



GDM 5 Final Project

Exploring the innovative use of generative AI to extract crucial information from geological documents, enhancing research and data accessibility.

Aicha EL BOU
Juan SIVORI
Samuel NAUTIN
Stephania ZAMBRANO

Tutors:
Henri BLONDELLE
David RHENALS

Evaluating LLM Accuracy in Document Processing

THIS PRESENTATION FOCUSES ON ASSESSING LLM RESPONSES IN EXTRACTING INFORMATION FROM GEOLOGICAL DOCUMENTS, EMPHASIZING ACCURACY AND MODEL COMPARISON.

COMPREHENSIVE PROJECT OVERVIEW

Exploring the Key Aspects of the Project

01. Introduction

02. Phase 1

03. Phase 2

04. Results and conclusions

08. Data Management

09. Project Management

10. Q&A

Our Timeline



Phase 2

DATA PREPROCESSING AND ENRICHMENT FOR LLMS USING RAG

- Definition of action plan for second phase, standardization of questions
- Definition of the documents manipulation and final configuration
- Evaluation of 5 documents

Defense

DEFENSE

- PROJECT SUSTENTATION



Phase 1

GETTING TO KNOW AGILE DD

- Project socialization
- Communicate and gather Geologist insight
- Platform recognition (Agile DD)
- Creation of initial prompt bank
- Benchmark among different LLMs

Phase 3

PROJECT DEVELOPING

- Data sets standardization
- Evaluation of 15 more documents
- Statistics
- Results and conclusions



PHASE 1: GETTING TO KNOW AGILE DD

- Initial Objective

Generate a benchmark to assess the accuracy of LLM responses in extracting information from technical documents.

- Expanded Scope

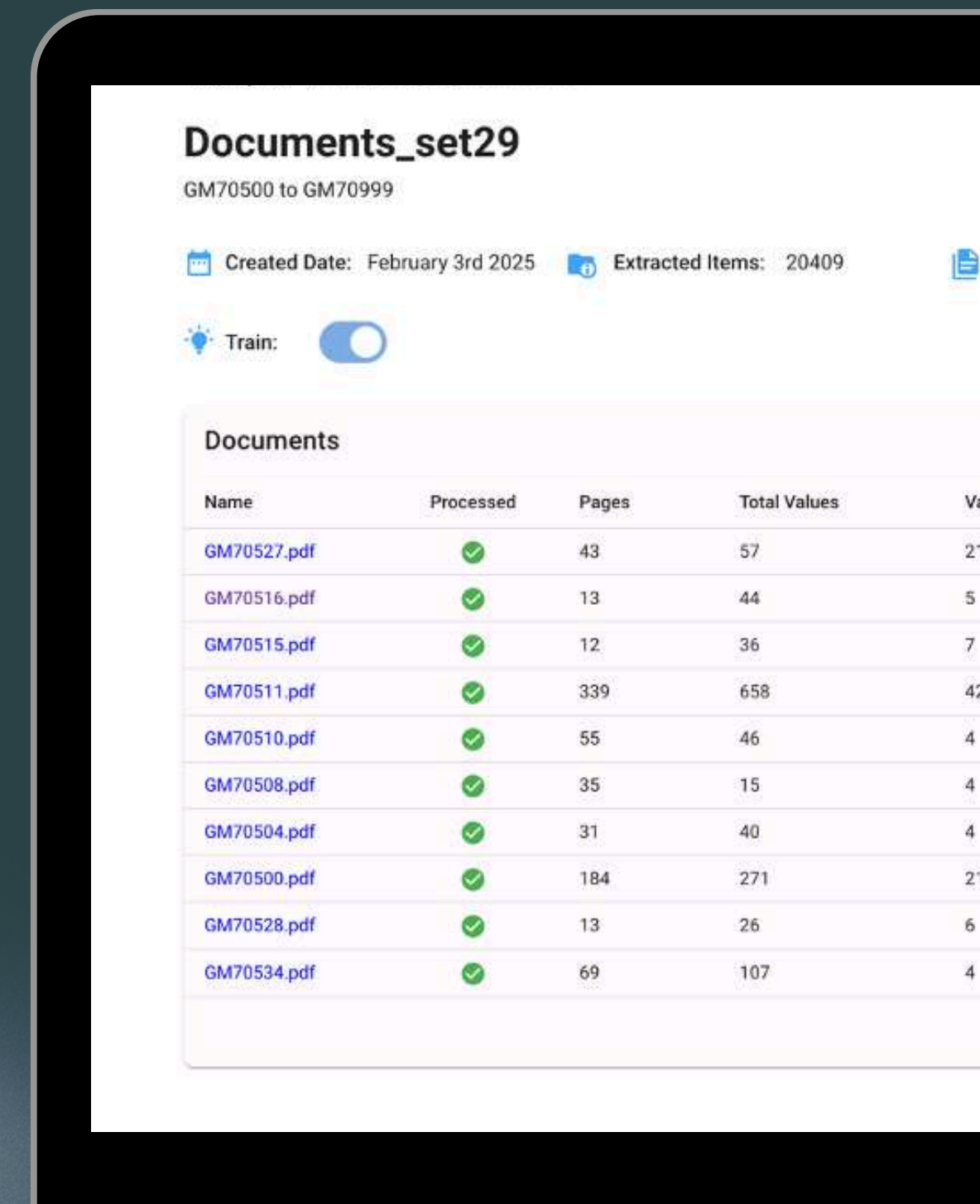
Evaluate the accuracy of Agile DD responses and compare them with other LLMs for performance insights.

- Importance of Evaluation

Identifying strengths and weaknesses in data comprehension aids in improving model selection for specific use cases.

- Relative Performance

Evaluating Agile DD against other LLMs showcases its accuracy and limitations in real-world applications.



