



# 41st IEEE International Conference on Data Engineering

HONG KONG SAR, CHINA | MAY 19 – 23, 2025

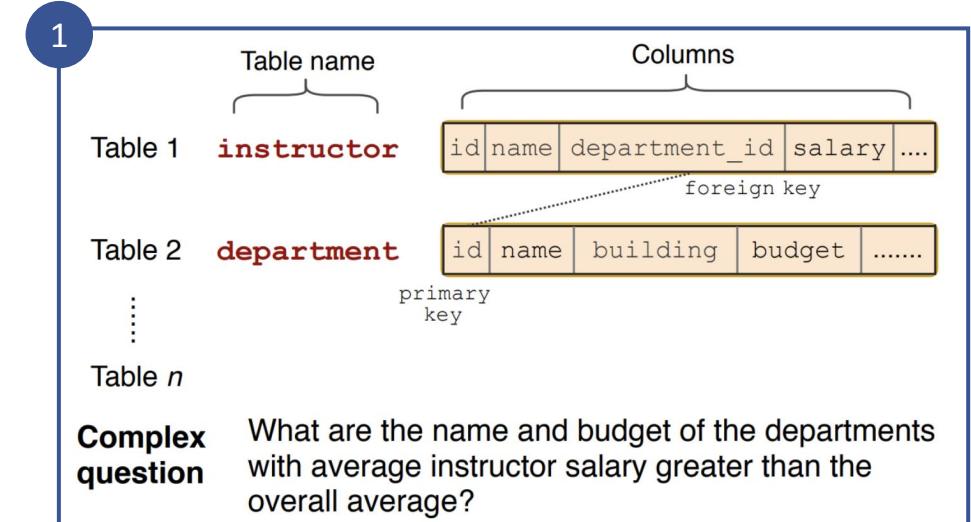


# Grounding Natural Language to SQL Translation with Data-Based Self-Explanations

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

# Natural Language to SQL (NL2SQL)

- Writing SQLs to query databases is NOT easy for non-SQL-savvies
- Natural Language to SQL (**NL2SQL**) comes in rescue:
  - ① Given a natural language query and the database schema
  - ② Generate the corresponding SQL query



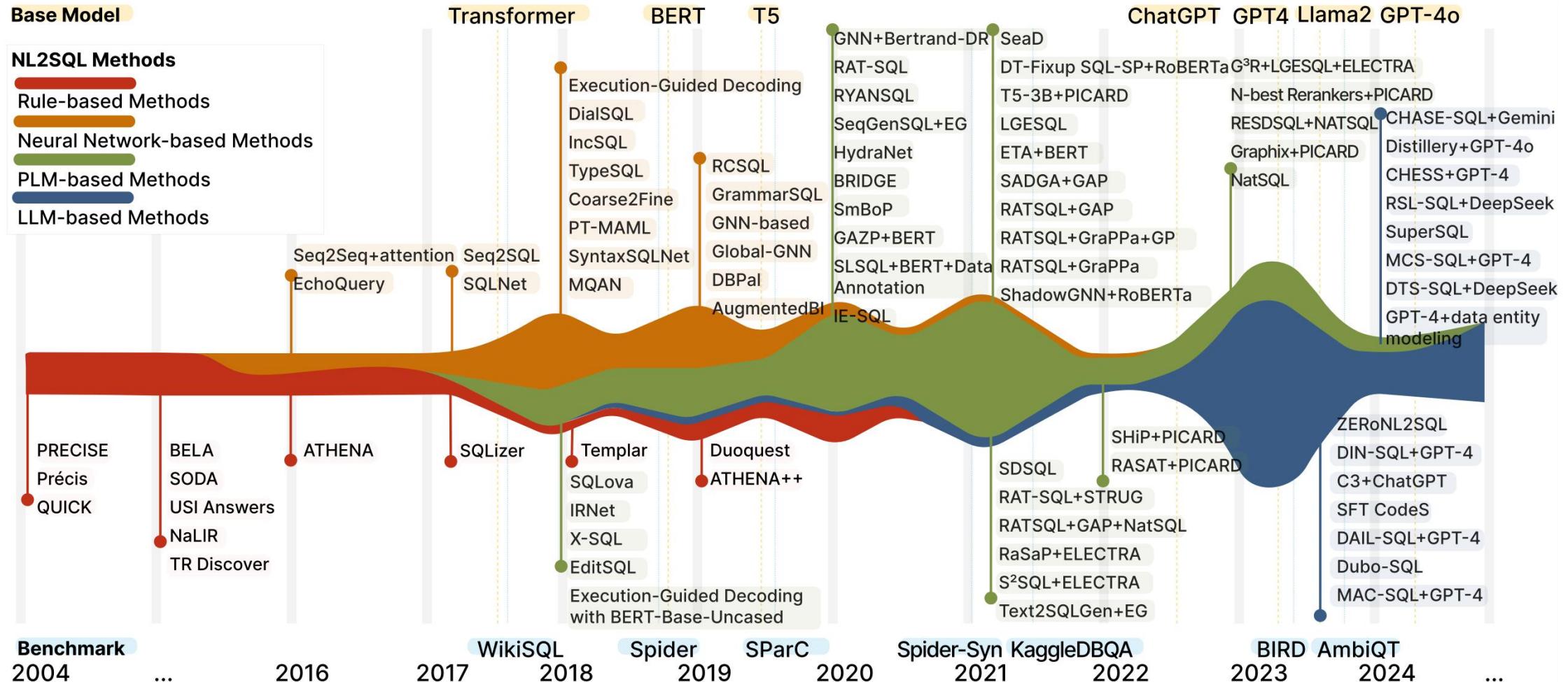
Complex SQL

```
SELECT T2.name, T2.budget
FROM instructor as T1 JOIN department as
T2 ON T1.department_id = T2.id
GROUP BY T1.department_id
HAVING avg(T1.salary) >
(SELECT avg(salary) FROM instructor)
```

