

Domain-Aware Prompting with LLMs for Data Science Notebooks

Sanjeet Chatterjee

Large Language Models (LLMs)



Transformer Architecture

Autoregressive neural networks that employ self-attention. This allows elements to interact directly and capture dependencies efficiently.



Emergent Capabilities

Scaling such models to billions of parameters has led to emergent behaviour such as in-context learning, multi-step reasoning and instruction following.



Coding Assistants

LLMs have demonstrated impressive results in translating natural language to code, giving rise to coding assistants such as GitHub Copilot.

Data Science Notebooks



Rapid Prototyping

Interactive Jupyter notebooks interleaving code and natural language markdown cells have become the standard for data science workflows.



Repetitive Workflows

Data scientists often go through repetitive workflows in data exploration and preparation.



Challenging for LLMs

Notebooks specifically offer additional challenges due to the multi-step nature of notebooks.

ARCADE: Exploratory Analysis

Table 2.1: Existing Datasets for Data Science Notebooks

Dataset	No. Notebooks	No. Tasks	Tasks / Notebook	Evaluation Method
JuICE	1457	3946	2.7	Exact Match + BLEU
MS DSP	305	1096	3.6	Unit Tests
ExeDS	277	534	1.9	Output Match
ARCADE	133	1082	8.55	Output Match

ARCADE: Exploratory Analysis

Table 2.1: Existing Datasets for Data Science Notebooks

Dataset	No. Notebooks	No. Tasks	Tasks / Notebook	Evaluation Method
JuICE	1457	3946	2.7	Exact Match + BLEU
MS DSP	305	1096	3.6	Unit Tests
ExeDS	277	534	1.9	Output Match
ARCADE	133	1082	8.55	Output Match

Table 2.4: Pandas Method Groups

Task Type	%	Met hods
Aggregation	42.7	groupby, agg ...
Transformation	27.8	apply, pivot_table, explode, cut, pct_change ...
Combination	3.0	concat, merge, append, join ...
Selection	68.5	loc, query, nlargest, sort_values, filter, isnull ...
Cleaning	22.4	fillna, rename, drop_duplicates, to_numeric ...
Strings	23.7	extract, startswith, str, replace, contains...
Computation	81.3	max, quantile, value_counts, div, corr ...
Datetime	2.7	to_datetime, strftime, period_range, to_timedelta ...
Visualization	5.9	plot, boxplot, barplot, scatter, hist ...

ARCADE: Exploratory Analysis

Table 2.1: Existing Datasets for Data Science Notebooks

Dataset	No. Notebooks	No. Tasks	Tasks / Notebook	Evaluation Method
JuICE	1457	3946	2.7	Exact Match + BLEU
MS DSP	305	1096	3.6	Unit Tests
ExeDS	277	534	1.9	Output Match
ARCADE	133	1082	8.55	Output Match

Table 2.4: Pandas Method Groups

Task Type	%	Met hods
Aggregation	42.7	groupby, agg ...
Transformation	27.8	apply, pivot_table, explode, cut, pct_change ...
Combination	3.0	concat, merge, append, join ...
Selection	68.5	loc, query, nlargest, sort_values, filter, isnull ...
Cleaning	22.4	fillna, rename, drop_duplicates, to_numeric ...
Strings	23.7	extract, startswith, str, replace, contains...
Computation	81.3	max, quantile, value_counts, div, corr ...
Datetime	2.7	to_datetime, strftime, period_range, to_timedelta ...
Visualization	5.9	plot, boxplot, barplot, scatter, hist ...

Figure 2.5: Tasks per Notebook

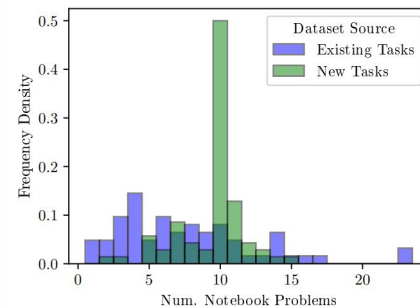
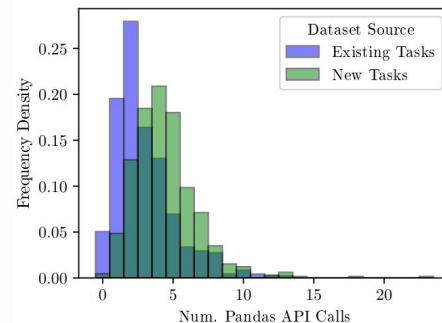


