

# Text to SQL



INSTITUT  
POLYTECHNIQUE  
DE PARIS

Mehwish Alam  
Associate Professor  
Télécom Paris  
Institut Polytechnique de Paris  
Winter Semester 2024-2025

# Outline

---

- What are the traditional approaches for T2S?
- Masked Language Models for T2S (TaBERT)
- Large Language Models for T2S
  - Preprocessing
  - Inference
  - Post-processing
  - Fine-tuning

# Outline

---

- What are the traditional approaches for T2S?
- Masked Language Models for T2S (TaBERT)
- Large Language Models for T2S
  - Preprocessing
  - Inference
  - Post-processing
  - Fine-tuning

# Parsing Approaches for T2S

---

## Single-turn Parsing (context independent).

- Given a NL question  $Q$  and the corresponding database schema  $S = \langle T, C \rangle$ , the goal is to generate a SQL query  $Y$ .
  - Question  $Q = \{q_1, q_2, \dots, q_{|Q|}\}$  is a sequence of  $|Q|$  tokens.
  - The database schema consists of  $|\mathcal{T}|$  tables  $\mathcal{T} = \{t_1, t_2, \dots, t_{|\mathcal{T}|}\}$  and  $|C|$  columns  $C = \{c_1, c_2, \dots, c_{|C|}\}$ .
- Each table  $t_i$  is described by its name that contains multiple words  $[t_{i,1}, t_{i,2}, \dots, t_{i,|t_i|}]$ .
- Each table  $c_j^{t_i}$  in the table  $t_i$  is represented by words (a phrase)  $[c_{j,1}^{t_i}, c_{j,2}^{t_i}, \dots, c_{j,|c_j^{t_i}|}^{t_i}]$
- The whole input is given as  $X = \langle Q, S \rangle$

# Parsing Approaches for T2S

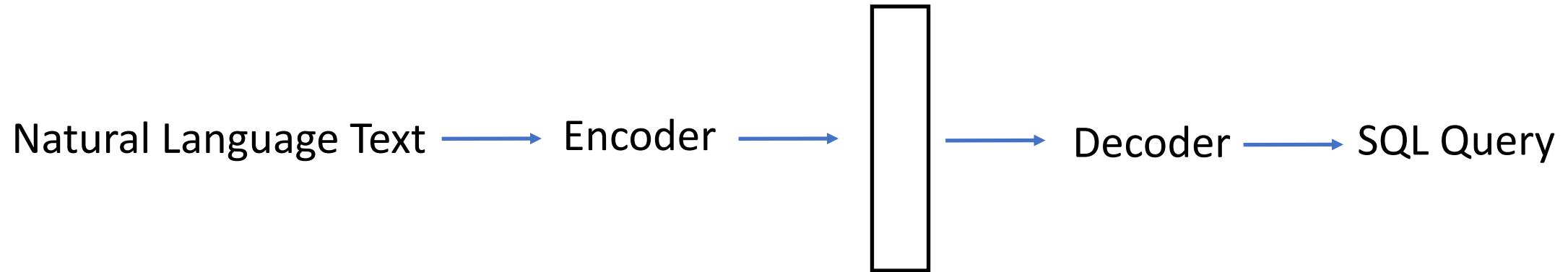
---

## Multi-turn Parsing (context dependent)

- A sequence of NL questions to the corresponding SQL queries
- NL questions may contain ellipsis and anaphora
- Let  $U = \{U_1, \dots, U_T\}$  denote a sequence of utterances with  $T$  turns
- $U_t = (X_t, Y_t)$  represents the  $t - th$  utterance
- $U_t$  is the combination of an NL question  $X_i$  and a SQL query  $Y_i$
- $S$  = corresponding database schema
- At the  $t - th$  turn, produce the SQL query  $Y_t$  conditioned on the current NL question  $X_t$
- The historical utterances  $\{U_i\}_{i=1}^{t-1}$ , and the database schema  $S$

# Deep Learning based Approaches – Single Turn

Seq2Seq task



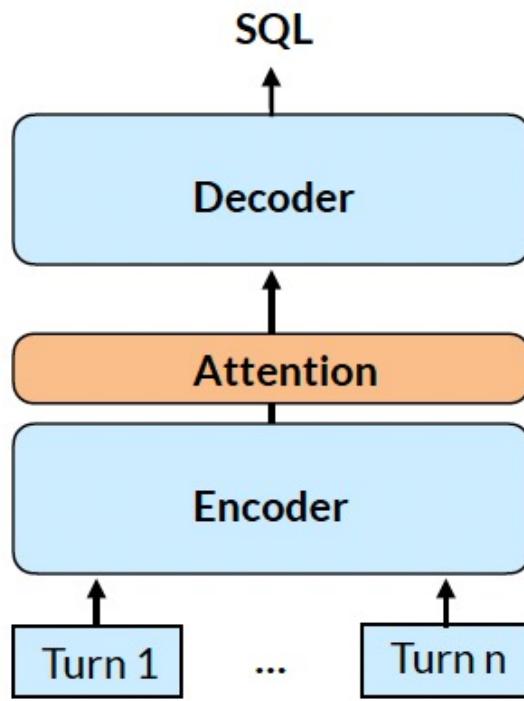
## Encoders.

- Schema Linking – SL
- LSTM
- Transformers
- Graph Neural Networks

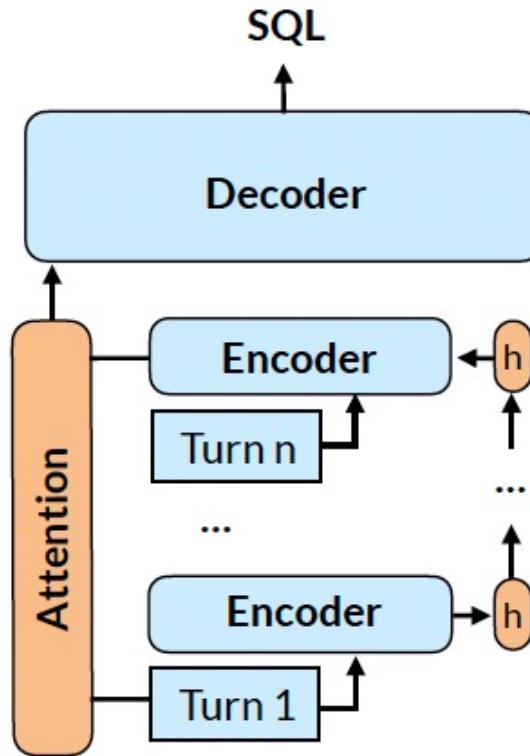
## Decoders.