2

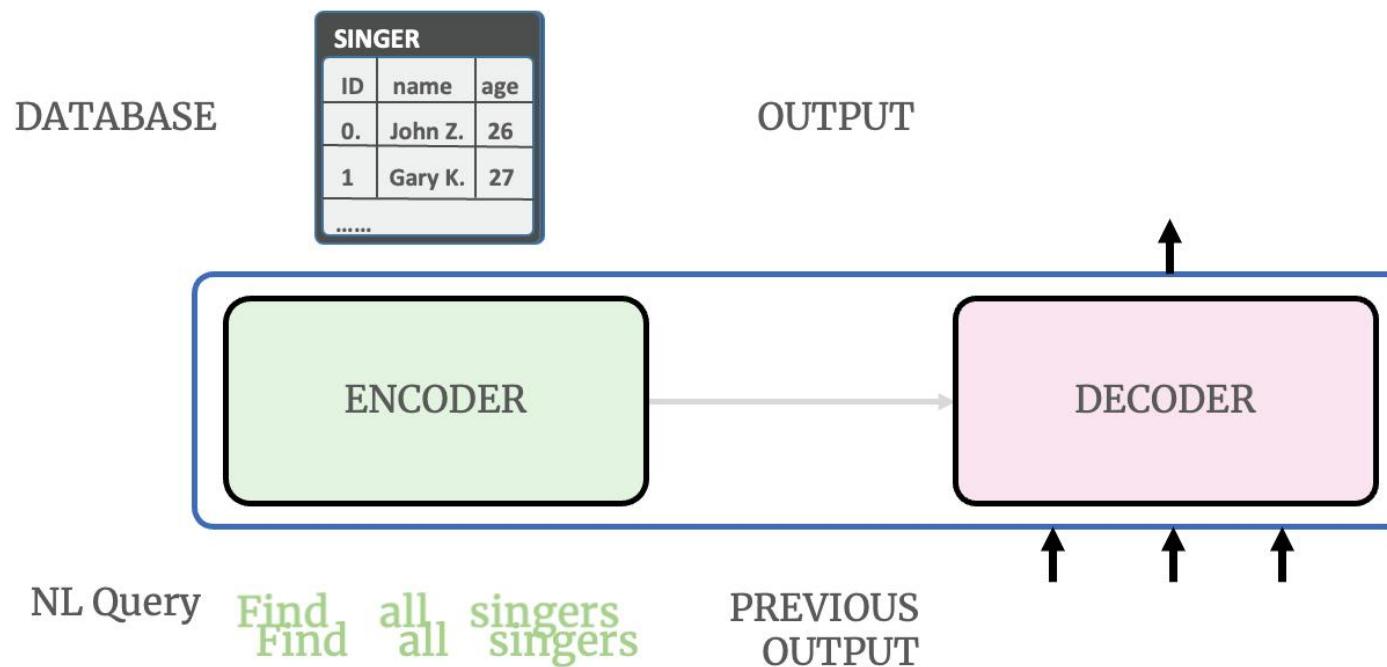
Generated SQL

# NL2SQL Evolution



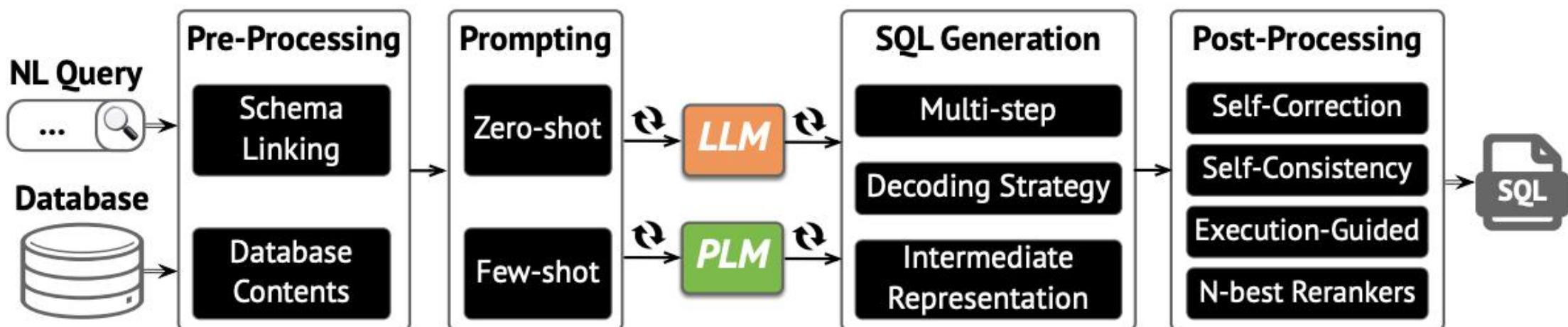
# Pre-LLM Era Approach - Seq2Seq

- Consider as machine translation problem
- Based on **Sequence-to-sequence** Encoder-Decoder framework



# Post-LLM Era Approach – LLMs

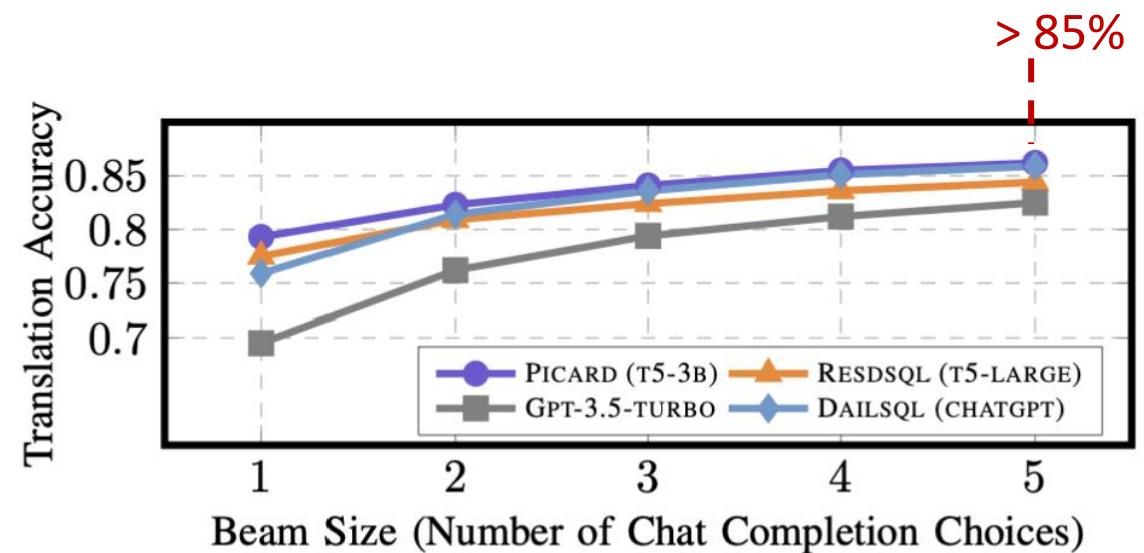
- Build NL2SQL system with **LLMs**
  - ① Pre-processing
  - ② SQL Generation
  - ③ Post-Processing



# Unsatisfied Results

---

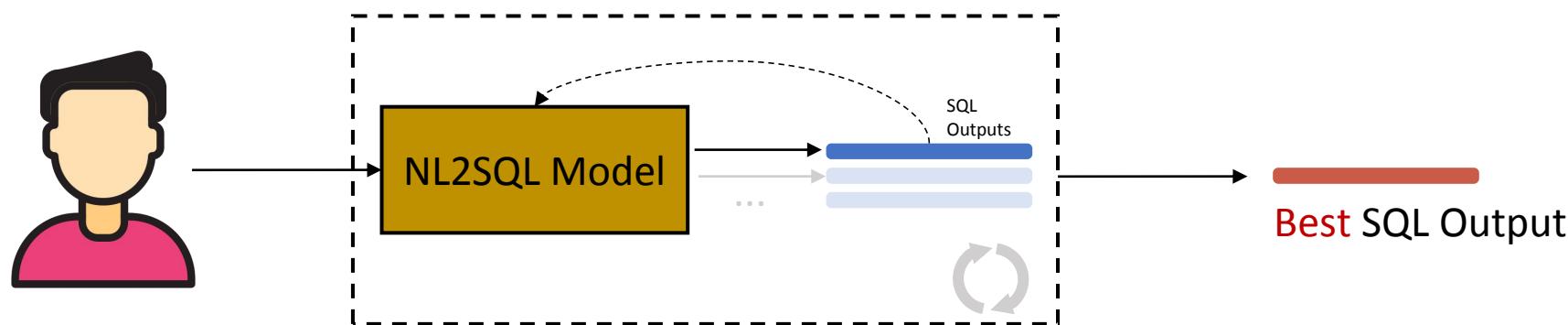
- Mainstream approaches (either Seq2seq models or LLMs) typically implement NL2SQL translation in **an end-to-end fashion**
- These models may benefit from broader exploration options over successive attempts