Figure 2.6: Pandas API Calls per Task



ARCADE: Sample Notebook

[1] `import pandas as pd`

c₁ `df = pd.read_csv('dataset/Gamepass_Games_v1.csv')`


[2] **u₁** Extract min and max hours as two columns 

```
def get_avg(x):  
    try: return float(x[0]) , float(x[1])  
    except: return 0, 0  
  
df['min'], df['max'] = zip(*df['TIME'].str.replace(  
    ' hours', '').str.split("-").apply(get_avg))
```

c₂


[3] `df['ADDED'] = pd.to_datetime(
 df['ADDED'], format="%d %b %y", errors='coerce')`

c₃

[4] **u₂** In which year was the most played game added? 


```
df['GAMERS']=df['GAMERS'].str.replace(  
    ',', ' ').astype(int)  
added_year=df[df['GAMERS'].idxmax()][ 'ADDED'].year
```

c₄

[5] **u₃** For each month in that year, how many games that has a rating of more than four? 


```
df[(df['ADDED'].dt.year== added_date.year) &  
(df['RATING']>4)].groupby(  
    df["ADDED"].dt.month)['GAME'].count()
```

c₅

[6] **u₄** What is the average maximum completion time for all fallout games added in 2021? 

```
fallout=df[df['GAME'].str.contains('Fallout')]  
fallout.groupby(fallout['ADDED'].dt.year).get_group(  
    2021)['max'].mean()
```