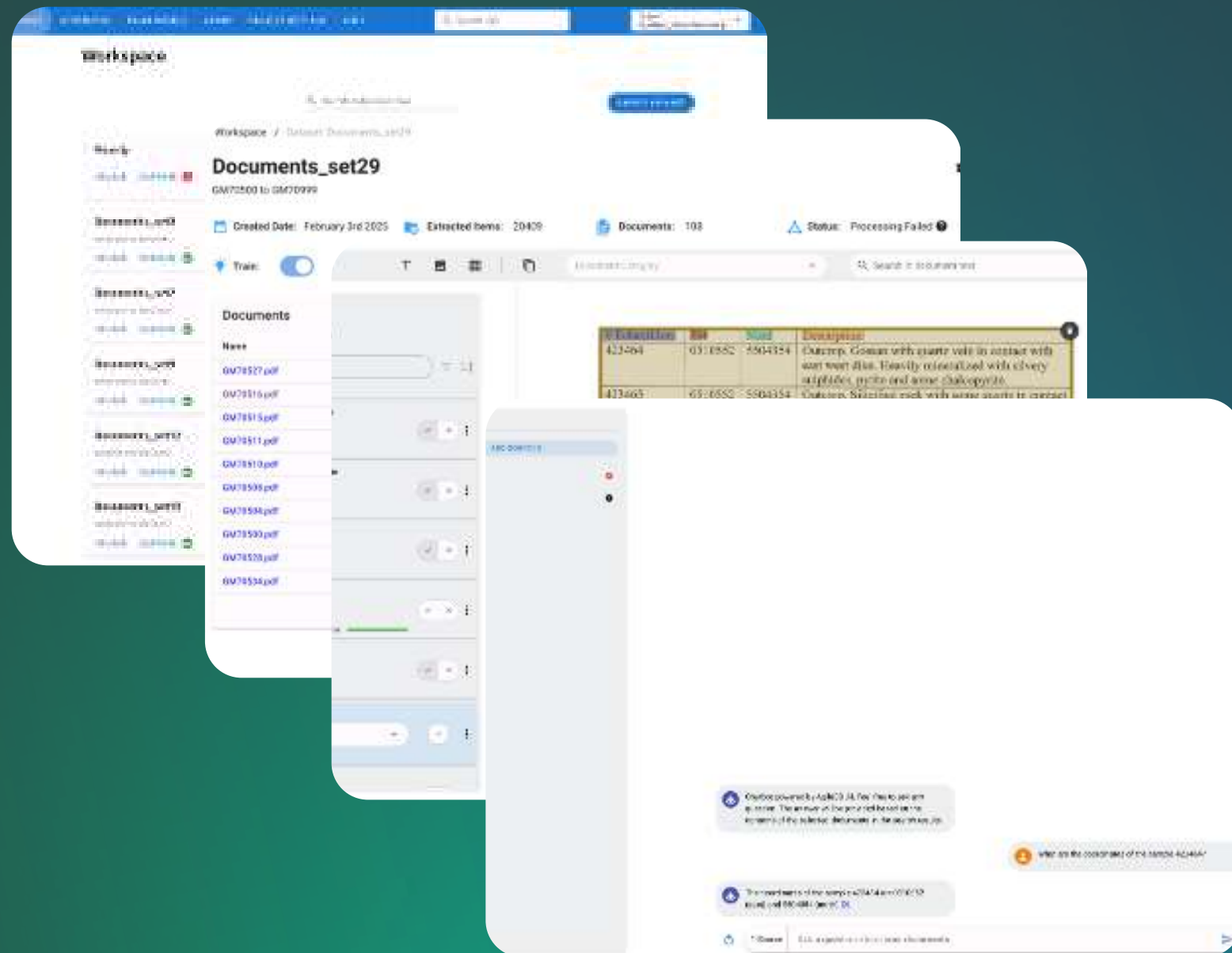
Documents_set29
GM70500 to GM70999

Created Date: February 3rd 2025 Extracted Items: 20409

Train: ☒

Documents				
Name	Processed	Pages	Total Values	Value
GM70527.pdf	✓	43	57	2
GM70516.pdf	✓	13	44	5
GM70515.pdf	✓	12	36	7
GM70511.pdf	✓	339	658	4
GM70510.pdf	✓	55	46	4
GM70508.pdf	✓	35	15	4
GM70504.pdf	✓	31	40	4
GM70500.pdf	✓	184	271	2
GM70528.pdf	✓	13	26	6
GM70534.pdf	✓	69	107	4

Agile DD



Services



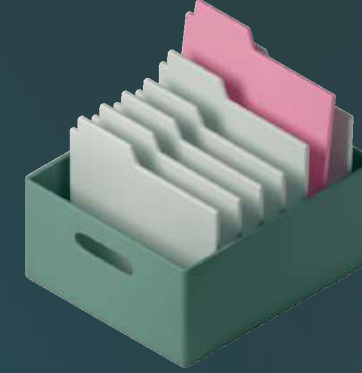
- AgileDD is an AI-powered platform that extracts and structures data from complex technical documents. It combines automation with human expertise to enhance data processing for industries like energy, mining, and defense.

Documents



- Geological Reports
- Assay Certificates (Geochemical Data)
- Sample tables (Mineral Composition)

Prompt Bank



43
Initial
Prompts

36
Refined
Prompts

6
Categories

4
Final
Prompts

✓
What are the coordinates of the sample #X

✓
Focusing on the Json/XML part, Which sample has the highest concentration of Zinc (Zn) ? excluding #X

✗
What are the geochemical anomalies, and do they indicate potential mineralization?

✗
What exploration recommendations were made based on the findings?

✓
Which sample has the highest concentration of Element (El) ? excluding #X

✗
Are there historical work results, and what were the assay values?

- Final Objective

Extract text from geological documents using generative AI with a focus on tables.

- Expanded Scope

Incorporate table recognition and structured data extraction to enhance AI responses.

- Importance of Table Data

Tables often contain critical information that AI models find challenging to interpret.

- Benefits of Structured Formats

Utilizing structured formats like JSON/XML can significantly improve accuracy in processing technical documents.

PHASE 2: DATA PREPROCESSING AND ENRICHMENT FOR LLMS USING RAG

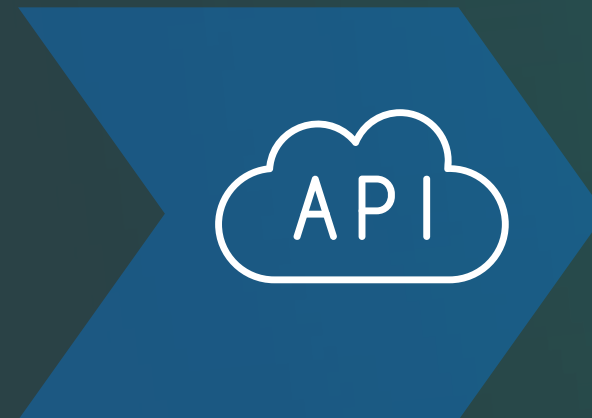
Data extraction and processing workflow

OCR Text Extraction



Utilizing OCR technology to extract text from PDF documents efficiently.

Table Identification



Recognizing and isolating tables from the extracted text for further processing.

Data Conversion



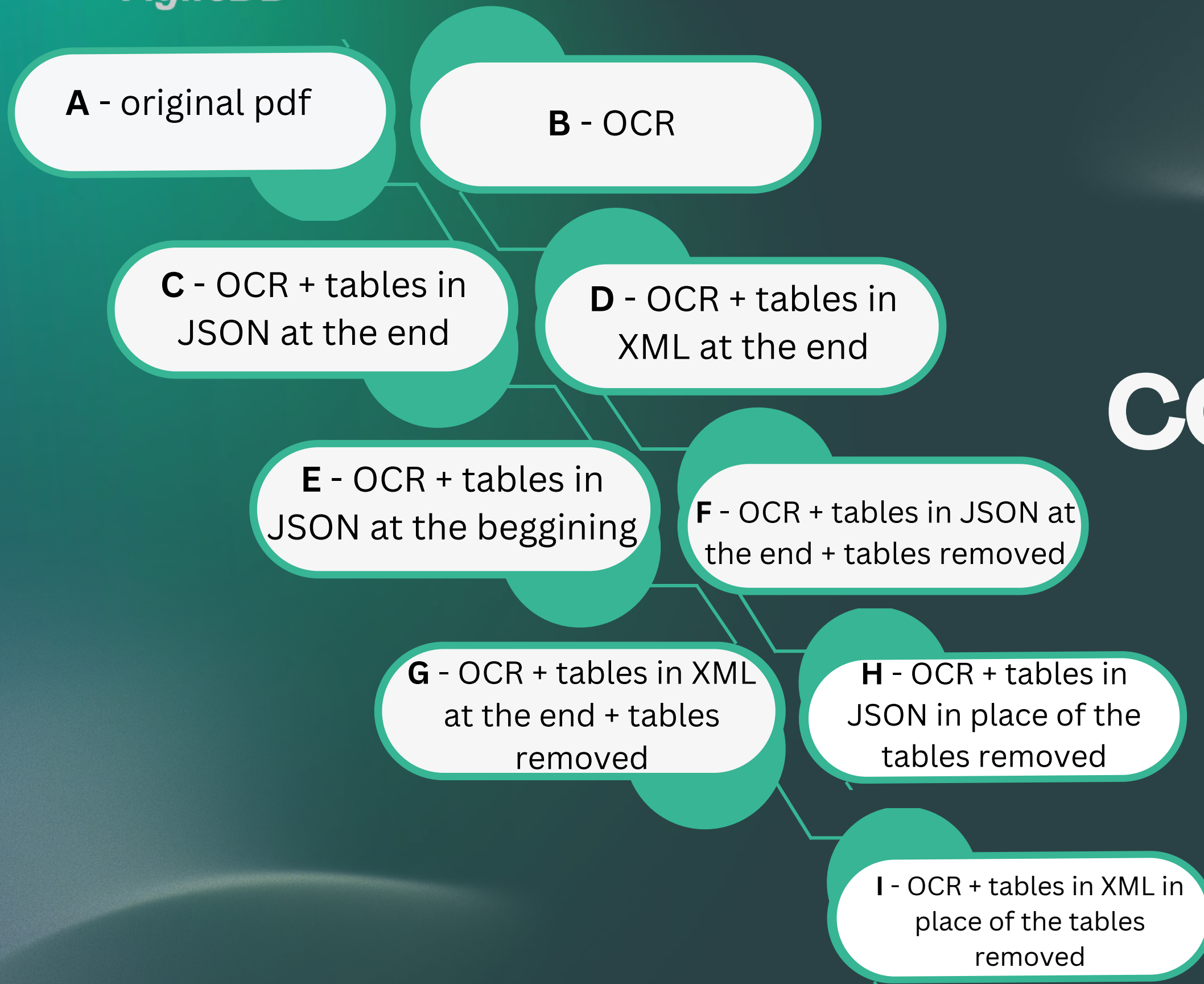
Transforming extracted tables into structured formats like JSON or XML.

Data Combination



Integrating OCR text with structured data formats in various configurations.

DOCUMENT CONFIGURATION



Q: table

XXXXXX, with a little more than 2,000 inhabitants. Room and board, general services and a workforce can be obtained there. The Quebec Ministry of Forest, Wildlife and Parks has an office in Label-sur-Quévillon. 5.0) HISTORY The historical work done on the property is summarized in Tables 2 and 3, below.