- LSTM
- Transformers
- Grammar
- Domain Specific Language - DSL
- Constrained Decoding
- Re-Ranking

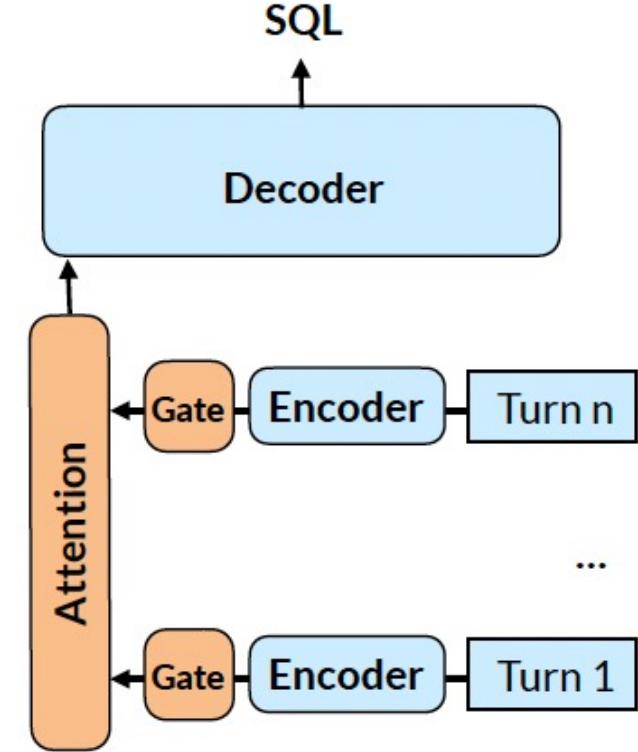
# Deep Learning based Approaches – Multi-turn



(a) Concat



(b) Turn



(c) Gate

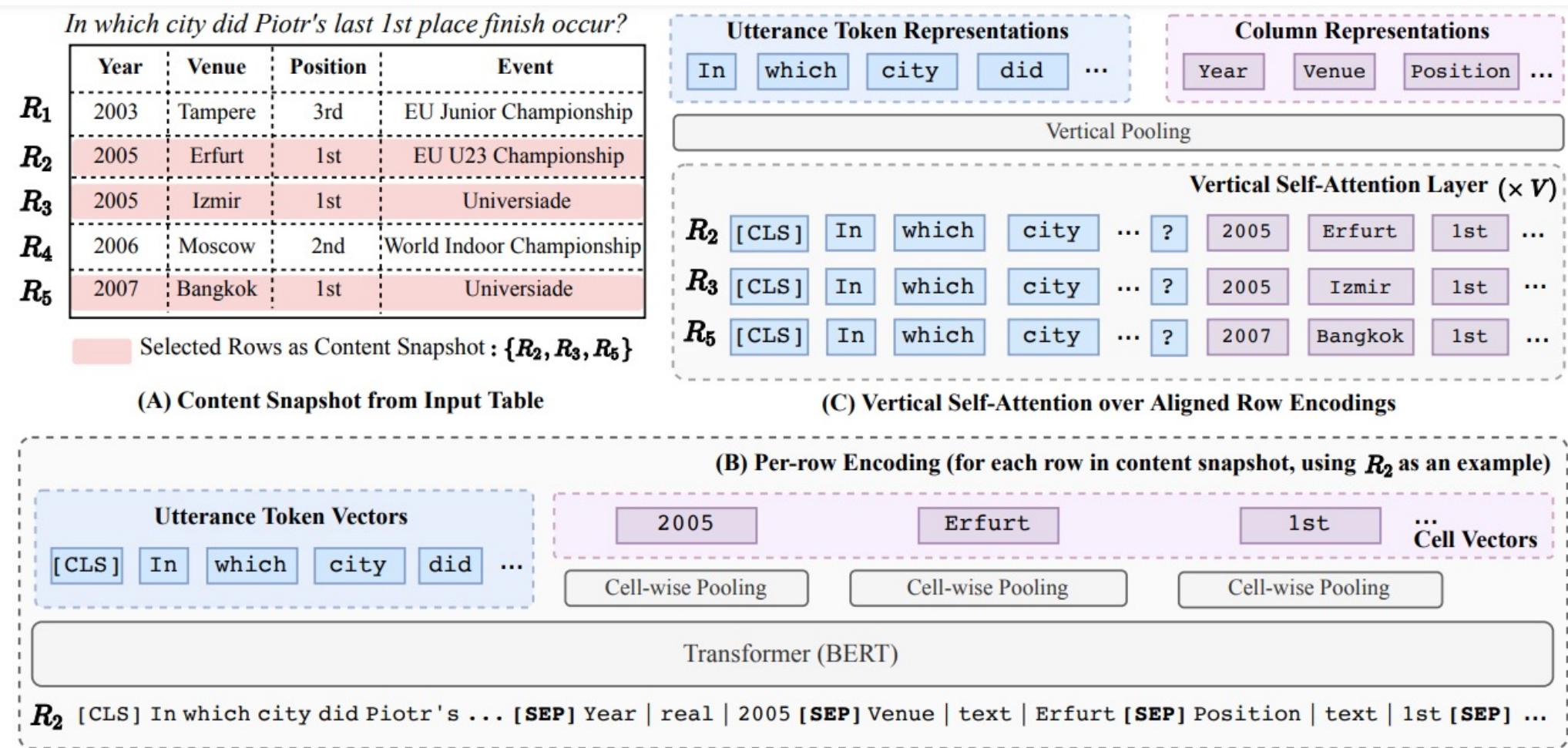
# Outline

---

- What are the traditional approaches for T2S?
- Masked Language Models for T2S (TaBERT)
- Large Language Models for T2S
  - Preprocessing
  - Inference
  - Post-processing
  - Fine-tuning

# Masked Language Models for T2S

## TaBERT – Learning Joint Representations over Textual and Tabular Data



# Masked Language Models for T2S - TaBERT

---

## Computing Representations for NL Utterances and Table Schemas

- Content Snapshot
- Row Linearization
- Vertical Self-Attention Mechanism
- Utterance and Column Representations

# Masked Language Models for T2S - TaBERT

---

## Computing Representations for NL Utterances and Table Schemas

### ➤ Content Snapshot

- TaBERT uses the **table contents** in addition to **column names**
- **Ambiguity** in column names
  - e.g., “Venue” can be ambiguous “City” may be more representative
- **Drawback.** Large number of rows in the table
- **Solution.**
  - Content Snapshot containing **only relevant rows** to input utterance
  - Calculate **content-sensitive column representations** from cell values
- **Strategy.**
  - Top-k based on **relevance between utterance and rows**
  - **Highest n-gram overlap ratio**
  - For k=1, synthetic row by selecting the cell values for each column with highest n-gram ratio

# Masked Language Models for T2S - TaBERT

## Computing Representations for NL Utterances and Table Schemas

### ➤ Row Linearization

- Creates a linearized sequence of each row in content snapshot
- Input to the transformer
- A concatenation of the utterance, columns, and their cell values
- Each cell is represented as follows:

Year		Real		2005
Column name		Column type		Cell Value

The diagram illustrates the structure of a table cell. It shows a row with four entries: 'Year', '|', 'Real', and '|', '2005'. Below this row, three red arrows point to the first three entries: 'Year' is labeled 'Column name', 'Real' is labeled 'Column type', and '2005' is labeled 'Cell Value'.

