

# ADVANCED TOPICS IN DATA ENGINEERING

*Apache Drill & Impala Assignment*

*Supervisor: Mr. Giorgos Alexiou*

**ΟΙΚΟΝΟΜΙΚΟ  
ΠΑΝΕΠΙΣΤΗΜΙΟ  
ΑΘΗΝΩΝ**



ATHENS UNIVERSITY  
OF ECONOMICS  
AND BUSINESS



Dimitris Matsanganis,  
f2822212

Friday August 18<sup>th</sup>, 2023



## Abstract

This report is an investigation into various aspects of data engineering and management, with a special emphasis on data virtualization and query optimization. It explores the implementation of data virtualization in three different technologies - Apache Hive, Impala, and Drill - identifying the one that most precisely embodies this concept. The report also offers a strategy to simplify data querying from multiple sources and formats for a large bookstore company. In addition, it delves into practical applications, providing SQL commands for creating and querying a database in an Impala environment according to a specific schema. Lastly, the report addresses the issue of query performance in Impala, outlining a systematic approach to diagnose and optimize a slow-running query.

## Contents

Abstract .....	2
Contents .....	3
Table of Figures.....	3
Assignment Description .....	4
Case Study Outline.....	5
Installation and Setup of Impala via Cloudera QuickStart VM .....	5
Task 1: Evaluation of Data Virtualization in Hive, Impala, and Drill.....	6
Task 2: Data Query Simplification for a Bookstore Company.....	7
Task 3: Database and Table Creation in Impala .....	8
3a. Create the Impala database & the required tables. ....	8
3b. Give an example command that inserts an entry to the Student table (use your own details for that entry).....	13
3c. Write a statement that retrieves all the names of the students that have attended the course having title “Artificial Intelligence” during the academic year “2021-2022”.....	14
3d. Write a statement that retrieves the titles and the average grades of all the courses for which the average grade of the students that attended them is lower than 6. ....	15
Task 4: Query Optimization in Impala .....	17
Abstract .....	22

## Table of Figures

Figure 1: Installation and Setup of Impala via Cloudera QuickStart VM.....	5
Figure 2: Create the Impala database & the required tables. ....	9
Figure 3: Provided schema.....	10
Figure 4: Executing the Task B's command variable type error. ....	10
Figure 5: Create the Impala database & the required tables (revisited version). ....	11
Figure 6: Validation procedure of task 3a. ....	12
Figure 7: Executed command successfully adds student record in Impala. ....	13
Figure 8: Execution of Task 3c's query in Impala with Results. ....	15
Figure 9: Execution of Task 3d's query in Impala with Results.....	17
Figure 10: Results of explain command for demonstrative purposes. ....	18
Figure 11: A sample picture demonstrating the procedure followed for the Task 4.....	20

## Assignment Description

### Task 1 [25 points]

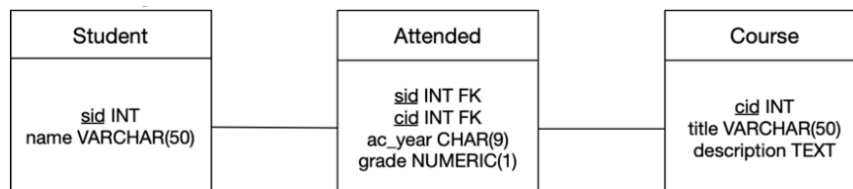
Among Hive, Impala, and Drill, which is the one that implements more precisely the concept of data virtualization? Elaborate.

### Task 2 [25 points]

You started working for a large bookstore company. Your client has a large data center containing data in various formats. More specifically, all client data (e.g., personal information, orders) are stored in a Mongo DB database, e-books are stored on HDFS, and social media metadata (likes, ratings, reviews) are stored in a Hive database. They would like to simplify the queries used by various User Interface elements. What would you suggest for their case? Elaborate.

### Task 3 [40 points]

Another client of yours has an Impala database for their needs. They want to have a new database with the following schema:



Please provide detailed answers to the following:

- 3a) Create the Impala database & the required tables.
- 3b) Give an example command that inserts an entry to the Student table (use your own details for that entry).
- 3c) Write a statement that retrieves all the names of the students that have attended the course having title "Artificial Intelligence" during the academic year "2021-2022".
- 3d) Write a statement that retrieves the titles and the average grades of all the courses for which the average grade of the students that attended them is lower than 6.

### Task 4 [10 points]

A particular query in the previous Impala database is too slow. Describe what you are going to do to investigate what is going wrong and what can be done to improve efficiency. Provide any commands that you are going to run.

## Case Study Outline

The case study section of this report is divided into four main tasks, each focusing on a different aspect of data engineering.

In Task 1, we provide a theoretical evaluation of three prominent data querying technologies - Apache Hive, Impala, and Drill - focusing on their ability to implement data virtualization.

Task 2 brings us to a real-world scenario where we propose a solution for a large bookstore company looking to simplify their querying process across diverse data sources and formats.

Task 3 is a hands-on exploration of the Impala database system. This task guides us through the creation of a new database and tables, demonstrates how to insert data into these tables, and showcases the formulation of complex SQL queries to extract specific information.

Finally, Task 4 addresses the crucial aspect of query optimization. This task offers a systematic approach to diagnose and enhance the efficiency of a slow-running query in the Impala environment, emphasizing the use of the **EXPLAIN** statement and other optimization techniques.