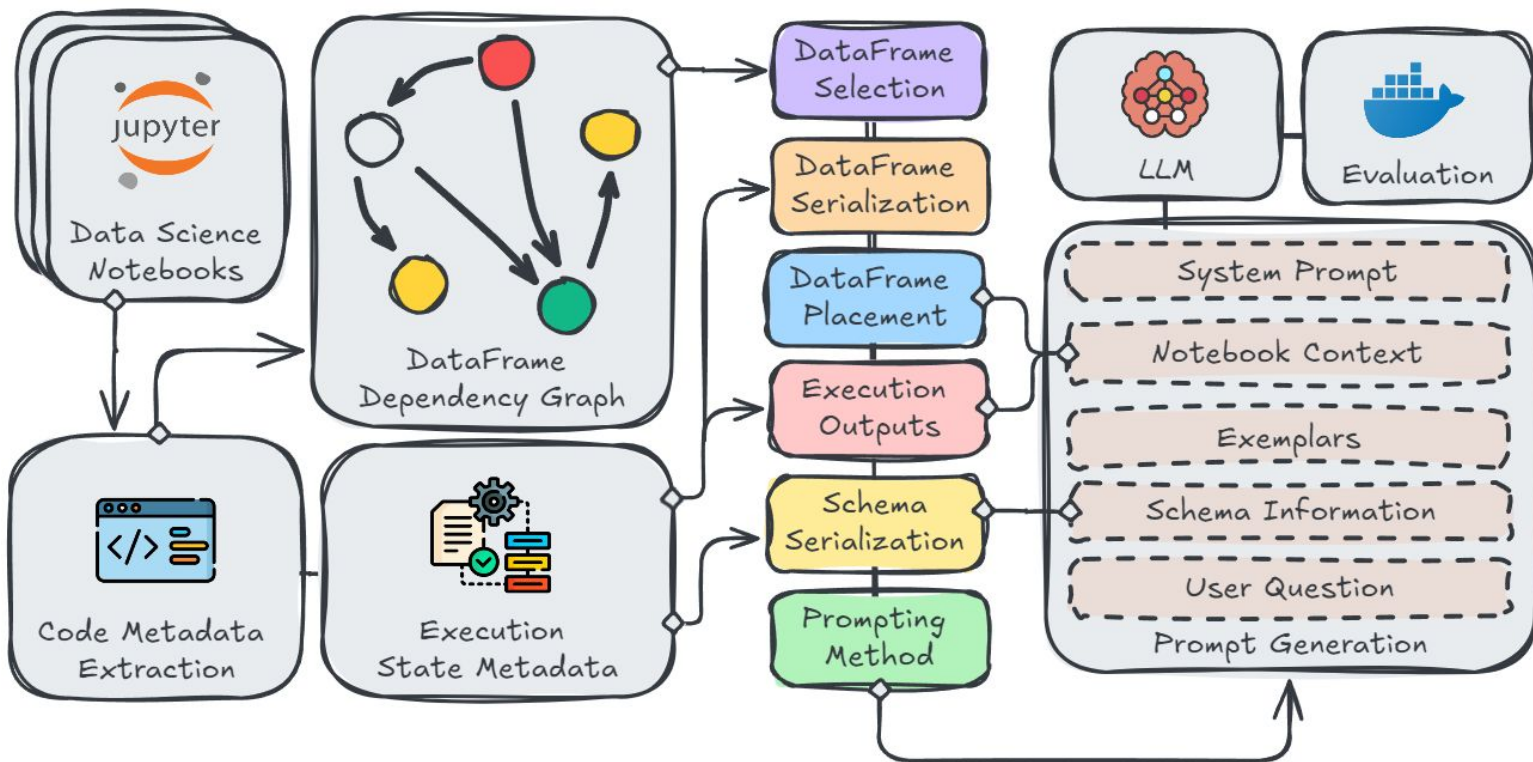
c₆

[7] **u₅** What is the amount of games added in each year for each month? (show a table with index as years, columns as months and fill null values with 0) 

```
pd.pivot_table(df, index=df['ADDED'].dt.year, ...,  
    aggfunc=np.count_nonzero,  
    fill_value='0').rename_axis(  
    index='Year', columns='Month')
```

c₇

System Overview



Execution State Metadata

- Imports
- Variables
- Functions



- Markdown
- Code Cells
- Outputs



Python Runtime

Notebook

Cell



Outputs

Cell execution results



DataFrame Variables

Pandas variables content

DataFrame Serialization

```
1 # In[:  
2 print(df_anim)  
3  
4 # Out[:  
5 [showing 3 sample rows out of 1812 rows]  
6  
7      theme (list)          name (str)  rating (float64)  
8 124    ["Historical", ..., "Gag Humor"]    "Gintama"          9.05  
9 1303          ["School"]    "Kareshi Jijou"          7.60  
1488          ["Mahou Shoujo"]    "Delicious Precure"          7.24
```

NEW

```
1 # Columns in df_anim with example values:  
2 # theme ([School]), name (Kareshi Jijou), rating (7.60)
```

