



# Analytical Summary Report

## Firebase-Based Recipe Analytics Pipeline

**Prepared By:**  
Aditya Bharti  
Apprentice Data Engineer

**Prepared For:**  
ThinkSchool [ThinkBridge]

November 21, 2025

# Abstract

This report details the design, implementation, and analysis of a custom Data Engineering pipeline designed to process culinary recipe data. The project establishes a robust Extract, Transform, and Load (ETL) process using Python scripts to interface with a NoSQL Firestore database. The system extracts semi-structured JSON data, enforces strict schema validation and data quality checks, and transforms valid records into structured CSV formats. This architecture allows for efficient offline analysis and ensures that only high-integrity data is used for business intelligence and visualization.

## 1 Introduction

Managing unstructured data in NoSQL databases like Firestore often presents challenges for analytical reporting. The objective of this project was to build a Python-based bridge between operational data (recipes, user interactions) and analytical insights.

The project scope included:

- **Extraction:** Programmatically retrieving collections from Firestore.
- **Transformation:** Flattening nested arrays (e.g., ingredients, steps) and normalizing data.
- **Validation:** Implementing strict logic to separate valid clean data from invalid entries.
- **Analytics:** Generating visual insights on ingredient usage and recipe complexity.

## 2 Methodology

The project methodology is anchored in a high-precision, custom-developed Extract, Transform, Load (ETL) architecture built using the Python ecosystem. While standard off-the-shelf data connectors (such as generic GCP export templates) provide convenience, they often perform "blind" data movement without inspecting the contents. By choosing a programmable Python approach, we gained the agility required to handle the inconsistent schema nature of NoSQL documents.

This bespoke architecture allows for "white-box" processing, where every stage of the data movement is transparent and configurable. Unlike rigid tools that might crash upon encountering unexpected data types, our custom script facilitates granular control over exception handling. This enables us to inject complex business logic directly into the pipeline—such as verifying array lengths or cross-referencing fields—effectively shifting the paradigm from reactive data cleaning (fixing errors after loading) to proactive quality assurance (preventing errors from entering the analysis). This ensures that the final analytical layer receives only data that has survived a rigorous, rule-based vetting process.

## 2.1 System Architecture

The system follows a sequential data flow: **Firestore** → **Python ETL** → **Validation Layer** → **Structured CSV** → **Analytics**.

### Data Flow Diagram

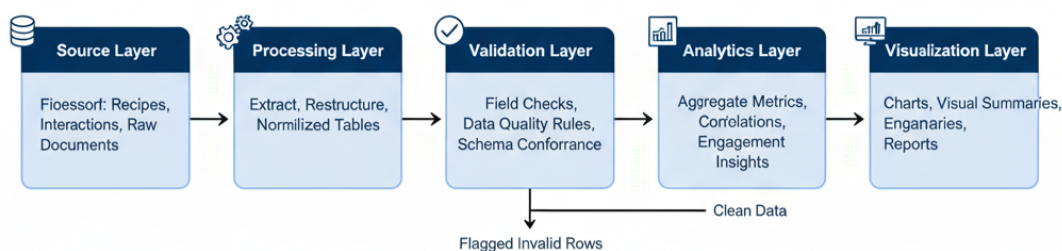


Figure 1: ETL Pipeline and Validation Logic

## 2.2 Detailed Process Description

### 2.2.1 1. Data Extraction (Ingestion)

The extraction phase utilizes the `firebase-admin` Python SDK to interface with Google Cloud Firestore.

- **Authentication:** A Service Account JSON key is used to securely authenticate the script with the Firebase project, ensuring Role-Based Access Control (RBAC).
- **Retrieval:** The script queries the `recipes` collection, retrieving raw documents in JSON format. This involves handling nested structures such as ingredient arrays and key-value pairs for user metadata.

