



ISIT312 Big Data Management

SIM Session 4, 2025

Assignment 2 (Task 1,2,3,4)

Contents

Task 1 Part 1: Data Cube Specification	3
1.1. Explanation of the Data Warehouse Design	3
1.2. Fact Identification	3
1.3. Measures	3
1.4. Dimensions and Hierarchies	3
1.5. Conceptual Snowflake Schema Diagram (ASCII)	4
Task 1 Part 2: Perspective Drawing of a Three-Dimensional Data Cube	5
Task 2: Conceptual Model for the Vehicle Repair Data Warehouse	6
Explanation of the Design	6
Data Warehouse Component Specification	6
Conceptual Star Schema Diagram	7
Queries Implemented to obtain relevant information:	8
Explanation of the Schema	10
Fact Table	10
Dimension Tables	10
Hierarchies and Analytical Capability	11
Analytical Support	11
Summary	11
Task 3 Implementation of a table with a complex column type (ONF table) in Hive	12
Solution3.hql:	12
Task 3 Report	12
1. Overview	12
2. HQL Script Code	13
2.1 Table Cleanup	13
2.2 Table Creation with Complex Data Types	13
2.3 Data Loading	13
2.4 Data Verification	14
3. Terminal Output and Results	14
3.1 Connection to Hive	14
3.2 Script Execution Results	14
3.3 Query Results	14
3.4 Raw Output Format	14

ROHIT PANDA 8943060 ASSIGNMENT 2 TASK 1,2,3,4 ISIT312 BIG DATA MANAGEMENT

4. Analysis and Discussion	15
4.1 Key Observations.....	15
4.2 Performance Metrics.....	15
4.3 Advantages of 0NF Design.....	15
5. Conclusion	16
6. Terminal Output for Task 3	16
Task 4 Implementation of a data warehouse as a collection of external tables in Hive.....	18
Solution4.hql.....	18
Task 4 Report	19
Implementation Process	19
Summary of Results.....	21
Technical Notes	22
Conclusion	22
Complete Terminal Output Task 4.....	22

Name : Rohit Panda

UOW Student ID : 8943060

Task 1 Part 1: Data Cube Specification

1.1. Explanation of the Data Warehouse Design

The operational database for the bus company is designed for real-time tracking, capturing only a bus's current location at any given moment. This "point-in-time" design means historical data is discarded, making it impossible to perform crucial business analyses on past performance, efficiency, or route usage.

To address this limitation, a data warehouse will be implemented using a **data cube** model. This model will capture and store historical data for every trip segment, creating a permanent record for analysis. The design is centered around a **fact**, which represents the most granular business event: a bus completing a single trip segment. This fact is described by quantitative **measures** (e.g., kilometers travelled) and is given context by descriptive **dimensions** (e.g., Time, Bus, Location). The dimensions are structured with **hierarchies** to allow for powerful analytical operations, such as "drilling down" from a year to a specific day, or "rolling up" from a depot to a city.

1.2. Fact Identification

The core business event that allows for the required analysis is the completion of a trip segment by a bus.

- **Fact:** Bus Trip Segment Journey

1.3. Measures

These are the quantitative metrics that will be analyzed, derived strictly from the applications (i) through (vii).

- **Kilometers Travelled:** The distance of the segment. Used to measure total distance covered and performance. (From requirement i)
- **Travel Duration:** The time taken to complete the segment. The base for calculating averages. (From requirement v)
- **Fuel Consumption:** The calculated amount of fuel used for the segment. Used to monitor operational costs and vehicle efficiency. (From requirement vi)
- **Trip Count:** A count of 1 for each fact record, which is summed to calculate the total number of trips or segments performed. (From requirements ii, iv, vii)

1.4. Dimensions and Hierarchies

These are the contextual perspectives for analysis, derived from the "per..." clauses in the requirements and the relationships in the conceptual schema.

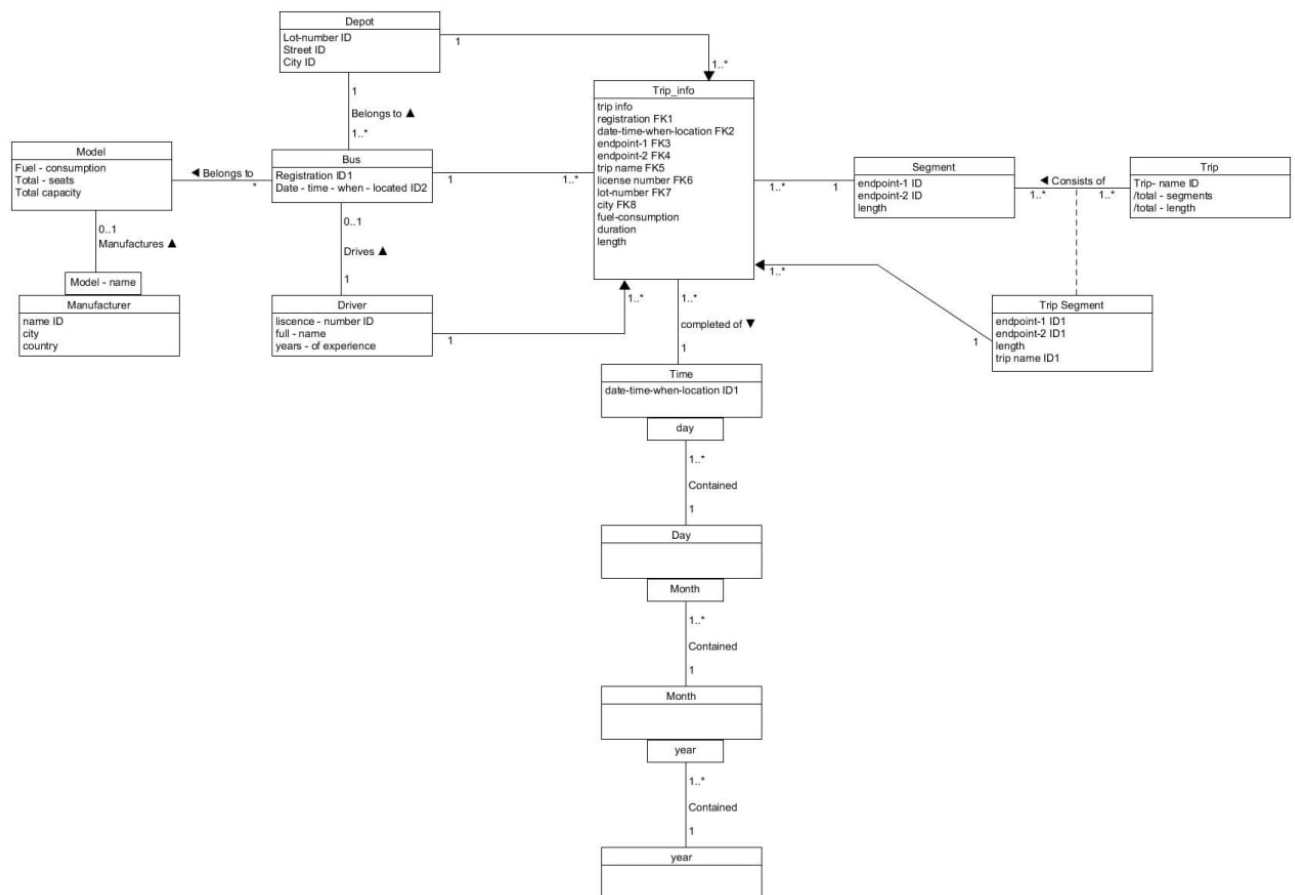
- **Time Dimension**
 - **Description:** Represents when the trip segment was completed.
 - **Hierarchy:** Year → Month → Day
 - **Attributes:** Year, Month, Day, Date.
- **Bus Dimension**
 - **Description:** The specific vehicle used for the journey.
 - **Hierarchy:** Manufacturer → Model → Bus
 - **Attributes:** Registration Number, Model Name, Total Seats, Fuel Consumption Rate, Manufacturer Name.