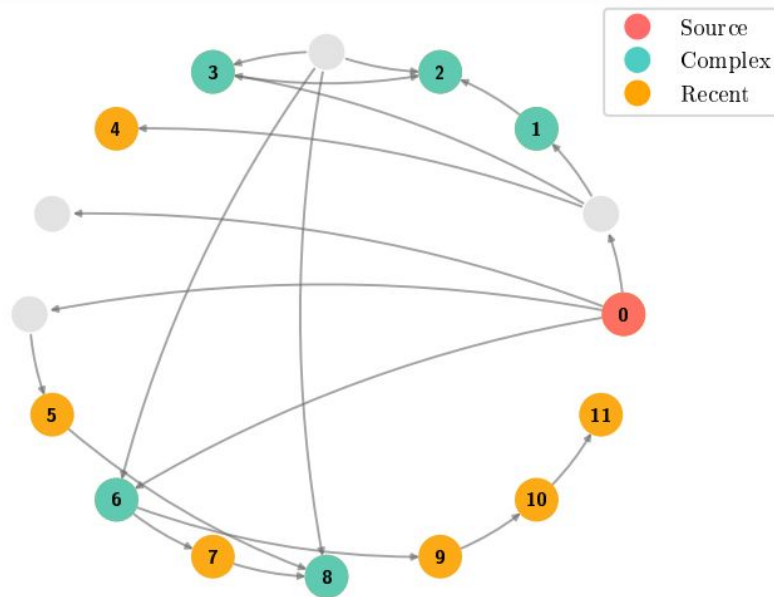
ORIGINAL

- ☐ Rounding Floats
- ☐ Truncating Strings
- ☐ Column Types

- ☐ Limit Rows
- ☐ String Quotes
- ☐ DataFrame Placement

DataFrame Dependency Graph

```
1 from pandas import Series, DataFrame
2 import pandas as pd
3
4 df = pd.read_csv('NYC_Restaurants.csv', dtype=str)
5
6 df_noduplicates = df.drop_duplicates(subset='RESTAURANT')
7 df_notchains = df_noduplicates.groupby("DBA").filter(lambda x: len(x) == 1)
8 boro_notchain_pivot = pd.pivot_table(df_notchains, index = 'BORO', values = '
    RESTAURANT', aggfunc = lambda x: len(x.unique()))
9 boro_restaurant_pivot = pd.pivot_table(df_noduplicates, index = 'BORO', values = '
    RESTAURANT', aggfunc = lambda x: len(x.unique()))
10 boro_notchain_pivot['TOTAL RESTAURANTS'] = boro_restaurant_pivot
11
12 cuis = df_noduplicates['CUISINE DESCRIPTION'].value_counts()
13 mask = (df['VIOLATION CODE']).isnull()
14 no_violations = df[mask]
15 no_violations[['CUISINE DESCRIPTION', 'RESTAURANT']]
16
17 cuisine_no_violations = no_violations['CUISINE DESCRIPTION'].value_counts()
18 mask = (df['VIOLATION CODE']).notnull()
19 violations = df[mask]
20 cuisine_violations = violations['CUISINE DESCRIPTION'].value_counts()
21
22 total_cuisine = pd.concat([cuisine_violations, cuisine_no_violations], axis = 1)
23 violations = pd.crosstab(df['BORO'], df['VIOLATION DESCRIPTION']).query('BORO != ["
    Missing"]')
24
25 vstack = violations.stack()
26 violations2 = vstack.unstack('BORO')
27 mostcommon = DataFrame({'Most Common Complaint':violations2.idxmax(),'Number of
    Complaints':violations2.max()})
```



Abstract Syntax Tree (AST) Parsing



Library Imports

Imported modules
and aliases.

pandas

numpy as np



Function Calls

Method call chains.
DataFrame variables identified
by last return method call.

df_agg

df
.groupby
.agg
.reset_index



DataFrame Graphs

Modified & used DataFrames.
Parent & derived variables
with creation methods.

df_agg

df

df_agg
— (groupby) —
df

Schema Serialization

```
1 DATAFRAMES: orders, cleaned_orders
2 COLUMNS: Country (str), Customer Name (str), Delivery Year (float64), Engine (str),
            Model Series (str), Order Month (str), Order Year (float64), Region (str),
            Delivery Total (float64), Order Total (float64), Unfilled Orders (float)
3
4 DATAFRAMES: df
5 COLUMNS: Order Total (float64)
6
7 DATAFRAMES: region
8 COLUMNS: region (str), year (float64)
```



Column Types



Aggregated Schemas

Prompting

```
1 You are a genius Python data science assistant.
2 Your task is to continue the notebook by answering the user question with the
   provided notebook context.
3 You must output Pandas Python code that will be parsed and executed in a stateful
   Jupyter notebook environment.
4 Think carefully about the DataFrames, columns, methods and use sound logical
   reasoning in your response.
5 You must output your answer in the requested format.
6 If no format is specified, your output should be a DataFrame or Series.
7 I will tip $10000000 if your code is clean and correct.
```

Listing 3.7: Default System Prompt

Prompting: Chain of Thought

- ❑ Detailed Instructions
- ❑ Message Tags
- ❑ Output Format

```
1 You are a genius AI programming data science expert.
2 Your goal is to output clear detailed steps such that a 5-year old can follow the
  steps to write Python code.
3 The current notebook code and markdown cells are provided in between the <
  notebook_context> tags.
4 Notebook variables with column information is provided in the <notebook_variables>
  tags.
5
6 <notebook_context>
7 ...
8 </notebook_context>
9
10 <notebook_variables>
11 ...
12 </notebook_variables>
13
14 My question is as follows:
15 <question>
16 ...
17 </question>
18
19 For each natural language step, think carefully about the following:
20 1. All dataframes, series, column names, functions and variables should be in
   quotation marks.
21 2. Think carefully about column data types. Observe dataframe contents and data
   formats. Consider if parsing is required, such as dates. Ensure null/empty
   values are handled correctly.
22 3. Think carefully about every column in the dataframe and its relationship with the
   user's question.
23 4. Do not make assumptions. You have been given all context required. Ensure your
   steps match the question requirements.
24
25 Write your steps in between <scratchpad> tags.
26 Write your code in between <python> tags.
27 Do not output any other text.
28
29 Take a deep breath and work on this problem step-by-step.
```