[CLS] In which city did Piotr's ... [SEP] Year | real | 2005 [SEP] Venue | text | Erfurt [SEP] Position | text | 1st [SEP]

# Masked Language Models for T2S - TaBERT

---

## Computing Representations for NL Utterances and Table Schemas

### ➤ Vertical Self-Attention Mechanism

- Allows for information flow across cell representations
- A fixed length initial vector for each cell at position  $< i, j >$  which is given by **mean pooling** over the sequence of the Transformer's output vectors
- The sequence of **word vectors** for the **NL utterance** are **concatenated** with the **cell vectors** as initial inputs to **the vertical attention layer**

# Masked Language Models for T2S - TaBERT

---

## Computing Representations for NL Utterances and Table Schemas

### ➤ Utterance and Column Representation

- A representation  $c_j$  is computed for each column  $c_j$  by mean-pooling over its vertically aligned cell vectors
- A representation for each utterance token,  $x_j$ , is computed similarly over the vertically aligned token vectors

# Masked Language Models for T2S - TaBERT

---

## Unsupervised Learning Objective

- **Masked Column Prediction (MCP)** objective encourages the model to recover the names and data types of masked columns
- **Cell Value Recovery (CVR)** objective to ensure information of representative cell values in content snapshots is retained after additional layers of vertical self-attention

# Outline

---

- What are the traditional approaches for T2S?
- Masked Language Models for T2S (TaBERT)
- Large Language Models for T2S
  - Preprocessing
  - Inference
  - Post-processing
  - Fine-tuning

# LLMs for Text2SQL - Preprocessing

## Question Representation.

- Natural Language Problem
- Related information from databases

## Layouts.

```
1  ### SQLite SQL tables, with their  
2      properties  
3  #  
4  # schools(SchID, city)  
5  # teachers(TechID, SchID)  
6  # students(StuID, SchID, TechID)  
7  # course(CourseID, SchID)  
8  # grade(StuID, CourseID, score)  
9  #  
10 #### {Question}  
11 SELECT
```

OpenAI Template

```
1  ### SQLite SQL tables, with their  
2      properties  
3  CREATE TABLE schools(  
4      SchID int primay key,  
5      city text  
6  );  
7  CREATE TABLE teachers(  
8      TechID int primay key,  
9      SchID int references schools(TechID)  
10 );  
11 #### {Question}  
12 SELECT
```

Create Table Template

# LLMs for Text2SQL - Preprocessing

---

- Add sample data to the prompt, i.e., the SQL query
- Combination of layout and sample data

## Observations on Question Formulations.

- **Layout.**
  - Significant **drop** in performance when **changing** from layout to unstructured format
- **Sample Data.**
  - With each of the layout **sample data** is **plug-able** and **beneficial**
  - More sample data is not always better
  - More sample data could be **risky** in case of **context length limitation**
- **Primary & Foreign Keys.**
  - **Removing foreign keys** significantly **decreases** the performance
  - SQLfuse adds **primary keys** as a **second contributor** to performance gain
  - Some **LLMs** see **drop** while adding the Fks

# LLMs for Text2SQL - Preprocessing

## Schema Linking.

Specify the tables and columns in the database that correspond to the phrases in the given query

Desired SQL:

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

Nature Language Question:

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

Schema:

cars_data						
id	mpg	cylinders	edispl	horspower	weight	year

**cars\_names**

make_id	model	make

**model\_list**

make_id	maker	model

**car\_makers**

id	maker	full_name

# LLMs for Text2SQL - Preprocessing

**Schema Linking.** Divide & Prompt. Clause by Clause DnP (CC-DnP).

## Input

These are tables(columns):

Products(Product\_ID, Product\_Type\_Code...),

...

This is the text: Show the names of products that are in at least two events.

Translate the text to SQL query base on tables(columns).

1. Write the **FROM clause**: FROM Products AS T1 JOIN

Products\_in\_Events AS T2 ON T1.Product\_ID = T2.Product\_ID ,

2. Write the **GROUP BY clause**: GROUP BY T1.Product\_Name,

3. Writer the **HAVING clause**: HAVING COUNT(\*) >= 2,

4. Write the **SELECT clause**: SELECT T1.Product\_Name

5. **Final SQL query**: SELECT T1.Product\_Name FROM Products

AS T1 JOIN Products\_in\_Events AS T2 ON T1.Product\_ID =

T2.Product\_ID GROUP BY T1.Product\_Name HAVING COUNT(\*)

>= 2;

These are tables(columns):

stadium(stadium\_id, location...),

singer(singer\_id, name, country...),

...

This is the text: Show the name and the release year of the song by the youngest singer.

Translate the text to SQL query base on tables and columns.

## Model output

1. Writer the **FROM clause**: FROM singer,

2. Writer the **subquery**: SELECT MIN(age) FROM singer,

3. Writer the **WHERE clause**: WHERE age = (SELECT MIN(age)  
FROM singer),

4. Writer the **SELECT clause**: SELECT song\_name,

song\_release\_year,

5. **The final SQL is:**

SELECT song\_name, song\_release\_year

FROM singer

WHERE age = (SELECT MIN(age) FROM singer)

# LLMs for Text2SQL - Preprocessing

Divide & Prompt. Schema Linking DnP (SL-DnP).

Input

These are tables(columns):

Products(Product\_ID, Product\_Type\_Code...),

...  
This is the text: Show the names of products that are in at least two events.

Translate the text to SQL query base on tables(columns).

**1. Identify the relevant tables:**

Products(Product\_ID, Product\_Type\_Code, Product\_Name,  
Product\_Price),  
Products\_in\_Events(Product\_in\_Event\_ID, Event\_ID, Product\_ID)

**2. Final SQL query:** SELECT T1.Product\_Name FROM Products  
AS T1 JOIN Products\_in\_Events AS T2 ON T1.Product\_ID =  
T2.Product\_ID GROUP BY T1.Product\_Name HAVING COUNT(\*)  
>= 2;

These are tables(columns):

stadium(stadium\_id, location...),  
singer(singer\_id, name, country...),

...  
This is the text: Show the name and the release year of the song by  
the youngest singer.

Translate the text to SQL query base on tables and columns.

