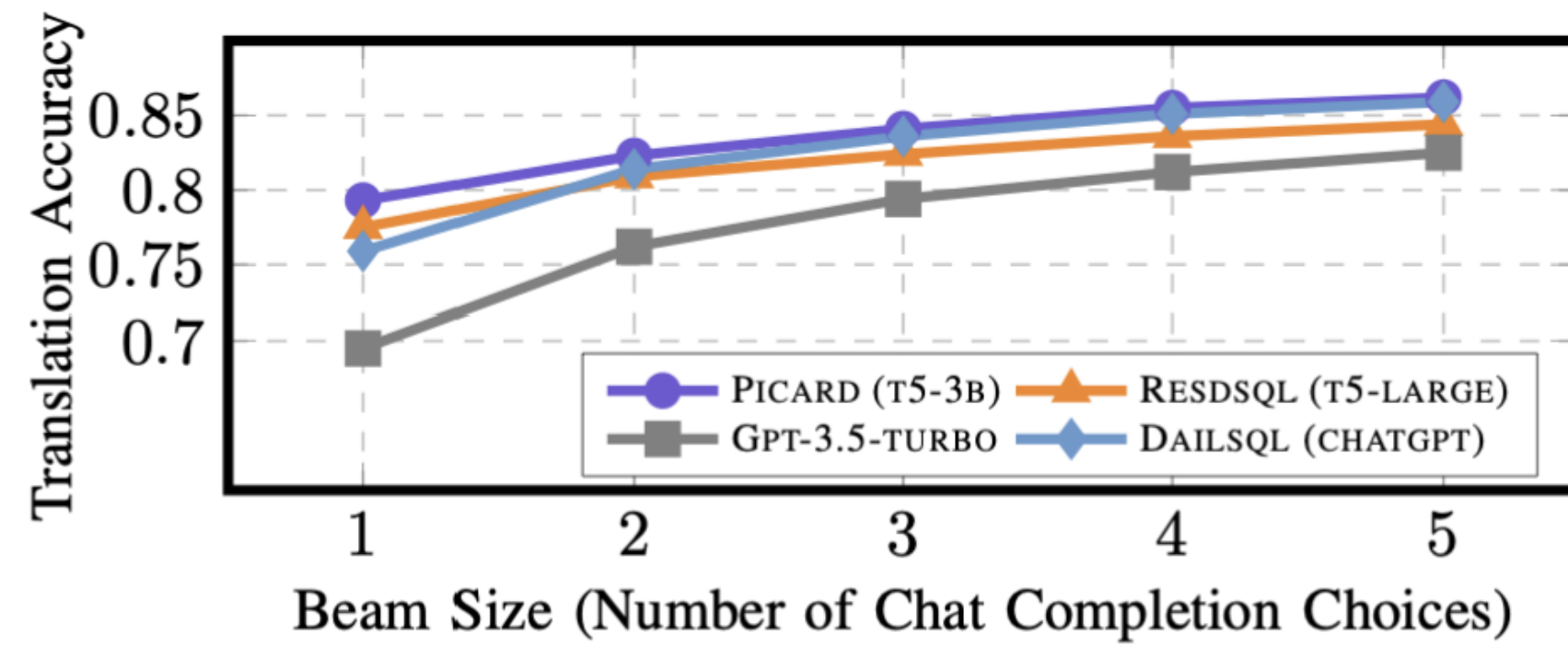
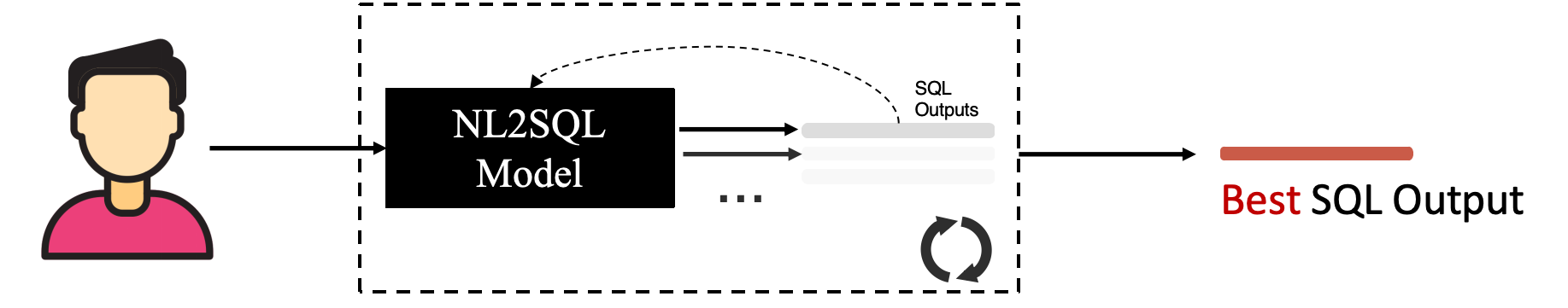