Table 2: Work by Mining Companies

GM # Yea Company Work Results
16333 1965 Sullico Mines Airborne Mag and EM Several weak EM surveys covering anomalies located in the approximately 75% of southern part of the Urban Thunder Thunder claims. claims.
38840 1979 Shell Canada Reconnaissance Airborne survey geological survey recommended.
44513 1987 Onyx Line cutting, Mag and Several VLF and Mag Resources, VLF-EM surveys anomalies located south of Sullivan Mines immediately south of the Urban Thunder claims. Inc. the current claims.
44514 1987 Sullivan Mines Airborne EM (Input) and Located immediately south Inc. Mag data interpretation of the claims.
45086 1987 Onyx Geological survey, Hole NO-86-22 cut basalt Resources stripping and sampling and felsic intrusive, no and 2 DDH totalling anomalous results
456.3 m, drilled on the obtained. Hole NO-86-24 southern edge of the cut basalt and gabbro property and to the units. Samples were taken south, but no assay results were provided.

11
Solumines
Label-sur-
Lac
Quévillon
Quetion
R1915
-
Great Thunder Urban Thunder

Raw

Q: table

XXXXXX, with a little more than 2,000 inhabitants. Room and board, general services and a workforce can be obtained there. The Quebec Ministry of Forest, Wildlife and Parks has an office in Label-sur-Quévillon. 5.0) HISTORY The historical work done on the property is summarized in Tables 2 and 3, below.

Table 2: Work by Mining Companies

{ " GM #":16333, " Year":1965, " Company":" Sullico Mines", " Work":" Airborne Mag and EM surveys covering approximately 75% of the Urban Thunder claims.", " Results":" Several weak EM anomalies located in the southern part of the Urban Thunder claims." }
{ " GM #":38840, " Year":1979, " Company":" Shell Canada", " Work":" Reconnaissance geological survey", " Results":" Airborne survey recommended." }
{ " GM #":44513, " Year":1987, " Company":" Onyx Resources, Sullivan Mines Inc.", " Work":" Line cutting, Mag and VLF-EM surveys immediately south of the current claims.", " Results":" Several VLF and Mag anomalies located south of the Urban Thunder claims." }
{ " GM #":44514, " Year":1987, " Company":" Sullivan Mines Inc.", " Work":" Airborne EM (Input) and Mag data interpretation", " Results":" Located immediately south of the claims." }
{ " GM #":45086, " Year":1987, " Company":" Onyx Resources", " Work":" Geological survey, stripping and sampling and 2 DDH totalling 456.3 m, drilled on the southern edge of the property and to the south.", " Results":" Hole NO-86-22 cut basalt and felsic intrusive, no anomalous results obtained. Hole NO-86-24 cut basalt and gabbro units. Samples were taken but no assay results were provided." }

11
Solumines
Label-sur-
Lac
Quévillon
Quetion
R1915
-
Great Thunder Urban Thunder
Scale Gold Corp.
0 10 km

JSON

Q: table

XXXXXX, with a little more than 2,000 inhabitants. Room and board, general services and a workforce can be obtained there. The Quebec Ministry of Forest, Wildlife and Parks has an office in Label-sur-Quévillon. 5.0) HISTORY The historical work done on the property is summarized in Tables 2 and 3, below.

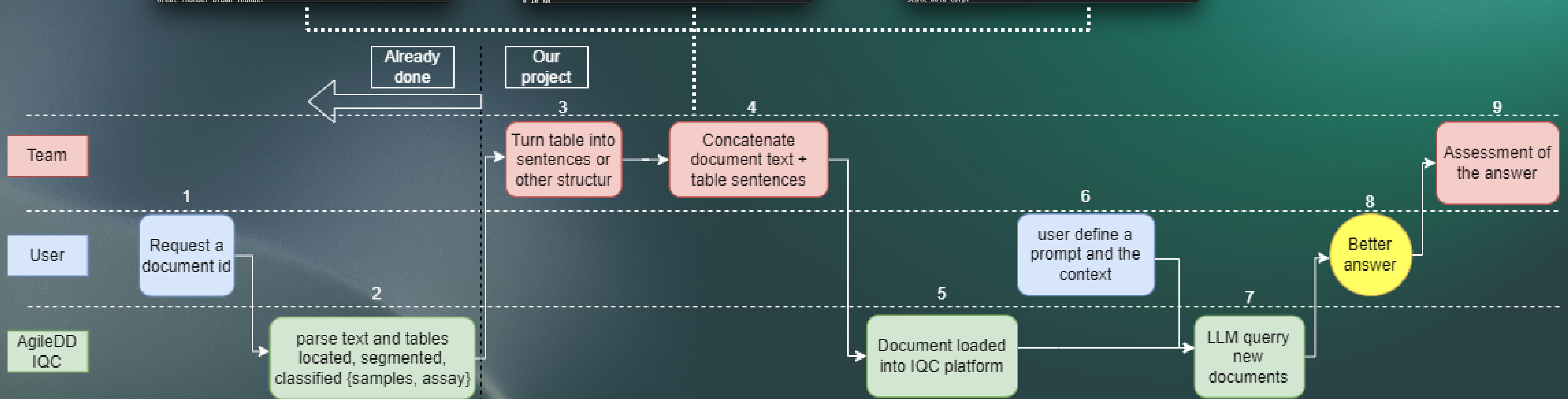
Table 2: Work by Mining Companies

<root><item>< GM #>16333</ GM #>< Year>1965</ Year>< Company> Sullico Mines</ Company>< Work> Airborne Mag and EM surveys covering approximately 75% of the Urban Thunder claims.</ Work>< Results> Several weak EM anomalies located in the southern part of the Urban Thunder claims.</ Results></item><item>< GM #>38840</ GM #>< Year>1979</ Year>< Company> Shell Canada</ Company>< Work> Reconnaissance geological survey</ Work>< Results> Airborne survey recommended.</ Results></item><item>< GM #>44513</ GM #>< Year>1987</ Year>< Company> Onyx Resources, Sullivan Mines Inc.</ Company>< Work> Line cutting, Mag and VLF-EM surveys immediately south of the current claims.</ Work>< Results> Several VLF and Mag anomalies located south of the Urban Thunder claims.</ Results></item><item>< GM #>44514</ GM #>< Year>1987</ Year>< Company> Sullivan Mines Inc.</ Company>< Work> Airborne EM (Input) and Mag data interpretation</ Work>< Results> Located immediately south of the claims.</ Results></item><item>< GM #>45086</ GM #>< Year>1987</ Year>< Company> Onyx Resources</ Company>< Work> Geological survey, stripping and sampling and 2 DDH totalling 456.3 m, drilled on the southern edge of the property and to the south.</ Work>< Results> Hole NO-86-22 cut basalt and felsic intrusive, no anomalous results obtained. Hole NO-86-24 cut basalt and gabbro units. Samples were taken but no assay results were provided.</ Results></item></root>

