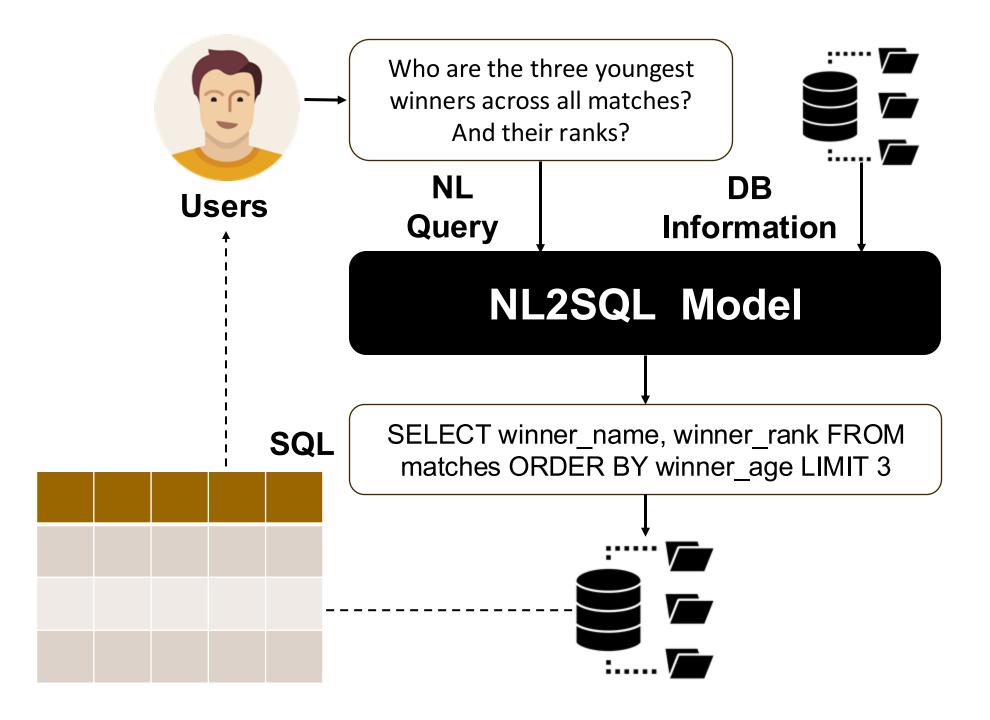


WHAT IS NATURAL LANGUAGE TO SQL (NL2SQL)?



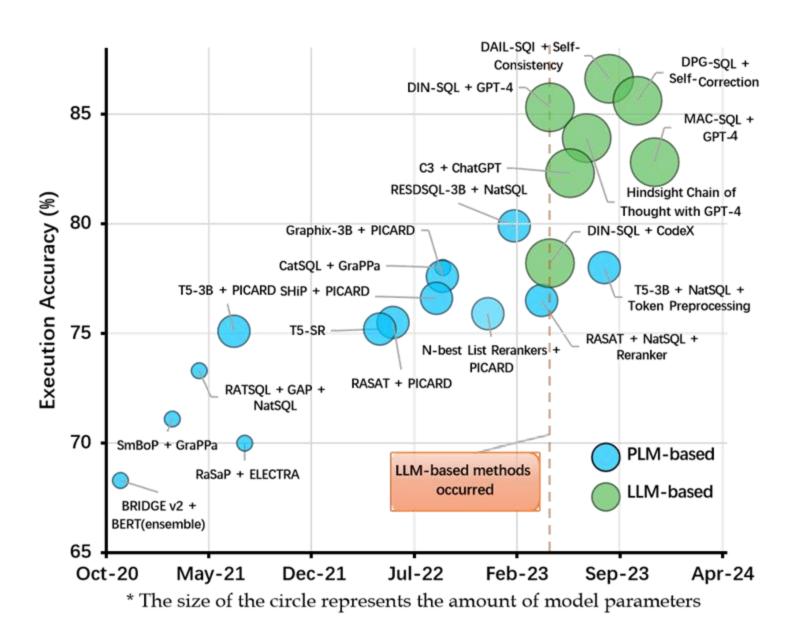




WHERE ARE WE NOW IN NL2SQL?



■ Recent advanced NL2SQL methods are mainly driven by LLM.



■ Can we conclude that **LLM-based**models are "the best choice" for any
NL2SQL application?



ONE MODEL DOES NOT FIT ALL!



