

DataFrame QA: A Universal LLM Framework on DataFrame Question Answering Without Data Exposure

Junyi Ye, Mengnan Du and Guiling Wang

New Jersey Institute of Technology, USA

December 5, 2024

Current LLM-Based Table QA Pipeline

Instruction: Using the table provided, answer the following questions based on the data.

Table:

Product	Price	Quantity Available	Rating
Laptop	\$1,000	15	4.5
Smartphone	\$800	30	4.7
Tablet	\$400	25	4.3
Headphones	\$200	50	4.1
Monitor	\$300	20	4.2

Question: What is the cheapest product available in the table?

Limitations of current LLM-Based Table QA pipeline:

- LLM's Context Window Limits.
- Inefficient Token Use & High Costs.
- Data Leakage Risk.
- Decreased Performance with Long Prompts.
- Weak Mathematical/Data Analysis Capabilities.
- Processing Unnecessary Data.
- Table miss alignment.

Figure 1: Typical prompt for table QA.

Introduction

- ❖ This paper introduces **DataFrame QA**, a novel task that utilizes **LLMs** to generate Pandas queries for information retrieval and data analysis on dataframes, emphasizing safe and non-revealing data handling.
- ❖ Emphasizes **data privacy** by relying on dataframe column names and types, not values.
- ❖ Addresses limitations of existing LLM-based data analysis methods.

System Prompt

You are a professional Python programming assistant. Write Pandas code to get the answer to the user's question.

- Assumptions:

- The Pandas library has been imported as `pd`. You can reference it directly.
- The dataframe `df` is loaded and available for use.
- All string values in the `df` have been converted to lowercase.

- Requirements:

- Use only Pandas operations for the solution.
- Store the answer in a variable named `result`.
- Do NOT include comments or explanations in your code.
- Place your code between the [PYTHON] and [/PYTHON] tags.

User Prompt

You are given a Pandas dataframe named `df`:

- Columns: ['Player', 'No.', 'Nationality', 'Position', 'Years in Toronto', 'School/Club Team']
- Data Types: ['str', 'str', 'str', 'str', 'str', 'str']
- User's Question: what is terrence ross' nationality

LLM Response

```
[PYTHON]
result = df.loc[df['Player'] == 'terrence ross', 'Nationality'].values[0]
[/PYTHON]
```

Ground Truth

"united states"

Pandas Query Execution Result

"united states"



Figure 2: An example where a LLM can generate a correct Pandas query to answer user question.

Methodology

Task: Generate executable Pandas queries for answering user's questions on dataframes.

Focus: Use column headers and metadata, ensuring **data security**.