Fig. 1: Translation accuracy on SPIDER validation set with varied beam sizes (or chat completion choices). Accuracy for beam sizes (or chat completions) > 1 is evaluated by matching any beam result.

**Motivation:** Can we build a **self-provided feedback loop** along with NL2SQL, to improve end-to-end performance?



**Question:** Is SQL2NL Back-Translation Sufficient to establish an NL-to-SQL-to-NL translation lifecycle?

**Answer:** ❌ NOT reliable!!!

**Fase positive feedback** for NL questions is observed, which leads to incorrect translation result ultimately.

aid	flno	origin	destination	aid	flno	name	distance
9	2	Los Angeles	Tokyo	1		Boeing 747-400	8430
3	7	Los Angeles	Sydney	2		Boeing 737-800	3383
13	13	Los Angeles	Chicago	3		Airbus A330-300	7170
10	68	Chicago	New York	4		British Aerospace Jetstream 41	1502
76	36	Chicago	Los Angeles	5		Embraer E30-145	1530
7	33	Los Angeles	Honolulu	6		SAAB 340	2128
5	34	Los Angeles	Honolulu	7		Piper Archer III	520
1	99	Los Angeles	Washington D.C.	8		Tupolev 154	4103
2	346	Los Angeles	Dallas	9		Lockheed L1011	6900
6	387	Los Angeles	Boston	10		Boeing 757-300	4010

**[Question]:** Show all flight numbers with aircraft Airbus A340-300  
**[Ground-truth SQL]:**  
 SELECT T1.fno FROM Flight AS T1  
 JOIN Aircraft AS T2 ON T1.aid = T2.aid  
 WHERE T2.name = 'Airbus A340-300'

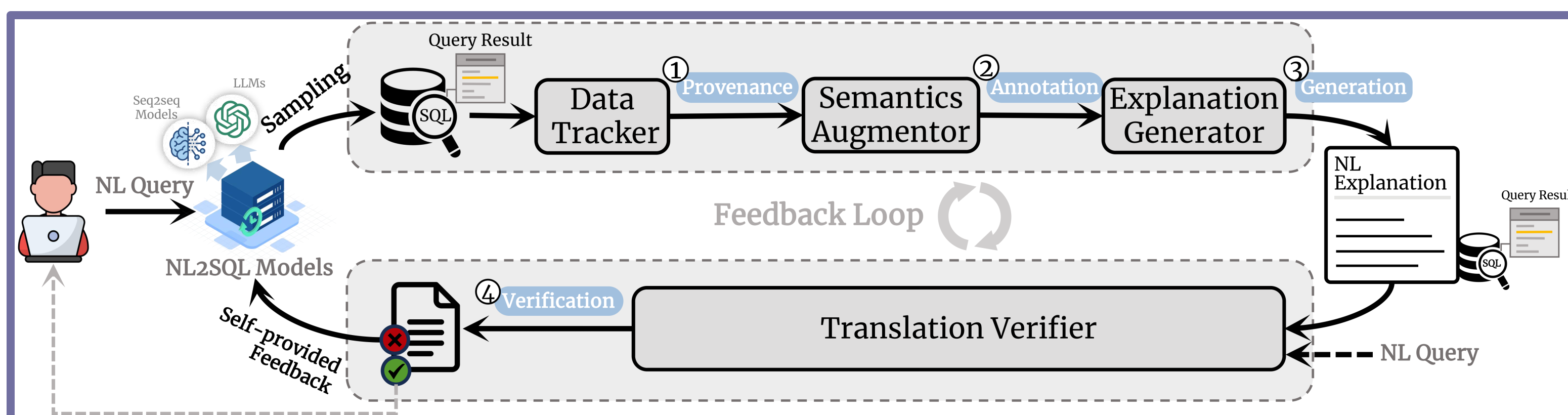
**[Translated SQL]:**  
 SELECT count(\*) FROM Flight AS T1  
 JOIN Aircraft AS T2 ON T1.aid = T2.aid  
 WHERE T2.name = 'Airbus A340-300'