SPECIFIC DOMAIN

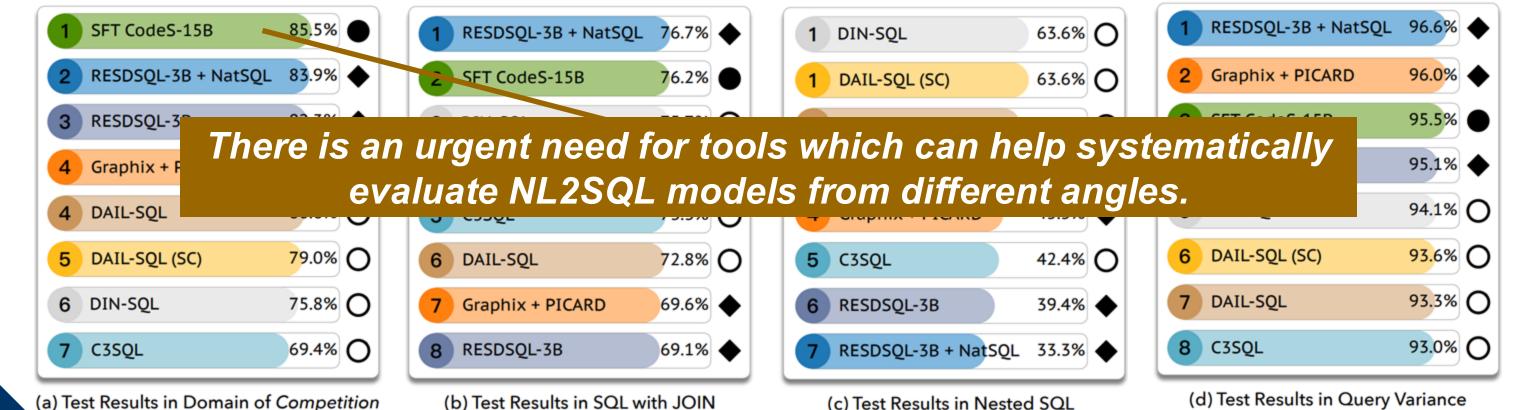
MULTIPLE JOINS

NESTED SQL QUERY QUESTION VARIANCE

- Domain-specific database
 - Competition
 - Medical
 - Finance

Query with multiple JOINs

- Query with nested SQL
- Different user questions for the same SQL query



NL2SQL360: A MULTI-ANGLE EVALUATION FRAMEWORK im





■ Dataset Filter

Case 1: SQL Complexity

Case 2: SQL Characteristics

– Case 3: NL Query Variance

Case 4: Database Domain

Case 5: Economy

Case 6: Efficiency



NL2SQL360: A MULTI-ANGLE EVALUATION FRAMEWORK







NL2SQL360: A MULTI-ANGLE EVALUATION FRAMEWORK







EXPERIMENTS WITH NL2SQL360



We utilized NL2SQL360 to evaluate 7 PLM-based and 13 LLM-based NL2SQL methods on 2 widely-used benchmarks, varing 15 settings and deriving a set of new findings.

DATASETS

- Spider Dataset
- BIRD Dataset

EVALUATION METRICS

- Execution Accuracy (EX)
- Exact Match Accuracy (EM)
- Valid Efficiency Score (VES)
- Query Variance Testing (QVT)

BASELINES

- PLM-based:
 - RESDSQL-3B, RESDSQL-3B + NatSQL,
 Graphix + PICARD
- LLM-based:
 - DAILSQL, DAILSQL(SC), DINSQL, C3SQL, SFT CodeS-1B/3B/7B/15B

EVALUATION ANGLES

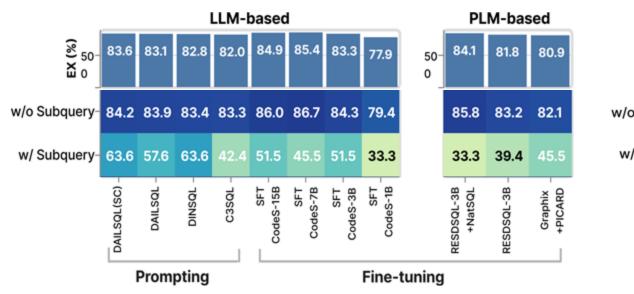
- SQL Characteristics
- Question Variance
- Database Domain Adaption
- Pre-training Corpus
- **■** Efficiency
- Economy (Cost)

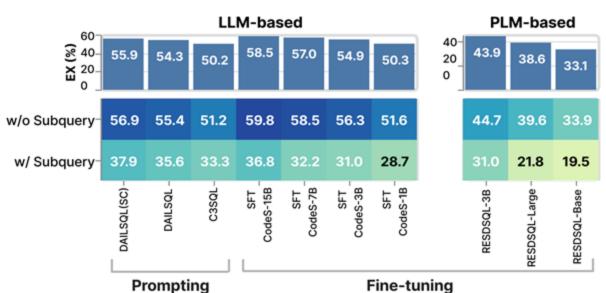


ACCURACY VS. SQL CHARACTERISTICS



 Motivation: Real-world applications often require generating SQL quries involving advanced operations like subquries.





■ In overall dataset

- SFT CodeS-15B > DAILSQL(SC) > DAILSQL > DINSQL
 - However, in subquery scenario
 - DAILSQL(SC) = DINSQL > DAILSQL > SFT CodeS-15B

In scenarios involving subqueries, LLM-based methods outperform PLMbased methods overall, with methods using GPT-4 showing particularly better performance.