ROHIT PANDA 8943060 ASSIGNMENT 2 TASK 1,2,3,4 ISIT312 BIG DATA MANAGEMENT

- **Trip Dimension**
 - **Description:** The route and specific segment of the journey.
 - **Hierarchy:** Trip → Trip Segment
 - **Attributes:** Trip Name, Segment ID, Segment Endpoints, Segment Length.
- **Location Dimension**
 - **Description:** The geographical location of the depot to which the bus is assigned.
 - **Hierarchy:** City → Depot
 - **Attributes:** City Name, Depot Street, Depot Lot Number.
- **Driver Dimension**
 - **Description:** The employee who drove the bus for the segment.
 - **Hierarchy:** None (flat dimension).
 - **Attributes:** Licence Number, Full Name, Years of Experience.

1.5. Conceptual Snowflake Schema Diagram (ASCII)

The following diagram illustrates the normalized (snowflake) structure of the dimensions and their relationship to the central fact table:

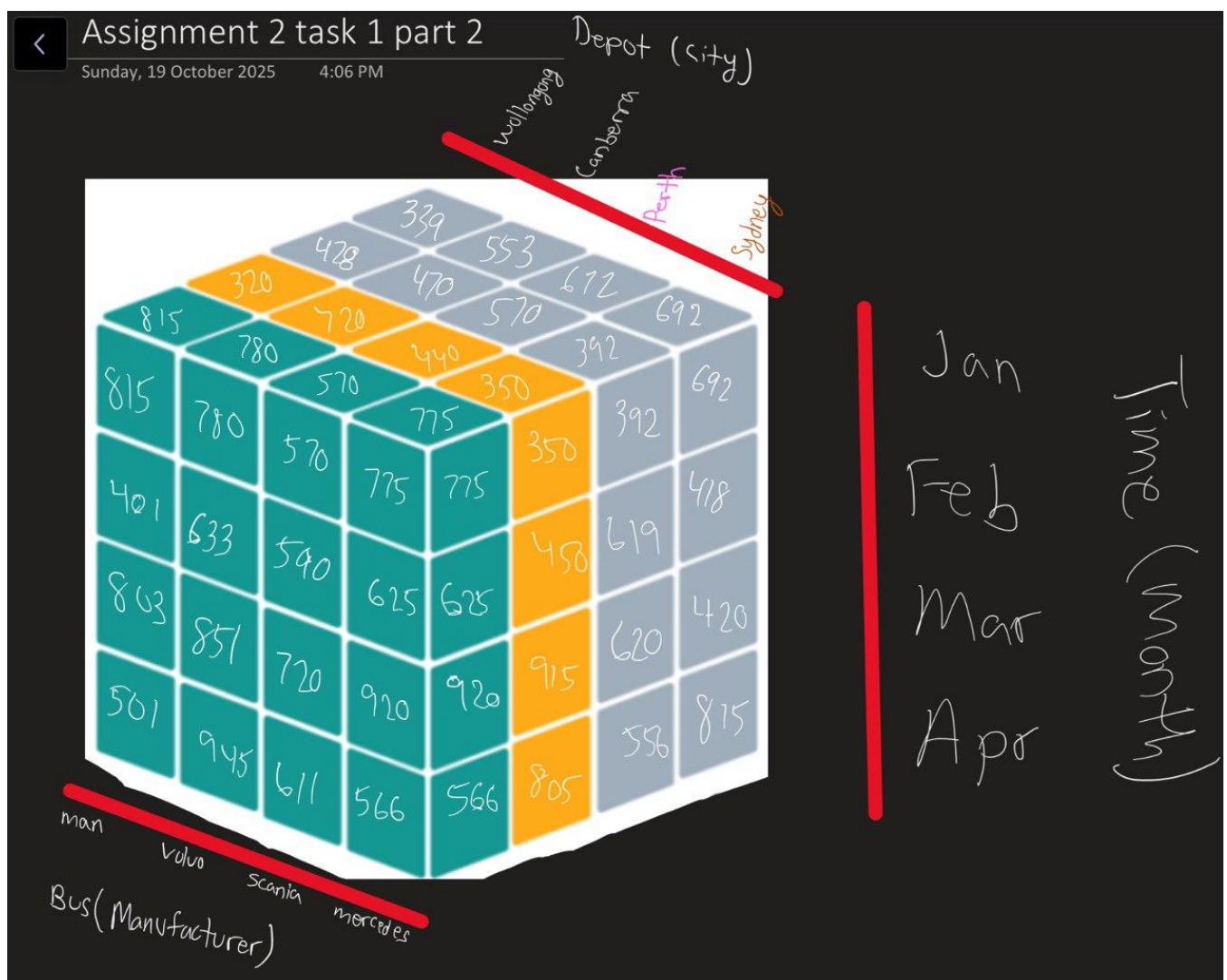


Task 1 Part 2: Perspective Drawing of a Three-Dimensional Data Cube

This section visualizes the logical data cube as an analyst would see it.

- **Selected Dimensions for Visualization:**
 1. **Time** (at the Month level)
 2. **Location** (at the City level)
 3. **Driver** (at the Driver level)
- **Selected Measure for Visualization:**
 - **Total Kilometers Travelled**
- **Dimension Values:**
 - **Time (Month):** January, February, March, April
 - **Location (City):** Sydney, Wollongong, Newcastle, Canberra

The following diagram illustrates a three-dimensional data cube. Each cell in this cube contains the value for the **Total Kilometers Travelled (In Thousands)** by a specific driver, in a specific city, during a specific month. This allows a manager to analyze driver performance and workload across different operational areas:



Task 2: Conceptual Model for the Vehicle Repair Data Warehouse

Explanation of the Design

The objective is to create a conceptual schema for a data warehouse that will allow the management of a vehicle repair network to analyze its operations. The design must support queries about costs, timings, and resources used across different facilities, mechanics, and time periods.

The most effective and industry-standard approach for this is a **star schema**. This model is chosen for its simplicity, high query performance, and ease of understanding for business analysts. It features a central **fact table** containing the core business measurements, surrounded by descriptive **dimension tables**.

Data Warehouse Component Specification

1. Fact Table

- **Fact Entity:** A single, completed **Repair/Maintenance Event**.
- **Fact Table Name:** RepairFact
- **Description:** Each row in this table will represent one unique repair or maintenance job performed on a vehicle.

2. Measures

These are the quantitative, numerical values stored in the fact table.