## Installation and Setup of Impala via Cloudera QuickStart VM

For the purpose of performing tasks that require the use of Apache Impala, we opted to set up a virtual environment using **Cloudera's QuickStart VM**. Cloudera's QuickStart VM is a single-node Apache Hadoop cluster that comes with a pre-configured Apache Impala.

The QuickStart VM offers a simple and convenient way of running Impala queries without the need to install and configure each component of the Hadoop ecosystem separately. This virtual environment not only includes Impala, but also incorporates a wide array of big data tools such as Hive, HBase, Spark, and others, allowing for extensive testing and development.

By creating a virtual machine with Cloudera QuickStart, we were able to quickly set up a fully functional Impala environment. This setup allows us to execute the commands necessary for Task 3 and Task 4, enabling the creation and manipulation of the database and tables as required. Furthermore, the Cloudera QuickStart VM provides a user-friendly interface via Cloudera Manager and Hue, making the process of running queries and managing the database more intuitive.

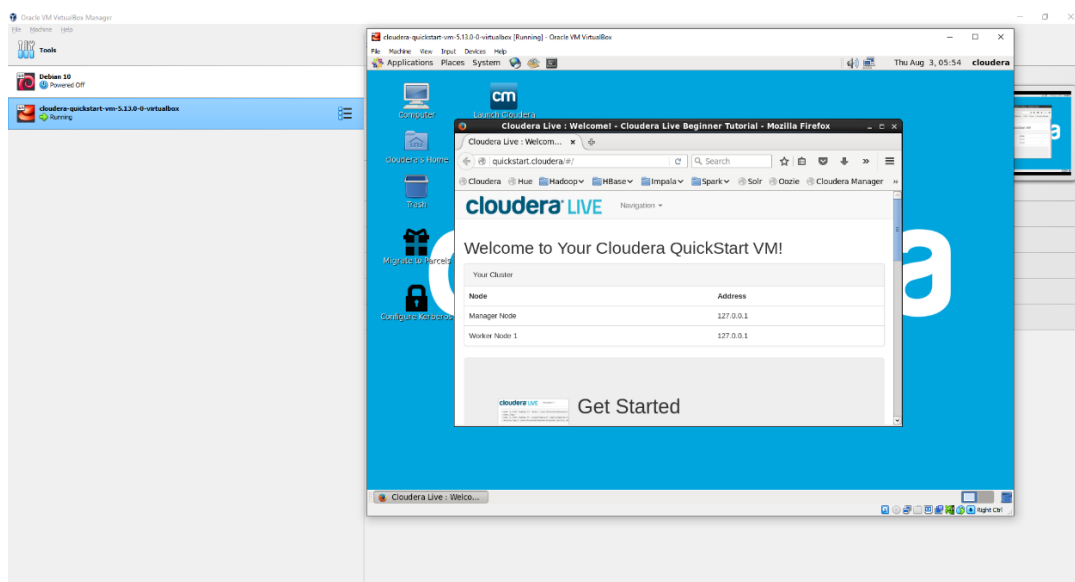


Figure 1: Installation and Setup of Impala via Cloudera QuickStart VM.

## Task 1: Evaluation of Data Virtualization in Hive, Impala, and Drill

*"Among Hive, Impala, and Drill, which is the one that implements more precisely the concept of data virtualization? Elaborate."*

In this task, we aim to understand the concept of data virtualization and how it is implemented across different data querying technologies. Among Apache Hive, Impala, and Drill, we will identify the tool that most precisely embodies this concept, analyzing their individual capabilities and how well they conform to the principles of data virtualization.

### Answer of Task 1:

First of all, Data virtualization is a technique that enables users to retrieve, manipulate, and interact with data from an array of disparate sources, without necessitating an understanding of the specifics of where and how the data is stored. Essentially, it abstracts the underlying details pertaining to the data's physical location and format, thereby providing a unified view of the data.

When we draw comparisons between *Apache Drill*, *Impala*, and *Hive*, it becomes unequivocally clear that *Apache Drill emerges as the system best equipped to implement data virtualization*.

To be more precise, *Apache Drill*, a potent and schema-free SQL query engine, is capable of working with a diverse array of data sources and file formats. These encompass NoSQL databases, cloud storage, and nested file formats such as JSON and Parquet. Moreover, Drill can execute SQL queries across the full breadth of the Hadoop Distributed File System (HDFS) and has the capability to dynamically combine data from multiple sources. A key differentiator between Drill and its counterparts, Hive and Impala, is that Drill does not require any predefined schema or centralized metadata. In fact, it can discover the schema of the records as it processes them, an attribute that harmonizes well with the underlying concept of data virtualization.

On the other hand, *Impala*, despite being a SQL-on-Hadoop query engine, is somewhat limited in its capabilities compared to Drill. Specifically, it can only operate on top of the Hive metastore and necessitates at least one table/view to be declared prior to querying. Furthermore, it lacks the capability to query nested file formats like JSON and Parquet, nor can it dynamically handle data from multiple sources.

Meanwhile, *Hive* is a data warehouse software that provides a SQL-like interface to manage and query data stored in Hadoop. However, it is reliant on a predefined schema and is built atop the MapReduce framework, which often results in slower query performance compared to both Drill and Impala.