Prompting: Chain of Thought

```
1 <scratchpad>
2 To find the average division population for different countries having events
  involving heavy rain, we need to follow these steps:
3 1. Extract the events involving heavy rain from the 'glc' dataframe.
4 2. Filter the events to only include those with a 'landslide_trigger' of 'Heavy Rain
  '.
5 3. Extract the 'country_name' and 'admin_division_population' columns from the
  filtered events.
6 4. Group the data by 'country_name' and calculate the average '
  admin_division_population' for each country.
7 </scratchpad>
8 <python>
9 heavy_rain_events = glc[glc.landslide_trigger == 'Heavy Rain']
10 avg_division_population = heavy_rain_events.groupby('country_name') \
11     ['admin_division_population'].mean()
12 print(avg_division_population)
13 </python>
```


Prompting: Multi-Step

REFLECTION + ERROR MESSAGES

```
1 #####
2 # Initial Prompt #
3 #####
4 [...]
5 Convert Levy and Mileage columns from object to integer.
6 [...]
7
8 #####
9 # Model Response #
10 #####
11
12 <scratchpad>
13 Step 1: Identify the columns that need to be converted from object to integer.
14 In this case, the columns are "Levy" and "Mileage".
15
16 Step 2: Remove any non-numeric characters from the "Levy" column.
17 We can use the "str.replace" method to remove any non-numeric characters.
18
19 Step 3: Convert the "Levy" column to integer.
20 We can use the "pd.to_numeric" function to convert the column to integer.
21
22 Step 4: Remove any non-numeric characters from the "Mileage" column.
23 We can use the "str.replace" method to remove any non-numeric characters.
24
25 Step 5: Convert the "Mileage" column to integer.
26 We can use the "pd.to_numeric" function to convert the column to integer.
27
28 Step 6: Assign the converted columns back to the original dataframe.
```

```
29 </scratchpad>
30 <python>
31 df['Levy'] = df['Levy'].str.replace('[^0-9]', '', regex=True).astype(int)
32 df['Mileage'] = df['Mileage'].str.replace('[^0-9]', '', regex=True).astype(int)
33 </python>
34
35 #####
36 # Follow-Up Prompt with Error Feedback #
37 #####
38 Your solution was executed in the stateful Jupyter notebook environment.
39 An error occurred during execution of the code you submitted:
40 <error>
41 RuntimeError: ValueError: invalid literal for int() with base 10: ''
42 </error>
43 [...]
44
45 #####
46 # Corrected Model Response #
47 #####
48 <analysis>
49 The error occurred because the code attempted to convert an empty string to an
50 integer, which is not a valid operation. This is likely due to the presence of
51 missing or empty values in the "Levy" or "Mileage" columns.
52
53 The question is asking to convert the "Levy" and "Mileage" columns from object type
54 to integer type. However, the original code did not properly handle the
55 possibility of missing or non-numeric values in these columns.
56
57 Upon re-examining the problem, it is clear that the code needs to handle these cases
58 to ensure correctness. One way to do this is to use the "errors" parameter of
59 the "pd.to_numeric" function to specify how to handle errors during the
60 conversion process.
61
62 </analysis>
63 <python>
64 df['Levy'] = df['Levy'].str.replace('[^0-9]', '', regex=True).apply(pd.to_numeric,
65 errors='coerce')
66 df['Mileage'] = df['Mileage'].str.replace('[^0-9]', '', regex=True).apply(pd.
67 to_numeric, errors='coerce')
68 </python>
```

- ❑ Reflection
- ❑ Reflection + Error Messages
- ❑ Reflection + Solution Outputs

Experiment Setup

Model	No. Params	No. Training Tokens	Context Length
PaLM	62B	3.6T	2048
Llama 3 70B	70B	15T	8192
Llama 3 8B	8B	15T	8192