# Idea: Feedback Loop in NL2SQL

---

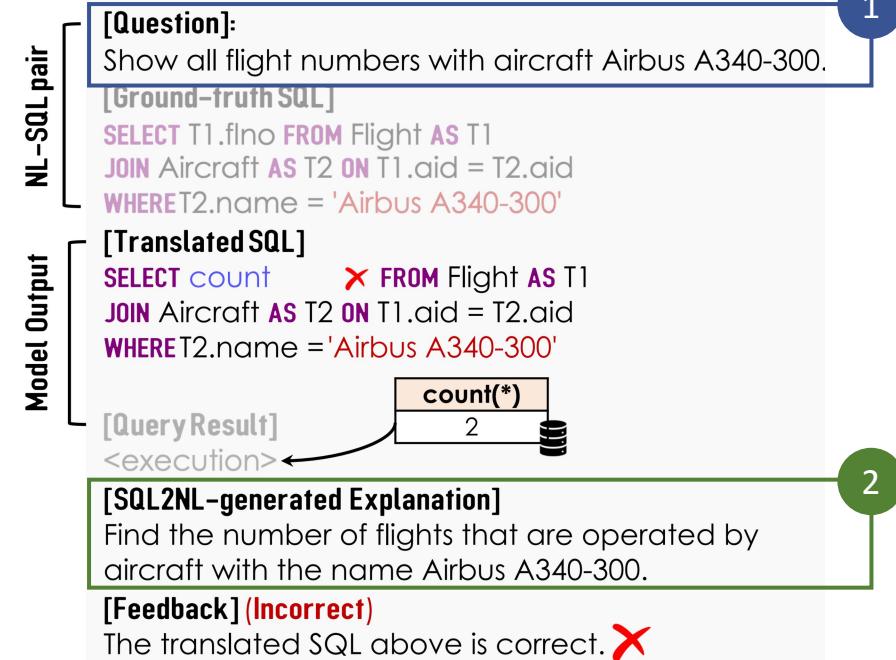
- Can we build a **self-provided feedback loop** along with NL2SQL, to improve end-to-end performance?



# Is SQL2NL Back-Translation Sufficient?

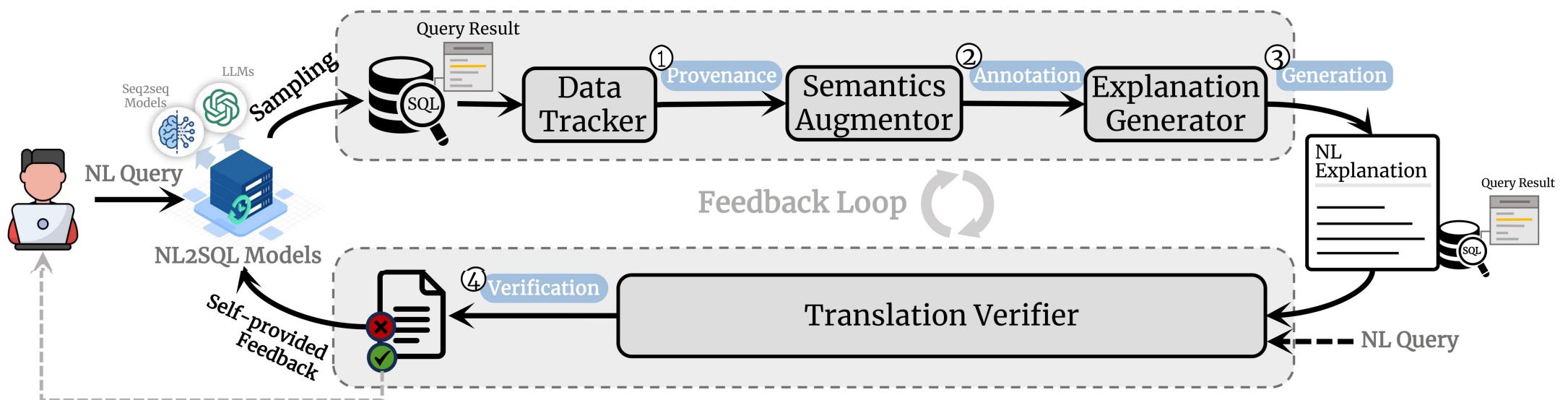
- Use **SQL2NL** to establish an **NL-to-SQL-to-NL** translation lifecycle?
  - **NOT reliable !!!**
- Fase Positive feedback ② for NL Question ①

aid	flno	origin	destination	aid	name	distance
9	2	Los Angeles	Tokyo	1	Boeing 747-400	8430
3	7	Los Angeles	Sydney	2	Boeing 737-800	3383
3	13	Los Angeles	Chicago	3	Airbus A340-300	7120
10	68	Chicago	New York	4	British Aerospace Jetstream 41	1502
9	76	Chicago	Los Angeles	5	Embraer ERJ-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

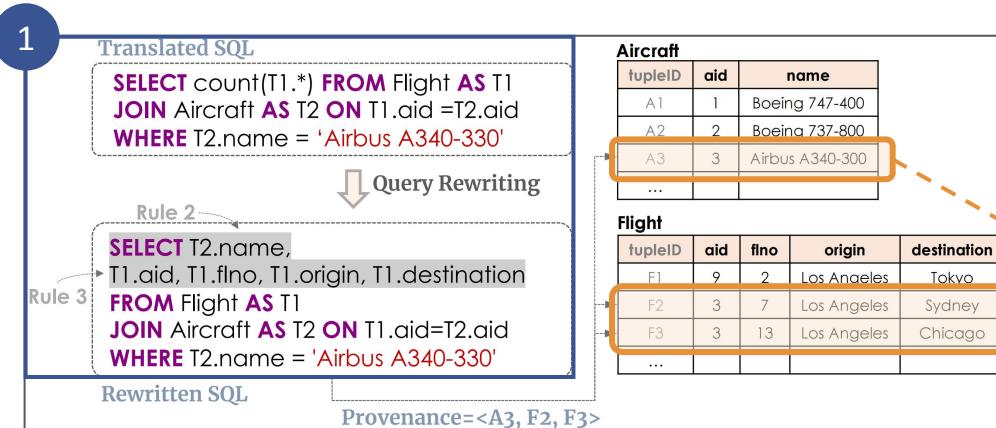


# CycleSQL: Self-Provided Feedback in NL2SQL

- ① Data Provenance
- ② Semantics Enrichment
- ③ Explanation Generation
- ④ Translation Verification

