




Internship Final Presentation



ABACUS digital
29 JULY 2025

PARINTARA (KAT)
WONGSRISOONTORN
DATA ENGINEER INTERN

Mentor: P'Tae

I N T R O



New York University (NYU)

BS Integrated Design and Media
MS Computer Science



A G E N D A



Project Setup &
Planning



Building the
Solution



Outcomes &
Reflections



O V E R V I E W



Context

Internal metadata is often **incomplete or unclear**, making it difficult for teams to understand unfamiliar tables and queries



Objective

Build scalable, maintainable tools to **automate metadata generation** for BigQuery tables and queries using LLMs, schema information, and internal documentation

OVERVIEW



Tech Stack

LangChain · OpenAI · BigQuery · Confluence API · Python



Scope

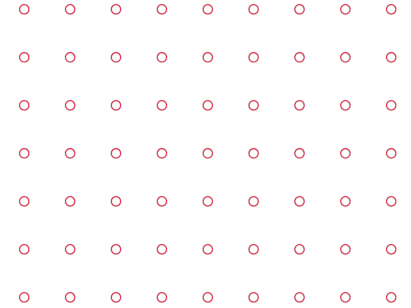
Generate metadata for [data-platform-prod.thunder.contract](#)
→ Column & table descriptions



Output

Structured JSON descriptions for column and table metadata

A G E N D A



Project Setup &
Planning



Building the
Solution

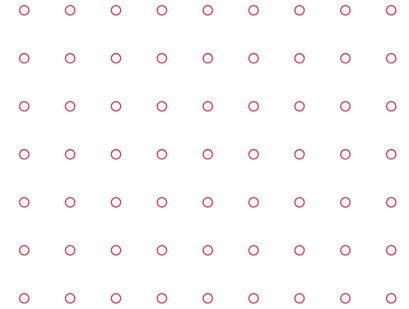


Outcomes &
Reflections



[Column Metadata]

W O R K F L O W



Input: table names



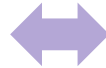
Loops through tables

get_table_schema
→ extracts column names



Loops through columns

LangChain Agent (ReAct)



Tools

Agent selects
and calls tools

Tool: fetch_confluence_chunked

Tool: query_from_bigquery_chunked



Input to LLM

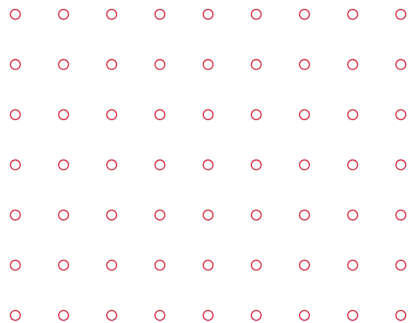
LLM generates
metadata descriptions



StructuredOutputParser
→ formats output as JSON



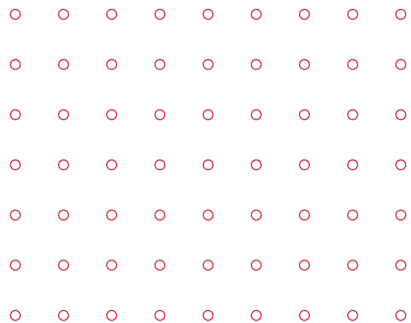
Saves per-column
JSON file



You are a **metadata assistant** that generates useful descriptions for database columns

Your task is to generate a ****{meta_key}**** for the column ****{column_name}**** in the BigQuery table ****{table}****

P
R
O
M
P
T



Follow this process:

1. Use the `query_from_bigquery_chunked`` tool to examine real data from the table

Only proceed to call the `fetch_confluence_chunked`` tool to retrieve internal documentation if the query results are **missing, unclear, or insufficient** to complete your task