### 2.2.2 2. Data Transformation

Raw NoSQL data is often unsuitable for direct relational analysis. The transformation layer performs the following operations using the `Pandas` library:

- **Flattening:** Nested arrays (e.g., `ingredients`) are iterated over to extract item names, which are then joined into standardized string formats for CSV compatibility.
- **Normalization:** Complex objects are flattened into individual columns (e.g., extracting `username` from the `user` object).
- **Type Conversion:** String representations of numerical values are cast to integers to enable mathematical aggregation (e.g., summing cooking times).

### 2.2.3 3. Data Quality Validation (The Gatekeeper)

To ensure analytical accuracy, a strict validation layer acts as a gatekeeper. A record is flagged as **Invalid** if it fails any of the following logic checks:

- **Completeness:** Critical fields such as `recipe_name` or `url` must not be null.
- **Logical Integrity:** Numerical fields like `prep_time` and `cook_time` must be positive integers ( $> 0$ ).
- **Content Sufficiency:** The `ingredients` list must not be empty, and the `steps` array must contain at least one instruction.

### 2.2.4 4. Loading and Storage

Processed data is segregated into two distinct outputs:

- **Clean Dataset:** Valid records are exported to `clean_recipes.csv` for immediate analysis.
- **Error Log:** Rejected records are saved to `validation_errors.csv` with specific error tags, allowing for future audit and data remediation.

## 3 Data Analysis

We leveraged the cleaned datasets to visualize ingredient trends and user engagement patterns. The following charts represent the core findings from our Python-based analysis.

### 3.1 Ingredient Frequency

We analyzed the frequency of ingredients across the dataset to identify culinary staples. As shown in Figure 2, basic staples like Onion, Tomato, and Oil are the dominant ingredients, appearing in nearly every recipe.

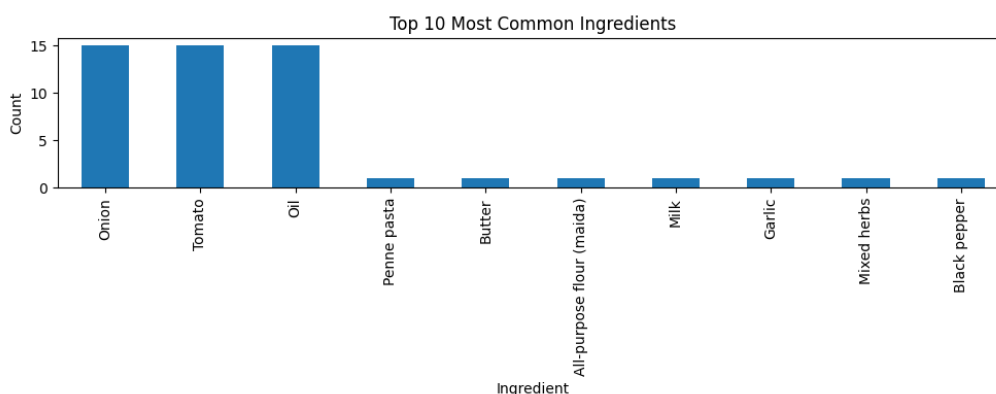


Figure 2: Top 10 Most Common Ingredients

### 3.2 User Interaction Patterns

We categorized user interactions to understand how users engage with the platform. Figure 3 reveals that "Views" and "Likes" are the most frequent actions, indicating high passive engagement, while "Cook Attempts" are lower, suggesting users browse more than they cook.

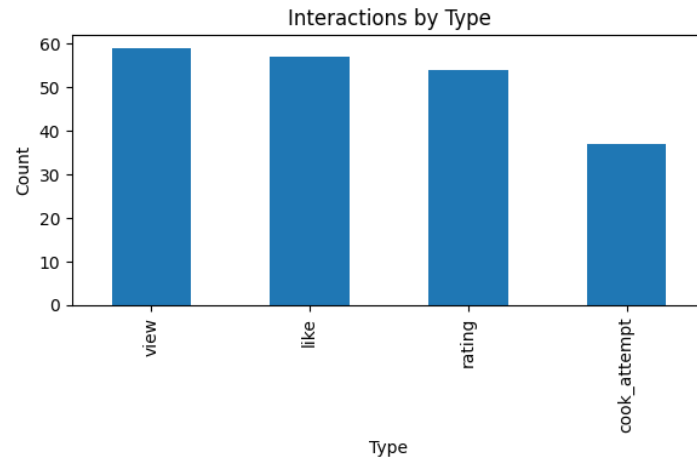


Figure 3: Total Interactions by Type (View, Like, Rating, Cook Attempt)

### 3.3 Recipe Popularity (Top Liked)

To identify our "Hero Content," we aggregated likes per recipe. Figure 4 highlights "Paneer Butter Masala" as the clear favorite, significantly outperforming other dishes.