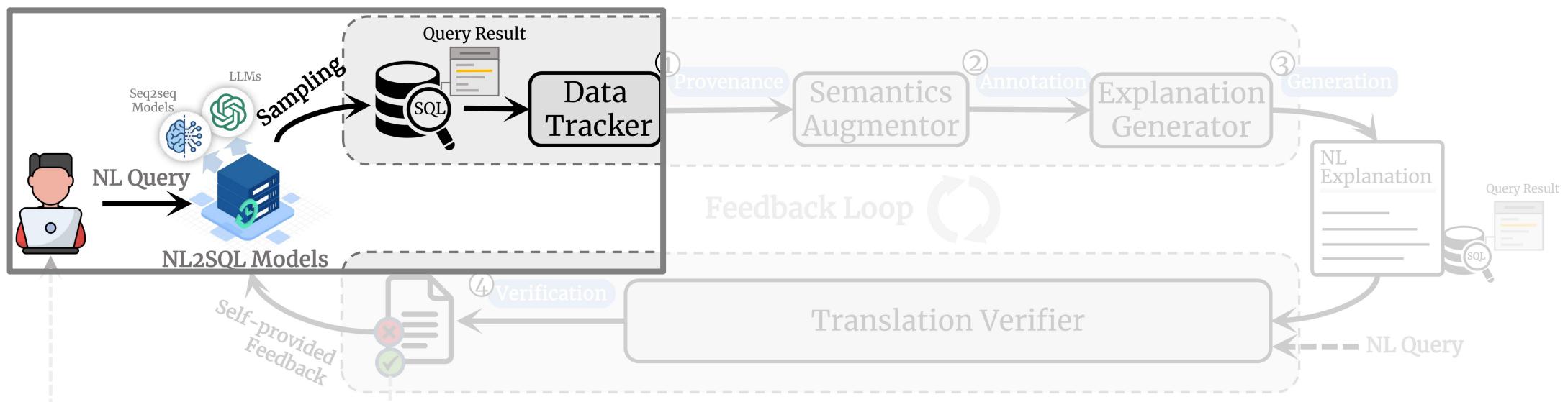


# Data Provenance



Capture provenance information using query rewriting 1

Provenance Tuples



# Semantics Enrichment

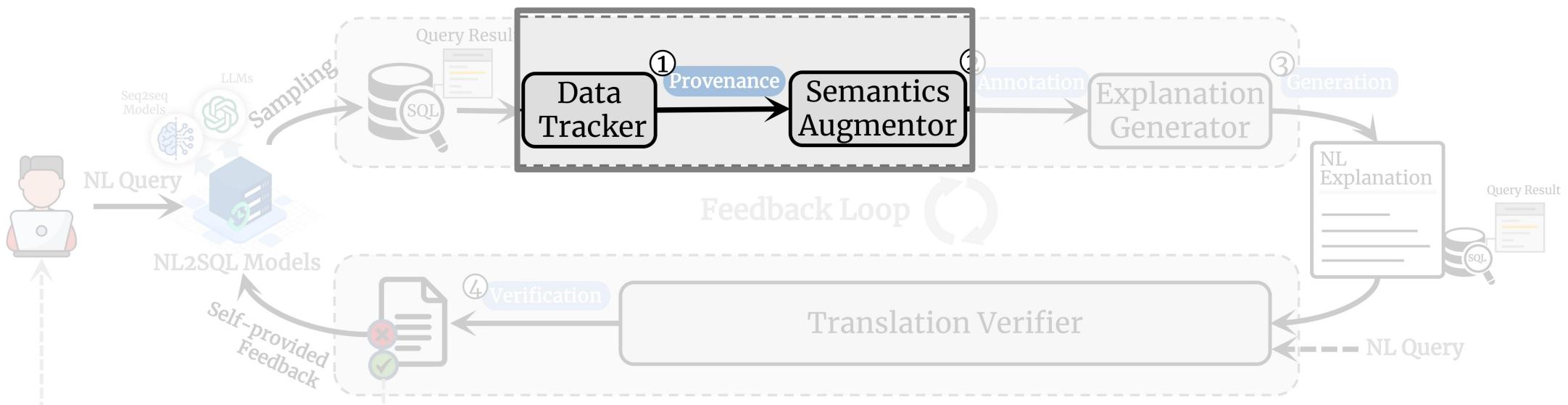
A screenshot of a database interface showing a query result for 'Flight-Aircraft'. The query is:

```
SELECT count(*)  
WHERE name = 'Airbus A340-300'
```

The result table has columns: tupleID, Flight.aid, Aircraft.name, Flight.flno, ...

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

Integrate **operation-level semantics** ① ② of the SQL queries to better reflect user query intent