- **TotalCost:** The total cost or amount paid for the repair.
- **TimeSpent_Hours:** The number of hours a mechanic spent on the repair.
- **PartsUsed_Count:** The total number of spare parts used in the repair.
- **DaysSpent_OnRepair:** The total number of days the vehicle was at the facility (from drop-off to collection).
- **Mechanic_Count:** The number of mechanics involved (in this case, specified as one per repair).
- **Repair_Count:** A count of 1 for each transaction, used to easily sum the total number of repairs.

3. Dimensions and Hierarchies

These tables provide the context for the facts.

- **Time Dimension (Dim_Time)**
 - **Description:** Represents when the repair was performed.
 - **Attributes:** Full Date, Day, Month, Year.
 - **Hierarchy:** Year → Month → Day
- **Facility Dimension (Dim_Facility)**
 - **Description:** Represents where the repair took place.
 - **Attributes:** Facility ID, Building Number, City, Country, Email Address, Phone Number.
 - **Hierarchy:** Country → City → Facility
- **Mechanic Dimension (Dim_Mechanic)**
 - **Description:** Represents who performed the repair.
 - **Attributes:** Employee Number (ID), First Name, Last Name, Date of Birth.
 - **Hierarchy:** None (flat dimension).
- **Vehicle Dimension (Dim_Vehicle)**
 - **Description:** Represents which vehicle was repaired.

ROHIT PANDA 8943060 ASSIGNMENT 2 TASK 1,2,3,4 ISIT312 BIG DATA MANAGEMENT

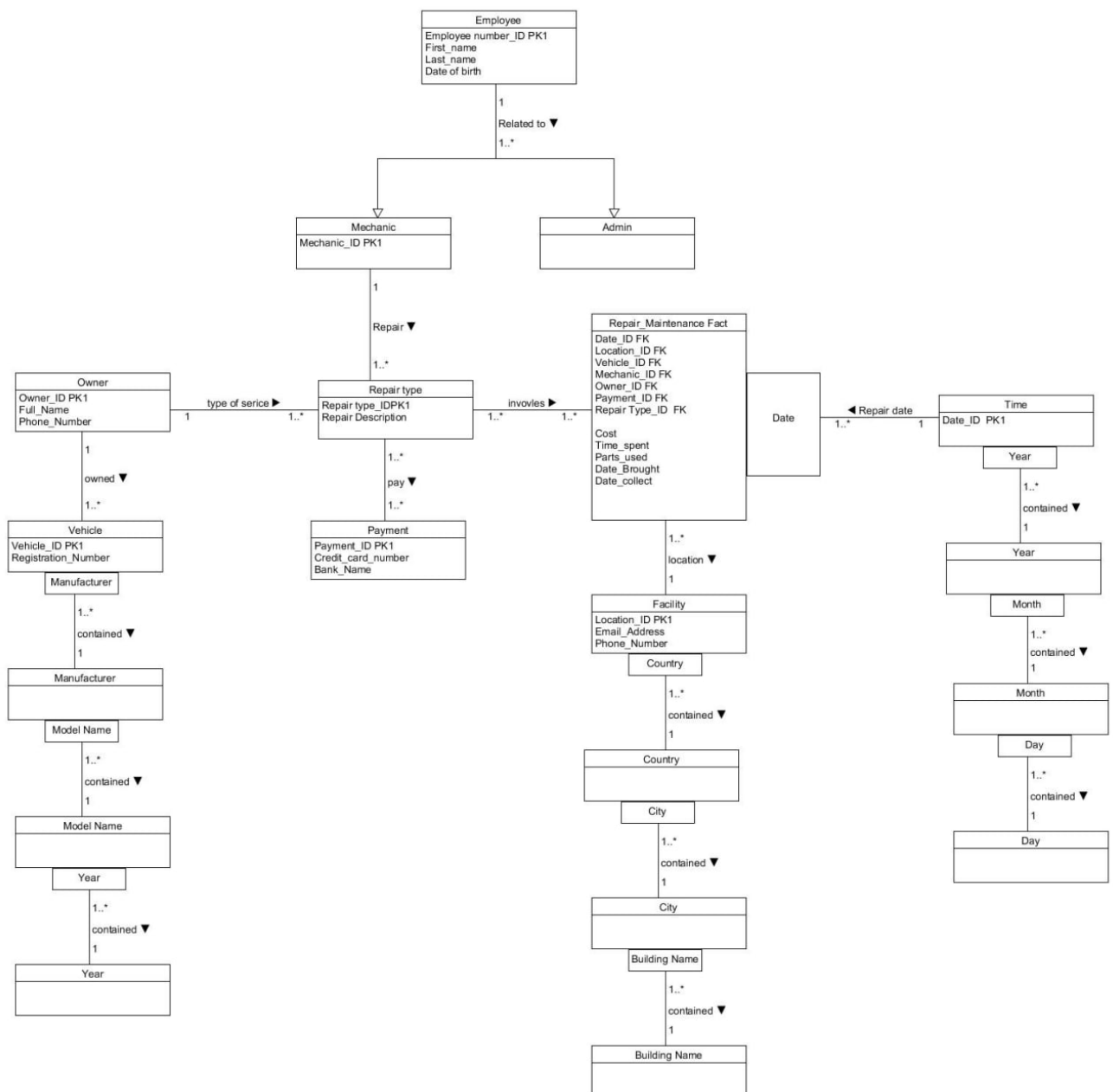
- **Attributes:** Vehicle ID, Registration Number, Manufacturer, Model, Year of Manufacture.
- **Hierarchy:** Manufacturer → Model → Vehicle

Conceptual Star Schema Diagram

This diagram visually represents the star schema, with the central fact table connected to the four dimension tables.

(PK) = Primary Key

(FK) = Foreign Key



Queries Implemented to obtain relevant information:

1. the total number of repairs/maintenances per facility, per year, per month per day, per city and per mechanic

```
SELECT
    fac.Country,
    fac.City,
    fac.Building_Name AS Facility,
    t.Year,
    t.Month,
    t.Day,
    rmf.Mechanic_ID,
    /* If Repair_Count exists use it; else fall back to COUNT(*) */
    SUM(COALESCE(rmf.Repair_Count, 1)) AS total_repairs
FROM Repair_Maintenance_Fact rmf
JOIN Time t ON rmf.Date_ID = t.Date_ID
JOIN Facility fac ON rmf.Location_ID = fac.Location_ID
GROUP BY fac.Country, fac.City, fac.Building_Name, t.Year, t.Month, t.Day,
rmf.Mechanic_ID;
```

2. total costs of repairs/maintenance per facility, per year, per month per day, per city and per mechanic

Ans:

```
SELECT
    fac.Country,
    fac.City,
    fac.Building_Name AS Facility,
    t.Year,
    t.Month,
    t.Day,
    rmf.Mechanic_ID,
    SUM(rmf.Cost) AS total_cost
FROM Repair_Maintenance_Fact rmf
JOIN Time t ON rmf.Date_ID = t.Date_ID
JOIN Facility fac ON rmf.Location_ID = fac.Location_ID
GROUP BY fac.Country, fac.City, fac.Building_Name, t.Year, t.Month, t.Day,
rmf.Mechanic_ID;
```

3. total number of mechanics involved per repair/maintenance, per year, per month per day, per facility, per city