11
Solumines
Label-sur-
Lac
Quévillon
Quetion
R1915
-
Great Thunder Urban Thunder
Scale Gold Corp.

XML

Data extraction and processing workflow



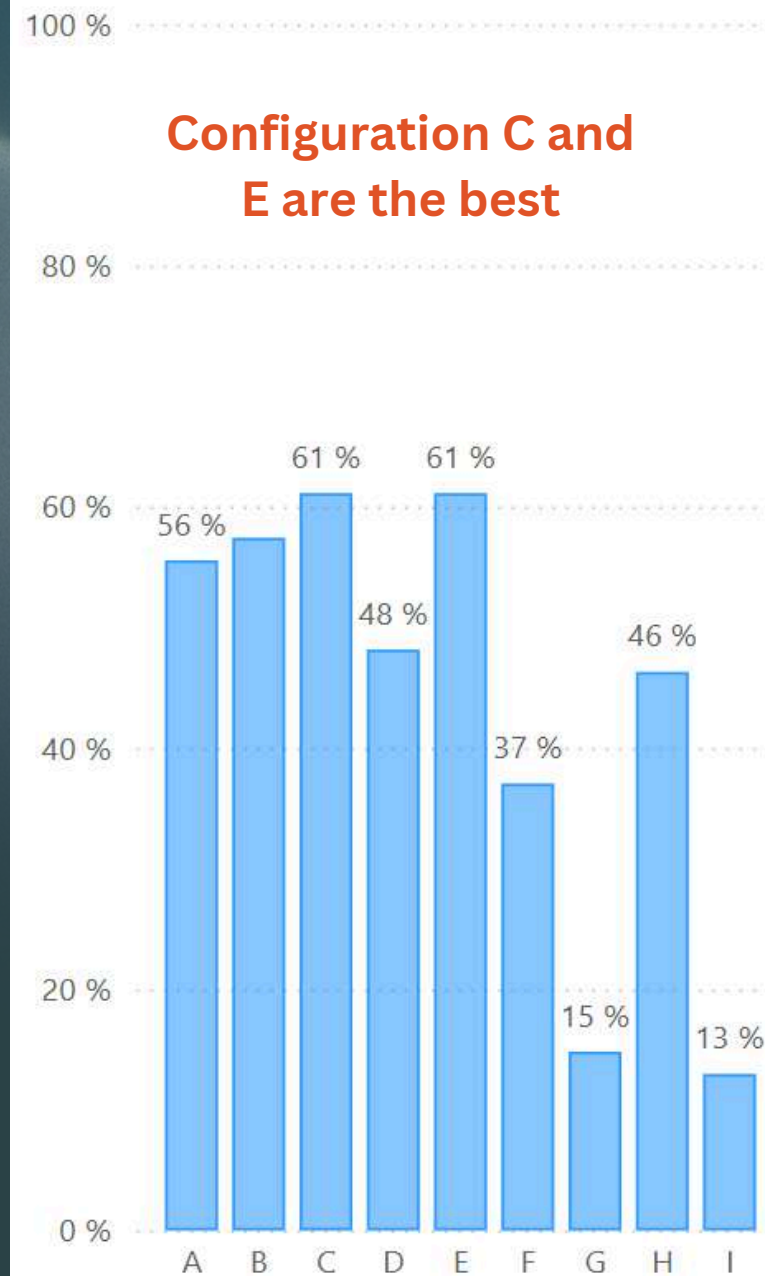
QUERYING & EVALUATION OF AI RESPONSES

Testing prompt strategies for AI response accuracy

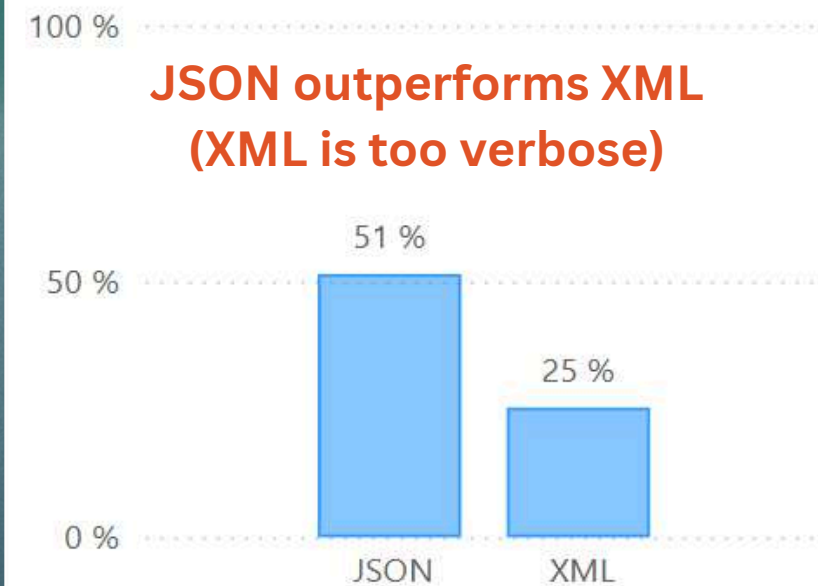
Item	prompt standard	prompt customized
1	What are the coordinates of the sample #X	What are the coordinates of the sample CV17-005
2	Which sample has the highest concentration of Element (El) ? excluding #X	Which sample has the highest concentration of Zinc (Zn) ? excluding CV17-005
3	Which sample has the highest concentration of Element (El), and what is its description .	Which sample has the highest concentration of Zinc (Zn), and what is its description.
4	Focusing on the JSON/XML part. Which sample has the highest concentration of Element (El) ? excluding #X	Focusing on the JSON/XML part. Which sample has the highest concentration of Zinc (Zn) ? excluding CV17-005

RESULTS AND CONCLUSIONS

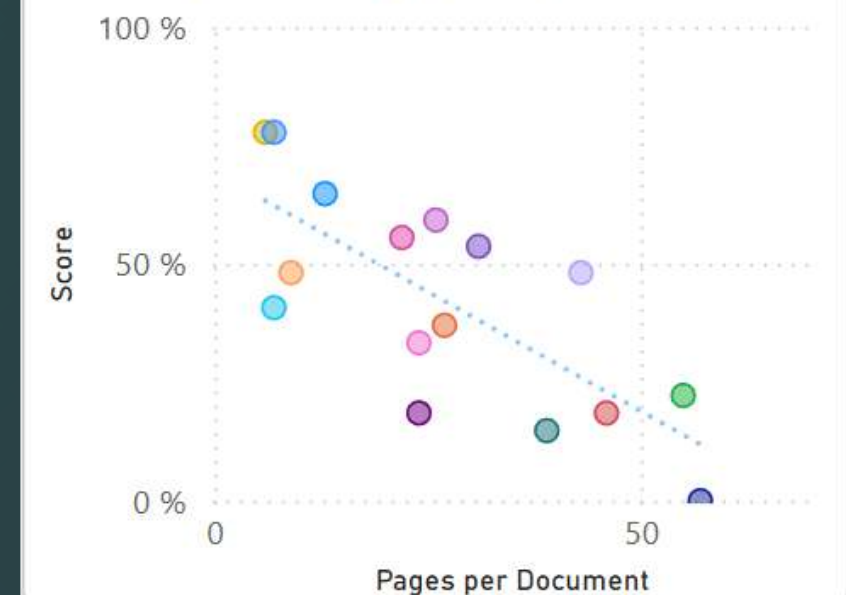
Score per Configuration



Score per Type of Table



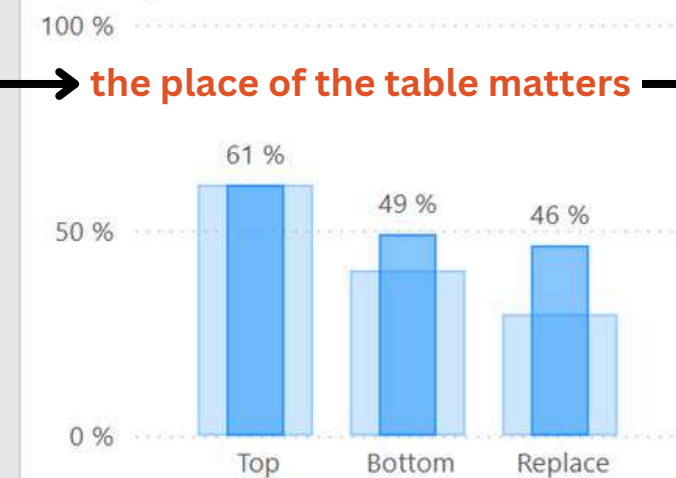
Correlation - Pages and Score



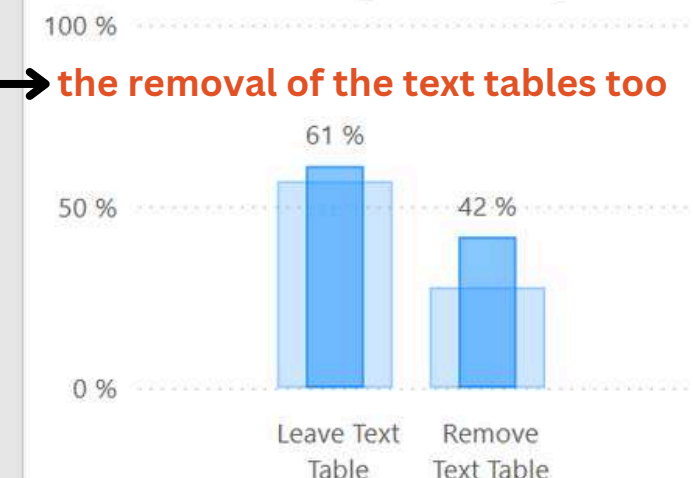
Score per Type of Table



Score per Location of Tables



Score for Removing or Leaving Text ...