ACCURACY VS. QUESTION VARIANCE



- Motivation: For the same SQL query, different users may pose different NL questions. The ability of handling question variance is crucial for NL2SQL applications.
- To this end, we proposed new metric called Question Variance Testing (QVT):

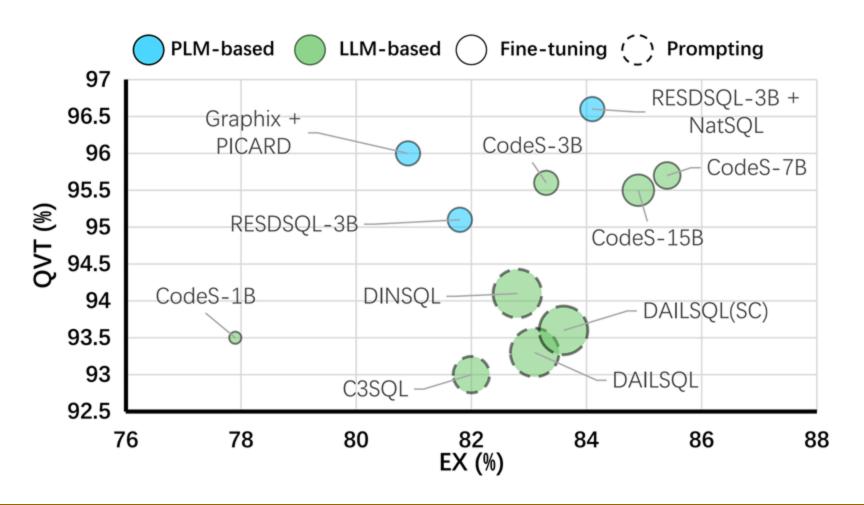
$$QVT = \frac{1}{M} \sum_{i=1}^{M} \left(\frac{\sum_{j=1}^{m_i} \mathbb{I}(\mathcal{F}(N_{ij}) = Q_i)}{m_i} \right)$$

- M is the total number of SQL quries in the test set.
- $-m_i$ is the number of natural language query variations corresponding to the SQL query Q_i .
- $\mathcal{F}(N_{ij})$ represents the SQL query generated by the NL2SQL model for the j-th natural language query variation of Q_i .
- $-\mathbb{I}(\cdot)$ is the indicator function that returns 1 if the query results inside are equal, and 0 otherwise.
- For NL2SQL models, higher QVT means higher stability of handing question variance.

ACCURACY VS. QUESTION VARIANCE



We evaluate QVT for multple NL2SQL methods on Spider development dataset.



Our finding:

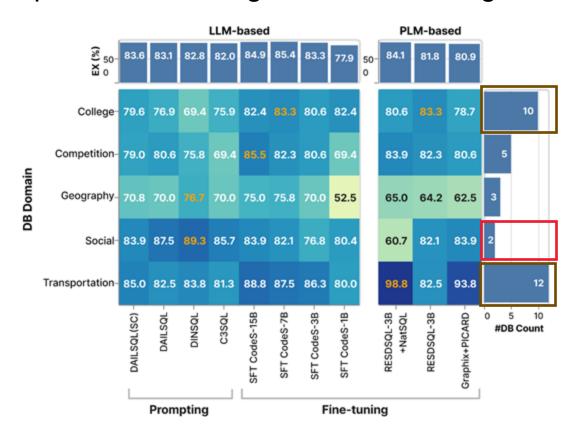
■ Fine-tuning the model with task-specific datasets may help stabilize its performance against NL variations.

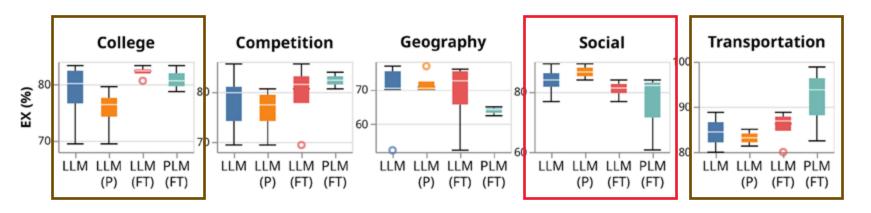


ACCURACY VS. DATABASE DOMAIN



 Motivation: In practical NL2SQL applications, scenarios usually involve domain-specific databases, each with unique schema designs and terminologies.





- LLM (P): Prompt-based LLMs
- LLM (FT): Fine-tuned LLMs
- PLM (FT): Fine-tuned PLMs

Our findings:

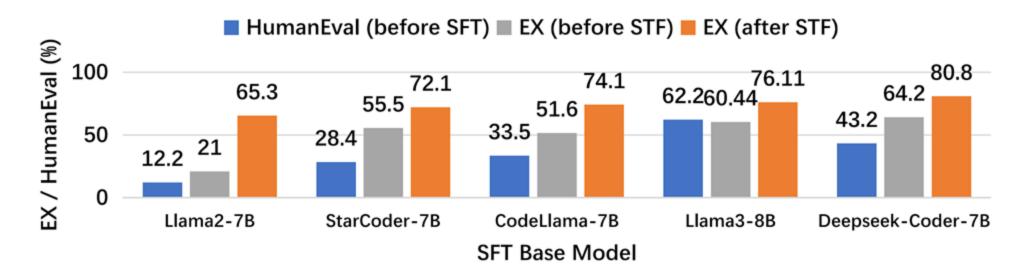
- Different methods exhibit varying biases towards different domains.
- However, in-domain training data during fine-tuning is crucial for model performance in specific domains.



ACCURACY VS. PRE-TRAINING CORPUS



- Motivation: A large number of LLMs with different pre-training corpus have emerged recently, such
 as Llama, StarCoder, etc. Which open-source LLMs are best suited for SFT in the NL2SQL task?
- We evaluate different LLMs which have different pre-training corpus or procedures, and evaluate them in Spider Dataset.