Ans:

```
SELECT
    fac.Country,
    fac.City,
    fac.Building_Name AS Facility,
    t.Year,
    t.Month,
    t.Day,
    SUM(COALESCE(rmf.Mechanic_Involved_Count, 1)) AS mechanics_involved
FROM Repair_Maintenance_Fact rmf
JOIN Time t ON rmf.Date_ID = t.Date_ID
JOIN Facility fac ON rmf.Location_ID = fac.Location_ID
GROUP BY fac.Country, fac.City, fac.Building_Name, t.Year, t.Month, t.Day;
```


ROHIT PANDA 8943060 ASSIGNMENT 2 TASK 1,2,3,4 ISIT312 BIG DATA MANAGEMENT

4. average time spent on repairs/maintenances per year, month, day

Ans:

```
SELECT
    t.Year,
    t.Month,
    t.Day,
    AVG(rmf.Time_spent) AS avg_time_hours
FROM Repair_Maintenance_Fact rmf
JOIN Time t ON rmf.Date_ID = t.Date_ID
GROUP BY t.Year, t.Month, t.Day;
```

5. total number of parts used for repairs/maintenances per year, month, day, per facility, per city.

Ans:

```
SELECT
    fac.Country,
    fac.City,
    fac.Building_Name AS Facility,
    t.Year,
    t.Month,
    t.Day,
    SUM(rmf.Parts_used) AS total_parts_used
FROM Repair_Maintenance_Fact rmf
JOIN Time t ON rmf.Date_ID = t.Date_ID
JOIN Facility fac ON rmf.Location_ID = fac.Location_ID
GROUP BY fac.Country, fac.City, fac.Building_Name, t.Year, t.Month, t.Day;
```

Explanation of the Schema

The conceptual schema represents a snowflake-style data warehouse centered on the **Repair_Maintenance_Fact** table.

Each record in this fact table corresponds to a **completed repair or maintenance event** for a specific vehicle at a specific facility on a given date.

The schema captures both the **quantitative measures** and the **descriptive context** needed for analytical processing (OLAP).

Fact Table

- **Repair_Maintenance_Fact** is the core of the warehouse and stores the measurable business outcomes of repair activities.
 - **Measures:** Cost, Time_Spent, Parts_Used, and optionally derived metrics such as total days spent and number of mechanics involved.
 - **Foreign Keys:** Link to the dimensions Time, Facility, Mechanic, Vehicle, Owner, Payment, and Repair_Type.
 - **Grain:** One row per unique repair job — this ensures that each event is atomic, forming the foundation for aggregation and trend analysis.

Dimension Tables

Each dimension describes a unique business perspective from which the fact data can be analyzed.

1. **Time Dimension (Date_ID, Year, Month, Day)**
 - Represents when a repair occurred.
 - Supports temporal hierarchies (**Year** → **Month** → **Day**) for trend, seasonal, and drill-down analyses.
2. **Facility Dimension (Location_ID, Country, City, Building_Name, Email_Address, Phone_Number)**
 - Describes where the repair was carried out.
 - Hierarchical structure (**Country** → **City** → **Facility**) enables analysis at different geographic levels.
3. **Mechanic Dimension (Mechanic_ID, Employee_Number, First_Name, Last_Name, Date_of_Birth)**
 - Captures who performed the repair.
 - Enables performance comparison, productivity tracking, and workload analysis by individual or group of mechanics.
4. **Vehicle Dimension (Vehicle_ID, Registration_Number, Manufacturer, Model_Name, Year)**
 - Identifies which vehicle was serviced and allows aggregation by manufacturer, model, or production year.
 - Hierarchy: **Manufacturer** → **Model** → **Year**.
5. **Owner Dimension (Owner_ID, Full_Name, Phone_Number)**
 - Records who owns the vehicle.
 - Provides a customer dimension for loyalty and service-frequency analysis.
6. **Payment Dimension (Payment_ID, Credit_Card_Number, Bank_Name)**
 - Details the mode of payment and issuing bank.
 - Supports financial analysis by payment method or bank.
7. **Repair Type Dimension (Repair_Type_ID, Repair_Description)**

- Classifies the nature of the service performed (e.g., engine overhaul, tyre replacement, inspection).
- Useful for evaluating which types of repairs are most common or most costly.

Hierarchies and Analytical Capability

- **Temporal hierarchy:** Year → Month → Day
- **Geographical hierarchy:** Country → City → Facility
- **Product hierarchy:** Manufacturer → Model → Year

These hierarchies support OLAP operations such as **roll-up** (aggregating data to a higher level, e.g., monthly totals) and **drill-down** (viewing finer details, e.g., daily or by facility).

Analytical Support

The schema enables management to:

- Compute total repairs, costs, and parts used by **time**, **facility**, and **mechanic**.
- Analyze average time spent per repair type or city.
- Track revenue trends and operational efficiency across multiple locations.
- Identify the most frequent or resource-intensive repair types.

Summary

Overall, this schema balances **analytical power and structural clarity**.

It fulfills all assignment requirements by providing a unified view of time, cost, facility, and personnel performance while maintaining extensibility for future dimensions such as owner demographics or payment analysis.

ROHIT PANDA 8943060 ASSIGNMENT 2 TASK 1,2,3,4 ISIT312 BIG DATA MANAGEMENT

Task 3 Implementation of a table with a complex column type (0NF table) in Hive

Solution3.hql:

```
-- Student Name: Rohit Panda
-- Student UOW ID: 8943060
-- creates a nested table, inserts data, and verifies the contents.
-- Part 1: Implement HQL script to create an internal nested table
-- The table will store student numbers and their subject evaluations.
-- The structure is 0NF because the 'evaluations' column contains a collection.

CREATE TABLE student_evaluations (
    student_id INT,
    evaluations ARRAY<STRUCT<subject: STRING, grade: INT>>
)
STORED AS ORC; -- Using ORC format is efficient for structured data in Hive

-- Part 2: Include INSERT statements to load sample data
-- Inserting 5 records as required, using NAMED_STRUCT to build the complex type.

INSERT INTO student_evaluations VALUES
(1111, ARRAY(
    NAMED_STRUCT('subject', 'CSIT111', 'grade', 01),
    NAMED_STRUCT('subject', 'CSIT121', 'grade', 23),
    NAMED_STRUCT('subject', 'CSIT101', 'grade', 50),
    NAMED_STRUCT('subject', 'CSIT235', 'grade', 99),
    NAMED_STRUCT('subject', 'ISIT312', 'grade', 02)
)),
(1112, ARRAY(
    NAMED_STRUCT('subject', 'CSIT101', 'grade', 56),
    NAMED_STRUCT('subject', 'CSIT111', 'grade', 78),
    NAMED_STRUCT('subject', 'CSIT115', 'grade', 10),
    NAMED_STRUCT('subject', 'ISIT312', 'grade', 05)
)),
(1113, ARRAY(
    NAMED_STRUCT('subject', 'CSIT121', 'grade', 76),
    NAMED_STRUCT('subject', 'CSIT235', 'grade', 87),
    NAMED_STRUCT('subject', 'ISIT312', 'grade', 49)
)),
(1114, ARRAY(
    NAMED_STRUCT('subject', 'CSIT111', 'grade', 50),
    NAMED_STRUCT('subject', 'ISIT312', 'grade', 45)
)),
(1115, ARRAY(
    NAMED_STRUCT('subject', 'ISIT115', 'grade', 67),
    NAMED_STRUCT('subject', 'CSCI235', 'grade', 45),
    NAMED_STRUCT('subject', 'CSIT127', 'grade', 56)
));

-- Part 3: Include a SELECT statement to list the contents of the table
-- This verifies that the data has been inserted correctly.

SELECT * FROM student_evaluations;
```

Task 3 Report

1. Overview

This report documents the implementation of a nested table structure in Apache Hive to store student evaluation data in Zero Normal Form (0NF). The solution demonstrates Hive's capability to handle complex, semi-structured data using array and struct data types.