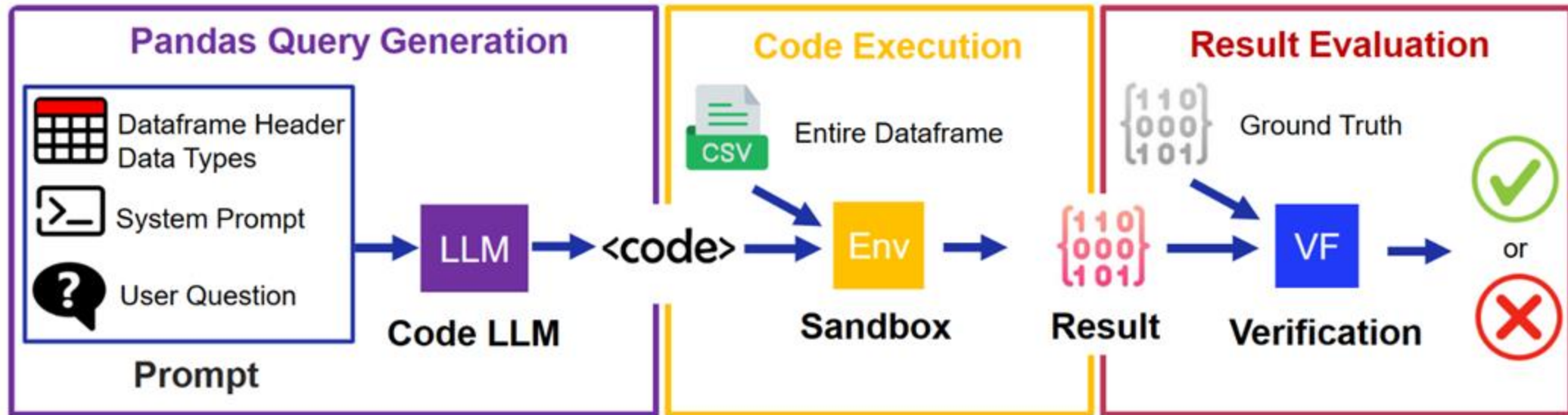


Figure 3: Framework of DataFrame QA. Note that, LLM in the figure can be replaced with any fine-tuned NLP model trained for the DataFrame QA task.

Datasets

1. WikiSQL (Simple Query Dataset) [1]

- ❖ Adapted for DataFrame QA, transforming tables into dataframes.
- ❖ Question types: 71% simple retrieval, 29% aggregation tasks.

2. UCI-DataFrameQA (Complex Dataset)

- ❖ Develop a DataFrame QA dataset reflecting **real-world scenarios** using GPT4.
- ❖ Data sourced from the **UCI dataset** [2], spanning various domains.
- ❖ Represent three real-life data interaction roles: **(1)Data Scientist; (2)General Users; (3)Data Owners**

User Question	Pandas Query	Types
which province is bay of islands in?	<pre>result = df.loc[df['Electorate']=='bay of islands', 'Province'].iloc[0]</pre>	Retrieval
how many combined days did go shiozaki have?	<pre>result = df.loc[df['Wrestler']=='go shiozaki', 'Combined days'].values[0]</pre>	Aggregation
how does the average shell weight vary across different numbers of rings?	<pre>result = df.groupby('Rings')['Shell_weight'].mean()</pre>	Data Analysis
can you create a new column 'volume' as a product of length, diameter, and height, then find the average volume for each sex?	<pre>df['Volume'] = df['Length'] * df['Diameter'] * df['Height'] result = df.groupby('Sex')['Volume'].mean()</pre>	Data Analysis

Table 1: Examples of Sample Questions and Corresponding Pandas Queries.