Our finding:

After SFT on open-source LLMs for the NL2SQL task, we found a positive correlation between the model's inherent coding ability and its performance in NL2SQL task.



RECAP: ONE MODEL DOES NOT FIT ALL!



SPECIFIC DOMAIN

MULTIPLE JOINS

NESTED SQL QUERY

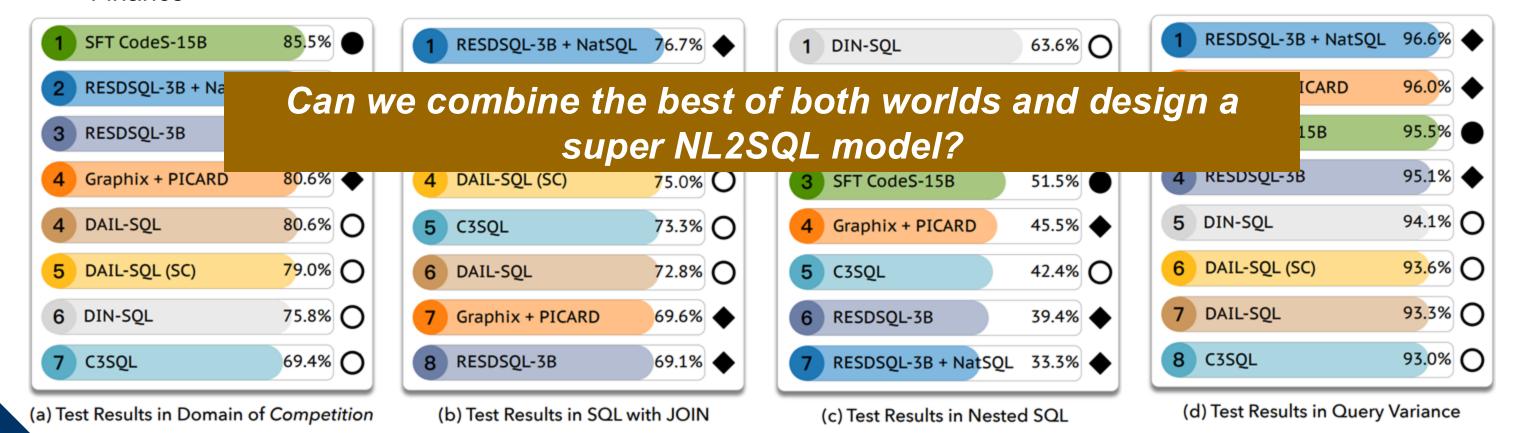
QUESTION VARIANCE

- Domain-specific database
 - Competition
 - Medical
 - Finance

■ Query with multiple JOINs

Query with nested SQL

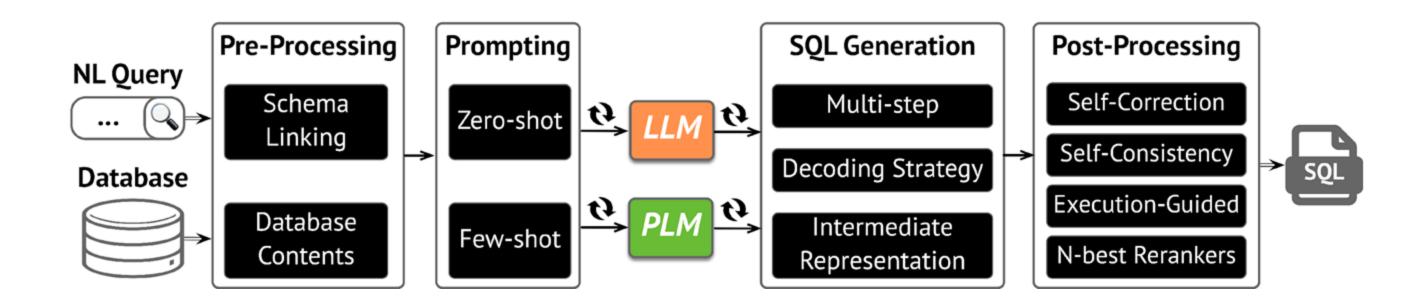
■ Different user questions for the same SQL query



WHAT IS NEXT FOR NL2SQL?



- Now, we have found that different NL2SQL models exhibit distinct advantages in specific scenarios.
 Next, how can we systematically design future NL2SQL models?
- We systematically categorized and analyzed recent NL2SQL modules based on LLMs and PLMs, highlighting their commonalities and distinct features.
- We proposed NL2SQL Design Space in a modular style.