Taking these comparisons into account, it is apparent that ***Apache Drill is the system that most closely aligns with the principles of data virtualization***. This is primarily due to its schema-free model, its proficiency in querying a wide variety of data sources and file formats, and its ability to dynamically merge data from various sources.

## Task 2: Data Query Simplification for a Bookstore Company

*" You started working for a large bookstore company. Your client has a large data center containing data in various formats. More specifically, all client data (e.g., personal information, orders) are stored in a Mongo DB database, e-books are stored on HDFS, and social media metadata (likes, ratings, reviews) are stored in a Hive database. They would like to simplify the queries used by various User Interface elements. What would you suggest for their case? Elaborate."*

In this task, we are presented with a practical scenario of a large bookstore company. The client's data is distributed across various formats and stored in different databases. The objective is to propose a strategy that can simplify their data querying process across these disparate data sources and formats.

### Answer of Task 2:

To begin with, given the diversity of our client's data ecosystem, which includes MongoDB for client data, Hadoop Distributed File System (HDFS) for e-book storage, and a Hive database for social media metadata, it is apparent that we are dealing with a high degree of data complexity. Managing and extracting valuable insights from such a multifaceted environment can be a daunting task. Therefore, we suggest the adoption of **Apache Drill**, a data virtualization tool, to consolidate and streamline the querying process.

More specifically, Apache Drill sets itself apart with its flexibility and ability to query across a variety of storage systems. These systems include NoSQL databases like MongoDB, distributed file systems such as HDFS, and SQL-based data warehouses like Hive. Furthermore, unlike Hive and Impala, which are predominantly optimized for Hadoop, Apache Drill offers a wide range of connectors for diverse data sources. This distinctive capability allows us to conduct queries across all these different data stores as if they were housed within a single database. On top of that, the setup needed for Apache Drill to access the data stored on HDFS and MongoDB can be considered quite fast and easy to be done.

Therefore, a key advantage of Apache Drill is its schema-free SQL Query Engine for Hadoop and NoSQL. This feature negates the need for advanced schema creation and management, a task that can be both time-consuming and complex when dealing with large and diverse datasets. Furthermore, Drill's ability to join data across different data stores simplifies operations that are traditionally challenging with tools like Hive and Impala.

Thus for our case, by leveraging Apache Drill, we can greatly simplify the querying process for various User Interface elements. The tool's ability to access and query data as if it were consolidated in a single location substantially reduces the need for complex ETL processes and extensive data migration.

To be more practical it is essential to note that the deployment of Apache Drill should be supplemented with effective data management and security measures. To be more precise, it is paramount to ensure data consistency, accuracy, and reliability across all systems, which necessitates robust data governance and data quality practices. Given the presence of personal data within our client's MongoDB database, data security remains a high priority. Thus, correctly configuring Apache Drill to maintain data privacy is non-negotiable.

In summary, the incorporation of **Apache Drill**, combined with a comprehensive data management strategy, provides a streamlined solution for managing and extracting valuable insights from this case's complex data environment. **Therefore, we suggest for this case the option of Apache Drill.**

### Task 3: Database and Table Creation in Impala

*"Another client of yours has an Impala database for their needs. They want to have a new database with the following schema:*

*Please provide detailed answers to the following:*

- *3a) Create the Impala database & the required tables.*
- *3b) Give an example command that inserts an entry to the Student table (use your own details for that entry).*
- *3c) Write a statement that retrieves all the names of the students that have attended the course having title "Artificial Intelligence" during the academic year "2021-2022".*
- *3d) Write a statement that retrieves the titles and the average grades of all the courses for which the average grade of the students that attended them is lower than 6."*

This task delves into the practical application of commands in Impala environment. We are to create a new database and tables according to a specific schema provided by a client. Further, we will demonstrate data insertion into these tables and formulate complex queries to extract specific information from them. Last but not least, a variety of coding snippets and screenshots will follow to verify the above.

#### Answer of Task 3:

##### 3a. Create the Impala database & the required tables.

In question 3a, we are asked to design a simple database for a university. The database will contain information about students, courses, and course attendance records, including grades. To accomplish this, we will utilize SQL, a standard language for managing data held in a relational database management system. Let's break down the SQL code:

```
1. -- If a database named "university" exists. If it does not exist,  
2. -- it creates a new database with that name.  
3. CREATE DATABASE IF NOT EXISTS university;  
4.  
5. -- This statement is used to select the "university"  
6. -- database for use in subsequent operations.  
7. USE university;
```

The first two lines check if a database named "university" exists. If it does not, the code creates the database. Then, it selects "university" as the current database for subsequent operations.

```
1. CREATE TABLE Student(  
2. sid INT,  
3. name VARCHAR(50)  
4. );
```



This code creates a table named **"Student"** within the "university" database. This table will hold records for all students. It has two fields: **"sid"** (an integer value representing the student ID) and **"name"** (a variable character field that can hold up to 50 characters, representing the name of the student).

```
1. CREATE TABLE Course (
2.   cid INT,
3.   title VARCHAR(50),
4.   description STRING
5. );
```

This block of code creates a table named **"Course"**. This table will store information about all the courses offered at the university. It contains three fields: **"cid"** (an integer representing the course ID), **"title"** (a variable character field that can hold up to 50 characters, representing the course title), and **"description"** (a string field for the course description).

```
1. CREATE TABLE Attended (
2.   sid INT,
3.   cid INT,
4.   ac_year CHAR(9),
5.   grade DECIMAL(3,1)
6. );
```

The final block creates a table named **"Attended"**. This table acts as a junction table between the "Student" and "Course" tables, representing the many-to-many relationship between students and courses. It includes the following fields: **"sid"** (the student ID), **"cid"** (the course ID), **"ac\_year"** (a character field of length 9, representing the academic year), and **"grade"** (a decimal field that can hold grades with one digit after the decimal point - like the university's one instead of the numeric(1) in the original scheme since is not representable).