**[Query Result]:**  
 count(\*)  
 2

**[SQL2NL-generated Explanation]:**  
 Find the number of flights that are operated by aircraft with the name Airbus A340-300.  
**[Feedback] (Incorrect):**  
 The translated SQL above is correct. ❌

## CycleSQL

### OVERVIEW



- ✓ **Data Tracker:** Track provenance information of the to-explained query result to retrieve **data-level information**
- ✓ **Semantics Augmentor:** Enhance provenance by associating it with **operation-level semantics**
- ✓ **Explanation Generator:** Interpret provenance to **explanation**
- ✓ **Translation Verifier:** Verify **correctness** of NL2SQL translation based on the generated NL explanation

### METHODOLOGIES

#### Data Provenance

Capture provenance information using **query rewriting**, with heuristic rules:

① **Result Transformation Rule:**

Translate query result into **WHERE clause** to specify as condition

② **Projection Enhancement Rule:**

Extract **all columns** used in the original query, as additional projection columns

③ **Aggregation Deconstructive Rule:**

Eliminate aggregations to unravel lineage

**Translated SQL:**  
 SELECT count(\*) FROM Flight AS T1  
 JOIN Aircraft AS T2 ON T1.aid = T2.aid  
 WHERE T2.name = 'Airbus A340-300'

**Query Rewriting:**

**Rule 1:**  
 SELECT T2.name  
 FROM Flight AS T1  
 JOIN Aircraft AS T2 ON T1.aid = T2.aid  
 WHERE T2.name = 'Airbus A340-300'

**Rule 2:**  
 SELECT count(\*)  
 FROM Flight AS T1  
 WHERE T1.aid = T2.aid

**Rule 3:**  
 SELECT count(\*)  
 FROM Flight AS T1  
 JOIN Aircraft AS T2 ON T1.aid = T2.aid  
 WHERE T2.name = 'Airbus A340-300'

**Rewritten SQL:**  
 SELECT count(\*)  
 FROM Flight AS T1  
 JOIN Aircraft AS T2 ON T1.aid = T2.aid  
 WHERE T2.name = 'Airbus A340-300'

**Provenance:** <A3, F2, F3>

Aircraft	aid	name
A1	1	Boeing 747-400
A2	2	Boeing 737-800
A3	3	Airbus A340-300

Flight	aid	flno	origin	destination
F1	9	2	Los Angeles	Tokyo
F2	3	7	Los Angeles	Sydney
F3	13	13	Los Angeles	Chicago

#### Semantics Enrichment

Main idea of integrating the operation-level semantics of the SQL queries is to better reflect **user query intent**.

As different parts of the query may **“contribute”** to different parts of the provenance information, we treat the SQL query as a text string and divides the string into **chunks** that correspond to each clause in the SQL query.

**Flight-Aircraft**

① WHERE name = 'Airbus A340-300'

② SELECT count(\*)

tupleID	Flight.aid	Aircraft.name	Flight.fno	...
<A3, F2>	3	Virgin America	7	
<A3, F3>	3	Virgin America	13	