P
R
O
M
P
T



P R O M P T

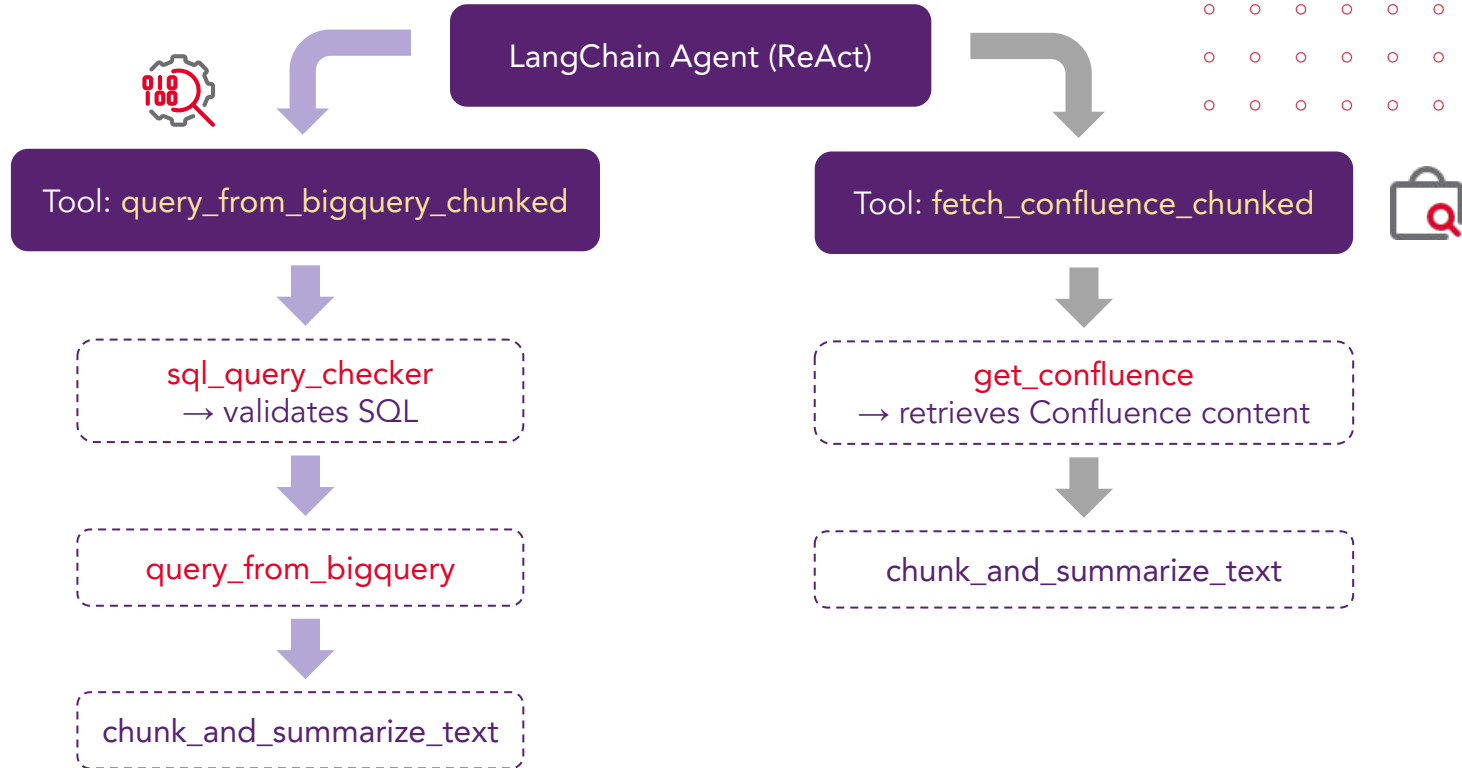
2. Based on the information you have gathered, write a meaningful and informative `**{meta_key}**` for the column

3. Finally, return your result in the `**exact JSON format**` below
`**Important: Output only the JSON object — no additional text, no explanations, no markdown formatting**`

`{format_instructions}`

[Agent Tool Execution]

W O R K F L O W



[LangChain Agent]