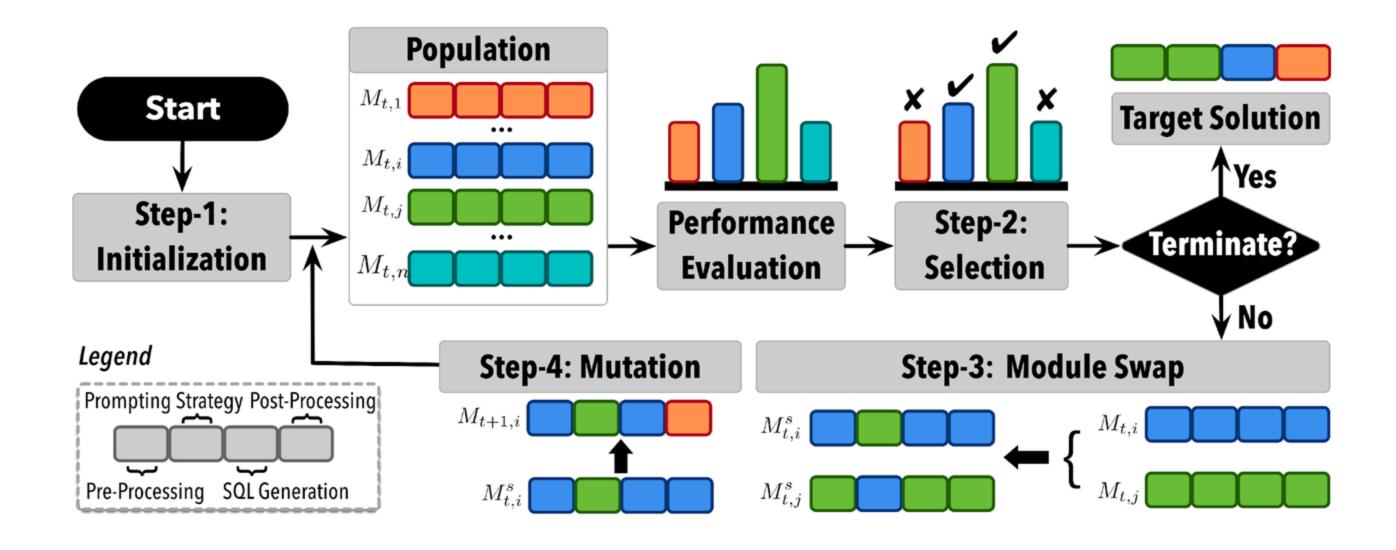
NL2SQL DESIGN SPACE



Types		Methods	Backbone	Example	Schema	DB	S	QL Generation Str	Post-processing	
		Wethous	Models	Selection (Few-shot)	Linking	Content	Multi-Step	Intermediate Representation	Decoding Strategy	Strategy
LLM-based	ting	DIN-SQL [38]	GPT-4	Manual	1	×	Classification Decomposition	NatSQL	Greedy Search	Self-Correction
		DAIL-SQL [11] (with Self-Consistency)	GPT-4	Similarity-based	Х	X	×	×	Greedy Search	Self-Consistency
	Prompting	MAC-SQL [52]	GPT-4	N/A	1	×	Sub-question	X	Greedy Search	Refiner
	Pr	Can we design an automatic algorithm to explore the NL2SQL design space and achieve stronger								
		performance?								
		RESDOQL + NatoQL [23]	13	14/21	· ·	•	Skeleton I arsing	NatoQL	Deant Scared	SQL Selector
	ng	Graphix + PICARD [25]	T5	N/A	✓	✓	X	X	PICARD	X
PLM-based	Fine-tuning	N-best Rerankers + PICARD [59]	T5	N/A	✓	✓	X	X	PICARD	N-best Rerankers
-ba		T5 + NatSQL + Token Preprocessing [42]	T5	N/A	✓	✓	X	NatSQL	Greedy Search	X
Ż		RASAT + PICARD [39]	T5	N/A	\	✓	X	X	PICARD	X
PI		SHiP + PICARD [16]	T5	N/A	X	✓	X	X	PICARD	X
		T5 + PICARD [45]	T5	N/A	X	✓	X	X	PICARD	X
		RATSQL + GAP + NatSQL [10]	BART	N/A	✓	✓	X	NatSQL	X	X
		BRIDGE v2 [28]	BERT	N/A	×	1	×	×	Schema-Consistency Guided Decoding	х

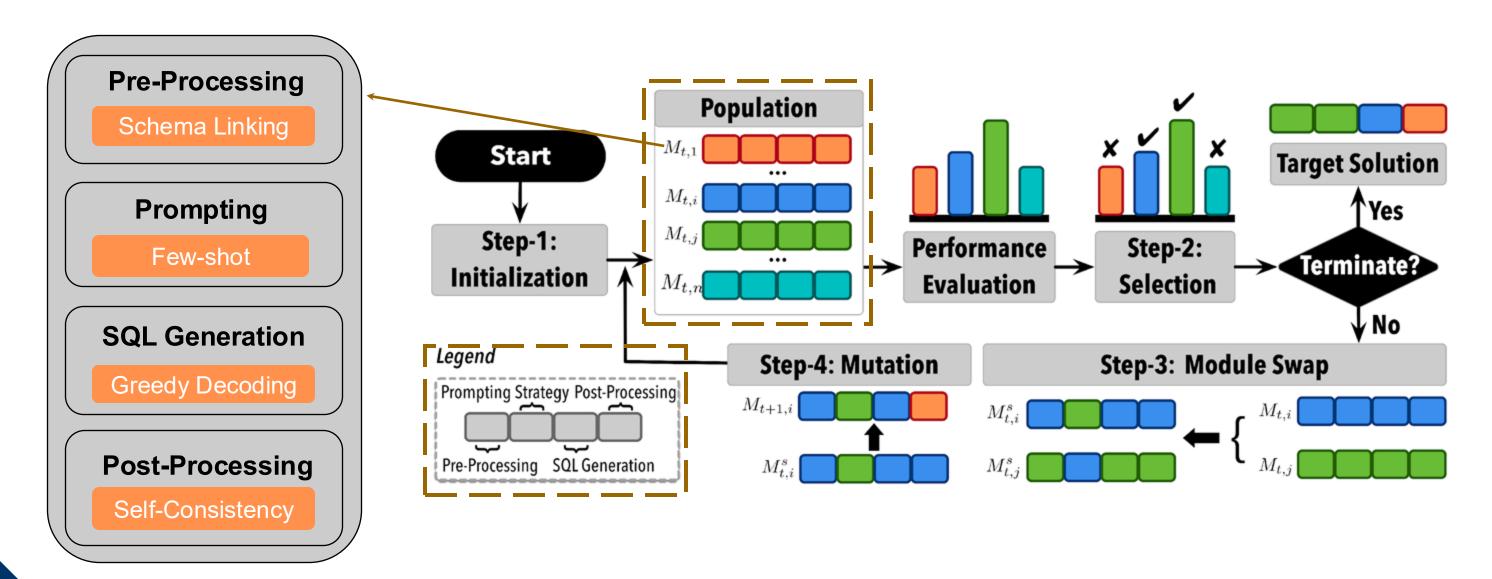
香港科技大学(广州) THE HONG KONG UNIVERSITY OF SCIENCE AND TECHNOLOGY (GUANGZHOU