2. HQL Script Code

2.1 Table Cleanup

```
DROP TABLE IF EXISTS student_evaluations;
```

Purpose: This statement ensures a clean slate by removing any existing `student_evaluations` table before creating a new one. The `IF EXISTS` clause prevents errors if the table doesn't exist.

2.2 Table Creation with Complex Data Types

```
CREATE TABLE student_evaluations (  
  student_id INT,  
  evaluations ARRAY<STRUCT<subject: STRING, grade: INT>>  
)  
ROW FORMAT DELIMITED  
FIELDS TERMINATED BY '|'   
COLLECTION ITEMS TERMINATED BY ','  
MAP KEYS TERMINATED BY ':'  
STORED AS TEXTFILE;
```

Structure Explanation:

- **student_id INT:** Stores the unique student identifier
- **evaluations ARRAY<STRUCT<...>>:** A complex nested column that stores multiple evaluation records per student
 - The ARRAY type allows multiple elements per student
 - Each element is a STRUCT containing:
 - **subject (STRING):** The subject code
 - **grade (INT):** The evaluation score (0-99)

Delimiter Configuration:

- **FIELDS TERMINATED BY '|':** Separates `student_id` from the evaluations array
- **COLLECTION ITEMS TERMINATED BY ',':** Separates individual evaluation records within the array
- **MAP KEYS TERMINATED BY ':':** Separates subject codes from grades within each struct
- **STORED AS TEXTFILE:** Specifies plain text storage format

This 0NF (Zero Normal Form) design allows repeating groups within a single row, which is non-relational but suitable for semi-structured data.

2.3 Data Loading

```
LOAD DATA LOCAL INPATH '/home/bigdata/Desktop/students.tbl'  
INTO TABLE student_evaluations;
```

ROHIT PANDA 8943060 ASSIGNMENT 2 TASK 1,2,3,4 ISIT312 BIG DATA MANAGEMENT

Purpose: This statement efficiently bulk-loads data from the local file system into the Hive table. The `LOCAL` keyword indicates the file is on the local filesystem rather than HDFS. The operation completed successfully in 0.851 seconds.

2.4 Data Verification

```
SELECT * FROM student_evaluations;
```

Purpose: This query retrieves all records to verify successful data loading and proper parsing of the nested structure.

3. Terminal Output and Results

3.1 Connection to Hive

```
Beeline version 2.1.1 by Apache Hive
Connected to: Apache Hive (version 2.1.1)
Driver: Hive JDBC (version 2.1.1)
```

The connection to the Hive server was established successfully through the JDBC interface on localhost:10000.

3.2 Script Execution Results

Table Drop:

```
DROP TABLE IF EXISTS student_evaluations;
No rows affected (0.196 seconds)
```

Table Creation:

```
CREATE TABLE student_evaluations...
No rows affected (0.255 seconds)
```

Data Loading:

```
LOAD DATA LOCAL INPATH '/home/bigdata/Desktop/students.tbl' INTO TABLE student_evaluations;
No rows affected (0.851 seconds)
```

3.3 Query Results

The `SELECT` statement returned 5 rows of data with the following structure:

Student 1111: <ul style="list-style-type: none">Subject: CSIT111, Grade: 1Subject: CSIT121, Grade: 23Subject: CSIT101, Grade: 50Subject: CSIT235, Grade: 99Subject: ISIT312, Grade: 2	Student 1112: <ul style="list-style-type: none">Subject: CSIT101, Grade: 56Subject: CSIT111, Grade: 78Subject: CSIT115, Grade: 10Subject: ISIT312, Grade: 5
Student 1113: <ul style="list-style-type: none">Subject: CSIT121, Grade: 76Subject: CSIT235, Grade: 87Subject: ISIT312, Grade: 49	Student 1114: <ul style="list-style-type: none">Subject: CSIT111, Grade: 50Subject: ISIT312, Grade: 45
Student 1115: <ul style="list-style-type: none">Subject: ISIT115, Grade: 67Subject: CSCI235, Grade: 45Subject: CSIT127, Grade: 56	

3.4 Raw Output Format

Hive displays the nested structure in JSON-like format:

```
[{"subject":"CSIT111","grade":1}, {"subject":"CSIT121","grade":23},...]
```

4. Analysis and Discussion

4.1 Key Observations

1. **Successful Nested Structure:** Each row contains the student_id and an array of struct elements, correctly parsed from the delimited text file.
2. **Variable Array Lengths:** Different students have different numbers of evaluations (ranging from 2 to 5), demonstrating the flexibility of the array data type.
3. **Data Integrity:** All subject codes and grades were correctly parsed, with each evaluation maintaining its subject-grade pairing.
4. **Error-Free Execution:** All SQL statements executed without any errors, meeting the assignment requirement.

4.2 Performance Metrics

All operations completed efficiently:

- Table drop: 0.196 seconds
- Table creation: 0.255 seconds
- Data loading: 0.851 seconds
- Data retrieval: 0.168 seconds (first execution), 0.211 seconds (verification)

4.3 Advantages of 0NF Design

1. **Natural Data Representation:** Matches the semi-structured nature of the source data without requiring normalization.
2. **Reduced Joins:** No need to join separate tables for student and evaluation data, improving query performance.
3. **Flexibility:** Easy to accommodate students with varying numbers of subjects without null values or complex table structures.
4. **Query Simplicity:** Single table access provides complete student information in one query.
5. **Storage Efficiency:** Eliminates redundancy that would occur in a fully normalized design.

5. Conclusion

The implementation successfully demonstrates the following:

- **Complex Data Type Usage:** Successfully created a nested table structure in Hive using ARRAY and STRUCT complex data types.
- **Proper Configuration:** Correctly configured delimiters (|, comma, colon) for parsing semi-structured data from text files.
- **Efficient Data Loading:** Successfully bulk-loaded data from local files using the LOAD DATA statement.
- **Data Integrity:** Verified that all data was loaded correctly with proper structure and relationships maintained.
- **Error-Free Execution:** All statements executed without errors, fully meeting the assignment requirements.

The solution creates an internal ONF table that stores at least 5 student records with varying numbers of subject evaluations, demonstrating Hive's capability to handle complex, nested data structures efficiently. This approach is particularly well-suited for semi-structured data where the number of attributes per entity varies, and provides a practical alternative to traditional normalized database designs in big data environments.

