

## Comprehensive Database Generation

### Schema Generation

#### Schema Augmentation

Add New Type  
Email: TEXT Info: JSON  
Phone: TEXT → { "Email": ..., "Phone": ... }

Using Reserved Keyword  
Table: FBUser → Table: User

#### Schema Translation

POINT → SDO\_GeOMETRY

#### Attribute Annotation

Birth: TEXT  
'21-JAN' ← Semantics: Date Format: 'DD-MON'  
'13-FEB'

### Data Generation

#### Filter Incompatible Data

'The sun was setting, casting...' (4000+letter)

#### Special Value Insertion

0 NULL '2024-02-29' ↗

### Input

Testing SQL

Schema

Database Setting

### Translator

Rule-based

LLM-based

### Execution-based Verification

Sample Database Generation  
SubQuery Tree

$q_1 \rightarrow q_2$  ask  $q_2$  to gen Data (1, 'a')

$q_2 \rightarrow q_3$  ask  $q_3$  to gen data (NULL, 'b')

$q_2 \rightarrow$  re-ask  $q_3$  to gen data: (15, 'c')

Sample DB

Result Comparison

Target DB

## Diverse SQL Pair Generation

### Translation Point Collection

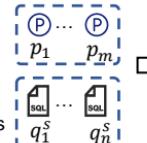
Extracting From Rule LLM-based Manual  
Translator Generation Generation

"SrcDialect": ... "TgtDialect": ...  
"SrcPattern": ... "TgtPattern": ...  
"Condition": ... "Category": ...

#### Points Written in TPDL

Translation Points  
 $p_1, p_m$

Real-world Reference SQLs  
 $q_1^s, q_n^s$



### SQL Pair Generation

Database Configuration

Reference SQL Selection

Reference SQL Set  
Iterative SQL Generation

Reference SQL  $p_i$

Instantiate Translation Point  
 $p_i$

Integrate to Reference SQL

Concatenate Expanded Reference SQLs

SQL Pair { "SrcDialect": ... "TgtDialect": ...  
"SrcSQL": ... "TgtSQL": ...  
"Condition": ... "Schema": ... }

### Translation Equivalence Verification

Manual Check

### Formal Verifier

SQL Standardization

str ILIKE fmt  $\Rightarrow$  lower(str) LIKE lower(fmt)

Unsupported Operator

Supported Operator

**SQLSolver** Formal Verification