



The Ohio State University



Dr.Spider: A Diagnostic Evaluation Benchmark towards Text-to-SQL Robustness

Shuaichen Chang¹, Jun Wang², Mingwen Dong², Lin Pan², Henghui Zhu², Alexander Hanbo Li²,
Wuwei Lan², Sheng Zhang², Jiarong Jiang², Joseph Lilien², Steve Ash², William Wang²,
Zhiguo Wang², Vittorio Castelli², Patrick Ng², Bing Xiang²

¹Ohio State University, ² AWS AI Labs

ICRL 2023

Selected as a notable-top-5% paper

Text-to-SQL

NLQ

Find the name and rank of the 3 youngest winners across all matches.

DB

Ranking	Date	Ranking
Player_ID	First_Name	Last_Name
Winner_Name	Winner_Rank	Age
Serena Williams	1	32
...

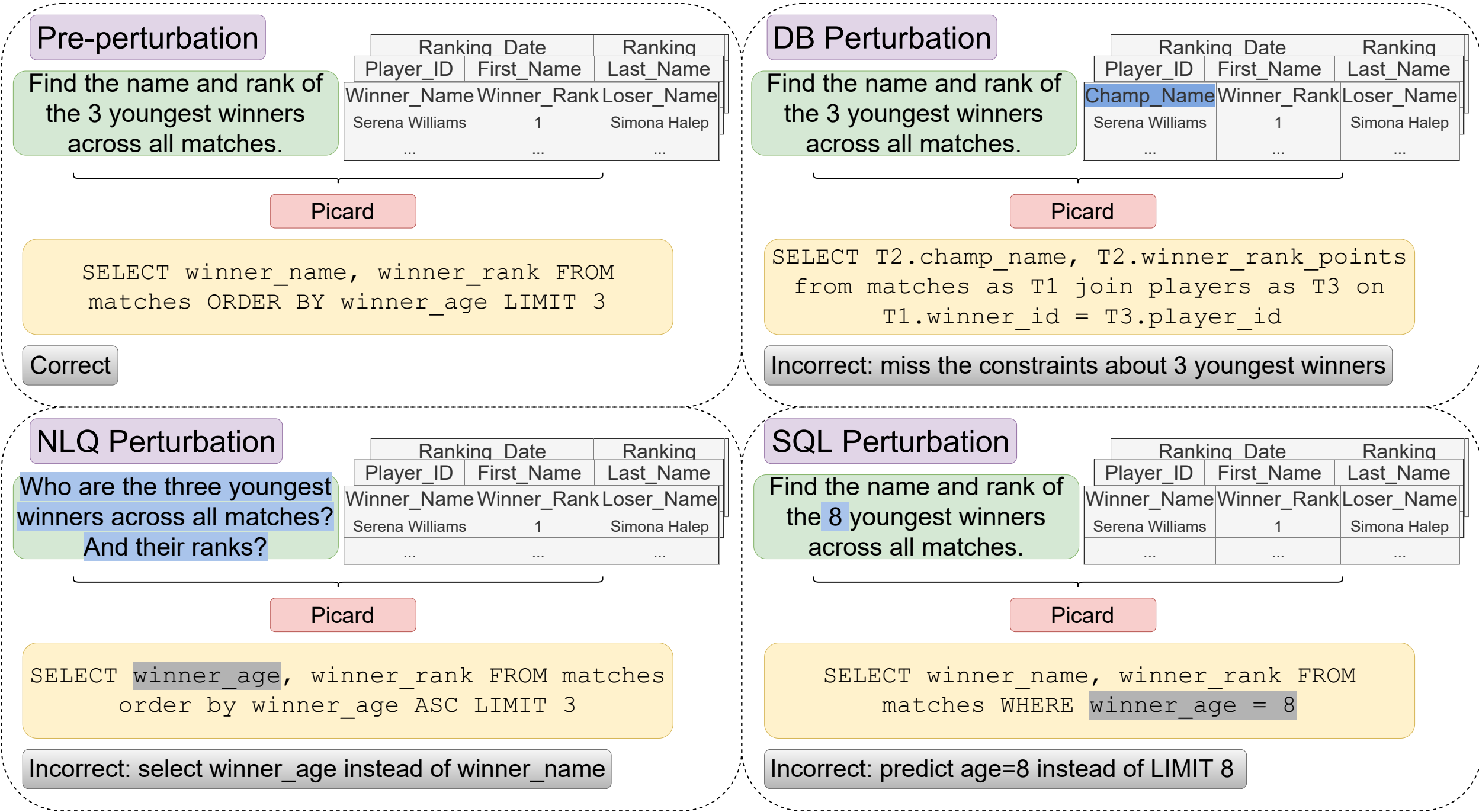
Text-to-SQL Model

SQL

```
SELECT winner_name, winner_rank FROM
matches ORDER BY winner_age ASC LIMIT 3
```

Does the high accuracy imply a robust model?

We perturb the data in Spider¹ to uncover the weaknesses in models that might not be evident when evaluated on the original data.



[1] Yu, Tao, et al. "Spider: A Large-Scale Human-Labeled Dataset for Complex and Cross-Domain Semantic Parsing and Text-to-SQL Task." *EMNLP*. 2018.

Our Evaluation Benchmark

Includes 17 perturbation types across all three text-to-SQL components

Diagnoses models with task-specific robustness phenomena

Provide insights for developing robust text-to-SQL models

Text-to-SQL Example

Pre-perturbation

Find the name and rank of the 3 youngest winners across all matches.

Ranking Date	Ranking	
Player_ID	First_Name	Last_Name
Winner_Name	Winner_Rank	Loser_Name
Serena Williams	1	Simona Halep
...

Picard

```
SELECT winner_name, winner_rank FROM
matches ORDER BY winner_age LIMIT 3
```

Correct

Database (DB) Perturbations

- Data in a database can be represented in various ways.

Pre-perturbation

Find the name and rank of the 3 youngest winners across all matches.