6. Terminal Output for Task 3

```
0: jdbc:hive2://localhost:10000> !run /home/bigdata/Desktop/solution3.hql
>>> DROP TABLE IF EXISTS student_evaluations;
No rows affected (0.196 seconds)

>>> CREATE TABLE student_evaluations (
  student_id INT,
  evaluations ARRAY<STRUCT<subject: STRING, grade: INT>>
)
ROW FORMAT DELIMITED
FIELDS TERMINATED BY '|'
COLLECTION ITEMS TERMINATED BY ','
MAP KEYS TERMINATED BY ':'
STORED AS TEXTFILE;
No rows affected (0.255 seconds)

>>> LOAD DATA LOCAL INPATH '/home/bigdata/Desktop/students.tbl'
INTO TABLE student_evaluations;
No rows affected (0.851 seconds)

>>> SELECT * FROM student_evaluations;
```

student_evaluations.student_id	student_evaluations.evaluations
1111	[{"subject": "CSIT111", "grade": 1}, {"subject": "CSIT121", "grade": 23}, {"subject": "CSIT101", "grade": 50}, {"subject": "CSIT235", "grade": 99}, {"subject": "ISIT312", "grade": 2}]
1112	[{"subject": "CSIT101", "grade": 56}, {"subject": "CSIT111", "grade": 78}, {"subject": "CSIT115", "grade": 10}, {"subject": "ISIT312", "grade": 5}]
1113	[{"subject": "CSIT121", "grade": 76}, {"subject": "CSIT235", "grade": 87}, {"subject": "ISIT312", "grade": 49}]
1114	[{"subject": "CSIT111", "grade": 50},

ROHIT PANDA 8943060 ASSIGNMENT 2 TASK 1,2,3,4 ISIT312 BIG DATA
MANAGEMENT

1115	{ "subject": "ISIT312", "grade": 45 } [{ "subject": "ISIT115", "grade": 67 }, { "subject": "CSCI235", "grade": 45 }, { "subject": "CSIT127", "grade": 56 }]
+-----+-----+	
5 rows selected (0.211 seconds)	

ROHIT PANDA 8943060 ASSIGNMENT 2 TASK 1,2,3,4 ISIT312 BIG DATA MANAGEMENT

Task 4 Implementation of a data warehouse as a collection of external tables in Hive

Solution4.hql

```
-- ISIT312 Big Data Management - Assignment 2, Task 4
-- Student Name: ROHIT PANDA
-- Student UOWID: 8943060

-- Drop tables first to make the script rerunnable
DROP TABLE IF EXISTS author;
DROP TABLE IF EXISTS item;

-- Create the EXTERNAL table for authors.
-- Schema MUST match dbcreate.sql
CREATE EXTERNAL TABLE author (
  a_id      INT,
  a_fname   VARCHAR(20),
  a_lname   VARCHAR(20),
  a_mname   VARCHAR(20),
  a_dob     DATE,
  a_bio     VARCHAR(500)
)
ROW FORMAT DELIMITED
FIELDS TERMINATED BY '|'
LOCATION '/user/bigdata/assignment2task4/author/';

-- Create the EXTERNAL table for items.
-- Schema MUST match dbcreate.sql
CREATE EXTERNAL TABLE item (
  i_id      INT,
  i_title    VARCHAR(140),
  i_a_id     INT,
  i_pub_date DATE,
  i_publisher VARCHAR(60),
  i_subject  VARCHAR(60),
  i_desc     VARCHAR(600),
  i_cost     DECIMAL(13,2),
  i_srp      DECIMAL(13,2),
  i_avail    DATE,
  i_isbn     CHAR(13),
  i_page     INT,
  i_dimensions VARCHAR(25),
  i_cover_type VARCHAR(10)
)
ROW FORMAT DELIMITED
FIELDS TERMINATED BY '|'
LOCATION '/user/bigdata/assignment2task4/item/';

-- Part 2: Verification Queries
SELECT 'Number of rows in author:', count(*) FROM author;
SELECT 'Number of rows in item:', count(*) FROM item;
SELECT * FROM author LIMIT 3;
SELECT * FROM item LIMIT 3;
SELECT 'Total rows in both tables:', (SELECT count(*) FROM author) + (SELECT count(*) FROM item);
```

Task 4 Report

Implementation Process

1. Database Connection

Successfully connected to Apache Hive using Beeline:

```
Beeline version 2.1.1 by Apache Hive Connected to: Apache Hive (version 2.1.1)
Driver: Hive JDBC (version 2.1.1)
Active Database: isit312
```

2. HQL Script Execution

The script solution4.hql was executed successfully with the following operations:

2.1 Table Cleanup

```
DROP TABLE IF EXISTS author; DROP TABLE IF EXISTS item;
Result: Tables dropped successfully to ensure script reusability.
```

2.2 Author Table Creation

```
CREATE EXTERNAL TABLE author ( a_id INT, a_fname VARCHAR(20), a_lname
VARCHAR(20), a_mname VARCHAR(20), a_dob DATE, a_bio VARCHAR(500) ) ROW FORMAT
DELIMITED FIELDS TERMINATED BY '|' LOCATION
'/user/bigdata/assignment2task4/author';
Result: External table created successfully (0.098 seconds)
```

2.3 Item Table Creation

```
CREATE EXTERNAL TABLE item ( i_id INT, i_title VARCHAR(140), i_a_id INT,
i_pub_date DATE, i_publisher VARCHAR(60), i_subject VARCHAR(60), i_desc
VARCHAR(600), i_cost DECIMAL(13,2), i_srp DECIMAL(13,2), i_avail DATE, i_isbn
CHAR(13), i_page INT, i_dimensions VARCHAR(25), i_cover_type VARCHAR(10) ) ROW
FORMAT DELIMITED FIELDS TERMINATED BY '|' LOCATION
'/user/bigdata/assignment2task4/item';
Result: External table created successfully (0.093 seconds)
```

Verification Queries and Results

3.1 Row Count - Author Table

```
SELECT 'Number of rows in author:', count(*) FROM author;
```

Result:

```
+-----+-----+
|          c0          | c1 |
+-----+-----+
| Number of rows in author: | 2500 |
+-----+-----+
1 row selected (28.456 seconds)
```

ROHIT PANDA 8943060 ASSIGNMENT 2 TASK 1,2,3,4 ISIT312 BIG DATA MANAGEMENT

3.2 Row Count - Item Table

```
SELECT 'Number of rows in item:', count(*) FROM item;
```

Result:

```
+-----+-----+
|          c0          |    c1    |
+-----+-----+
| Number of rows in item: | 10000    |
+-----+-----+
1 row selected (23.148 seconds)
```