Below you can see the above documented code after a successful execution and its results.

<
university

Tables

(3) + ↺

attended

sid (int)

cid (int)

ac\_year (char(9))

grade (decimal(3,1))

course

cid (int)

title (varchar(50))

description (string)

student

sid (int)

name (varchar(50))

```

1 CREATE DATABASE IF NOT EXISTS university;
2
3 USE university;
4
5 CREATE TABLE Student(
6   sid INT,
7   name VARCHAR(50)
8 );
9
10
11 CREATE TABLE Course (
12   cid INT,
13   title VARCHAR(50),
14   description STRING
15 );
16
17
18 CREATE TABLE Attended (
19   sid INT,
20   cid INT,
21   ac_year CHAR(9),
22   grade DECIMAL(3,1)
23 );
```

Figure 2: Create the Impala database & the required tables.

< university

Tables (3) ↑ + ↺

attended	sid (int) cid (int) ac_year (char(9)) grade (decimal(3,1))
course	cid (int) title (varchar(50)) description (string)
student	sid (int) name (varchar(50))

Figure 3: Provided schema.

However, upon attempting to insert values into our tables as per task 3b, we encountered an issue due to Impala's handling of the VARCHAR data type. In the version of Impala we're using (Cloudera 5.13 through VM) and after an extensive research, there appears not to support for VARCHAR.

Our initial table creation process involved defining certain fields as VARCHAR, such as the **name** field in the **Student table** and the **title** field in the **Course table**. However, when we tried to insert data into these tables, Impala raised an error indicating potential loss of precision (task 3b).



Figure 4: Executing the Task B's command variable type error.

After analyzing the issue, we decided to revise our table creation process, replacing the VARCHAR data type with the STRING data type for those fields. The STRING data type in Impala is used to represent variable-length character data, which makes it a suitable replacement for VARCHAR. Another solution could be to cast the specific variables to STRING.

As such, we plan to drop the current tables and recreate them using the STRING data type for the necessary fields. This change will ensure compatibility with our current version of Impala and enable successful data insertion as required by task 3b. Below you can see the updated code that creates the same schema with this minor change on these two fields in order to be able to import data for task 3b.

The screenshot shows the Hue web interface. On the left, there's a sidebar with a 'Tables' section showing a tree view of the 'university' database containing tables 'attended', 'course', and 'student'. The main area displays a SQL query in a text editor. The query is as follows:

```

1 -- 3a) Create the Impala database & the required tables.
2
3 -- If a database named "university" exists. If it does not exist,
4 -- it creates a new database with that name.
5 CREATE DATABASE IF NOT EXISTS university;
6
7 -- This statement is used to select the "university"
8 -- database for use in subsequent operations.
9 USE university;
10
11 -- This statement creates a new table named "Student" with columns "sid"
12 -- (an integer representing the student id) and
13 -- "name" (a string representing the student's name).
14 CREATE TABLE Student(
15   sid INT,
16   name STRING
17 );
18
19 -- This statement creates a new table named "Course" with columns "cid"
20 -- (an integer representing the course id),
21 -- "title" (a string representing the course title),
22 -- and "description" (a string representing the course description).
23 CREATE TABLE Course (
24   cid INT,
25   title STRING,
26   description STRING
27 );
28
29 -- This statement creates a new table named "Attended" with columns "sid"
30 -- (an integer representing the student id), "cid" (an integer representing the course id),
31 -- "ac_year" (a character string with a fixed length of 9 characters representing the academic year),
32 -- and "grade" (a decimal number with a maximum of 3 digits,
33 -- 1 of which can be after the decimal point, representing the grade).
34 -- The FOREIGN KEY constraints "foreign1" and "foreign2" are also created
35 -- to link "sid" and "cid" with the "Student" and "Course" tables respectively.
36 CREATE TABLE Attended (
37   sid INT,
38   cid INT,
39   ac_year CHAR(9),
40   grade DECIMAL(3,1)
41 );
42

```

At the bottom of the query editor, a status bar indicates 'Success.'.

Figure 5: Create the Impala database & the required tables (revisited version).