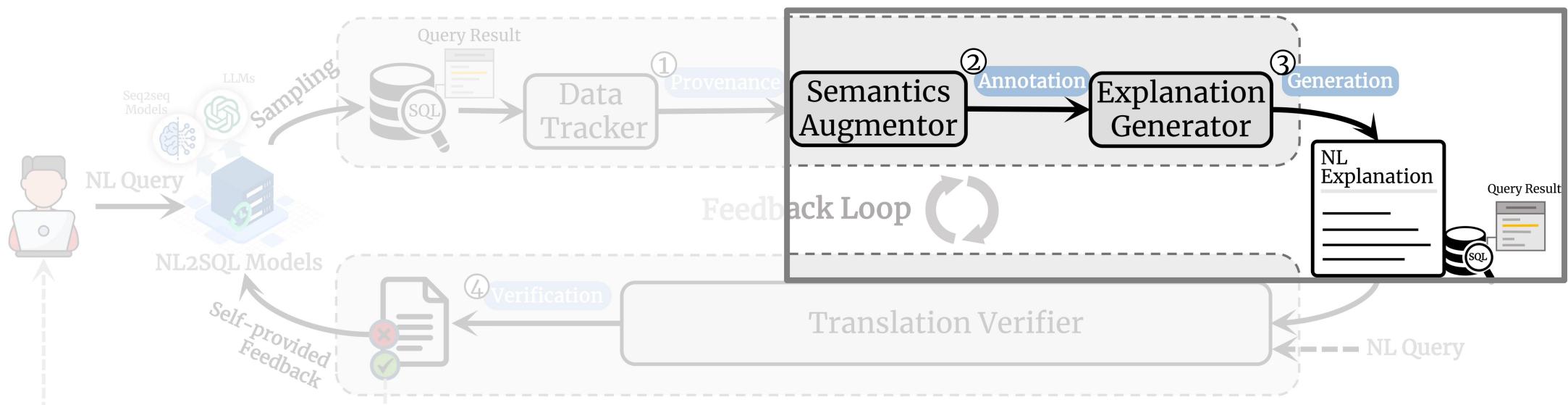
# NL Explanation Text Generation

Synthesize **NL explanation** based on a rule-based method

## Synthesized Explanation Text Example:

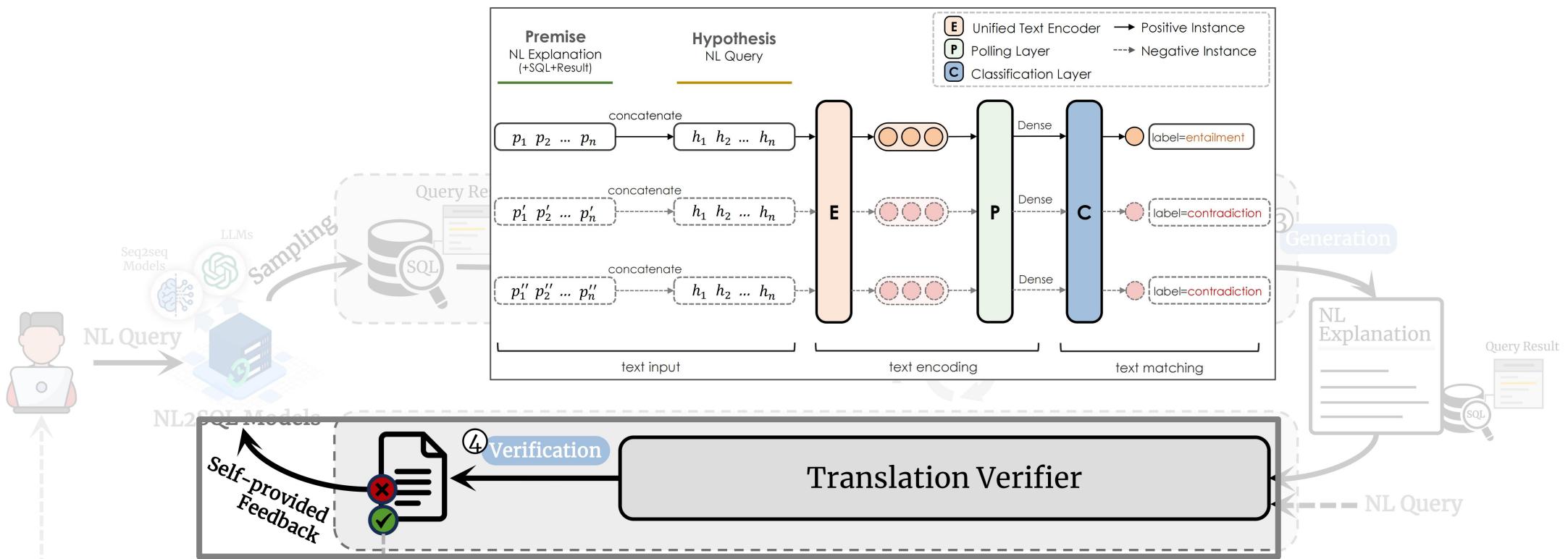
The query returns a result with one column of aggregation type (count) and one row.

For lights with aircraft, named Airbus A340-300, there are 2 flights in total.



# NL2SQL Validation

- Formulate the translation validation problem as a **textual entailment task**
- Use textual entailment model to determine if the translation is correct or not



# Experimental Settings

---

## ■ Benchmarks

- Spider Spider and its variants (Spider-DK/Spider-Syn/Spider-Realistic)
- ScienceBenchmark

## ■ Baselines

- **Seq2seq-based:** SmBoP/PICARD/RESDSQL
- **LLM-based:** GPT-3.5-Turbo/GPT-4/CHESS/DAILSQ

## ■ Metrics

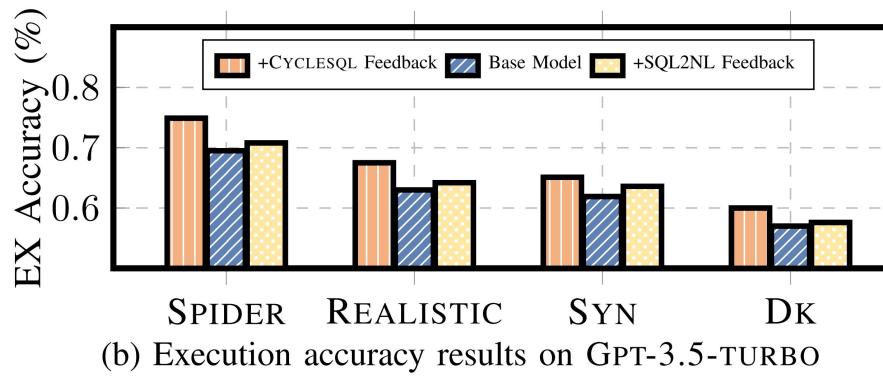
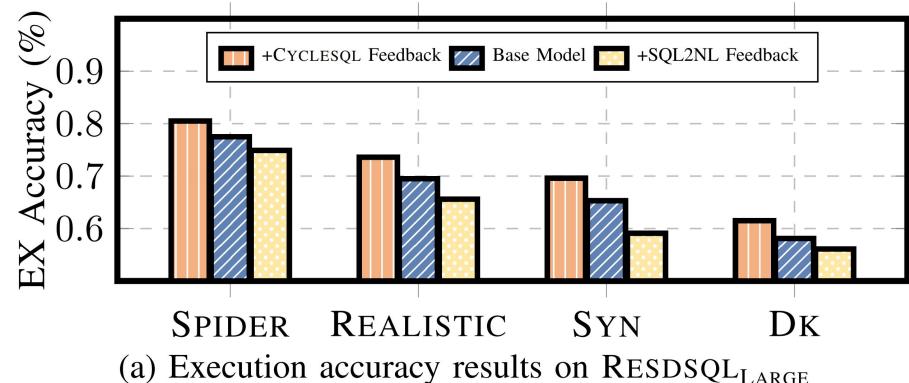
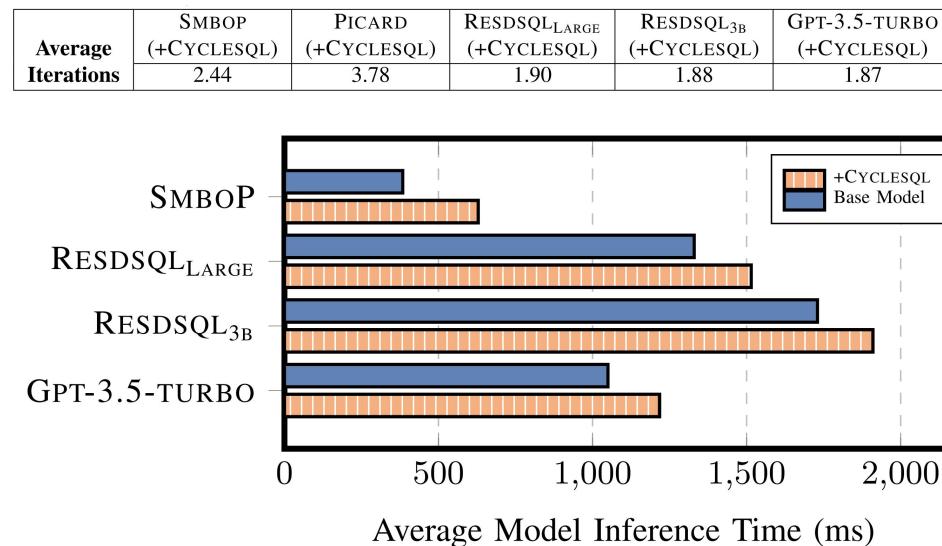
- Syntactic Accuracy (EM): *exact match* the ground truth
- **Execution Accuracy (EX):** *execution result match*
- Test Suit Accuracy (TS): Similar to EX, but more robust