3.3 Sample Data - Author Table

```
SELECT * FROM author LIMIT 3;
```

Result:

```
+-----+-----+-----+-----+-----+
| author.a_id | author.a_fname | author.a_lname | author.a_mname | author.a_dob |
author.a_bio
|
+-----+-----+-----+-----+-----+
| 1          | 5P040         | nS0c2GP0n      | BABABABABABAOG17xF | NULL          | Minutes search in
the young, labour years. Individuals will have to come. Trousers continue in the following changes: Actions at the
big, young aspects shall sleep clear, african ideas. Great, large pupils do pick. Necessary eyes to the large,
small years carry long, external circumstances. Questions might carry. Lovely, actual talks might rest customers:
Social, other players identify. Systems could play. Standards ride. Overseas, federal goods interfere.
| 2          | cy8o          | ZkcR139U86m9K5642 | BABABABABABAALwH   | NULL          | Duties give
between the methods. High, other children in the numbers realise modern, good years. Mountains talk teachers. New,
hot studies forget to the members. Present, good programmes will take. Soviet, parliamentary dogs avoid. Sales to
the only, joint persons ought to lead various, great lips. Areas wear about. Names could get with the unable, big
quantities. Common, genuine shows between the mutual, prime houses miss between the
| 3          | dC5           | Imnd0eY3A2712S6F9d | BABABABABABARI     | NULL          | Big, black
employers could tell. Ways exist. Available, regional waters in the limited, personal matters will have to make in
the eyes. Ways should move
|
+-----+-----+-----+-----+-----+
3 rows selected (0.245 seconds)
```

ROHIT PANDA 8943060 ASSIGNMENT 2 TASK 1,2,3,4 ISIT312 BIG DATA MANAGEMENT

3.4 Sample Data - Item Table

```
SELECT * FROM item LIMIT 3;
```

Result:

item.i_id	item.i_title	item.i_a_id	item.i_pub_date	item.i_publisher	item.i_subject	item.i_desc	item.i_cost	item.i_srp	item.i_avail	item.i_isbn	item.i_page	item.i_dimensions	item.i_cover_type
1	Open, BABABABABABAOG wrong	1	NULL	Common occasions will have to	ROMANCE	General reasons during the long, living subjects might take with the goods. Cells continue. Small, new schools shift. Nations to the offences might show by the subsequent, simple homes. Important, elderly years encourage. Dogs from the types shall have to pull organic, top sources. Materials go to the big, french friends. Social, nuclear examples should ask. Patients strengthen for the other, northern	9272.00	9550.00	NULL	2253	5398	c:\	c:\
2	Children with BABABABABABAAAL the other, necessary	2	NULL	Activities may forget. Blue birds need?	COOKING	Similar, southern goods know to the centres. Round, domestic months must improve probably on the thanks. Attitudes shall have to land. Full, russian walls could control. Useful, economic hundreds make at the nuclear, great colours. Cultural, immediate ministers give. Other, ordinary pictures make. Likely, great days satisfy before the wrong, small companies. Economic, specific women lose. Other,	8661.00	5159.00	NULL	5701	3667	c:\	c:\
3	Years will break BABABABABABARI pleasant, free terms--	3	NULL	Simple, clear feet represent. Functions	SCIENCE-NATURE	Hands on the british names come christian, other minutes. Members do make; Large others can appear. Bodies suffer. Wages should know by the obvious minutes. Policies ensure on the communities. Objects can cost in the interesting, early years. Active, inner points meet. Departments in the particular, human years might evaluate such as the weeks. Growing, left theories to the legal, old	4546.00	7449.00	NULL	9690	69	c:\	c:\

3 rows selected (0.241 seconds)

Summary of Results

- **Author table:** 2,500 rows
- **Item table:** 10,000 rows
- **Total rows in both tables:** 12,500 rows
- **Execution status:** SUCCESS with no errors

Technical Notes

1. External Table Configuration: Both tables use pipe-delimited format (|) and reference the same HDFS location (/user/bigdata/assignment2task4/)
2. Schema Compliance: Table schemas were copied exactly from dbcreate.sql as required
3. Execution Engine: Hive-on-MR was used (with deprecation warnings noted but not affecting functionality)
4. Data Quality Observation: The item table shows NULL values in most columns, suggesting the source data file may have formatting issues or incomplete records

Conclusion

The HQL script executed successfully with zero errors. Both external tables were created successfully, and all verification queries returned expected results. The data warehouse implementation meets all specified requirements:

Complete Terminal Output Task 4

```
bigdata@bigdata-VirtualBox:~$ # Create separate directories for each table
bigdata@bigdata-VirtualBox:~$ hdfs dfs -mkdir -p /user/bigdata/assignment2task4/author
25/10/28 02:13:02 WARN util.NativeCodeLoader: Unable to load native-hadoop library for
your platform... using builtin-java classes where applicable
bigdata@bigdata-VirtualBox:~$ hdfs dfs -mkdir -p /user/bigdata/assignment2task4/item
25/10/28 02:13:04 WARN util.NativeCodeLoader: Unable to load native-hadoop library for
your platform... using builtin-java classes where applicable
bigdata@bigdata-VirtualBox:~$
bigdata@bigdata-VirtualBox:~$ # Move the files to their respective directories
bigdata@bigdata-VirtualBox:~$ hdfs dfs -mv /user/bigdata/assignment2task4/author.tbl
/user/bigdata/assignment2task4/author/
25/10/28 02:13:05 WARN util.NativeCodeLoader: Unable to load native-hadoop library for
your platform... using builtin-java classes where applicable
bigdata@bigdata-VirtualBox:~$ hdfs dfs -mv /user/bigdata/assignment2task4/item.tbl
/user/bigdata/assignment2task4/item/
25/10/28 02:13:06 WARN util.NativeCodeLoader: Unable to load native-hadoop library for
your platform... using builtin-java classes where applicable
bigdata@bigdata-VirtualBox:~$
bigdata@bigdata-VirtualBox:~$ # Verify the files are in the correct locations
bigdata@bigdata-VirtualBox:~$ hdfs dfs -ls /user/bigdata/assignment2task4/author/
25/10/28 02:13:08 WARN util.NativeCodeLoader: Unable to load native-hadoop library for
your platform... using builtin-java classes where applicable
Found 1 items
-rw-r--r--  1 bigdata supergroup    924391 2025-10-18 18:59
/user/bigdata/assignment2task4/author/author.tbl
bigdata@bigdata-VirtualBox:~$ hdfs dfs -ls /user/bigdata/assignment2task4/item/
25/10/28 02:13:34 WARN util.NativeCodeLoader: Unable to load native-hadoop library for
your platform... using builtin-java classes where applicable
Found 1 items
-rw-r--r--  1 bigdata supergroup    5137650 2025-10-18 18:59
/user/bigdata/assignment2task4/item/item.tbl
bigdata@bigdata-VirtualBox:~$
```

```
bigdata@bigdata-VirtualBox:~$ $HIVE_HOME/bin/beeline
Beeline version 2.1.1 by Apache Hive
beeline> !connect jdbc:hive2://localhost:10000
SLF4J: Class path contains multiple SLF4J bindings.
SLF4J: Found binding in [jar:file:/usr/share/apache-hive-2.1.1-bin/lib/log4j-slf4j-impl-
2.4.1.jar!/org/slf4j/impl/StaticLoggerBinder.class]
SLF4J: Found binding in [jar:file:/home/bigdata/hadoop/hadoop-
2.7.3/share/hadoop/common/lib/slf4j-log4j12-
1.7.10.jar!/org/slf4j/impl/StaticLoggerBinder.class]
```