NL2SQL AUTOMATED ARCHITECTURE SEARCH ALGORITHM



香港科技大学(广州) THE HONG KONG UNIVERSITY OF SCIENCE AND TECHNOLOGY (GUANGZHOU)

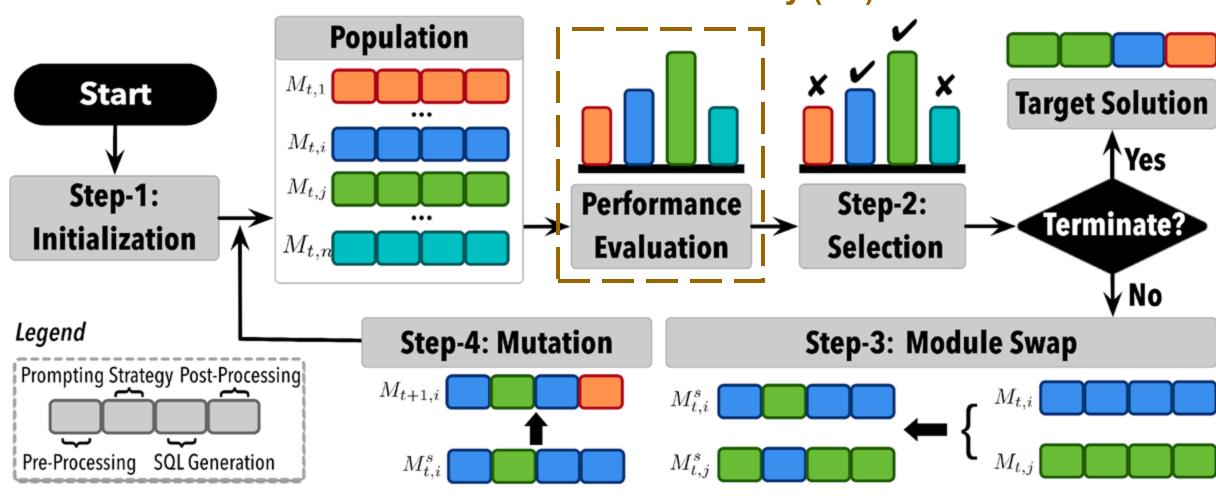
NL2SQL AUTOMATED ARCHITECTURE SEARCH ALGORITHM





NL2SQL AUTOMATED ARCHITECTURE SEARCH ALGORITHM

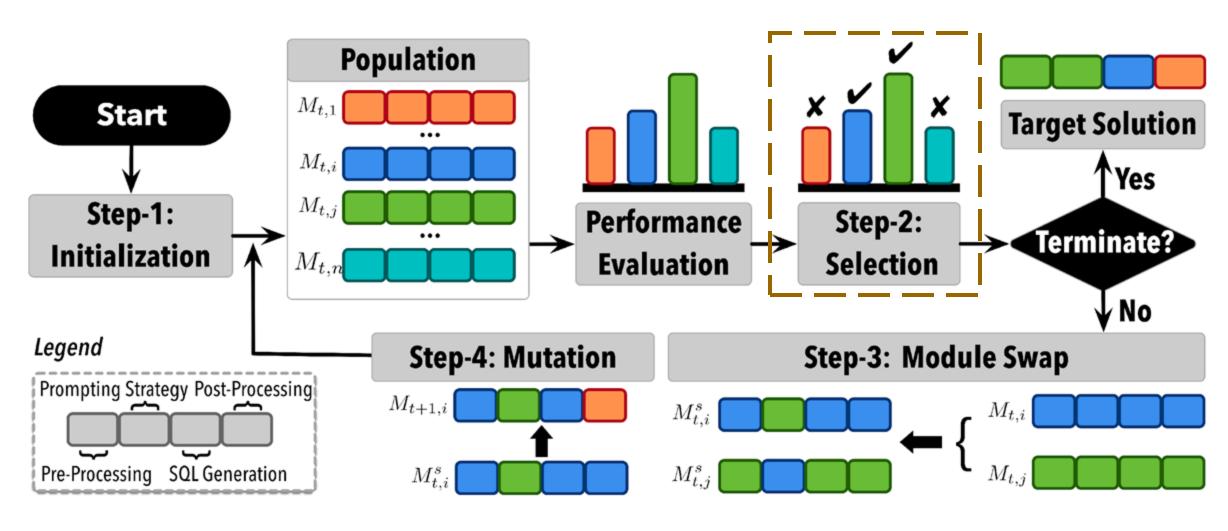
Target Metric, e.g., Execution Accuracy (EX)