# Overall Results

- CycleSQL **consistently** improves over all base models

# Break-Down Results

Model		Easy	Medium	Hard	Extra Hard
<b>SMBOP</b>	Base	90.7	82.7	70.7	52.4
	+CYCLESQ	90.7	<b>84.1</b> ( $\uparrow 1.4$ )	69.5( $\downarrow 1.2$ )	53.0( $\uparrow 0.6$ )
<b>PICARD<sub>3B</sub></b>	Base	95.6	85.4	67.8	50.6
	+CYCLESQ	95.6	86.1( $\uparrow 0.7$ )	69.5( $\uparrow 1.7$ )	50.6
<b>RESDSQL<sub>LARGE</sub></b>	Base	92.3	83.4	66.1	51.2
	+CYCLESQ	93.5( $\uparrow 1.2$ )	86.1( $\uparrow 0.7$ )	<b>73.0</b> ( $\uparrow 6.9$ )	53.6( $\uparrow 2.4$ )
<b>RESDSQL<sub>3B</sub></b>	Base	94.0	85.7	65.5	55.4
	+CYCLESQ	94.0	<b>89.0</b> ( $\uparrow 3.3$ )	<b>74.7</b> ( $\uparrow 9.2$ )	53.0( $\downarrow 0.4$ )
<b>GPT-3.5-TURBO</b> (5-SHOTS)	Base	84.3	78.5	65.5	48.2
	+CYCLESQ	86.3( $\uparrow 2.0$ )	<b>83.0</b> ( $\uparrow 4.5$ )	<b>73.0</b> ( $\uparrow 7.5$ )	<b>56.0</b> ( $\uparrow 7.8$ )
<b>GPT-4</b> (5-SHOTS)	Base	90.3	84.3	63.8	56.6
	+CYCLESQ	90.7( $\uparrow 0.4$ )	85.4( $\uparrow 1.1$ )	<b>66.7</b> ( $\uparrow 2.9$ )	<b>59.6</b> ( $\uparrow 3.0$ )
<b>CHESS</b>	Base	70.2	25.3	39.7	19.3
	+CYCLESQ	<b>71.6</b> ( $\uparrow 1.4$ )	25.6( $\uparrow 0.3$ )	<b>41.1</b> ( $\uparrow 1.4$ )	<b>21.6</b> ( $\uparrow 2.3$ )
<b>DAILSQ</b> <sub>3.5</sub>	Base	91.1	86.5	77.0	57.2
	+CYCLESQ	91.1	<b>86.8</b> ( $\uparrow 0.3$ )	<b>77.6</b> ( $\uparrow 0.6$ )	<b>59.0</b> ( $\uparrow 1.8$ )



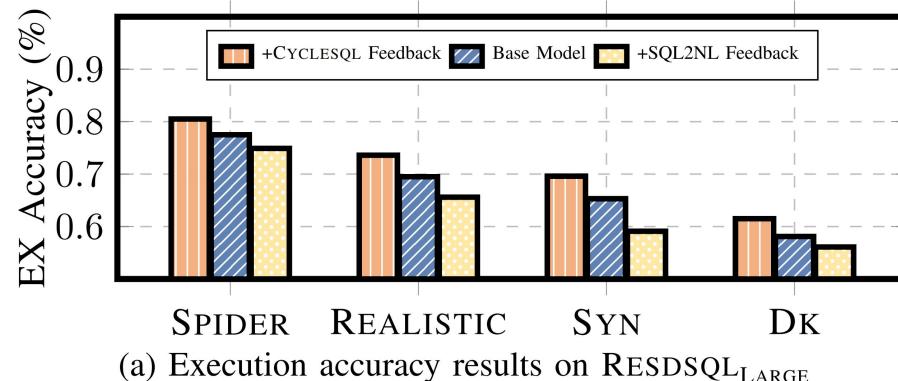
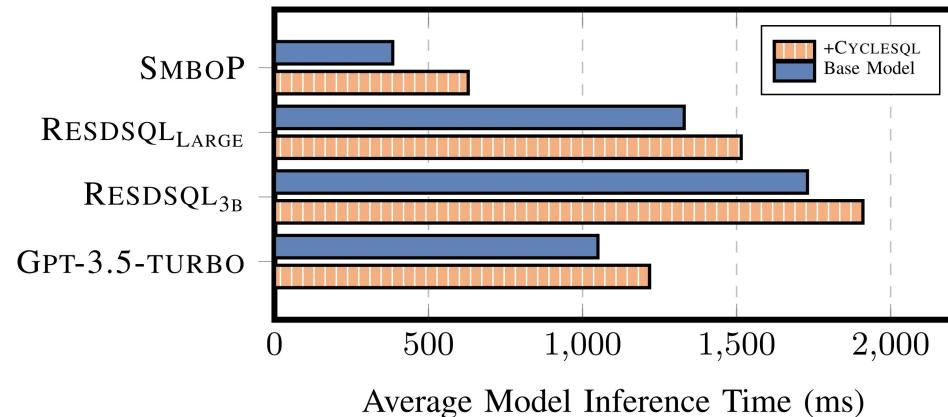
# Break-Down Results

Model		Easy	Medium	Hard	Extra Hard
SMBOP	Base	90.7	82.7	70.7	52.4
	+CYCLESQ	90.7	<b>84.1</b> ( $\uparrow 1.4$ )	69.5( $\downarrow 1.2$ )	53.0( $\uparrow 0.6$ )
PICARD <sub>3B</sub>	Base	95.6	85.4	67.8	50.6
	+CYCLESQ	95.6	86.1( $\uparrow 0.7$ )	69.5( $\uparrow 1.7$ )	50.6
	Base	92.3	83.4	66.1	51.2

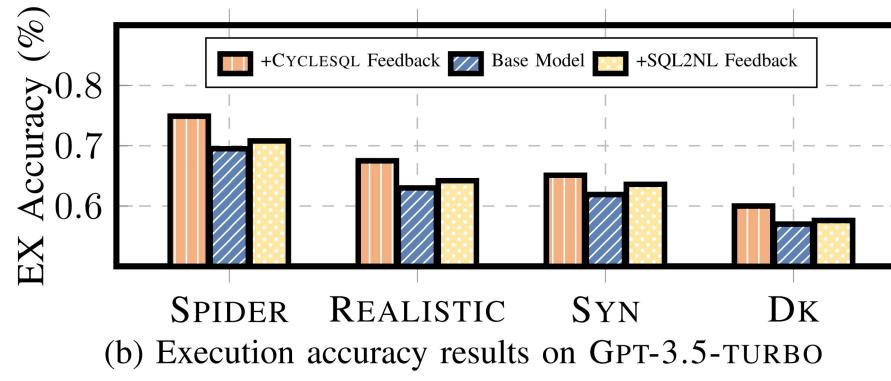
Significant Improvement over Extra-Hard-Queries on LLM Models