INSIGHTS

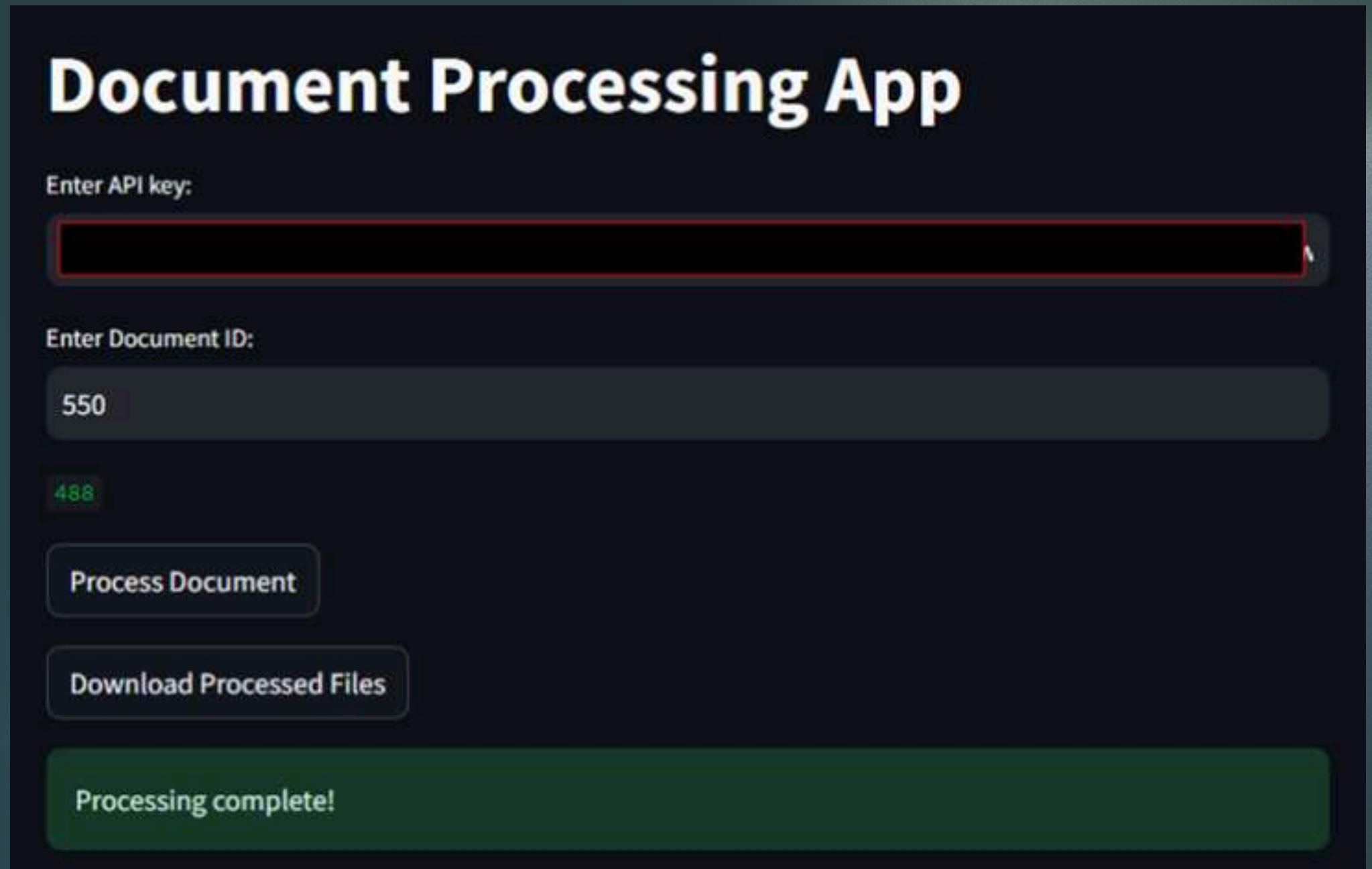
MOST COMMON ERRORS	SOME UNCOMMON ERRORS
“Information not found”	The LLM read the text backward
Error for the application reading inside tables in pdf format and txt ocr format	LLM failed extracting good samples from txt/Json/XML for the highest concentration value
OCR misplacing values inside tables from unstructured data	The LLM described the correct sample without giving its ID
LLM confuses the description from one sample to another	
Provides correct answers but incomplete description	
Errors reading the Json/XML due to the application splitting formats	
If there is a "0" in front of a coordinate number the LLM takes it out	

WORKFLOW AUTOMATION

Custom and time saving web application to increase the speed of the process.
Accessible for everyone in the team.

Characteristics:

- Selection of one GM document.
- Used Agile DD's APIs:
 - Extraction of text from the PDF.
 - Extraction of tables from the PDF.
- Conversion from csv tables to JSON and XML.
- Compression of the text and tables files into a zip file.
- Download of the zip file containing all the inputs for the creation of configuration files.



Document Processing App

Enter API key:

Enter Document ID:

488

Process Document

Download Processed Files

Processing complete!

STANDARDIZATION

Template for data entry:

We created a template in Excel so each member was able to test the GM documents and load the results in a consistent way.

Also, reference data was created to limit the different ways of entering the results from the tests.

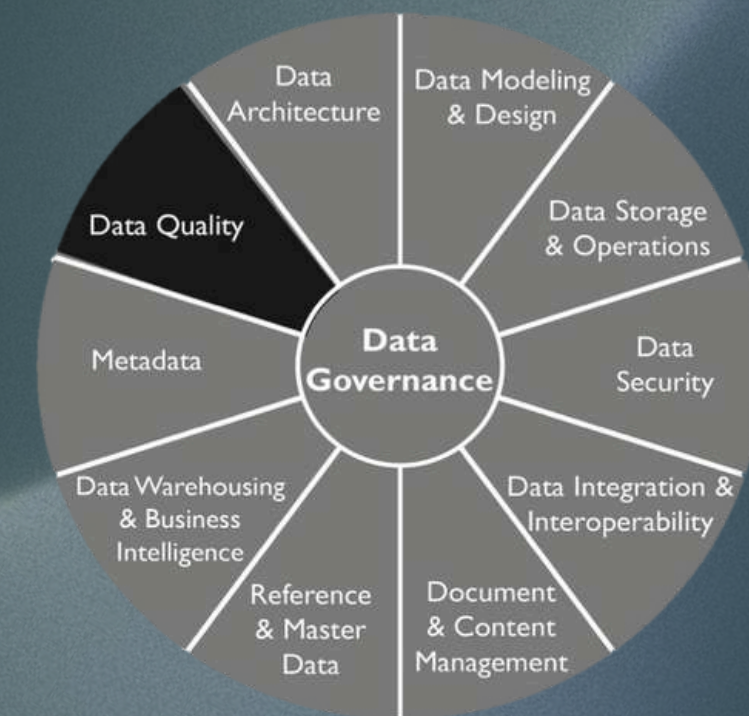
This improved the quality of the gathered data in several dimensions:

- Consistency (same template for all)
- Completeness (the template is like a form, so it is difficult to forget values)
- Validity (reference lists)

Standardized File Naming :

Structured naming conventions applied systematically to track document versions and configurations clearly:

- GMXXXXX_A_ocr_pdf
- GMXXXXX_B_extracted_text
- GMXXXXX_xml or GMXXXXX_json
- GMXXXXX_D_extracted_text_plus_xml
- GMXXXXX_C_extracted_text_plus_json



DATA MODEL

Data Dictionary

Field	Type	Description	Constraints
author_id	int	Unique identifier for the author	Primary Key
author_name	str	Name of the author	Not Null

2. Document

Field	Type	Description	Constraints
document_id	int	Unique identifier for each document	Primary Key
document_name	str	Name or title of the document	Not Null
number_of_pages	int	Total number of pages in the document	Not Null, >= 1
number_of_tables	int	Number of tables present in the document	Not Null, >= 0

3. Prompt

Field	Type	Description	Constraints
prompt_id	int	Unique identifier for the prompt	Primary Key
prompt_std_msg	str	Standardized message used in the prompt	Not Null

4. Configuration

Field	Type	Description	Constraints
configuration_id	str	Unique identifier for a configuration setup	Primary Key
contains_text_tables	bool	Indicates if text-based tables are included	Not Null
contains_json_tables	bool	Indicates if JSON-based tables are included	Not Null
contains_xml_tables	bool	Indicates if XML-based tables are included	Not Null
location_of_added_tables	str	Where additional tables are added to the extracted text	Not Null

