# NoSQL Assignment 3

IMT2022001 Siddharth Menon
IMT2022006 Shreyas Arun Saggere
IMT2022082 Sathvik S. Rao
IMT2022507 Vishruth Vijay

April 14, 2025

## Question 1:

The objective in this question was to define a Hive schema for the data files given, and create Hive tables accordingly. We initially created a staging table in order to load the data as it is. Then we loaded the data following our schema. We also defined an error schema for erroneous records. Our Hive schema looks as follows:

1. **Course_Attendance.csv** The **course_attendance_staging** table was used for staging.

| Column Name | Data Type | Description |
|---|---|---|
| course | TEXT | Course identifier or name |
| name_hash | TEXT | Hashed student name |
| email_hash | TEXT | Hashed email address |
| member_id_hash | TEXT | Hashed member ID |
| classes_attended | INTEGER | Number of classes attended |
| classes_absent | INTEGER | Number of classes missed |
| instructors | ARRAY$\langle String \rangle$ | List of instructors (parsed from string) |
| avg_attendance_percent | DECIMAL(5,2) | Average attendance percentage |

Table 1: Schema of the `course_attendance` Table

2. **Enrollment_Data.csv** The **enrollment_data_staging** table was used for staging.

| Column Name | Data Type | Description |
|---|---|---|
| serial_no | INTEGER | Serial number |
| course | TEXT | Course name |
| status | TEXT | Enrollment status |
| course_type | TEXT | Type of course |
| course_variant | TEXT | Course variant |
| academia_lms | TEXT | NA |
| student_id | TEXT | Student ID |
| student_name | TEXT | Name of the student |
| program | TEXT | Program enrolled in |
| batch | TEXT | Batch or cohort |
| period | TEXT | Academic period or semester |
| enrollment_date | DATE | Date of enrollment |
| primary_faculty | ARRAY⟨TEXT⟩ | List of faculty members for the course |

Table 2: Schema of the `enrollment_data` Table

3. **grade_roster_report.csv**

| Column Name | Data Type | Description |
|---|---|---|
| academy_location | STRING | Location of the academy or campus |
| student_id | STRING | Unique identifier for the student |
| student_status | STRING | Status of the student |
| admission_id | STRING | Admission ID associated with the student |
| admission_status | STRING | Status of the admission |
| student_name | STRING | Full name of the student |
| program_code_name | STRING | Program code and name |
| batch | STRING | Batch or cohort information |
| period | STRING | Academic period or term |
| subject_code_name | STRING | Subject code and name |
| course_type | STRING | Type of course |
| section | STRING | Section identifier |
| faculty_name | STRING | Name of the faculty/instructor |
| course_credit | INT | Credit value of the course |
| obtained_marks_grade | STRING | Marks or grade obtained by the student |
| out_of_marks_grade | STRING | Maximum marks or grade possible |
| exam_result | STRING | Exam result |

Table 3: Schema of the `grade_roster_report` Table

4. **Erroneous record schema**

| Column Name | Data Type | Description |
|---|---|---|
| source_table | STRING | Name of the source table where the error originated |
| original_row | STRING | Raw data row containing the error |
| column_name | STRING | Name of the column that caused the error |
| issue_type | STRING | Type or category of the issue (e.g., NULL, format error) |
| error_description | STRING | Detailed description of the error |

Table 4: Schema of the `error_records` Table

# Question 2

The **course_attendance table, enrollment_data table, grade_roster_report** table were joined on the fields **student_id and course field** to create a warehouse filtering out erroneous data. This warehouse was used in order to perform our analytical queries.
The analytical queries were:

1. **Query 1**: Comparing the mode grade and average attendance of different students of different programs for every course

2. **Query 2**: Percentage of students with a below 50% attendance record, per program, per course.

3. **Query 3**: Identification of top 3 courses by average attendance for each program.

# Question 3

We set up a new data warehouse that is partitioned by **batch/program** and bucketed into 4 buckets by **course**. The goal was to evaluate the impact of partitioning on query performance. We chose to partition by **program** because our queries frequently include a 'GROUP BY' operation on the program field. For the same reason we chose to bucket by **course**.

# Question 4: Converting Hive queries to Pig

**Important note:** A difference in schema between Hive and Pig: in the Hive schema, we have two fields `batch` and `period` – which would contain values like (Master of Technology, 2024–2026 CSE), while these same fields are renamed to `program_name` and `batch`, where `batch` is analogous to the `period` field in the Hive schema. For this section, those two fields will be referred to by the names they are referred to in the Pig schema.

To analyse the runtime performance of pig queries, after performing transformations on the data, we had to then use the `dump` command, as Pig uses lazy evaluation, and the actual operations aren't done until some kind of output operation is done.

Furthermore, since the code for each of the queries was put in a different file and executed one after the other from the terminal, the load had to be done in all the three query files, causing an additional overhead for a load every time. If the 3 queries are executed together, and not timed separately, the time taken is lesser, as we load only once. That analysis, too, has been done.