GPT-3.5-TURBO (5-SHOTS)	Base	84.3	78.5	65.5	48.2
	+CYCLESQ	86.3( $\uparrow 2.0$ )	<b>83.0</b> ( $\uparrow 4.5$ )	<b>73.0</b> ( $\uparrow 7.5$ )	<b>56.0</b> ( $\uparrow 7.8$ )
GPT-4 (5-SHOTS)	Base	90.3	84.3	63.8	56.6
	+CYCLESQ	90.7( $\uparrow 0.4$ )	85.4( $\uparrow 1.1$ )	<b>66.7</b> ( $\uparrow 2.9$ )	<b>59.6</b> ( $\uparrow 3.0$ )
CHESS	Base	70.2	25.3	39.7	19.3
	+CYCLESQ	<b>71.6</b> ( $\uparrow 1.4$ )	25.6( $\uparrow 0.3$ )	<b>41.1</b> ( $\uparrow 1.4$ )	<b>21.6</b> ( $\uparrow 2.3$ )
DAILSQ <sub>L3.5</sub>	Base	91.1	86.5	77.0	57.2
	+CYCLESQ	91.1	<b>86.8</b> ( $\uparrow 0.3$ )	<b>77.6</b> ( $\uparrow 0.6$ )	<b>59.0</b> ( $\uparrow 1.8$ )

Average Iterations	SMBOP (+CYCLESQ)	PICARD (+CYCLESQ)	RESDSQL <sub>LARGE</sub> (+CYCLESQ)	RESDSQL <sub>3B</sub> (+CYCLESQ)	GPT-3.5-TURBO (+CYCLESQ)
	2.44	3.78	1.90	1.88	1.87



(a) Execution accuracy results on RESDSQL<sub>LARGE</sub>



(b) Execution accuracy results on GPT-3.5-TURBO

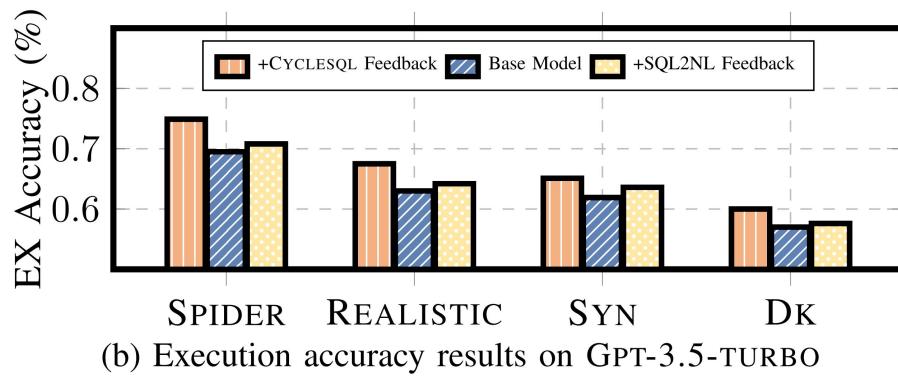
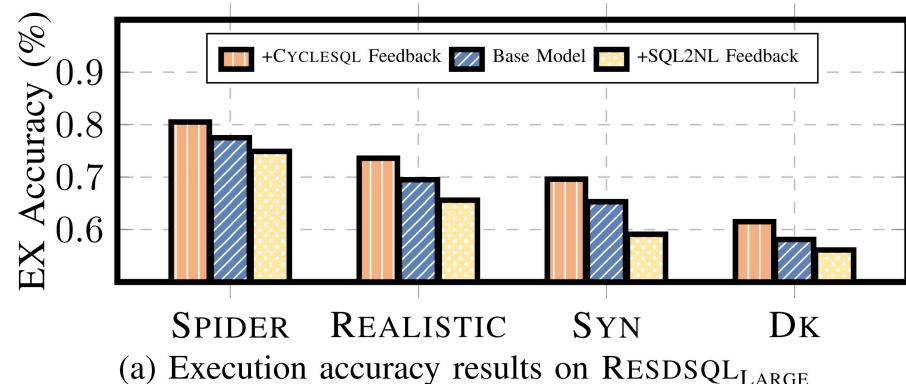
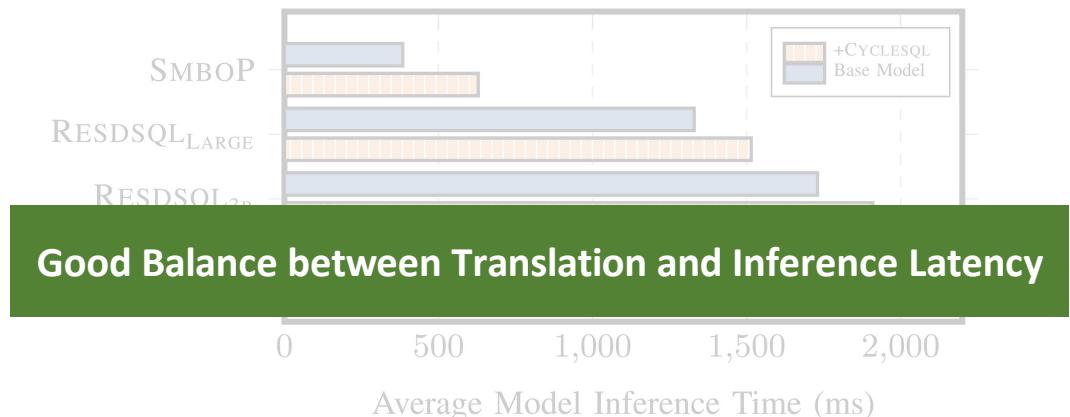
# Break-Down Results

Model		Easy	Medium	Hard	Extra Hard
<b>SMBOP</b>	Base	90.7	82.7	70.7	52.4
	+CYCLESQ	90.7	<b>84.1</b> ( $\uparrow 1.4$ )	69.5( $\downarrow 1.2$ )	53.0( $\uparrow 0.6$ )
<b>PICARD<sub>3B</sub></b>	Base	95.6	85.4	67.8	50.6
	+CYCLESQ	95.6	86.1( $\uparrow 0.7$ )	69.5( $\uparrow 1.7$ )	50.6
	Base	92.3	83.4	66.1	51.2

Significant Improvement over Extra-Hard-Queries on LLM Models

<b>GPT-3.5-TURBO</b> (5-SHOTS)	Base	84.3	78.5	65.5	48.2
	+CYCLESQ	86.3( $\uparrow 2.0$ )	<b>83.0</b> ( $\uparrow 4.5$ )	<b>73.0</b> ( $\uparrow 7.5$ )	<b>56.0</b> ( $\uparrow 7.8$ )
<b>GPT-4</b> (5-SHOTS)	Base	90.3	84.3	63.8	56.6
	+CYCLESQ	90.7( $\uparrow 0.4$ )	85.4( $\uparrow 1.1$ )	<b>66.7</b> ( $\uparrow 2.9$ )	<b>59.6</b> ( $\uparrow 3.0$ )
<b>CHESS</b>	Base	70.2	25.3	39.7	19.3
	+CYCLESQ	<b>71.6</b> ( $\uparrow 1.4$ )	25.6( $\uparrow 0.3$ )	<b>41.1</b> ( $\uparrow 1.4$ )	<b>21.6</b> ( $\uparrow 2.3$ )
<b>DAILSQ</b> <sub>3.5</sub>	Base	91.1	86.5	77.0	57.2
	+CYCLESQ	91.1	<b>86.8</b> ( $\uparrow 0.3$ )	<b>77.6</b> ( $\uparrow 0.6$ )	<b>59.0</b> ( $\uparrow 1.8$ )