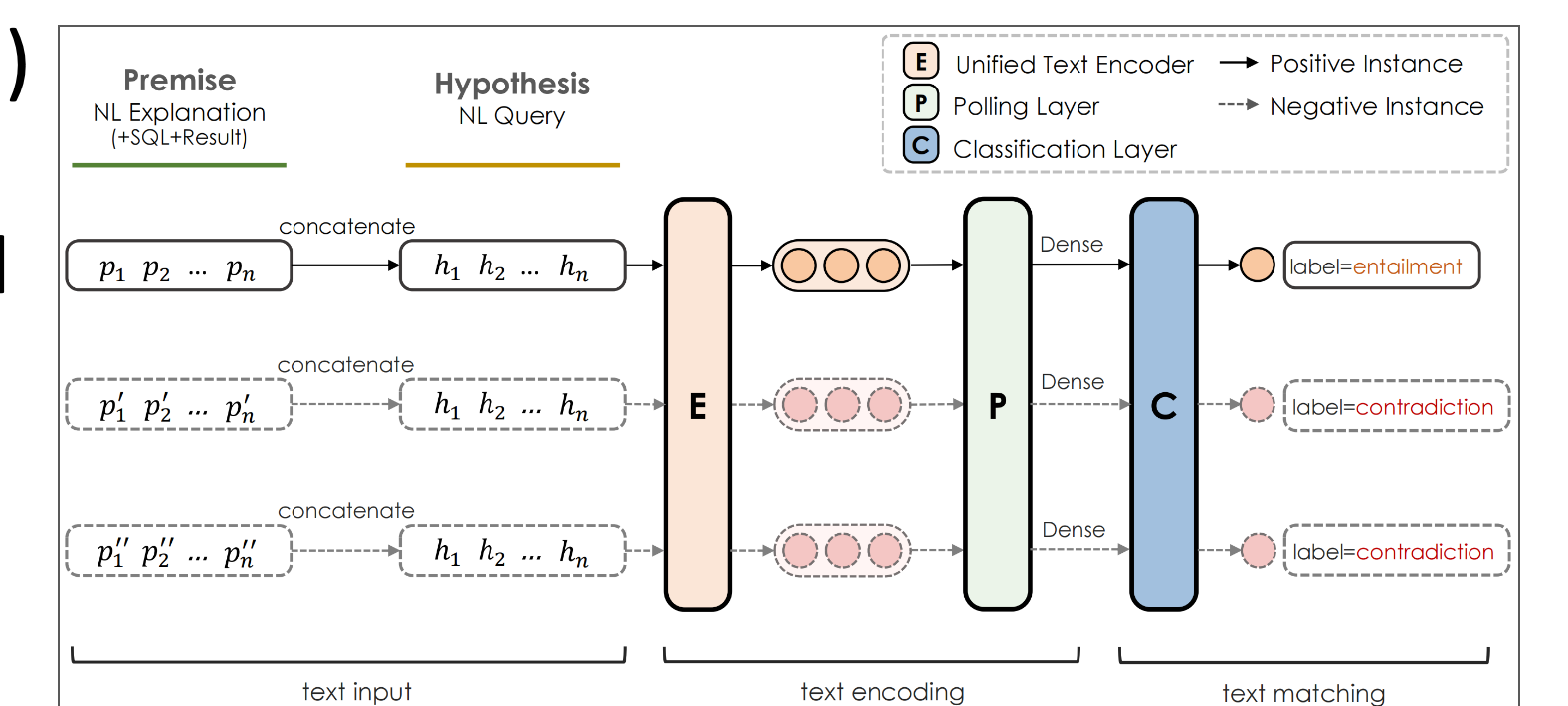
#### Explanation Text Generation

To interpret the enriched provenance to a textual expression in free-text forms, we employ a **rule-based method** to synthesize the NL explanation using the enriched provenance information:

*The query returns a result with one column of aggregation type (count) and one row. For flights with aircraft, named Airbus A340-300, there are 2 flights in total.*

#### NL2SQL Translation Validation

Drawing inspiration from recent achievements in the NLP field, we thus formulate the translation validation problem as a **textual entailment task (NLI)** and leverage deep learning techs to build the **NLI model** for this purpose.