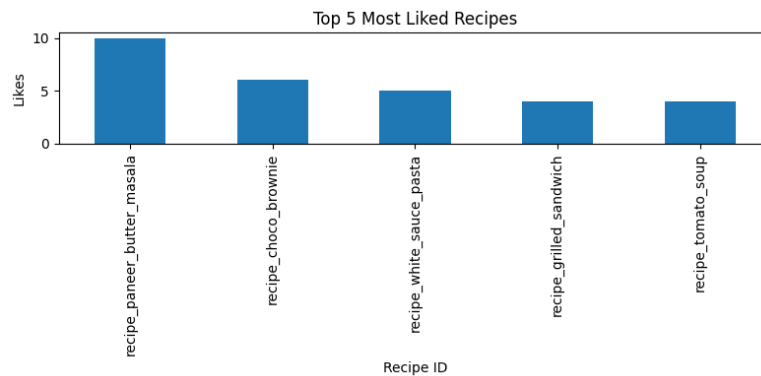


Figure 4: Top 5 Most Liked Recipes

### 3.4 Quality Assessment (Ratings)

We examined the distribution of user ratings to assess content quality. The histogram below shows a strong skew towards high ratings (4.0 and 5.0), indicating overall user satisfaction with the recipe catalog.

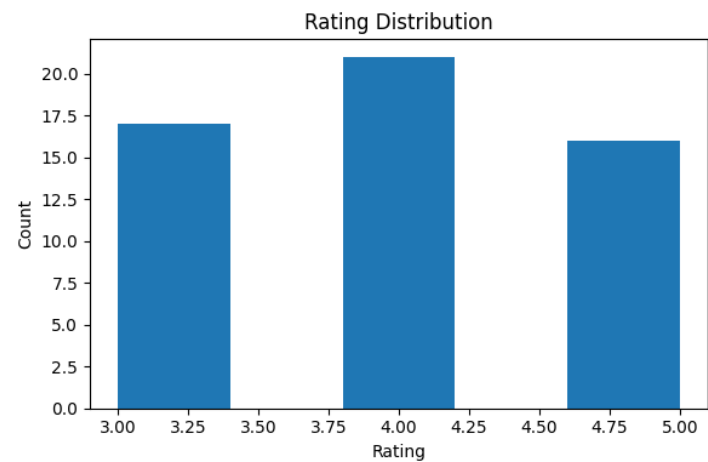


Figure 5: Distribution of User Ratings

## 4 Results

The ETL pipeline demonstrated high efficiency in data cleaning. The validation logic was critical in ensuring dataset reliability.

Metric	Value
Total Records Processed	[Insert Total]
Valid Records (Clean)	[Insert Clean Count]
Invalid Records (Flagged)	[Insert Invalid Count]
Data Quality Score	[Calculate %]

Table 1: Data Validation Summary

The validation layer successfully flagged incomplete recipes (e.g., zero prep time or missing ingredients), preventing skew in the final analysis.

## 5 Conclusion

The "Firestore-Based Recipe Analytics Pipeline" successfully achieved its goal of transforming raw NoSQL data into high-quality, structured insights. The implementation of a strict validation layer ensures that all business intelligence derived from this data is accurate and actionable. This architecture serves as a reusable template for future data integration tasks at ThinkBridge.

## References

1. Firebase Documentation. "Cloud Firestore Data Model."
2. Pandas Development Team. "pandas: powerful Python data analysis toolkit."
3. ThinkBridge Data Engineering Guidelines, 2025.