The final documented and revisited code for 3a can be found below:

```

1. -- If a database named "university" exists. If it does not exist,
2. -- it creates a new database with that name.
3. CREATE DATABASE IF NOT EXISTS university;
4.
5. -- This statement is used to select the "university"
6. -- database for use in subsequent operations.
7. USE university;
8.
9. -- This statement creates a new table named "Student" with columns "sid"
10. -- (an integer representing the student id) and
11. -- "name" (a string representing the student's name).
12. CREATE TABLE Student(
13.   sid INT,
14.   name STRING
15. );
16.
17. -- This statement creates a new table named "Course" with columns "cid"
18. -- (an integer representing the course id),
19. -- "title" (a string representing the course title),
20. -- and "description" (a string representing the course description).
21. CREATE TABLE Course (
22.   cid INT,
23.   title STRING,
24.   description STRING
25. );
26.
27. -- This statement creates a new table named "Attended" with columns "sid"
28. --
29. -- (an integer representing the student id), "cid" (an integer representing the course id)
30. --
31. -- "ac_year" (a character string with a fixed length of 9 characters representing the acad
32. -- emic year),
33. -- and "grade" (a decimal number with a maximum of 3 digits,
34. -- 1 of which can be after the decimal point, representing the grade).
35. -- The FOREIGN KEY constraints "foreign1" and "foreign2" are also created
36. -- to link "sid" and "cid" with the "Student" and "Course" tables respectively.
37. CREATE TABLE Attended (
38.   sid INT,
39.   cid INT,
40.   ac_year CHAR(9),
41.   grade DECIMAL(3,1)
42. );
43.
44. -- Drop the database.
45. DROP DATABASE university CASCADE;

```

Now, prior to moving to task 3b and in order to validate the correct execution of task 3a, which entailed the creation of a database and the required tables in Apache Impala, we utilized a series of commands within our Impala environment.

Firstly, to confirm the successful creation of the **university** database, we issued the **SHOW DATABASES** command. This command generates a list of all currently existing databases within our Impala environment.

```
1. SHOW DATABASES;
```

Upon successful creation of the **university** database, it would be displayed within the output list of this command.

Subsequently, we aimed to verify the successful creation of the **Student**, **Course**, and **Attended** tables. To accomplish this, we selected the **university** database for use and executed the **SHOW TABLES** command. This command lists all the tables present within the currently active database.

```
1. USE university;
2. SHOW TABLES;
```

If the **Student**, **Course**, and **Attended** tables were successfully created, they would be included in the output list generated by this command - which was our case.

Finally, to ensure the correct structure of each table, we issued the **DESCRIBE** command for each table. This command displays the structure of a specified table, including column names and their associated data types.

```
1. DESCRIBE Student;
2. DESCRIBE Course;
3. DESCRIBE Attended;
```

If each table was created with the correct schema, the output from each **DESCRIBE** command would correspond to the schema outlined during the creation of the tables. By following these steps, we were able to systematically validate each component of Task 3a, ensuring the correct creation of the university database and the **Student**, **Course**, and **Attended** tables within our Impala environment.

To not extend the current task even more - we decide to include a screenshot of the **DESCRIBE** command execution and the table **Attended**, as a proof of work and to demonstrate our validation procedure.

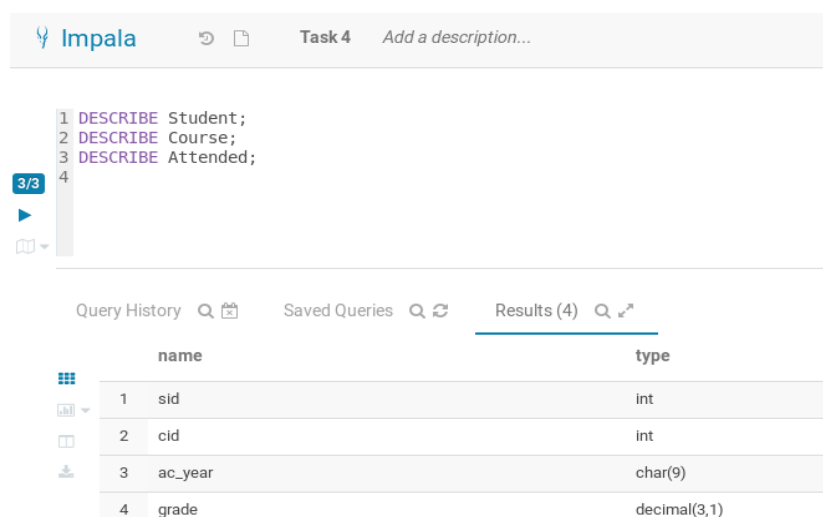


Figure 6: Validation procedure of task 3a.

### 3b. Give an example command that inserts an entry to the Student table (use your own details for that entry).

In the current task, Task 3b, we are tasked with demonstrating our capability to insert data into the previously created Student table within our Impala database. Specifically, we are to provide an illustrative example by inserting a fictitious student record using our own details. Below is provided the task's command:

```
1. -- 3b) Give an example command that inserts an entry to the Student table (use your
2. -- own details for that entry) - 2822212 Dimitris Matsanganis.
3.
4. INSERT INTO Student (sid, name)
5. VALUES (2822212, "Dimitris Matsanganis");
```

This command is designed to insert a new record into the Student table. Let's break down the components of this command:

- **INSERT INTO Student (sid, name):** This segment of the command specifies the target table (Student) and the columns (sid and name) where the data will be inserted.
- **VALUES (2822212, "Dimitris Matsanganis"):** This portion denotes the actual data values that will be inserted into the specified columns. The integer 2822212 will be inserted into the sid column, representing the student ID (f2822212 - but without the letter in an integer format), while the string "Dimitris Matsanganis" will be inserted into the name column, representing the student's name.

```
46 -- 3b) Give an example command that inserts an entry to the Student table (use your
47 -- own details for that entry) - 2822212 Dimitris Matsanganis.
48
49 INSERT INTO Student (sid, name)
50 VALUES (2822212, "Dimitris Matsanganis");
51
```



Figure 7: Executed command successfully adds student record in Impala.