## Evaluation

### Overall Results

Models	SPIDER				SPIDER VARIANTS				SCIENCEBENCHMARK								
	VALIDATION		TEST		REALISTIC		SYN		DK		ONCOMX		CORDIS		SDSS		
	EM	EX	TS	EM	EX	EM	EX	TS	EM	EX	EM	EX	EX	EX	EX	EX	
Seq2seq-based Models																	
SMBOP	Base	74.5	77.8	72.5	69.5	71.1	60.2	61.4	56.1	60.2	64.9	58.3	54.0	59.1	16.2	16.0	7.0
	+CYCLESQ	74.1	78.8	73.1	-	-	59.8	62.2	56.7	59.1	65.5	58.6	53.3	59.8	17.2	16.0	8.0
PICARD <sub>3b</sub>	Base	75.9	79.3	72.8	72.4	75.1	68.9	71.7	64.6	65.6	69.8	63.2	49.7	63.2	32.3	26.0	8.0
	+CYCLESQ	76.1	79.9	72.8	-	-	68.7	71.9	64.6	65.9	70.8	63.6	50.3	63.6	33.3	24.0	8.0
RESDSQL <sub>LARGE</sub>	Base	74.0	77.5	71.6	66.8	73.5	67.3	69.5	61.6	58.8	65.3	59.6	50.1	58.1	33.0	29.0	4.0
	+CYCLESQ	75.1	80.5	74.0	68.9	77.1	69.1	73.6	64.2	61.6	69.6	63.2	50.5	61.5	35.0	32.0	8.0
RESDSQL <sub>3b</sub>	Base	76.0	79.4	73.5	70.2	78.4	68.9	72.2	64.8	62.9	69.5	63.2	52.0	59.3	34.1	28.0	6.0
	+CYCLESQ	76.8	82.0	76.0	72.5	81.6	71.3	76.5	67.5	65.5	73.3	66.6	52.5	61.3	37.4	31.0	8.0
LLM-based Models																	
GPT-3.5-TURBO (5-SHOTS)	Base	43.8	72.8	61.7	44.7	69.4	39.4	65.4	52.8	34.0	61.0	49.8	39.1	60.6	34.3	30.0	10.0
	+CYCLESQ	49.2	77.8	66.2	50.6	75.1	46.5	71.1	56.7	40.0	66.7	55.4	43.7	65.6	43.4	34.0	12.0
GPT-4 (5-SHOTS)	Base	51.9	77.9	71.2	-	-	46.7	68.6	55.3	45.2	75.0	64.4	53.7	68.5	51.5	40.0	14.0
	+CYCLESQ	58.7	79.8	73.1	-	-	49.0	70.6	56.9	46.2	76.0	66.3	57.4	68.5	56.6	43.0	15.0
CHESS	Base	23.4	41.1	37.5	-	-	21.2	39.7	36.6	17.4	37.0	32.5	19.5	35.6	64.6	46.0	40.0
	+CYCLESQ	25.6	42.3	38.2	-	-	23.1	41.1	37.9	19.3	38.1	33.0	20.2	36.1	65.7	52.0	40.0
DAISQL <sub>3.5</sub>	Base	65.1	81.3	74.1	-	-	-	-	-	-	-	-	-	-	-	-	-
	+CYCLESQ	67.2	81.8	74.3	-	-	-	-	-	-	-	-	-	-	-	-	-

- ✓ CycleSQL **consistently** improves over all base models, particularly evident in the EX and TS metrics.
- ✓ Remarkable outcome: CycleSQL+RESDSQL-3B on EX metric **surpasses** the best-reported result among leading Seq2seq models on Spider

### Analysis

Model	Easy	Medium	Hard	Extra Hard	
SMBOP	Base	90.7	82.7	70.7	52.4
	+CYCLESQ	90.7	84.1(+1.4)	69.5(+1.2)	53.0(+0.6)
PICARD <sub>3b</sub>	Base	95.6	85.4	67.8	50.6
	+CYCLESQ	95.6	86.1(+0.7)	69.5(+1.7)	50.6
RESDSQL <sub>LARGE</sub>	Base	92.3	83.4	66.1	51.2
	+CYCLESQ	93.5(+1.2)	86.1(+2.7)	73.0(+6.9)	53.0(+1.8)
RESDSQL <sub>3b</sub>	Base	94.0	85.7	65.5	55.4
	+CYCLESQ	94.0	89.0(+3.3)	74.7(+9.2)	53.0(+0.4)
GPT-3.5-TURBO (5-SHOTS)	Base	84.3	78.5	65.5	48.2
	+CYCLESQ	86.3(+2.0)	83.0(+4.5)	73.0(+7.5)	56.0(+7.8)
GPT-4 (5-SHOTS)	Base	90.3	84.3	63.8	56.6
	+CYCLESQ	90.7(+0.4)	85.4(+1.1)	66.7(+2.9)	59.6(+3.0)
CHESS	Base	70.2	25.3	39.7	19.3
	+CYCLESQ	71.6(+1.4)	25.6(+0.3)	41.1(+1.4)	21.6(+2.3)
DAISQL <sub>3.5</sub>	Base	91.1	86.5	77.0	57.2
	+CYCLESQ	91.1	86.8(+0.3)	77.6(+0.6)	59.0(+1.8)

- ✓ Significant Improvement over **Extra-Hard-Queries** on LLM Models
- ✓ Good **Balance** between Translation and Inference Latency
- ✓ CycleSQL Feedback is **Better** than SQL2NL Feedback