香港科技大学(广州) THE HONG KONG UNIVERSITY OF SCIENCE AND TECHNOLOGY (GUANGZHOU

NL2SQL AUTOMATED ARCHITECTURE SEARCH ALGORITHM

Russian Roulette Selection



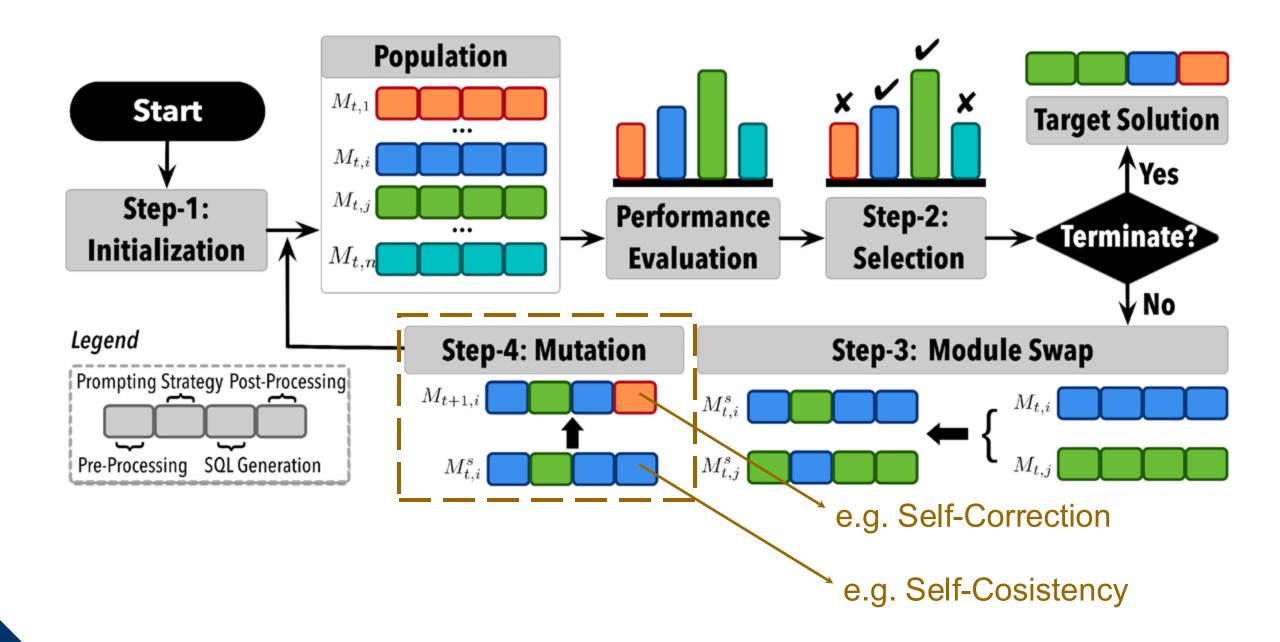
香港科技大学(广州) THE HONG KONG UNIVERSITY OF SCIENCE AND TECHNOLOGY (GUANGZHOU

NL2SQL AUTOMATED ARCHITECTURE SEARCH ALGORITHM



香港科技大学(广州) THE HONG KONG UNIVERSITY OF SCIENCE AND TECHNOLOGY (GUANGZHOU)