Table 4.1: Model Comparison

```
1 def extract_code_from_response(text):
2     pattern = r'``(?:python\s+)?(.*?)``(?:``|$)|<code>(.*?)</>|<python>(.*?)`(?:</|$)`'
3     match = re.search(pattern, text, re.IGNORECASE | re.DOTALL)
4     if match:
5         code_block = next(group for group in match.groups() if group)
6         return code_block.strip(), text
7     return text, text
```

Listing 4.1: Code Extraction

$$\text{pass@k} := \mathbb{E} \left[1 - \left(\frac{C(n-c, k)}{C(n, k)} \right) \right]$$

Figure 2.10: The pass@k metric

Results: pass@5

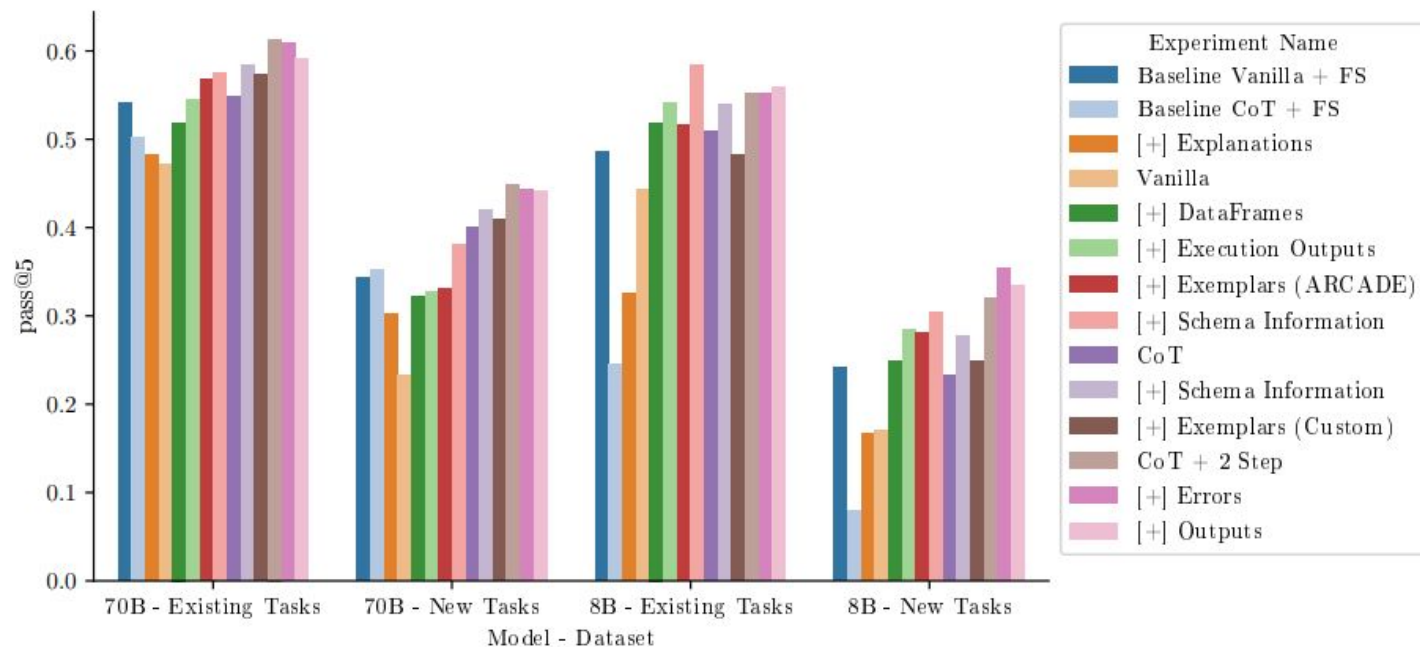


Figure 4.1: pass@5 evaluation results on the ARCADE dataset

Results: Errors

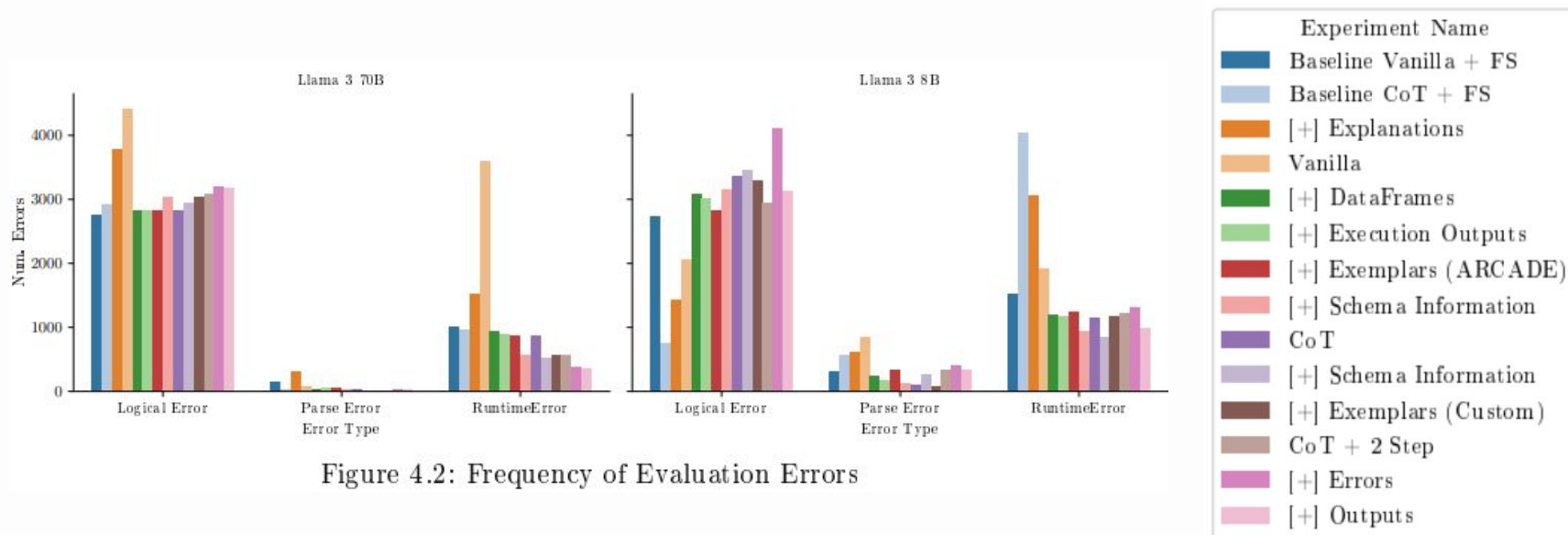


Figure 4.2: Frequency of Evaluation Errors

Ablation Studies

Table 4.4: pass@5 evaluation results comparing DataFrame placement

Experiment Name	Existing Tasks		New Tasks	
	Llama 3 70B	Llama 3 8B	Llama 3 70B	Llama 3 8B
CoT (Inline DataFrames)	58.8	50.8	40.1	23.3
CoT (Appended DataFrames)	56.5	50.6	37.2	24.8