D E C I S I O N



```
> Entering new AgentExecutor chain...
```

```
I need to gather data from the BigQuery table `data-platform-prod.thunder.contract` to understand the `contract_id` column. I will use the `query_from_bigquery_chunked` tool from this table.
```

```
Action: query_from_bigquery_chunked
```

```
Action Input: "SELECT contract_id FROM `data-platform-prod.thunder.contract` LIMIT 100"
```

```
The «contract_id» column contains unique numerical identifiers assigned to contracts. These IDs are primarily large integers, typically starting with "100" followed by six digits, though some smaller numeric values also appear. Each contract_id serves as a distinct reference key for individual contracts within the dataset. I have gathered sufficient information about the `contract_id` column from the query results. Now, I will formulate a concise summary and a detailed explanation for the `long_description`.
```

```
```json
{
 "short_description": "The contract_id is a unique identifier for each contract in the dataset, primarily represented as large integers.",
 "long_description": "The contract_id column in the data-platform-prod.thunder.contract table serves as a unique identifier for each contract recorded in the dataset. This identifier is crucial for distinguishing between different contracts and ensuring data integrity. The values in this column are primarily large integers, typically formatted to start with '100' followed by six additional digits, although there are instances of smaller numeric values. Each contract_id is assigned at the time of contract creation and remains constant throughout the contract's lifecycle. This consistency allows for reliable tracking and referencing of contracts across various systems and reports. The uniqueness of the contract_id ensures that no two contracts can share the same identifier, which is essential for maintaining accurate records and facilitating efficient data retrieval and analysis."
}...

Thought: The task is to generate a long description for the column contract_id in the specified BigQuery table. The provided information is already sufficient to create the required JSON output.
Final Answer: {
 "short_description": "The contract_id is a unique identifier for each contract in the dataset, primarily represented as large integers.",
 "long_description": "The contract_id column in the data-platform-prod.thunder.contract table serves as a unique identifier for each contract recorded in the dataset. This identifier is crucial for distinguishing between different contracts and ensuring data integrity. The values in this column are primarily large integers, typically formatted to start with '100' followed by six additional digits, although there are instances of smaller numeric values. Each contract_id is assigned at the time of contract creation and remains constant throughout the contract's lifecycle. This consistency allows for reliable tracking and referencing of contracts across various systems and reports. The uniqueness of the contract_id ensures that no two contracts can share the same identifier, which is essential for maintaining accurate records and facilitating efficient data retrieval and analysis."
}...