Model output

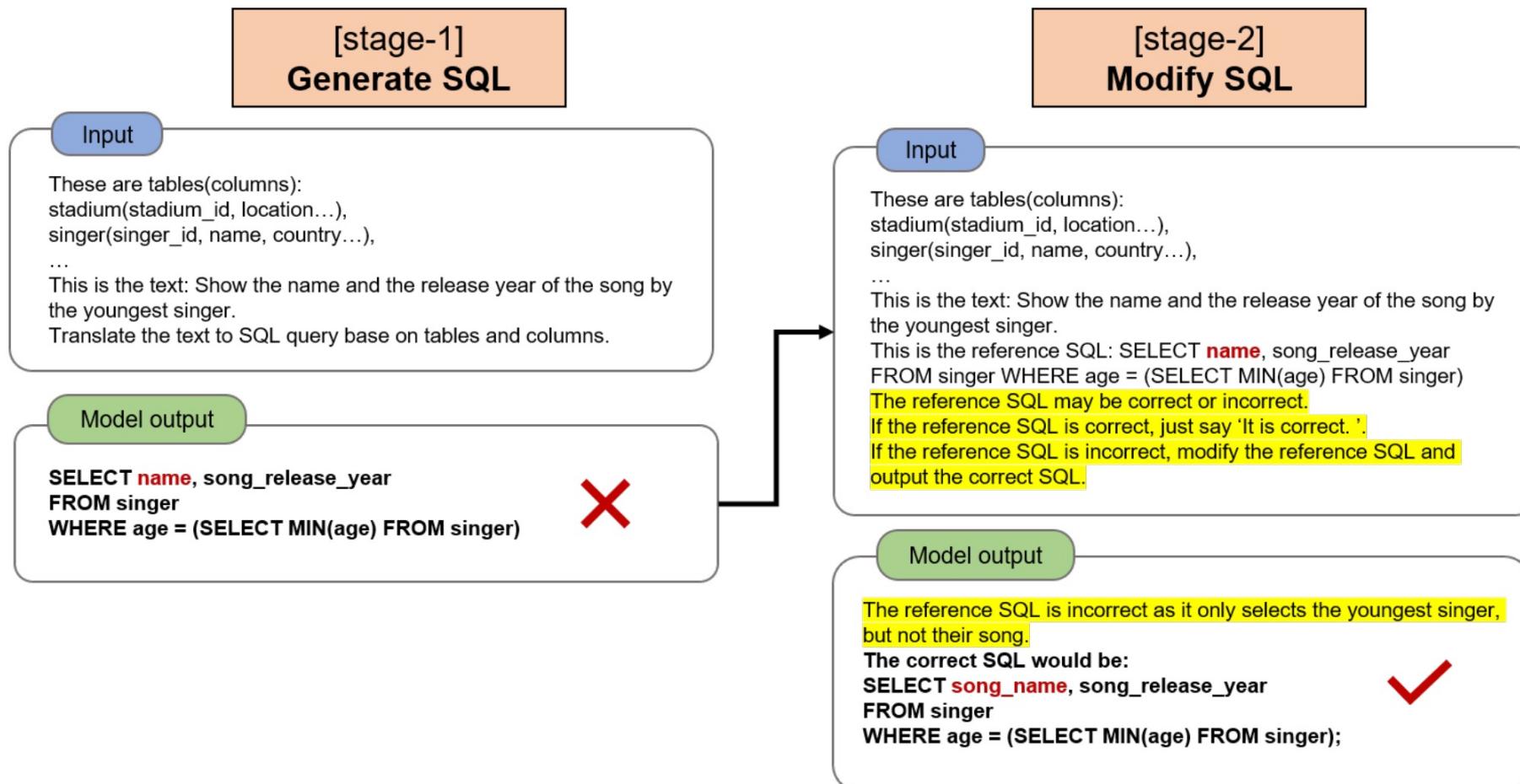
**1. identify the relevant tables:** singer(singer\_id, name, country,  
song\_name, song\_release\_year, age, is\_male).

**2. The final SQL is:**

SELECT song\_name, song\_release\_year  
FROM singer  
WHERE age = (SELECT MIN(age) FROM singer)

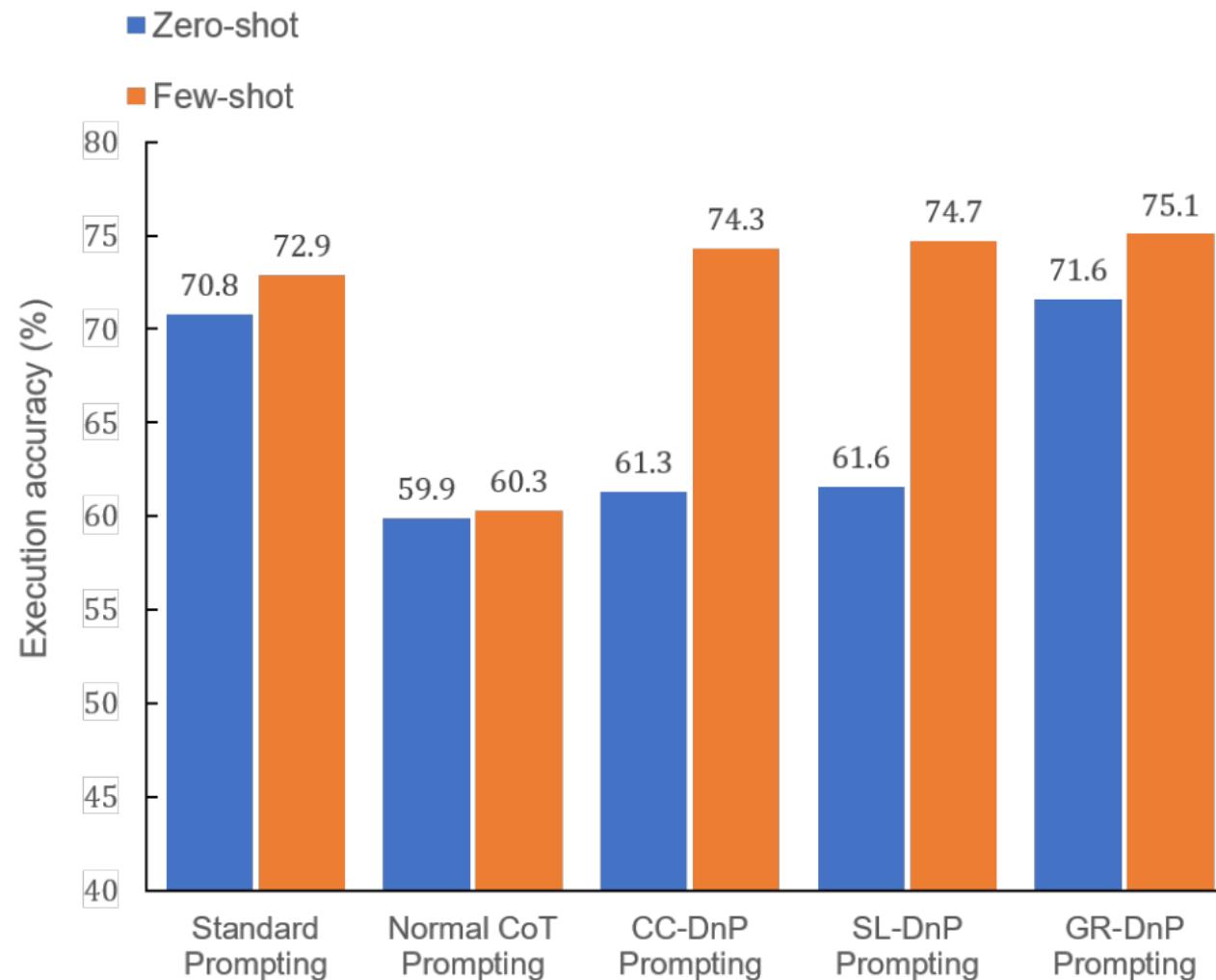
# LLMs for Text2SQL - Preprocessing

Divide & Prompt. Generate and Refine (GR-DnP).