Average Iterations	SMBOP (+CYCLESQ)	PICARD (+CYCLESQ)	RESDSQL <sub>LARGE</sub> (+CYCLESQ)	RESDSQL <sub>3B</sub> (+CYCLESQ)	GPT-3.5-TURBO (+CYCLESQ)
	2.44	3.78	1.90	1.88	1.87



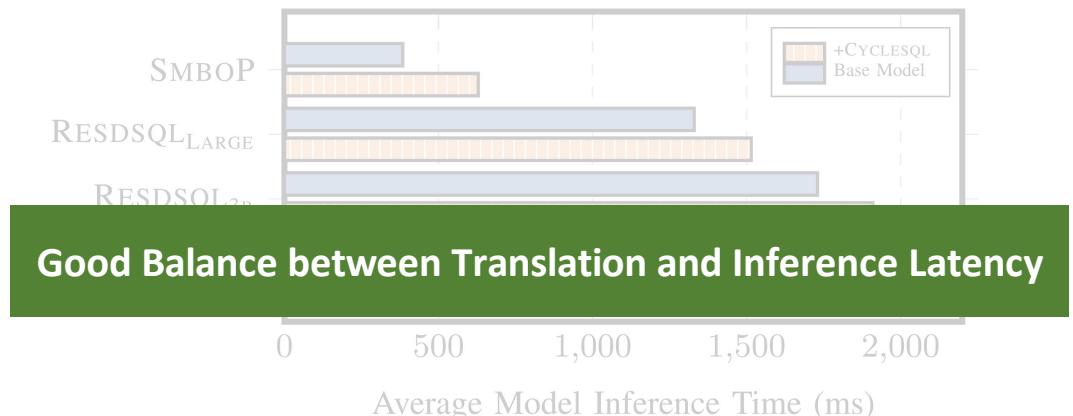
# Break-Down Results

Model		Easy	Medium	Hard	Extra Hard
<b>SMBOP</b>	Base	90.7	82.7	70.7	52.4
	+CYCLESQ	90.7	<b>84.1</b> ( $\uparrow 1.4$ )	69.5( $\downarrow 1.2$ )	53.0( $\uparrow 0.6$ )
<b>PICARD<sub>3B</sub></b>	Base	95.6	85.4	67.8	50.6
	+CYCLESQ	95.6	86.1( $\uparrow 0.7$ )	69.5( $\uparrow 1.7$ )	50.6
	Base	92.3	83.4	66.1	51.2

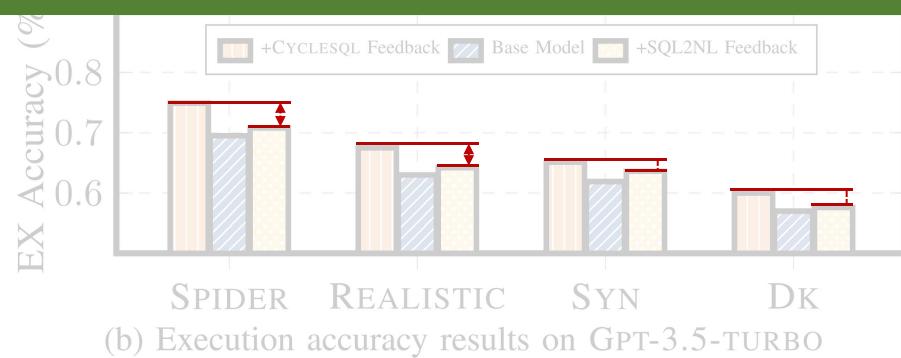
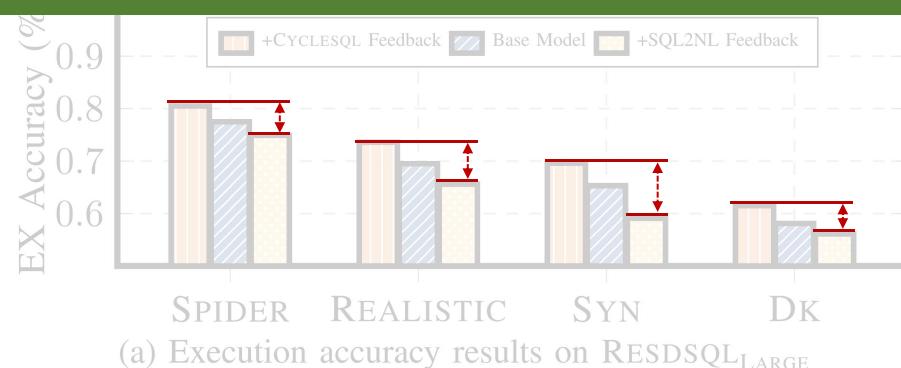
Significant Improvement over Extra-Hard-Queries on LLM Models

<b>GPT-3.5-TURBO</b> (5-SHOTS)	Base	84.3	78.5	65.5	48.2
	+CYCLESQ	86.3( $\uparrow 2.0$ )	<b>83.0</b> ( $\uparrow 4.5$ )	<b>73.0</b> ( $\uparrow 7.5$ )	<b>56.0</b> ( $\uparrow 7.8$ )
<b>GPT-4</b> (5-SHOTS)	Base	90.3	84.3	63.8	56.6
	+CYCLESQ	90.7( $\uparrow 0.4$ )	85.4( $\uparrow 1.1$ )	<b>66.7</b> ( $\uparrow 2.9$ )	<b>59.6</b> ( $\uparrow 3.0$ )
<b>CHESS</b>	Base	70.2	25.3	39.7	19.3
	+CYCLESQ	<b>71.6</b> ( $\uparrow 1.4$ )	25.6( $\uparrow 0.3$ )	<b>41.1</b> ( $\uparrow 1.4$ )	<b>21.6</b> ( $\uparrow 2.3$ )
<b>DAILSQ</b> <sub>3.5</sub>	Base	91.1	86.5	77.0	57.2
	+CYCLESQ	91.1	<b>86.8</b> ( $\uparrow 0.3$ )	<b>77.6</b> ( $\uparrow 0.6$ )	<b>59.0</b> ( $\uparrow 1.8$ )

Average Iterations	SMBOP (+CYCLESQ)	PICARD (+CYCLESQ)	RESDSQL <sub>LARGE</sub> (+CYCLESQ)	RESDSQL <sub>3B</sub> (+CYCLESQ)	GPT-3.5-TURBO (+CYCLESQ)
	2.44	3.78	1.90	1.88	1.87



CycleSQL Feedback is Better than SQL2NL Feedback!



# More Results in the Paper

---

- Comparison of different translation verifier selection
- Qualitative evaluation results

# Conclusion

---

- **Closing the feedback loop** for NL2SQL brings good benefits
- Good feedback improves not only **accuracy**, but **explanability**

Yuankai Fan <kaimary1221@gmail.com>

<https://github.com/Kaimary/CycleSQL>