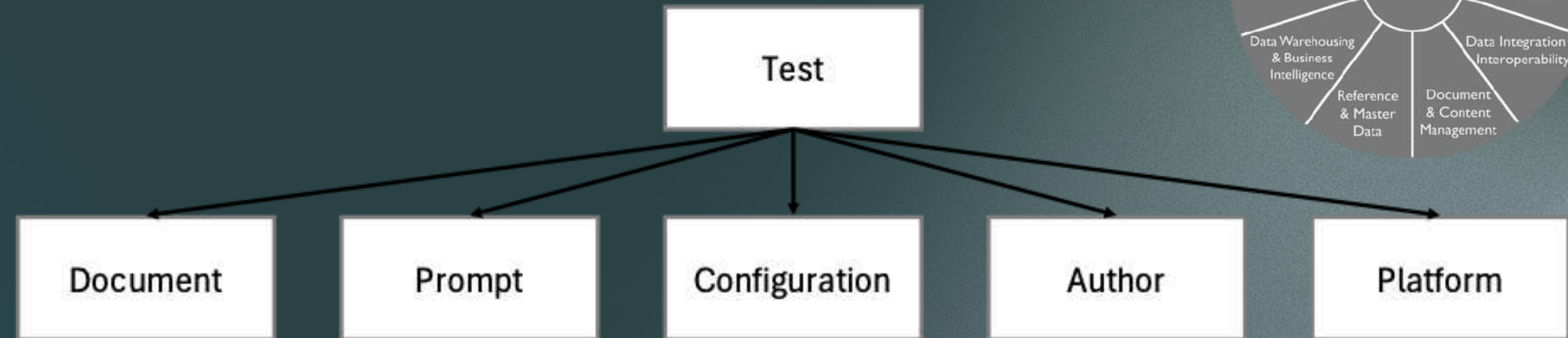
5. Test

Field	Type	Description	Constraints
test_id	int	Unique identifier for a test instance	Primary Key
platform_id	int	Platform used for the test	Foreign Key (Platform)
document_id	int	Document being evaluated	Foreign Key (Document)
author_id	int	Author conducting the test	Foreign Key (Author)
configuration_id	str	Configuration used for the test	Foreign Key (Configuration)
prompt_id	int	Standard prompt used in the test	Foreign Key (Prompt)
prompt_custom_msg	str	Customized version of the prompt	Not Null
answer	str	Extracted value from the document	Not Null
pertinence	bool	Indicates if the answer is relevant	Not Null
correct_coordinates	bool	Whether extracted coordinates were correct	Not Null
correct_sample	bool	Whether the correct sample was identified	Not Null
correct_concentration	bool	Whether the concentration value is correct	Not Null
correct_description	bool	Whether the extracted description is correct	Not Null
remarks	str	Additional comments on test results	Optional

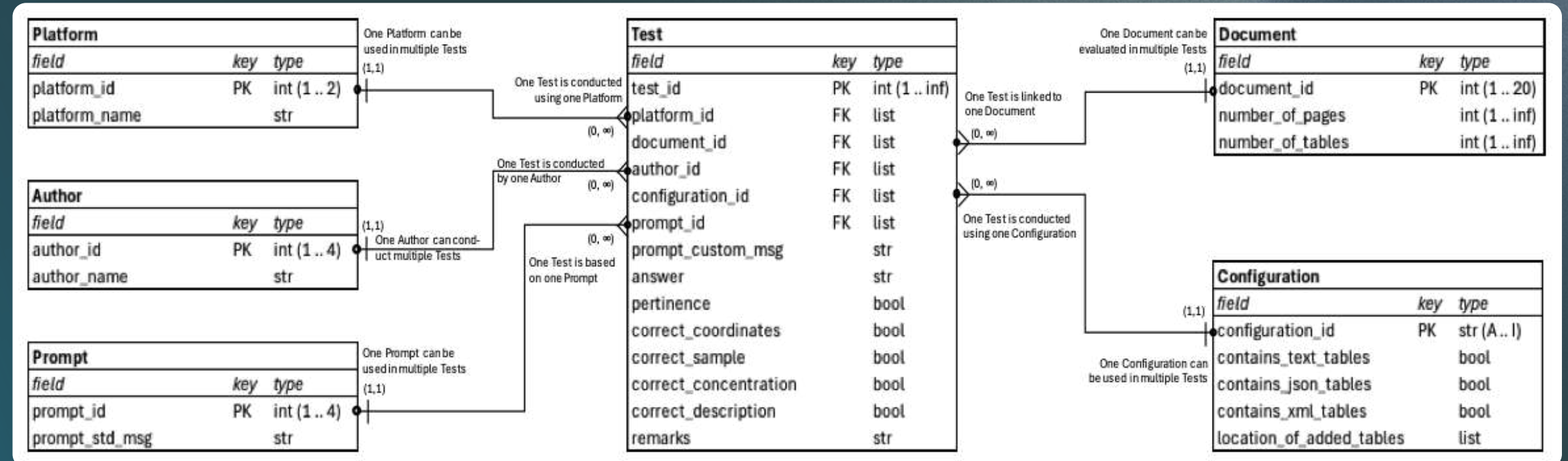
6. Platform

Field	Type	Description	Constraints
platform_id	int	Unique identifier for the platform	Primary Key
platform_name	str	Standardized message used in the prompt	Not Null