Table 4.5: pass@5 evaluation results on the ARCADE dataset

Experiment Name	Existing Tasks		New Tasks	
	Llama 3 70B	Llama 3 8B	Llama 3 70B	Llama 3 8B
CoT (*)	58.8	50.8	40.1	23.3
+ Schema Information	58.4	54.0	42.1	27.7
CoT + Schema Information	55.0	49.8	39.3	23.6

(*) includes DataFrames and Execution Outputs

Conclusion

Table 5.1: pass@30 baseline comparisons on the New Tasks dataset

Experiment Name	New Tasks (pass@30)			
	PaLM 62B ([†])	PaChiNCo 62B ([‡])	Llama 3 70B	Llama 3 8B
Baseline Vanilla + FS	39.8	48.6	-	-
Baseline CoT + FS	-	52.9	-	-
CoT (*)	-	-	52.5	43.3
CoT (*) [Temp. Sampling]	-	-	54.3	42.7
(*) includes DataFrames, Execution Outputs and Schemas				
([†]) fine-tuned on Python				
([‡]) fine-tuned on Python and Notebooks				

Future Work



Exemplars

Although adding exemplars to prompt has been shown to significantly boost results, our static crafted approach did not give good results.



Agentic Systems

Multi-step prompting naturally extends to agentic systems, unlocking potential for a more interactive, collaborative environment.