> Finished chain.
```

Sample Execution of Agent Decision Making and Tool Selection

[ Table Metadata ]

# W O R K F L O W



Input: table names



Loops through tables

LangChain Agent (ReAct)



get\_column\_descriptions  
→ reads column metadata



get\_table\_schema  
→ extracts column names

Loops through columns



Looks up column  
descriptions



LLM generates  
metadata descriptions

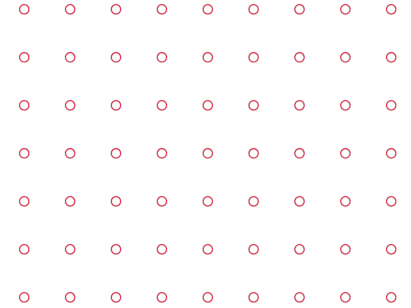


StructuredOutputParser  
→ formats output as JSON



Saves table-level  
JSON file

# A G E N D A



Project Setup &  
Planning



Building the  
Solution

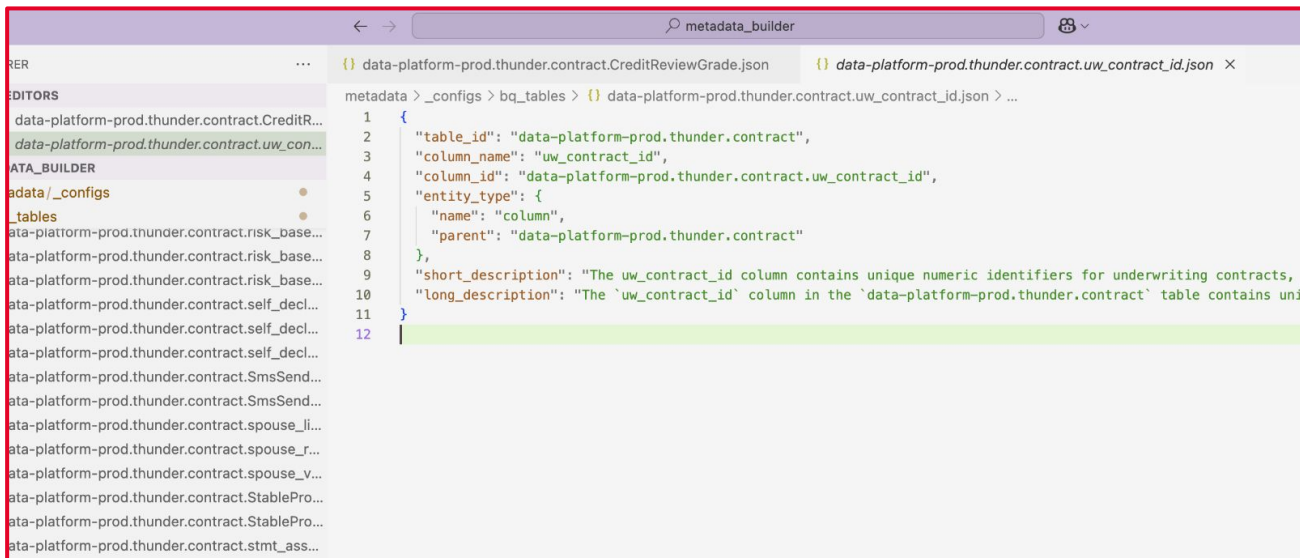


Outcomes &  
Reflections



## [ Column Metadata ]

# O U T P U T



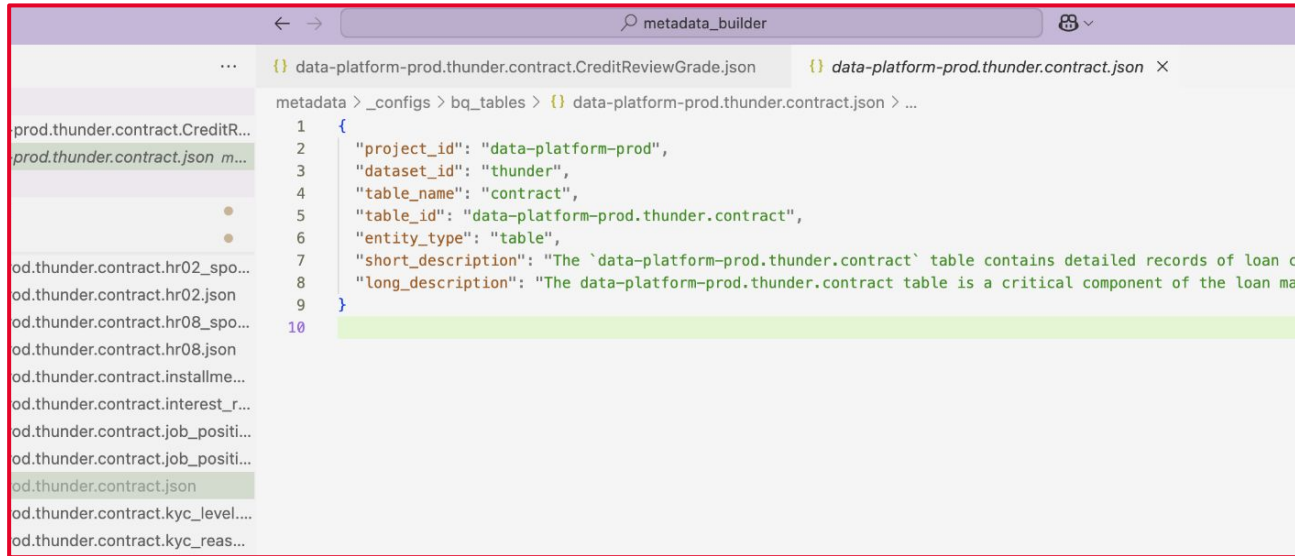
The screenshot shows a web application titled 'metadata\_builder'. The left sidebar contains a list of file paths, including 'data-platform-prod.thunder.contract.CreditReviewGrade.json' and 'data-platform-prod.thunder.contract.uw\_contract\_id.json'. The main area displays the JSON metadata for the 'uw\_contract\_id' column. The JSON structure includes 'table\_id', 'column\_name', 'column\_id', 'entity\_type' (with 'name' and 'parent' fields), 'short\_description', and 'long\_description'.