Conceptual Data Model

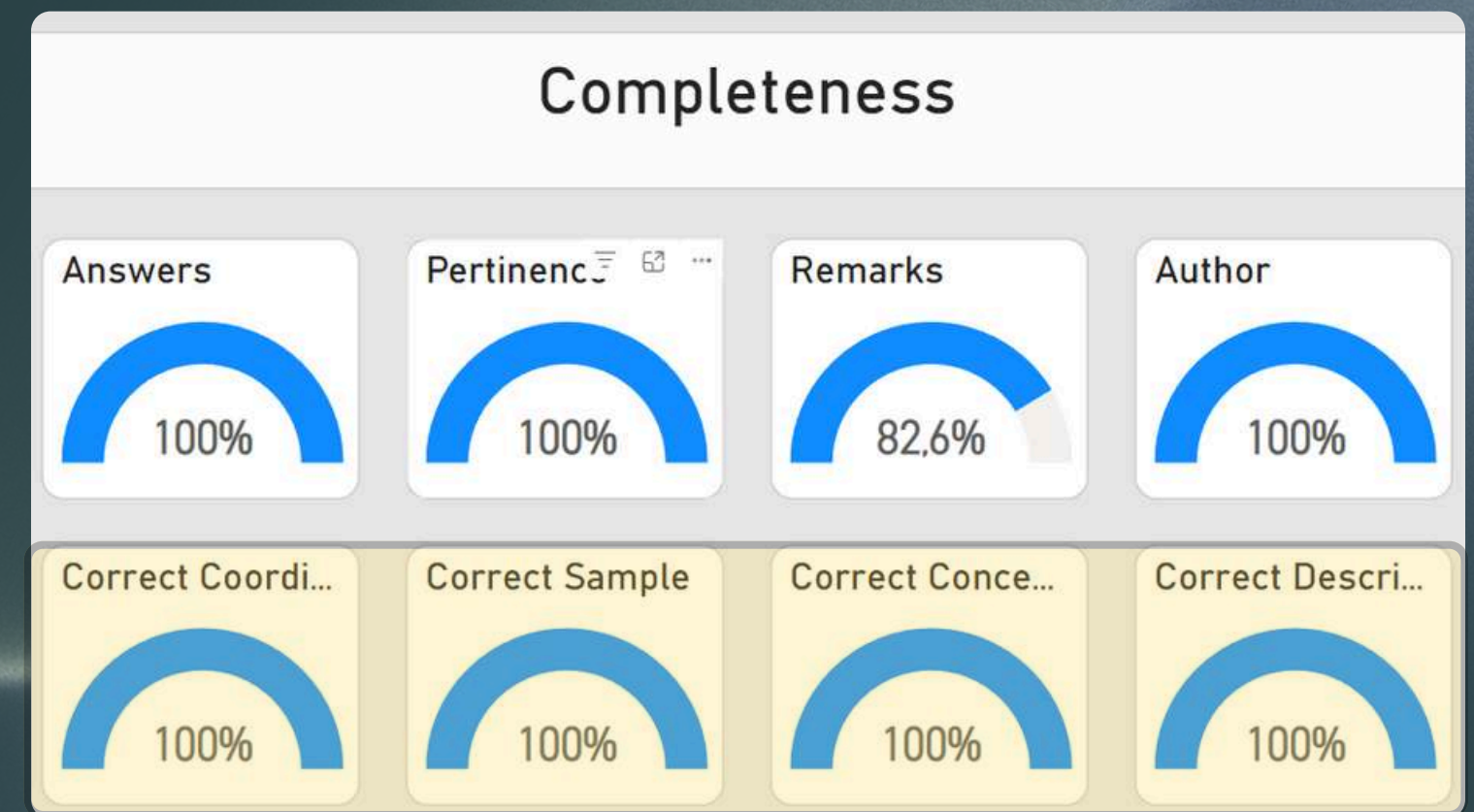
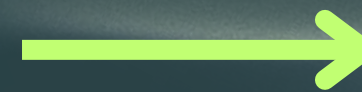
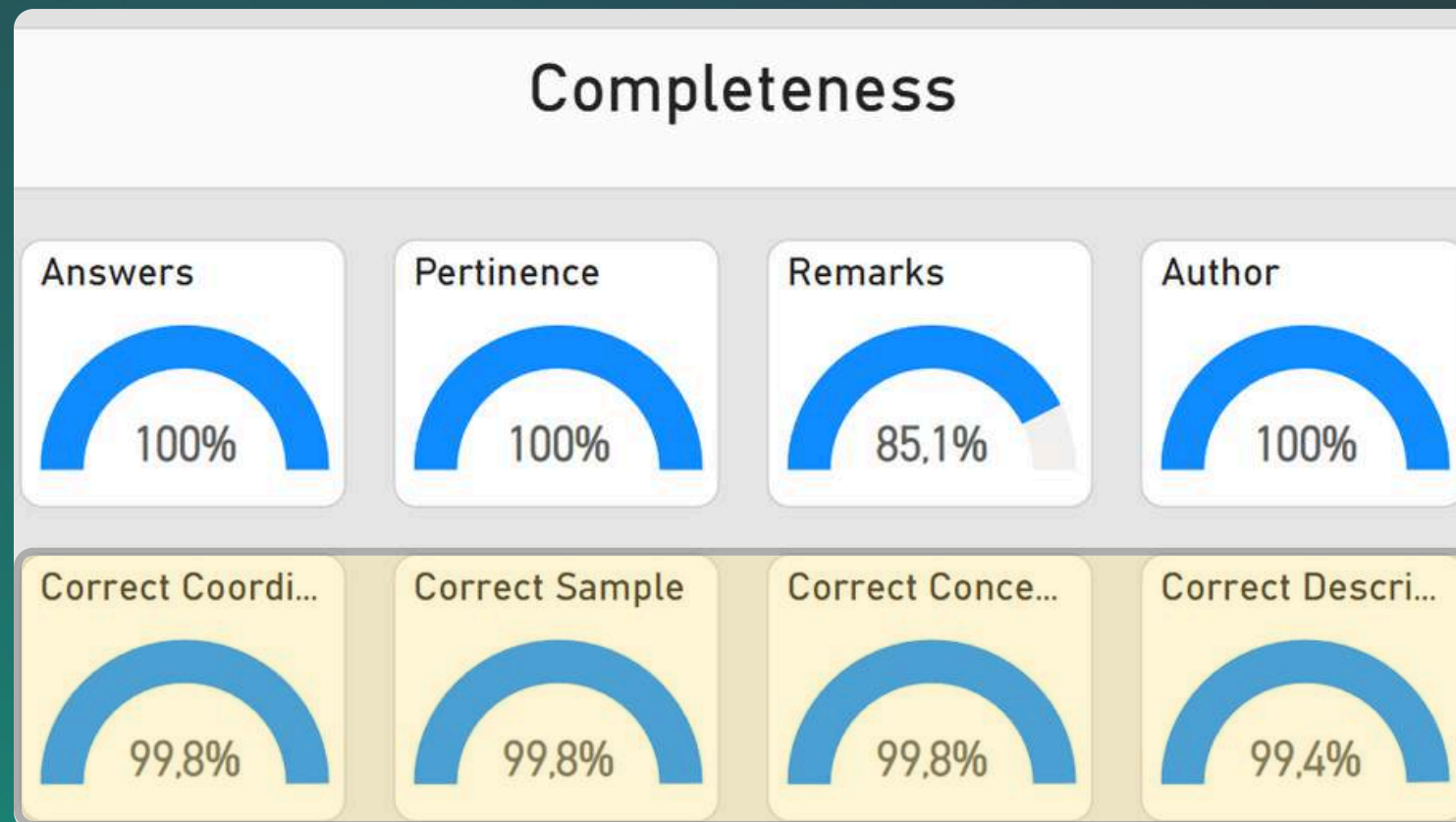
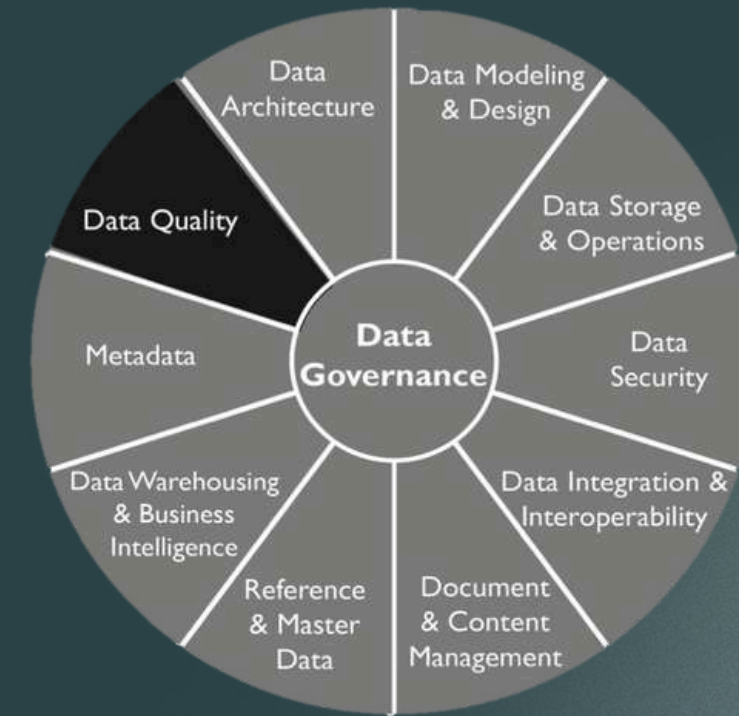


Logical Data Model



DATA QUALITY

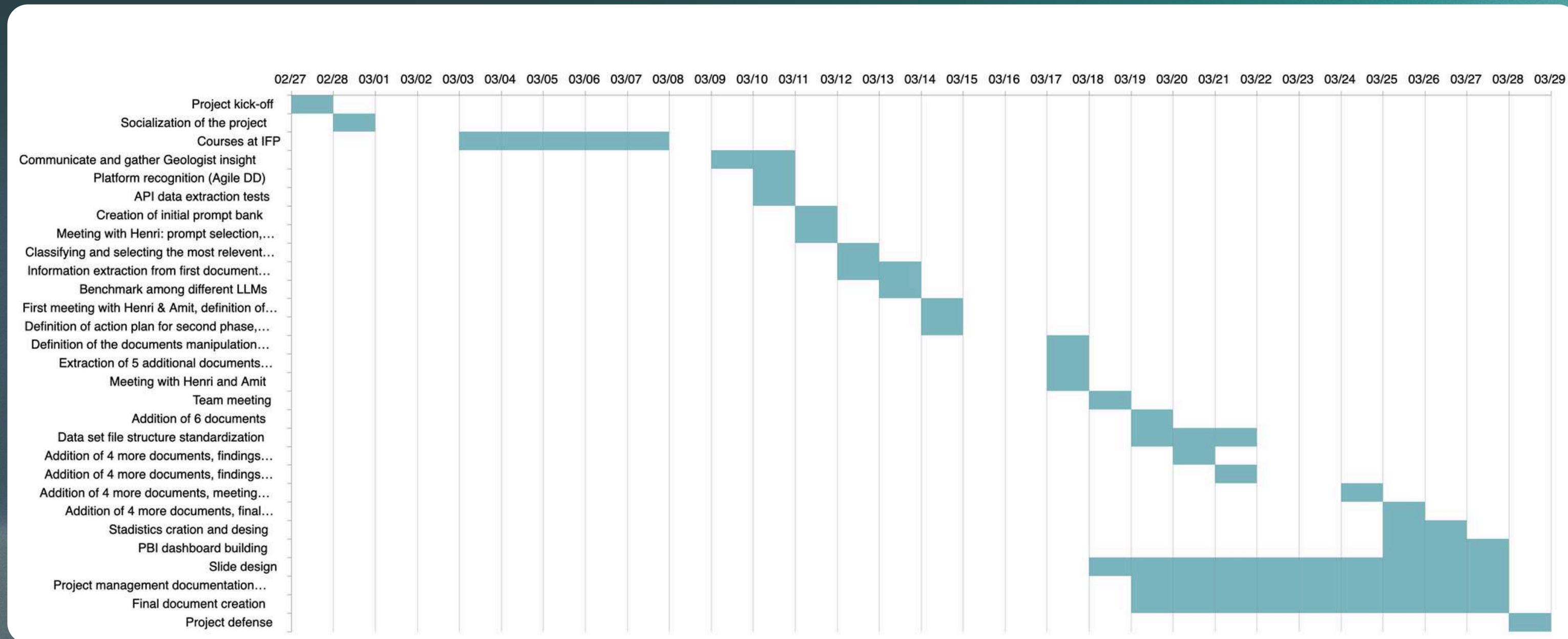
- Data quality dashboard
- Easily detect errors
- Mostly regarding the completeness dimension