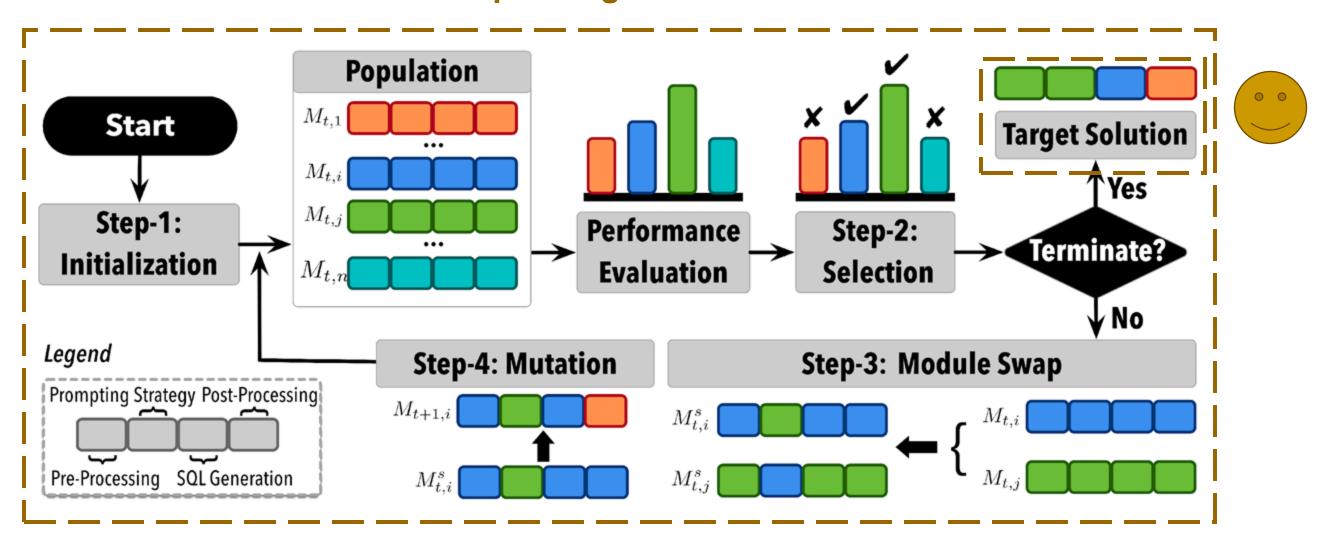
NL2SQL AUTOMATED ARCHITECTURE SEARCH ALGORITHM





NL2SQL AUTOMATED ARCHITECTURE SEARCH ALGORITHM

Loop for T generations



SUPERSQL



• To validate the effectiveness of the NL2SQL360-AAS algorithm, we conduct a case study.

■ Dataset: Spider

■ Target Metric: Execution Accuracy (EX)

■ Model Backend: GPT-3.5-Turbo

■ Predefined Search Space:

NL2SQL Pipeline	Module Candidates					
Pre-processing	Schema Linking Module (Source: RESDSQL), DB Content Module (Source: BRIDGE v2)					
Prompting	Similarity-based Few-shot Module (Source: DAILSQL)					
SQL Generation	Greedy Search Decoding (Limited by OpenAI)					
Post-Processing	Self-Correction Module (Source: DINSQL), Self-Consistency Module (Source: C3SQL) Execution-Guided SQL Selector Module (Source: RESDSQL)					

SUPERSQL



- From the final generation produced by NL2SQL360-AAS, we select the individual with the highest EX metric as our final NL2SQL solution, i.e., **SuperSQL**. Note that we replace the backbone model with GPT-4 for more powerful performance.
- Take a deep look at prompt improvement compared to DAILSQL.

```
/* Given the following database schema: */
CREATE TABLE airports (
    City text,
    AirportCode text primary key,
    AirportName text
)

CREATE TABLE airlines (
    uid int primary key,
    Airline text
)

/* Answer the following: What are airport names at City
'Aberdeen'?*/
SELECT AirportName FROM airports WHERE City = "Aberdeen";
```

DAILSQL Prompt

```
/* Given the following database schema (with linked
value list after column definition): */
CREATE TABLE airports (
    City text,
    AirportCode text primary key,
    AirportName text -- ["Aberdeen"]
)

CREATE TABLE airlines (
    uid int primary key,
    Airline text
)

/* Answer the following: What are airport names at City
'Aberdeen'?*/
SELECT AirportName FROM airports WHERE City = "Aberdeen";
```

SuperSQL Prompt



SUPERSQL



Effectiveness of SuperSQL.

Model	Spider-Dev(%)	Spider-Test(%)	BIRD-Dev(%)	BIRD-Test(%)
BRIDGE v2	70.3	-	-	-
RESDSQL-3B	81.8	-	43.9	-
C3SQL	82.0	82.3	50.2	-
DINSQL	82.8	85.3	50.7	55.9
DAILSQL(SC)	83.6	86.6	55.9	57.4
SuperSQL	87.0 (3.4↑) (17.7↑)	87.0 (0.4↑) (4.7↑)	58.5 (2.6↑) (14.6↑)	62.7 (5.3↑) (6.8↑)

Table 5: Accuracy vs. LLM Economy on Spider/BIRD Dev Set.

Methods	LLMs	Avg. Tokens / Query		Avg. Cost / Query		EX(%)		EX / Avg. Cost	
Methods		Spider	BIRD	Spider	BIRD	Spider	BIRD	Spider	BIRD
C3SQL	GPT-3.5	5702	5890	0.0103	0.0104	82.0	50.2	7961	4825
DINSQL	GPT-4	9571	-	0.2988	-	82.8	-	277	-
DAILSQL	GPT-4	930	1559	0.0288	0.0486	83.1	54.3	2885	1117
DAILSQL(SC)	GPT-4	1063	1886	0.0377	0.0683	83.6	55.9	2218	819
SuperSQL	GPT-4	942	1412	0.0354	0.0555	87.0	58.5	2458	1053

CONCLUSION



- We proposed a multi-angle testbed, named **NL2SQL360**, for evaluating NL2SQL methods from different perspectives.
- We utilized our NL2SQL360 to evaluate **13** LLM-based and **7** PLM-based NL2SQL methods on **2** widely-used benchmarks, varying **15** settings and deriving **a set of new findings**.
- We employed our NL2SQL360 to analyze the NL2SQL Design Space for NL2SQL solutions and automatically search (NL2SQL360-AAS) for one of the best solutions, named SuperSQL.

THANK YOU & QUESTIONS?