# LLMs for Text2SQL - Preprocessing

Divide & Prompt. Results.



# LLMs for Text2SQL - Inference

## Decomposition. QDecomp

### Model Input (API Doc Format)

```
### SQLite SQL tables, with their properties:  
#  
# allergy_type (allergy, allergytype)  
# has_allergy (stuid, allergy)  
# student (stuid, lname, fname, age, sex, major, advisor, city_code)  
#  
### Show first name, last name, age for all female students. Their sex is F.
```

### Single-Pass Prompting

#### A. Chain-of-Thought

# This query chooses the records from the student table, followed by a WHERE clause that selects records where the sex column is equal to F. It then selects the fname, lname, and age columns

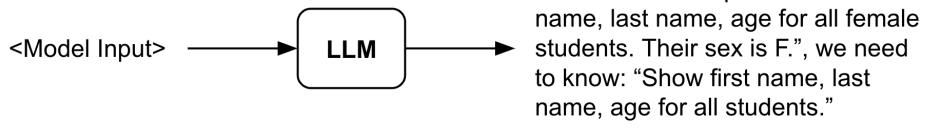
# Thus, the answer for the question is: Show first name, last name, age for all female students. Their sex is F.

```
SELECT fname, lname, age FROM student WHERE sex = 'F'
```

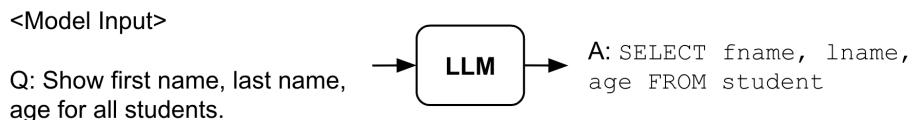
### Iterative Prompting

#### B. Least-to-Most Prompting

##### Problem Reduction:



##### Problem Solving:



##### <Model Input>

Q: Show first name, last name, age for all students.

A: SELECT fname, lname, age FROM student

Q: Show first name, last name, age for all female students.  
Their sex is F.

A: SELECT fname, lname, age FROM student WHERE sex = 'F'

# LLMs for Text2SQL - Inference

## Decomposition. QDecomp

### Model Input (API Doc Format)

```
### SQLite SQL tables, with their properties:  
#  
# allergy_type (allergy, allergytype)  
# has_allergy (stuid, allergy)  
# student (stuid, lname, fname, age, sex, major, advisor, city_code)  
#  
### Show first name, last name, age for all female students. Their sex is F.
```

### C. QDecomp

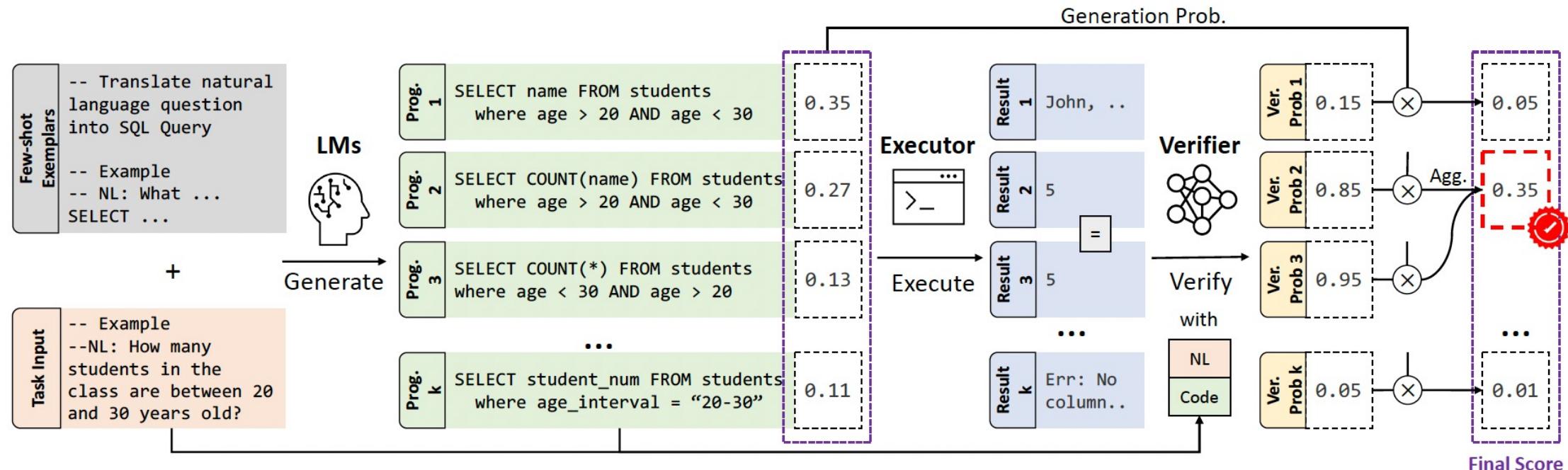
1. Show first name, last name, age for all students.
  2. Show first name, last name, age for all female students. Their sex is F.
- # Thus, the answer for the question is: Show first name, last name, age for all female students. Their sex is F.  
SELECT fname, lname, age FROM student WHERE sex = 'F'

### D. QDecomp + InterCOL

1. Show first name, last name, age for all students.  
SQL table (column): student (fname, lname, age)
  2. Show first name, last name, age for all female students. Their sex is F.  
SQL table (column): student (sex)
- # Thus, the answer for the question is: Show first name, last name, age for all female students. Their sex is F.  
SELECT fname, lname, age FROM student WHERE sex = 'F'

# LLMs for Text2SQL - Postprocessing

**Consistency.** Adopts a majority voting strategy, allowing the same LLM to generate multiple answers to the same question by setting a certain temperature.



# LLMs for Text2SQL - Postprocessing

---

LEVER.

Language to Code Generation

$$P_{LM} (y|x) = P(y \mid \text{prompt}(x, \{(x_i, y_i)\}_{i < m}))$$

$$\hat{y}_{greedy} \approx \arg \max_y P_{LM} (y|x)$$

# LLMs for Text2SQL - Postprocessing

---

LEVER.

Reranking. Reranking of Program Candidates

$$\hat{y}_{rerank} = \arg \max_{\hat{y} \in S} R(x, \hat{y})$$

# LLMs for Text2SQL - Postprocessing

---

LEVER.