```
{
 "table_id": "data-platform-prod.thunder.contract",
 "column_name": "uw_contract_id",
 "column_id": "data-platform-prod.thunder.contract.uw_contract_id",
 "entity_type": {
 "name": "column",
 "parent": "data-platform-prod.thunder.contract"
 },
 "short_description": "The uw_contract_id column contains unique numeric identifiers for underwriting contracts,",
 "long_description": "The 'uw_contract_id' column in the 'data-platform-prod.thunder.contract' table contains uni..."
}
```

Sample JSON output for one of 249 columns  
in the data-platform-prod.thunder.contract table

[ Table Metadata ]

O U T P U T



The screenshot shows a web browser window with the address bar displaying 'metadata\_builder'. The page content is divided into a left sidebar and a main panel. The sidebar contains a list of file names, including 'prod.thunder.contract.CreditReviewGrade.json', 'prod.thunder.contract.json m...', 'od.thunder.contract.hr02\_spo...', 'od.thunder.contract.hr02.json', 'od.thunder.contract.hr08\_spo...', 'od.thunder.contract.hr08.json', 'od.thunder.contract.installme...', 'od.thunder.contract.interest\_r...', 'od.thunder.contract.job\_positi...', 'od.thunder.contract.job\_positi...', 'od.thunder.contract.json', 'od.thunder.contract.kyc\_level...', and 'od.thunder.contract.kyc\_reas...'. The main panel displays a JSON object for the selected file 'data-platform-prod.thunder.contract.json'. The JSON object is as follows:

```
1 {
2 "project_id": "data-platform-prod",
3 "dataset_id": "thunder",
4 "table_name": "contract",
5 "table_id": "data-platform-prod.thunder.contract",
6 "entity_type": "table",
7 "short_description": "The `data-platform-prod.thunder.contract` table contains detailed records of loan co",
8 "long_description": "The data-platform-prod.thunder.contract table is a critical component of the loan mar",
9 }
10
```

Table-level JSON output for the  
`data-platform-prod.thunder.contract` table



# CHALLENGES



## Token Limit

**Issue:** Large Confluence pages and queries exceeded model limits

**Solution:** Implemented `chunk_and_summarize_text` to reduce content size before sending to LLM



## SQL Format Error

**Issue:** Agent sometimes generated invalid queries

**Solution:** Wrapped SQL tool in `try-except` and looped until a valid response was returned

# CHALLENGES



## Parsing Failure

**Issue:** LLM output did not always match the expected JSON format

**Solution:** Refined the prompt and used `StructuredOutputParser` with fallback handling



## Rate Limiting

**Issue:** Too many requests triggered API limits

**Solution:** Paused after every 50 columns for 60 seconds to avoid exceeding limits

# T A K E A W A Y S



Gained exposure to the **data engineering** industry and practical tools used in the field



Discovered how **data, AI, and engineering** intersect in real-world and future innovation



Learned about fintech and how technology supports **money lending and digital finance**



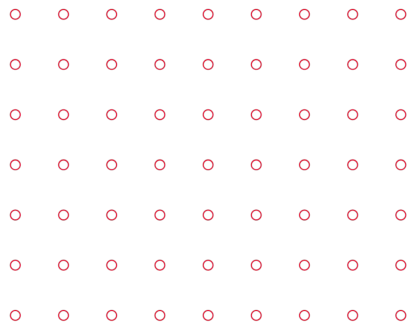
Enjoyed working on **hands-on projects** that go beyond academic concepts



Inspired by the smart, collaborative people and the **open culture of a modern company**



# Internal Tools Shape Company Culture ”



# Special Thanks to the Data Engineering (DE) Team



NATTH  
BEJRABURNIN  
Chief Data & AI  
Officer



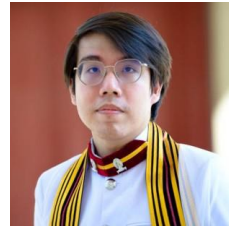
VARUNYA  
THAVORNUN  
Head of Data  
Engineering



WISAROOT  
LERTTHAWEEDECH  
Senior Data Engineer



THAYANEE  
RUENNARK  
Data Engineer

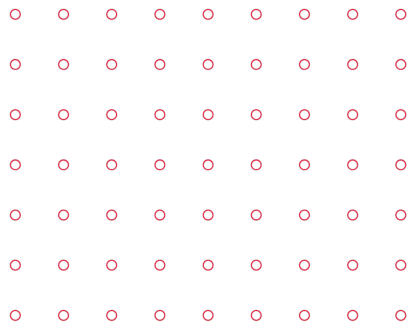


WARIS  
PUNNAHITANOND  
Data Engineer



PONTHAI  
KLINSAMPAS  
Senior Data Engineer

# Q & A



THANK YOU