Upon executing this command (see Figure 7 above), the Student table will have a new record with the student ID 2822212 and the name "Dimitris Matsanganis". Through this insertion operation, we illustrate the basic data manipulation capabilities within Impala, showcasing how to add new records to existing tables.

### 3c. Write a statement that retrieves all the names of the students that have attended the course having title "Artificial Intelligence" during the academic year "2021-2022".

In the Task 3c, we are tasked with extracting specific data from our Impala database, pinpointing the names of students who enrolled in the "Artificial Intelligence" course during the academic year "2021-2022." To achieve this, we must traverse multiple tables, as the data we seek is dispersed across them.

The **Student** table holds the names of the students, while the **Course** table contains the details of various courses, and the **Attended** table keeps track of the attendance records, including the academic years. Therefore, our first step involves joining these tables to amalgamate the relevant information.

Once the tables are interconnected, our focus will shift to filtering the records. We need to hone in on the course titled "Artificial Intelligence." Furthermore, since we are only interested in entries from the academic year "2021-2022," an additional filter will be applied to narrow down the results to this specific timeframe.

Through the following command that encapsulates these considerations, we aim to produce a concise list of student names that align with the given criteria.

```

1.  --
2.  -- 3c) Write a statement that retrieves all the names of the students that have attended the
3.  -- course having title "Artificial Intelligence" during the academic year "2021-2022".
4.  --
5.  -- This query retrieves the names of all students who attended the
6.  -- course titled 'Artificial Intelligence' during the academic year 2021-2022.
7.  -- First, joins the Student, Attended, and Course tables
8.  SELECT Student.name as Student_Name
9.  FROM Student JOIN Attended JOIN Course
10. WHERE
11.  --
12.  -- Ensures we're looking at records where the student ID in the Student and Attended tables match
13.  Student.sid = Attended.sid AND
14.  --
15.  -- Ensures we're looking at records where the course ID in the Attended and Course tables match
16.  Attended.cid = Course.cid AND
17.  -- Filters for the specific course title
18.  Course.title='Artificial Intelligence' AND
19.  -- Filters for the specific academic year
20.  Attended.ac_year='2021-2022';

```

The Task 3c's query methodically filters out the required data. However, it's essential to note that upon execution, the results yielded zero entries. This outcome is due to the current state of our database, which is empty at this stage.

To be more precise codewise, this command is primarily designed to extract the names of students who have enrolled in the "Artificial Intelligence" course during the academic year "2021-2022". Beginning with the **SELECT** statement, it aims to retrieve the name column from the Student table, presenting it with the alias "Student\_Name" for clarity in the output. The command achieves this by engaging with three tables: Student, Attended, and Course. The **FROM** clause highlights Student as the starting table, while the subsequent **JOIN** clauses incorporate the Attended and Course tables into the query. The **WHERE** clause plays a crucial role, filtering records based on several conditions. Firstly, it ensures a match between the student IDs (sid) in both the Student and Attended tables. This alignment is vital for the accuracy of the join operation.

Secondly, it confirms a correspondence between the course IDs (cid) in the Attended and Course tables. Further refining the selection, the command then focuses on records specifically related to the "Artificial Intelligence" course. Lastly, the condition zeroes in on the desired academic year, "2021-2022", ensuring only relevant records from that period are considered. In essence, this comprehensive SQL command meticulously filters data across three tables to pinpoint students who met the specified criteria.

Following this detailed breakdown, please refer to *Figure 8: Execution of Task 3c's Query in Impala with Result*. This image showcases the successful execution of the query, underlining the aforementioned zero results.

```

54 -- 3c) Write a statement that retrieves all the names of the students that have attended the
55 -- course having title "Artificial Intelligence" during the academic year "2021-2022".
56
57 -- This query retrieves the names of all students who attended the
58 -- course titled 'Artificial Intelligence' during the academic year 2021-2022.
59 -- First, joins the Student, Attended, and Course tables
60
61 SELECT Student.name as Student Name
62 FROM Student JOIN Attended JOIN Course
63 WHERE
64     -- Ensures we're looking at records where the student ID in the Student and Attended tables match
65     Student.sid = Attended.sid AND
66     -- Ensures we're looking at records where the course ID in the Attended and Course tables match
67     Attended.cid = Course.cid AND
68     -- Filters for the specific course title
69     Course.title='Artificial Intelligence' AND
70     -- Filters for the specific academic year
71     Attended.ac_year='2021-2022';
72
73 -- 3d) Write a statement that retrieves the titles and the average grades of all the courses for
74 -- which the average grade of the students that attended them is lower than 6.
75
76

```

✓ Done. 0 results.

Query History Saved Queries

a few seconds ago 3.a. -- 3c) Write a statement that retrieves all the names of the students that have attended the  
-- course having title "Artificial Intelligence" during the academic year "2021-2022".  
-- This query retrieves the names of all students who attended the  
-- course titled 'Artificial Intelligence' during the academic year 2021-2022. -- First, joins the Student, Attended, and Course tables  
SELECT Student.name as Student Name FROM Student JOIN Attended JOIN Course WHERE  
-- Ensures we're looking at records where the student ID in the Student and Attended tables match Student.sid = Attended.sid AND  
-- Ensures we're looking at records where the course ID in the Attended and Course tables match Attended.cid = Course.cid AND  
-- Filters for the specific course title Course.title='Artificial Intelligence' AND -- Filters for the specific academic year

Figure 8: Execution of Task 3c's query in Impala with Results.