The three queries that were converted from Hive to Pig Latin were:

1. Comparing the mode grade and average attendance of different students of different programs for every course.



Figure 1: Query 1 execution



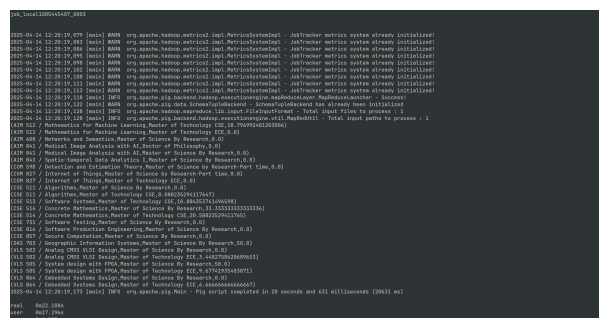Figure 2: Time taken for query 1

2. Percentage of students with a below 50% attendance record, per program, per course.



Figure 3: Query 2 execution



Figure 4: Time taken for query 2

3. Identification of top 3 courses by average attendance for each program.

Figure 5: Query 3 execution



Figure 6: Time taken for query 3

Executing all the 3 queries together and then timing it:



Figure 7: All queries together - execution



Figure 8: All queries together - time taken

| Query | Execution Time (seconds) |
|---|---|
| Query 1 | 21.174 |
| Query 2 | 22.108 |
| Query 3 | 15.485 |
| All Queries Executed Together | 50.720 |

Table 5: Execution times of individual queries and combined execution

# 1 Comparing the performance of Hive(Partitioned and non-partitioned) and Pig

**Query 1: Average attendance and mode grade per course per batch**



Figure 9: Query 1 - Non-Partitioned Execution



Figure 10: Query 1 - Partitioned Execution

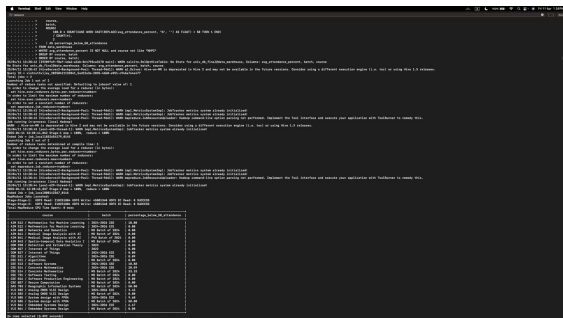**Query 2: Percentage of students with below 50% attendance per course and program**



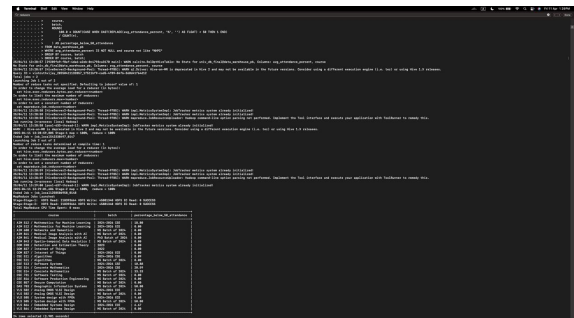Figure 11: Query 2 - Non-Partitioned Execution



Figure 12: Query 2 - Partitioned Execution
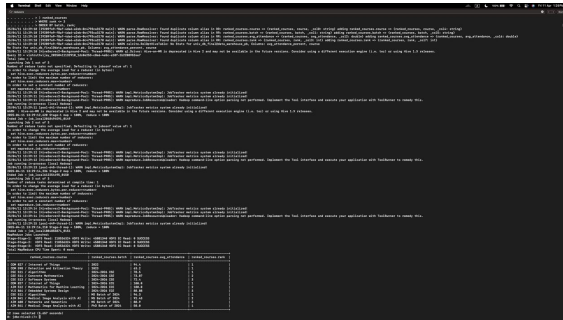
# Query 3: Top 3 courses by average attendance per program
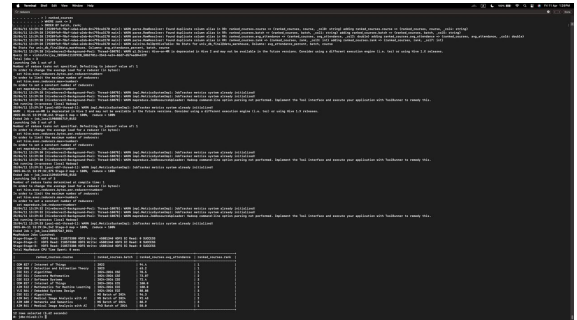


Figure 13: Query 3 - Non-Partitioned Execution



Figure 14: Query 3 - Partitioned Execution

# Execution Time Comparison Table

| Query | Partitioned + Bucketed Time (s) | Non-Partitioned Time (s) | Pig (s) |
|-------|--------------------------------|--------------------------|---------|
| Query 1 | 9.356 | 9.62 | 21.174 |
| Query 2 | 3.892 | 3.901 | 22.108 |
| Query 3 | 5.62 | 5.657 | 15.485 |

Table 6: Execution Time Comparison: Partitioned vs Non-Partitioned Data Warehouse