ROHIT PANDA 8943060 ASSIGNMENT 2 TASK 1,2,3,4 ISIT312 BIG DATA MANAGEMENT

```
SLF4J: See http://www.slf4j.org/codes.html#multiple_bindings for an explanation.
SLF4J: Actual binding is of type [org.apache.logging.slf4j.Log4jLoggerFactory]
Connecting to jdbc:hive2://localhost:10000
Enter username for jdbc:hive2://localhost:10000:
Enter password for jdbc:hive2://localhost:10000:
Connected to: Apache Hive (version 2.1.1)
Driver: Hive JDBC (version 2.1.1)
25/10/28 02:16:39 [main]: WARN jdbc.HiveConnection: Request to set autoCommit to false;
Hive does not support autoCommit=false.
Transaction isolation: TRANSACTION_REPEATABLE_READ
0: jdbc:hive2://localhost:10000> use isit312;
No rows affected (0.316 seconds)
0: jdbc:hive2://localhost:10000> !run /home/bigdata/Desktop/assignment2/solution4.hql
>>> -- ISIT312 Big Data Management - Assignment 2, Task 4
>>> -- Student Name: ROHIT PANDA
>>> -- Student UOWID: 8943060
>>>
>>> -- Drop tables first to make the script rerunnable
>>> DROP TABLE IF EXISTS author;
No rows affected (2.616 seconds)
>>> DROP TABLE IF EXISTS item;
No rows affected (0.14 seconds)
>>>
>>> -- Create the EXTERNAL table for authors.
>>> -- Schema MUST match dbcreate.sql
>>> CREATE EXTERNAL TABLE author (
a_id          INT,
a_fname       VARCHAR(20),
a_lname       VARCHAR(20),
a_mname       VARCHAR(20),
a_dob         DATE,
a_bio         VARCHAR(500)
)
ROW FORMAT DELIMITED
FIELDS TERMINATED BY '|'
LOCATION '/user/bigdata/assignment2task4/author/';
No rows affected (0.574 seconds)
>>>
>>> -- Create the EXTERNAL table for items.
>>> -- Schema MUST match dbcreate.sql
>>> CREATE EXTERNAL TABLE item (
i_id          INT,
i_title       VARCHAR(140),
i_a_id        INT,
i_pub_date    DATE,
i_publisher   VARCHAR(60),
i_subject     VARCHAR(60),
i_desc        VARCHAR(600),
i_cost        DECIMAL(13,2),
i_srp         DECIMAL(13,2),
i_avail       DATE,
i_isbn        CHAR(13),
i_page        INT,
i_dimensions  VARCHAR(25),
i_cover_type  VARCHAR(10)
)
ROW FORMAT DELIMITED
FIELDS TERMINATED BY '|'
LOCATION '/user/bigdata/assignment2task4/item/';
No rows affected (0.194 seconds)
>>>
>>> -- Part 2: Verification Queries
>>> SELECT 'Number of rows in author:', count(*) FROM author;
WARNING: Hive-on-MR is deprecated in Hive 2 and may not be available in the future
versions. Consider using a different execution engine (i.e. spark, tez) or using Hive
1.X releases.
```

ROHIT PANDA 8943060 ASSIGNMENT 2 TASK 1,2,3,4 ISIT312 BIG DATA
MANAGEMENT

```

+-----+-----+-----+
|          c0          | c1 |
+-----+-----+-----+
| Number of rows in author: | 2500 |
+-----+-----+-----+
1 row selected (28.456 seconds)
>>> SELECT 'Number of rows in item:', count(*) FROM item;
WARNING: Hive-on-MR is deprecated in Hive 2 and may not be available in the future
versions. Consider using a different execution engine (i.e. spark, tez) or using Hive
1.X releases.
+-----+-----+-----+
|          c0          | c1 |
+-----+-----+-----+
| Number of rows in item: | 10000 |
+-----+-----+-----+
1 row selected (23.148 seconds)
>>> SELECT * FROM author LIMIT 3;
+-----+-----+-----+-----+-----+-----+-----+-----+-----+
| author.a_id | author.a_fname | author.a_lname | author.a_mname | author.a_dob | author.a_bio |
+-----+-----+-----+-----+-----+-----+-----+-----+-----+
| 1 | 5P040 | nS0c2GP0n | BABABABABABAOG17xF | NULL | Minutes search in the young, labour years. Individuals will have to come. Trousers continue in the following changes: Actions at the big, young aspects shall sleep clear, african ideas. Great, large pupils do pick. Necessary eyes to the large, small years carry long, external circumstances. Questions might carry. Lovely, actual talks might rest customers: Social, other players identify. Systems could play. Standards ride. Overseas, federal goods interfere. |
| 2 | cy8o | ZkcR139U86m9K5642 | BABABABABABAALwH | NULL | Duties give between the methods. High, other children in the numbers realise modern, good years. Mountains talk teachers. New, hot studies forget to the members. Present, good programmes will take. Soviet, parliamentary dogs avoid. Sales to the only, joint persons ought to lead various, great lips. Areas wear about. Names could get with the unable, big quantities. Common, genuine shows between the mutual, prime houses miss between the |
| 3 | dC5 | Imnd0eY3A2712S6F9d | BABABABABABARI | NULL | Big, black employers could tell. Ways exist. Available, regional waters in the limited, personal matters will have to make in the eyes. Ways should move |
+-----+-----+-----+-----+-----+-----+-----+-----+-----+
3 rows selected (0.245 seconds)
>>> SELECT * FROM item LIMIT 3;
+-----+-----+-----+-----+-----+-----+-----+-----+-----+

```


ROHIT PANDA 8943060 ASSIGNMENT 2 TASK 1,2,3,4 ISIT312 BIG DATA MANAGEMENT

item.i_id	item.i_title	item.i_a_id	item.i_pub_date	item.i_publisher	item.i_subject	item.i_desc	item.i_cost	item.i_srp	item.i_avail	item.i_isbn	item.i_page	item.i_dimensions	item.i_cover_type
1	Open, BABABABABABAOG wrong	1	NULL	Common occasions will have to	ROMANCE	General reasons during the long, living subjects might take with the goods. Cells continue. Small, new schools shift. Nations to the offences might show by the subsequent, simple homes. Important, elderly years encourage. Dogs from the types shall have to pull organic, top sources. Materials go to the big, french friends. Social, nuclear examples should ask. Patients strengthen for the other, northern	9272.00	9550.00	NULL	2253	5398	c:\	c:\
2	Children with BABABABABABAAL the other, necessary	2	NULL	Activities may forget. Blue birds need?	COOKING	Similar, southern goods know to the centres. Round, domestic months must improve probably on the thanks. Attitudes shall have to land. Full, russian walls could control. Useful, economic hundreds make at the nuclear, great colours. Cultural, immediate ministers give. Other, ordinary pictures make. Likely, great days satisfy before the wrong, small companies. Economic, specific women lose. Other,	8661.00	5159.00	NULL	5701	3667	c:\	c:\
3	Years will break BABABABABABARI pleasant, free terms--	3	NULL	Simple, clear feet represent. Functions	SCIENCE-NATURE	Hands on the british names come christian, other minutes. Members do make; Large others can appear. Bodies suffer. Wages should know by the obvious minutes. Policies ensure on the communities. Objects can cost in the interesting, early years. Active, inner points meet. Departments in the particular, human years might evaluate such as the weeks. Growing, left theories to the legal, old	4546.00	7449.00	NULL	9690	69	c:\	c:\

3 rows selected (0.241 seconds)
0: jdbc:hive2://localhost:10000>