Player_ID	First_Name	Last_Name
Winner_Name	Winner_Rank	Loser_Name
Serena Williams	1	Simona Halep
...

Picard

```
SELECT winner_name, winner_rank FROM matches ORDER BY winner_age LIMIT 3
```

Correct

DB Perturbation

Find the name and rank of the 3 youngest winners across all matches.

Player_ID	First_Name	Last_Name
Champ_Name	Winner_Rank	Loser_Name
Serena Williams	1	Simona Halep
...

Picard

```
SELECT T2.champ_name, T2.winner_rank_points from matches as T1 join players as T3 on T1.winner_id = T3.player_id
```

Incorrect: miss the constraints about 3 youngest winners

Natural Language Question (NLQ) Perturbations

- The same query intent can be expressed with different phrasings.

Pre-perturbation

Find the name and rank of the 3 youngest winners across all matches.

Player_ID	First_Name	Last_Name
Winner_Name	Winner_Rank	Loser_Name
Serena Williams	1	Simona Halep
...

Picard

```
SELECT winner_name, winner_rank FROM matches ORDER BY winner_age LIMIT 3
```

Correct

NLQ Perturbation

Who are the three youngest winners across all matches? And their ranks?

Player_ID	First_Name	Last_Name
Winner_Name	Winner_Rank	Loser_Name
Serena Williams	1	Simona Halep
...

Picard

```
SELECT winner_age, winner_rank FROM matches order by winner_age ASC LIMIT 3
```

Incorrect: select winner_age instead of winner_name

SQL Perturbations

- Modifying logical and symbolic units in a question should not impact understanding of other components.

Pre-perturbation

Find the name and rank of the 3 youngest winners across all matches.

Player_ID	First_Name	Last_Name
Winner_Name	Winner_Rank	Loser_Name
Serena Williams	1	Simona Halep
...

Picard

```
SELECT winner_name, winner_rank FROM matches ORDER BY winner_age LIMIT 3
```

Correct

SQL Perturbation

Find the name and rank of the 8 youngest winners across all matches.