Reranking. Program Sampling from Code LLMs

- $k$  programs from  $P_{LM}(y|x)$  with **temperature sampling**, i.e.,  $\{\hat{y}_i\}_{i=1}^k \sim P_{LM}(y|x)$
- **Deduplication** is performed since the same program can be sampled more than once

$$S = \{\hat{y}_i\}_{i=1}^n \text{ where } n \leq k$$

- Sampling is performed instead of beam search because
  - Beam search for code generation results in **worst performance** due to **degenerated programs**
  - Beam search is **not available or efficiently implemented** for all **LLMs** used in the work

# LLMs for Text2SQL - Postprocessing

LEVER.

Reranking. Verification with Execution

Reranking probability

Joint probability

$$P_R(\hat{y}, v_{=1} | x) = P_{\text{LM}}(\hat{y} | x) \cdot P_{\theta}(v_{=1} | x, \hat{y}, \mathcal{E}(\hat{y}))$$

generation

Passing the verification

$x$  = Input

$\hat{y}$  = Candidate Program

$\mathcal{E}(\hat{y})$  = Execution Results

# LLMs for Text2SQL - Postprocessing

LEVER.

Reranking. Execution Result Aggregation

$$R(x, \hat{y}) = P_R(\mathcal{E}(\hat{y}), v_{=1} | x) = \sum_{y \in S, \mathcal{E}(y) = \mathcal{E}(\hat{y})} P_R(y, v_{=1} | x)$$



Likelihood of the best output for the input

# LLMs for Text2SQL - Postprocessing

---

LEVER.

Learning the Verifiers. Learning Objective

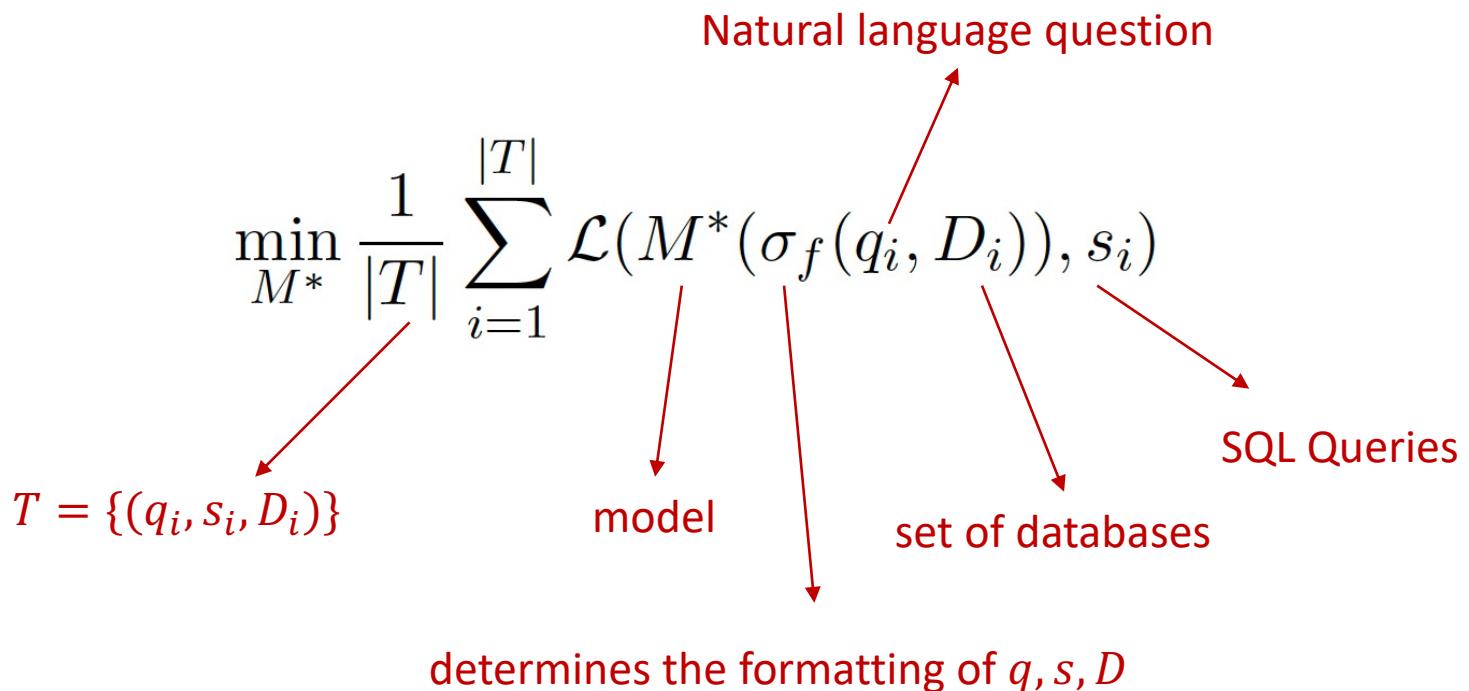
Negative log-likelihood

$$\mathcal{L}_\theta(x, S) = -\frac{1}{|S|} \cdot \sum_{\hat{y}_i \in S} \log P_\theta(v_i | x, \hat{y}_i, \hat{z}_i)$$