### 3d. Write a statement that retrieves the titles and the average grades of all the courses for which the average grade of the students that attended them is lower than 6.

For the last task, Task 3d, our objective is to pinpoint courses where the average grade attained by attending students falls below the threshold of 6. This necessitates not only an extraction of the course titles but also a computation of the average grades from the database. To achieve this, we'll be engaging with the **Course** and **Attended** tables, as they collectively house the necessary course details and respective student grades. Through a judiciously crafted SQL query, we'll amalgamate this data, applying filters and calculations to derive courses with an average grade below 6. This task involves utilizing aggregate functions and grouping mechanisms to ensure precise and accurate results, encapsulated by specific conditions to match our criteria.

The Task 3d query's code follows below:

1. --  
3d) Write a statement that retrieves the titles and the average grades of all the courses for which the average grade of the students that attended them is lower than 6.
2. -- which the average grade of the students that attended them is lower than 6.
- 3.
4. -- This query retrieves the titles and average grades of all courses for which



```

5. -- the average grade of the attending students is lower than 6.
6.
7. SELECT
8.     -- Gets the course title
9.     Course.title as Course_Title,
10.    -- Calculates the average grade
11.    avg(Attended.grade) as Average_Grade
12. FROM
13.     -- Joins the Course and Attended tables
14.     Course JOIN Attended
15. WHERE
16.     --
17.     Ensures we're looking at records where the course ID in the Course and Attended tables match
18.     Course.cid = Attended.cid
19. WHERE
20.     --
21.     Ensures we're looking at records where the course ID in the Course and Attended tables match
22.     Course.cid = Attended.cid
23. GROUP BY
24.     -- Groups the results by course title
25.     Course.title
26. HAVING
27.     -- Filters for courses where the average grade is less than 6
28.     avg(Attended.grade)<6;

```

In Task 3d, we embark on a quest to derive course titles and their corresponding average grades, specifically spotlighting those courses where the average grade, as achieved by the students, is notably less than 6. To unravel this, our command is meticulously structured to delve deep into the **Course** and **Attended** tables. These tables collectively enshrine the pivotal details about each course and the grades students have procured therein.

Initiating our exploration with the **SELECT** statement, we explicitly call for the title from the **Course** table, simultaneously invoking the SQL aggregate function **AVG()** to compute the average of the **grade** column from the **Attended** table. For clarity in our output, this average is christened as "Average\_Grade".

As we weave our way through the query, the **FROM** clause designates the **Course** table as our primary point of departure. This is seamlessly integrated with the Attended table using the **JOIN** clause, ensuring a harmonious merger based on matching course IDs (**cid**). This joining operation is pivotal, as it facilitates a holistic view of each course alongside its corresponding student grades.

However, our expedition doesn't halt here. The **WHERE** clause instills precision, ensuring that we're solely inspecting records where the course IDs in both the **Course** and **Attended** tables resonate in harmony. Venturing further, the **GROUP BY** clause clusters our results by course title, ensuring each course is distinctly represented, preventing any overlap of data.

Then, unlike the traditional **WHERE** clause which filters rows, the **HAVING** clause filters groups. In our scenario, it meticulously sieves through the aggregated average grades, spotlighting only those groups where this average dips below the 6 mark.

In culmination, through this sophisticated command, we have adeptly navigated the database's depths, emerging with a list of courses and their average grades, with a special emphasis on those averaging less than 6.



```

74 -- 3d) Write a statement that retrieves the titles and the average grades of all the courses for
75 -- which the average grade of the students that attended them is lower than 6.
76
77 -- This query retrieves the titles and average grades of all courses for which
78 -- the average grade of the attending students is lower than 6.
79 SELECT
80     -- Gets the course title
81     Course.title as Course_Title,
82     -- Calculates the average grade
83     avg(Attended.grade) as Average_Grade
84 FROM
85     -- Joins the Course and Attended tables
86     Course JOIN Attended
87 WHERE
88     -- Ensures we're looking at records where the course ID in the Course and Attended tables match
89     Course.cid = Attended.cid
90 GROUP BY
91     -- Groups the results by course title
92     Course.title
93 HAVING
94     -- Filters for courses where the average grade is less than 6
95     avg(Attended.grade)<6;

```

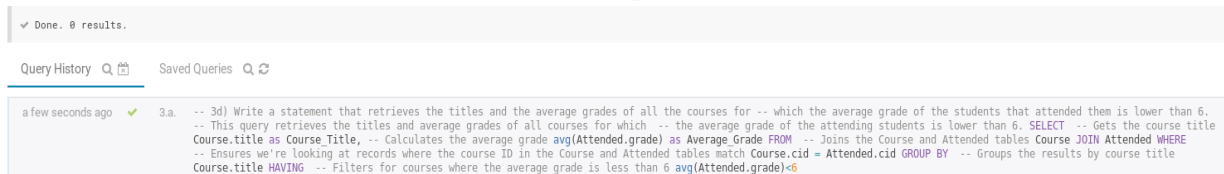


Figure 9: Execution of Task 3d's query in Impala with Results.

For a comprehensive examination of the query associated with **Task 3**, we kindly direct you to the [Abstract](#) section of the current document.