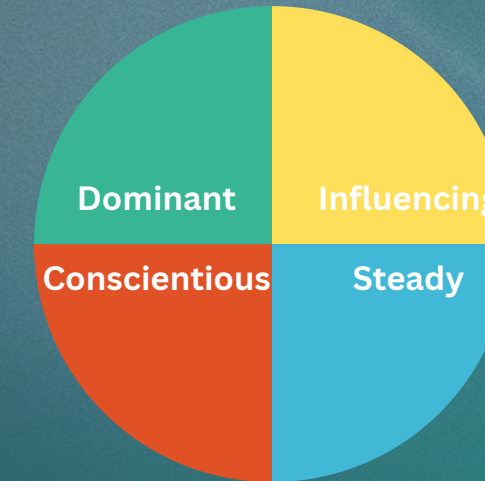
OUR TEAM



WORKING SCHEDULE



WORK BY TASK



	Full Team	Juan SIVORI	Stephania ZAMBRANO	Aicha LE BOU	Samuel NAUTIN
Phase 1	<ul style="list-style-type: none"> Recognized the platform. Tested API data extraction. Created the initial prompt bank. Extracted data from the first document (GM70036). Held daily team meetings. Held weekly meetings with tutors. 	<ul style="list-style-type: none"> Automated the workflow for downloading input files through a web-based app. Defined configuration settings for documents. 	<ul style="list-style-type: none"> Project Manager Defined the workflow inspired by agile methodology. Facilitated communication with tutors and the team. Documented processes. Established the document structure for meeting minutes. 	<ul style="list-style-type: none"> Conducted a benchmark among different LLMs. 	<ul style="list-style-type: none"> Communicating and gathering Geologist insight Compiled and Classifying and selecting the most technical questions out of the bank Presentation on document classification
Phase 2	<ul style="list-style-type: none"> Defined new objectives. Designed an action plan. Standardized document manipulation processes. Extracted data from the first five documents. Held daily team meetings. Held weekly meetings with tutors. 	<ul style="list-style-type: none"> Automated and standardized data entry in order to store consistent prompt results. Calculated statistics to assess configuration performance. 	<ul style="list-style-type: none"> Developed a new workflow and delegated tasks. Created the project schedule. Supervised task completion to ensure progress. 	<ul style="list-style-type: none"> Organized Excel files daily according to new versions. Extracted and categorized remarks. 	<ul style="list-style-type: none"> Configuration definition for documents Define naming convention for documents Standardized data entry in order to store consistent prompt results.
Phase 3	<ul style="list-style-type: none"> Standardized the dataset file structure. Added 15 more documents to complete the representative sample. Discussed and presented the results to Agile DD CEO and tutor Built visuals and schemas for processes Created the deliverables 	<ul style="list-style-type: none"> Created the data model to represent the data gathered. Developed the final dashboard with performance and data quality metrics. 	<ul style="list-style-type: none"> Refined the workflow, delegated tasks, and supervised execution. Ensured the quality of deliverables. Created the final presentation. Finalized and adjusted the project document. 	<ul style="list-style-type: none"> Created the first version of the presentation. Adjusted the final presentation. Adjusted the final report 	<ul style="list-style-type: none"> Creation of final document (structure, visuals, table, redaction)

STRATEGIES

AGILE-INSPIRED WORKFLOW

Adaptive Project Management

- Adopted an Agile-inspired methodology for flexibility and responsiveness.
- Worked iteratively, adjusting tasks and objectives based on evolving insights.

Daily Meetings

- Morning sessions:
Defined daily goals and assigned responsibilities.
- Evening sessions:
Reviewed progress, resolved issues, and outlined next steps.

Weekly Progress Reviews

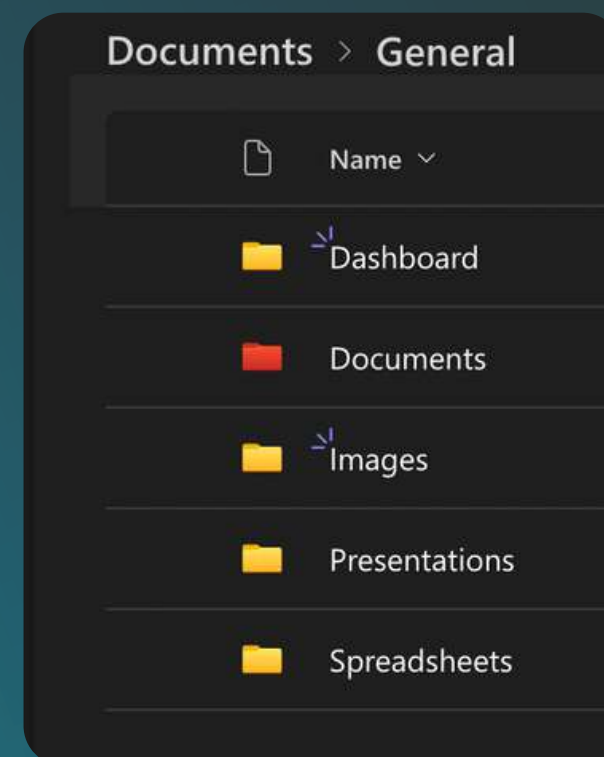
- Regular meetings with our tutor (client role) to validate results.
- Continuous refinement of strategies informed by weekly feedback.

STRATEGIES

FILE ORGANIZATION & DOCUMENT TRACKING

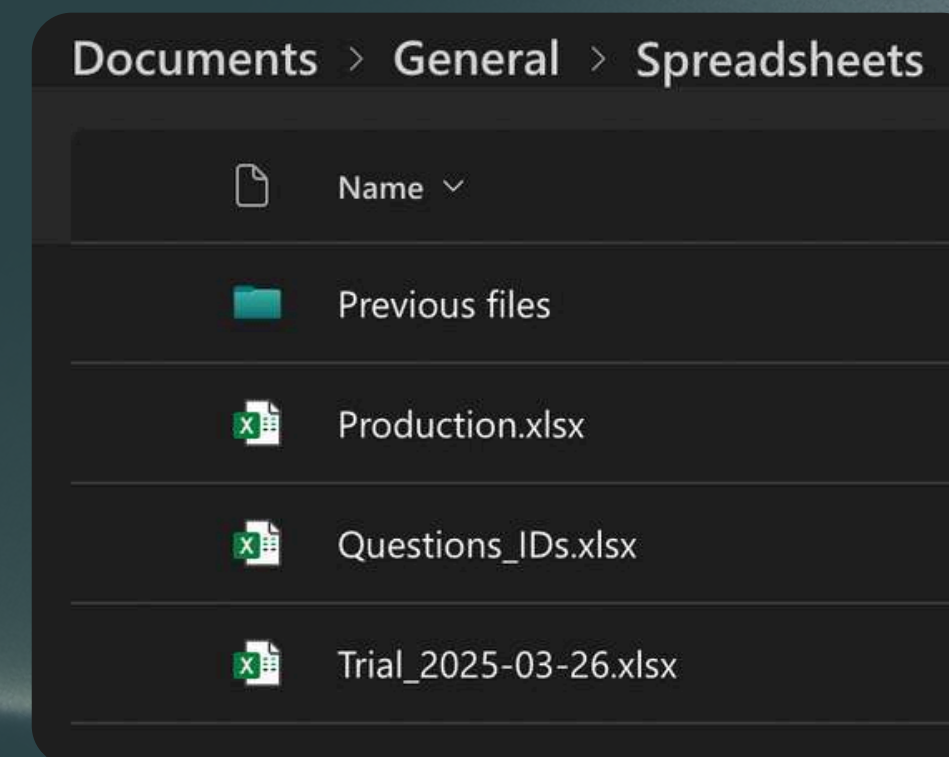
Organized Collaboration with Microsoft Teams

All files were sorted into shared folders
(Spreadsheets, Documents, Presentations,
Dashboards, Images)



Daily evaluations performed using a standardized Excel template:

New tracking sheets created daily (e.g Trial_2025-03-24.xlsx), with previous versions archived systematically for version control and data traceability.



Thank You