# LLMs for Text2SQL - Finetuning

## DTS-SQL. Supervised Fine-tuning

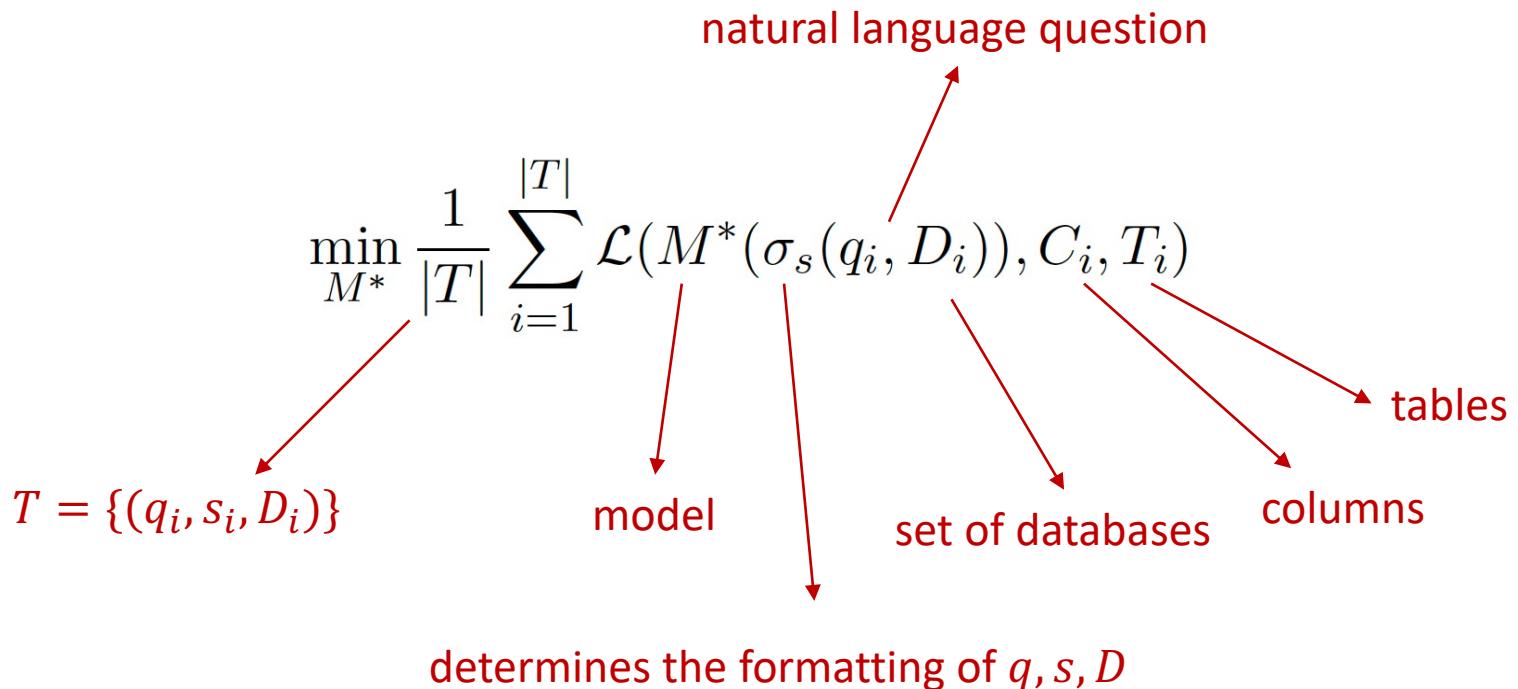
Next token prediction loss



# LLMs for Text2SQL - Finetuning

## DTS-SQL. Decomposed Supervised Fine-tuning

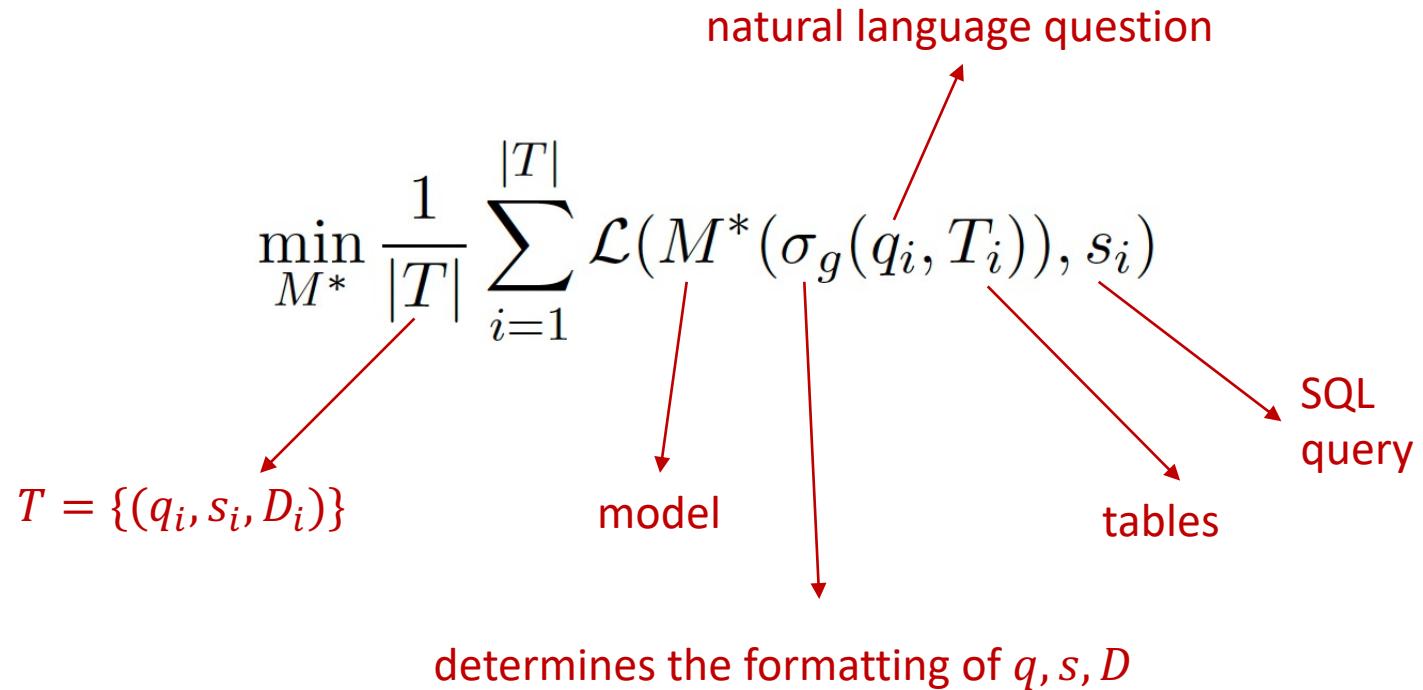
### Schema Linking Fine-Tuning



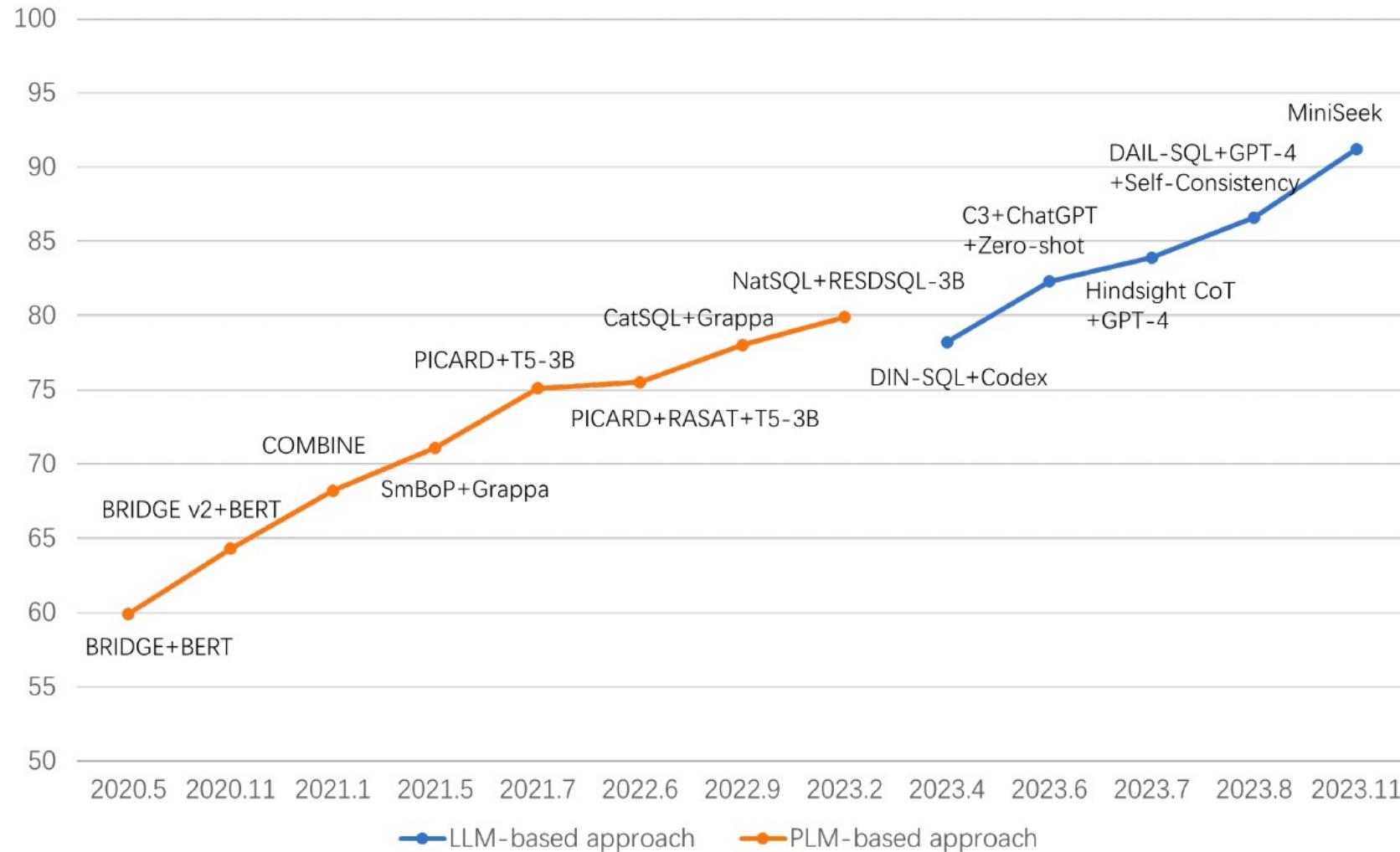
# LLMs for Text2SQL - Finetuning

## DTS-SQL. Decomposed Supervised Fine-tuning

### SQL Generation Fine-Tuning



# Accuracy for T2S over time



# Benchmark Datasets for T2S - Traditional

---

## WikiSQL.

- A large crowd-sourced dataset for developing NL interfaces for relational databases.
- 80654 hand-annotated questions and SQL queries across 24241 tables from Wikipedia.

## Spider 1.0.

- Parallel data of utterances and queries
- 10,181 examples across 200 DBs
- Each example consists of an occurrence and a DB with one or more tables and an annotated query

## KaggleDBQA.

- Cross-domain and complex evaluation dataset of real Web databases,
- with domain-specific data types, original formatting, and unrestricted questions
- Databases are pulled from real-world data sources and not normalized
- Provides database documentation

# Benchmark Datasets for T2S - Traditional

---

## Spider-Realistic.

- Creates scenarios where the explicit appearance of the database schema in question is eliminated

## Spider-SYN.

- Replaces schema-related words in questions with manually selected synonyms

## SPIDER-CG.

- Splits the sentences in Spider into sub-sentences,
- Annotates each sub-sentence with the corresponding SQL clause,

## Cspider.

- Translates Spider into Chinese to test the performance of Chinese text-to-SQL.

## ADVETA.

- Adds adversarial table perturbation to deceive Text-to-SQL parsers

# Benchmark Datasets for T2S – LLM Era

---

## ScienceBenchmark

- Adaptation of spider dataset to domain specific queries is limited
- Sciencebenchmark is created with the help of SQL experts and domain specific researchers

## BIRD

- Existing datasets focus on the database schema with few rows of database values
- Dataset of 37 professional domains
- Dirty and noisy database values
- External knowledge grounding between NL questions and database values
- Efficiency (in the context of massive databases)
- GPT-4 achieved 54.89% accuracy, human result 92.96%

# Benchmark Datasets for T2S – LLM Era

---