Player_ID	First_Name	Last_Name
Winner_Name	Winner_Rank	Loser_Name
Serena Williams	1	Simona Halep
...

Picard

```
SELECT winner_name, winner_rank FROM matches WHERE winner_age = 8
```

Incorrect: predict age=8 instead of LIMIT 8

Perturbation Data Creation

- For DBs and SQL queries, we perturb them programmatically, using their inherent structure.
- For NLQ perturbations, we propose to use LLMs to simulate human paraphrasing with various task-specific phenomena that we find in crowdsourcing annotations.

Paraphrase Category	Example
Keyword-synonym (Replace SQL keyword indicators with synonyms)	Question: What is the code of airport that has the highest number of flights? Paraphrases: Show me the code for the airport that currently has the most flights.
Keyword-carrier (Imply SQL keyword indicators by carrier phrases)	Question: Show the name and theme for all concerts and the number of singers in each concert. Paraphrase: List the names and themes for all concerts and how many singers are in each.
Column-synonym (Replace column indicators with synonyms)	Question: List the name of teachers whose hometown is not Little Lever Urban District. Paraphrases: Find the name of teachers who were not born in Little Lever Urban District.
Column-carrier (Imply column indicators by carrier phrases)	Question: Show the name of teachers aged either 32 or 33? Paraphrases: Which teachers are aged either 32 or 33.
Column-attribute (Imply column indicators by aggregated attributes)	Question: What is the name of the conductor who has worked the greatest number of year ? Paraphrases: Who has worked the longest as conductor?
Column-value (Imply column indicators by values)	Question: What are the ids of the students who do not own cats as pets ? Paraphrases: Find the IDs of students who don't own cats .
Value-synonym (Replace value indicators with synonyms)	Question: Find all airlines that have at least 10 flights. Paraphrases: Show the number of airlines that have at least ten flights.
Multitype (Contain multiple phenomena above)	Question: Find number of pets owned by students who are older than 20 ? Paraphrases: How many pets are owned by students over 20 ?
Others	Question: what is the name and nation of the singer who have a song having 'Hey' in its name? Paraphrases: Which singers have 'Hey' in their song's name? List their name and nation.

Our Main Findings

- SOTA supervised learning model (Picard) suffers
 - **14%** overall performance drop
 - **50%** performance drop to the most challenging perturbation type `DBcontent-equivalence`

Our Main Findings

- SOTA supervised learning model (Picard) suffers
 - 14% overall performance drop
 - 50% performance drop to the most challenging perturbation type `DBcontent-equivalence`
- Code-pretrained LLM (GPT-3 Codex) with in-context learning is
 - **more robust** than supervised learning models to **DB and SQL perturbations**
 - **less robust** than supervised learning models to **NLQ perturbations**

Our Main Findings

- SOTA supervised learning model (Picard) suffers
 - 14% overall performance drop
 - 50% performance drop to the most challenging perturbation type DBcontent-equivalence
- LLM (GPT-3 Codex) with in-context learning is
 - more robust than supervised learning models to DB and SQL perturbations
 - less robust than supervised learning models to NLQ perturbations
- Some model designs are beneficial to model robustness, e.g. **larger model size**
- Some model designs create a **reversed performance** to certain perturbations
 - strengthen the model performance on the original data but weaken the model robustness

Our Main Findings

- SOTA supervised learning model (Picard) suffers
 - 14% overall performance drop
 - 50% performance drop to the most challenging perturbation type DBcontent-equivalence
- LLM (GPT-3 Codex) with in-context learning is
 - more robust than supervised learning models to DB and SQL perturbations
 - less robust than supervised learning models to NLQ perturbations
- Some model designs are beneficial to model robustness, e.g. larger model size
- Some model designs create a reversed performance to certain perturbations
 - strengthen the model performance on the original data but weaken the model robustness
- More findings can be found in our paper <https://openreview.net/pdf?id=Wc5bmZZU9cy>