Experiment and Results

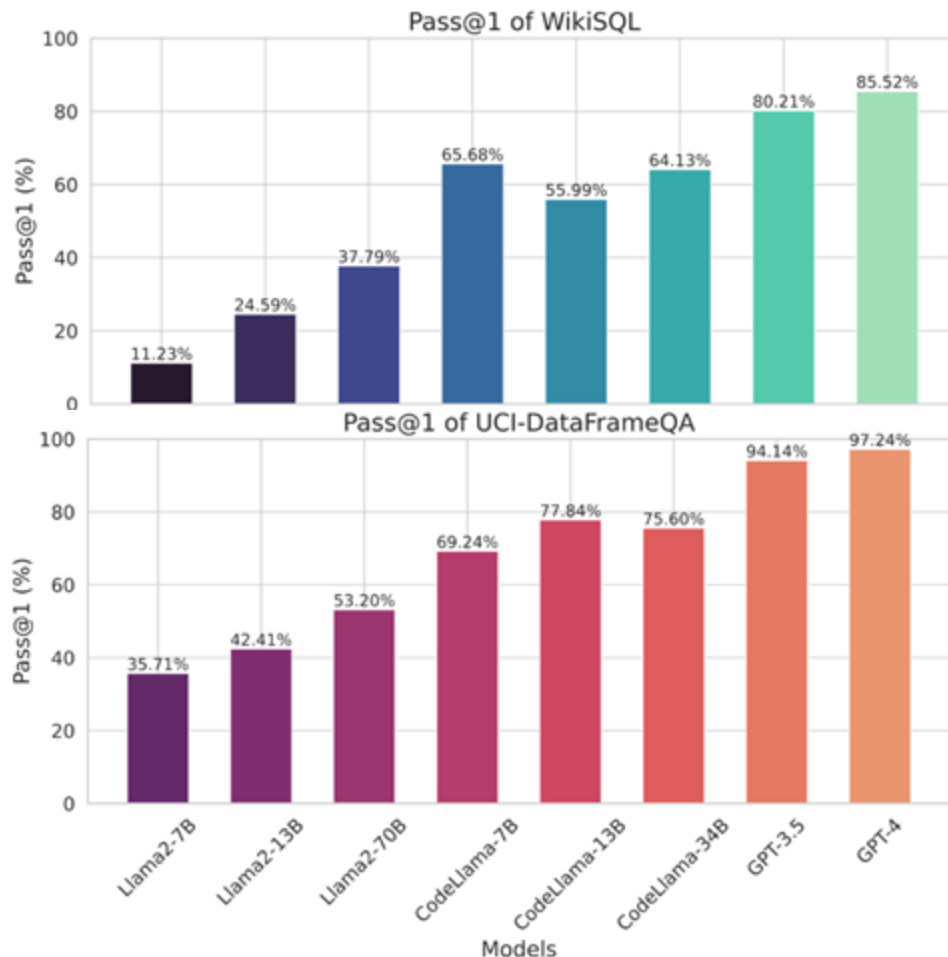


Figure 4: Performance of LLMs on Two Datasets.

Findings:

- ❖ **GPT-4 Dominance:** GPT-4 demonstrated exceptional performance on both datasets, achieving high accuracy.
- ❖ **Scaling Laws [3]:** Larger models generally perform better, but CodeLlama models are an exception.
- ❖ **Comparison to Text-to-SQL:** GPT-4's zero-shot approach and the complexity of Pandas queries slightly hinder its performance compared to specialized Text-to-SQL models.

Error Analysis

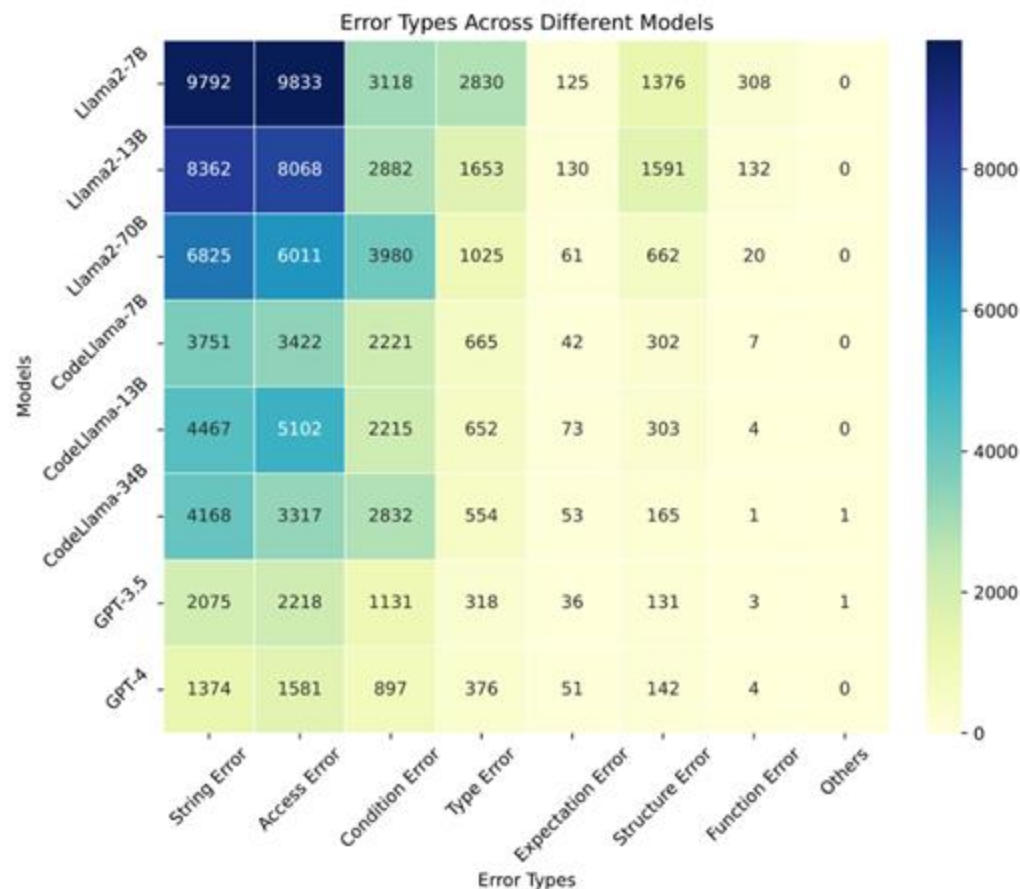


Figure 5: Distribution of Error Types Among Different LLMs on WikiSQL.