## Dr.Spider

- Existing models are sensitive to task specific perturbations
- 17 perturbations in T2S tasks
- Leveraged PLMs to simulate human behaviors in creating natural questions

## Spider2.0

- Real scenarios involve complex cloud or local data across DB systems
- Multiple SQL queries in various dialects
- Operations from data transformation to analytics

# Evaluation Metrics – Single Turn T2S

## Exact Set Match Accuracy (EM)

- Compare ground truth SQL query with predicted SQL query
- SQL clauses: SELECT, GROUP BY, WHERE, ORDER BY, KEYWORDS

$$score(\hat{Y}, Y) = \begin{cases} 1, & \hat{Y} = Y \\ 0, & \hat{Y} \neq Y \end{cases}$$

Where  $\hat{Y} = \{(\hat{k}^i, \hat{v}^i), i \in (1, m)\}$  and  $Y = \{(k^i, v^i), i \in (1, m)\}$

$$EX = \frac{\sum_{n=1}^N score(\hat{Y}_n, Y_n)}{N}$$

**Drawback.** Often underestimates the prediction accuracy of the model

k = SQL clause

v = corresponding value

# Evaluation Metrics – Single Turn T2S

---

## Execution Accuracy (EX).

- Comparison of the execution results
- the ground-truth SQL query and the predicted SQL query

$$score(\hat{V}, V) = \begin{cases} 1, & \hat{V} = V \\ 0, & \hat{V} \neq V \end{cases}$$

$$EX = \frac{\sum_{n=1}^N score(\hat{Y}_n, Y_n)}{N}$$

## Drawback.

- SQL with different logic may yield the same results when executed
- EX overestimates the prediction accuracy of the model

# Evaluation Metrics – Multiturn T2S

---

## Question Match Accuracy (QM).

- EM score over all questions
- Its value is 1 for each question only if all predicted SQL clauses are correct.

$$score(\hat{Y}, Y) = \begin{cases} 1, & \hat{Y} = Y \\ 0, & \hat{Y} \neq Y \end{cases}$$

$$QM = \frac{\sum_{m=1}^M score(\hat{Y}_m, Y_m)}{M}$$

$M$  = total number of question

# Evaluation Metrics

---

## Valid Efficiency Score (VES).

- SQL execution efficiency in the evaluation scope
- The higher the correctness rate of the generated SQL, the higher the execution efficiency of the generated SQL, and the higher the VES value

$$VES = \frac{\sum_{n=1}^N 1(V_n, \hat{V}_n) \cdot R(Y_n, \hat{Y}_n)}{N}, R(Y_n, \hat{Y}_n) = \sqrt{\frac{E(Y_n)}{E(\hat{Y}_n)}}$$

# Summary

---

- What are the traditional approaches for T2S?
- Masked Language Models for T2S (TaBERT)
- Large Language Models for T2S
  - Preprocessing
  - Inference
  - Post-processing
  - Fine-tuning