Error Type	Example
Value Retrieval Error	<p><i>Q</i> : Which province is grey and bell electorate in?</p> <p><i>P'</i> : <code>result = df[(df['Electorate']=='grey') (df['Electorate']=='bell)]['Province'].unique()</code></p> <p><i>P</i> : <code>result=df[(df['Electorate']=='grey and bell)]['Province'].unique()</code></p>
Column Reference Error	<p><i>Q</i> : What is the variance in resting blood pressure (trestbps) among different heart disease severity levels?</p> <p><i>P'</i> : <code>df.groupby('ca')['trestbps'].var()</code></p> <p><i>P</i> : <code>df.groupby('num')['trestbps'].var()</code></p>
Instruction Misalignment	<p><i>Q</i> : What are the mean and median lengths of abalone by each sex category?</p> <p><i>P'</i> : <code>import pandas as pd</code> <code># Group the dataframe by sex</code> <code>...</code></p> <p><i>P</i> : <code>result = df.groupby('Sex')['Length'].agg(['mean', 'median'])</code></p>
Aggregation Error	<p><i>Q</i> : what is the total amount of allied-unrelated where the component is human capital?</p> <p><i>P'</i> <code>result = df[df['Component']=='human capital']['Allied-Unrelated'].sum()</code></p> <p><i>P</i> <code>result = df[df['Component']=='human capital']['Allied-Unrelated'].count()</code></p>
Function-Column Ambiguity Error	<p><i>Q</i> : what is the average score when the swimsuit score is 8.503?</p> <p><i>P'</i> : <code>result = df[df['Swimsuit']=='8.503']['Average'].mean()</code></p> <p><i>P</i> : <code>result = df[df['Swimsuit']=='8.503']['Average']</code></p>
Insufficient Column Data/Format Information	<p><i>Q</i> : which team played on december 5?</p> <p><i>P'</i> <code>df['Date'] = pd.to_datetime(df['Date'])</code> <code>result = df[df['Date'] == 'december 5']['Team']</code></p> <p><i>P</i> : <code>result = df[df['Date'] == 'december 5']['Team']</code></p>
Coding Syntax Error	<p><i>Q</i> : how does the average shell weight vary across different numbers of rings?</p> <p><i>P'</i> : <code>result = df.groupby('Rings').mean()['Shell.weight']</code></p> <p><i>P</i> : <code>result = df.groupby('Rings')['Shell.weight'].mean()</code></p>
Hallucination Error	<p><i>Q</i> : i'm interested in knowing the most common age of abalone. can you find that for me?</p> <p><i>P'</i> : <code>result = df['Age'].mode()</code></p> <p><i>P</i> : <code>result = df['Rings'].mode()</code></p>

Table 2: Typical Failure Cases in DataFrame QA Task.

Conclusion

- ❖ We introduce **DataFrame QA**, a **secure, zero-shot LLM** framework that leverages dataframe headers to address **data privacy** and **minimize extraneous prompts**.
- ❖ By enriching prompts with dataset descriptions, the framework improves performance.
- ❖ The success of this task relies on both **coding abilities** and **query comprehension**, with GPT-4 achieving high practical accuracy.

Reference

1. Zhong, Victor, Caiming Xiong, and Richard Socher. "Seq2sql: Generating structured queries from natural language using reinforcement learning." *arXiv preprint arXiv:1709.00103* (2017).
1. Blake, Catherine L. "UCI repository of machine learning databases." <http://www.ics.uci.edu/~mlern/MLRepository.html> (1998).
1. Kaplan, Jared, et al. "Scaling laws for neural language models." *arXiv preprint arXiv:2001.08361* (2020).

Thank you!
Q&A



GitHub