## Task 4: Query Optimization in Impala

*"A particular query in the previous Impala database is too slow. Describe what you are going to do to investigate what is going wrong and what can be done to improve efficiency. Provide any commands that you are going to run."*

The final task is concerned with improving the performance of a slow-running query in an Impala database. We will describe a systematic approach to diagnose the cause of the inefficiency and propose solutions that can enhance the query's performance. This task emphasizes the importance of efficient data querying in maintaining optimal system performance.

### Answer of Task 4:

To be more precise in addressing task 4, our primary goal is to diagnose the inefficiencies within a particular query that is exhibiting slow performance in our Impala database. To achieve this, we adopt a multi-pronged approach, applying diagnostic commands and then implementing potential solutions. **This procedure has been systematically applied to all queries up to this point, in order to achieve our goal.**

Firstly, we utilize the **EXPLAIN** command:

```

1. EXPLAIN SELECT Student.name as Student_Name
2. FROM Student
3. JOIN Attended
4. JOIN Course
5. WHERE (Student.sid = Attended.sid AND Attended.cid = Course.cid AND
6. Course.title='Artificial Intelligence' AND Attended.ac_year='2021-2022');

```

This command provides a detailed overview of how Impala plans to execute the given query. By analyzing the execution plan, we can identify potential bottlenecks such as suboptimal join orders or full table scans that could be contributing to the query's slow performance.

A demonstrative figure from the results of an execution of the **EXPLAIN** command follows below (due to space and time constraints we will stick only to this results of only one query - but in the procedure we execute them for all the queries):

Query History		Saved Queries	Results (33)
Explain String			
1	Max Per-Host Resource Reservation: Memory=3.88MB	19	
2	Per-Host Resource Estimates: Memory=19.88MB	20	03:HASH JOIN [INNER JOIN, BROADCAST]
3	WARNING: The following tables are missing relevant table and/or column statistics.	21	hash predicates: Student.sid = Attended.sid
4	university.attended, university.course, university.student	22	runtime filters: RF001 <- Attended.sid
5		23	
6	PLAN-ROOT SINK	24	-05:EXCHANGE [BROADCAST]
7		25	
8	07:EXCHANGE [UNPARTITIONED]	26	01:SCAN HDFS [university.attended]
9		27	partitions=1/1 files=0 size=0B
10	04:HASH JOIN [INNER JOIN, BROADCAST]	28	predicates: Attended.ac_year = '2021-2022'
11	hash predicates: Attended.cid = Course.cid	29	runtime filters: RF000 -> Attended.cid
12	runtime filters: RF000 <- Course.cid	30	
13		31	00:SCAN HDFS [university.student]
14	-06:EXCHANGE [BROADCAST]	32	partitions=1/1 files=1 size=29B
15		33	runtime filters: RF001 -> Student.sid
16	02:SCAN HDFS [university.course]		
17	partitions=1/1 files=0 size=0B		
18	predicates: Course.title = 'Artificial Intelligence'		

Figure 10: Results of explain command for demonstrative purposes.

Next, to gather more granular insights into the runtime behavior of the query, we employ the **PROFILE** command:

```
1. SELECT Student.name as Student_Name
2. FROM Student
3. JOIN Attended
4. JOIN Course
5. WHERE (Student.sid = Attended.sid AND Attended.cid = Course.cid AND
6. Course.title='Artificial Intelligence' AND Attended.ac_year='2021-2022');
7. PROFILE;
```

The **PROFILE** command offers a comprehensive breakdown of the query's execution, revealing how much time is spent in each phase. Such insights are invaluable in pinpointing exact stages that might be causing the slowdown.

For a more concise representation of the query's execution phases, and for more complete research we turn to the **SUMMARY** command for each query of Task 3:

```
1. SELECT Student.name as Student_Name
2. FROM Student
3. JOIN Attended
4. JOIN Course
5. WHERE (Student.sid = Attended.sid AND Attended.cid = Course.cid AND
6. Course.title='Artificial Intelligence' AND Attended.ac_year='2021-2022');
7. SUMMARY;
```

The **SUMMARY** command provides a succinct overview of the timings for the different phases of query execution. By contrasting this summary with the detailed profile, we can further narrow down areas requiring optimization.

Then, since in Impala, the process of query planning relies on statistics to generate efficient execution plans. Without accurate statistics, Impala might make suboptimal decisions, such as choosing an inefficient join order or misallocating memory resources. This can lead to slower query execution times. Therefore, we utilize the **COMPUTE STATS** command.

To be more precise the **COMPUTE STATS** command gathers statistics on tables, including row counts, data size, and distinct values for columns. This information allows Impala to better estimate the costs associated with different execution plans and choose the most efficient one.

```
1. -- To further optimize, statistics can be collected for the tables involved,  
2. -- aiding Impala in making more informed decisions during query planning.  
3. -- This can potentially lead to more efficient execution plans.  
4. COMPUTE STATS Student;  
5. COMPUTE STATS Attended;  
6. COMPUTE STATS Course;
```

By running the above commands, we are essentially providing Impala with a deeper understanding of the data distribution and characteristics in the **Student**, **Attended**, and **Course** tables. Once these statistics are computed, subsequent queries on these tables are likely to benefit from more efficient execution plans.

Furthermore, as data evolves (new data gets inserted, old data gets deleted), it's a good practice to periodically recompute statistics to ensure that they remain accurate and reflective of the current state of the data - *a complete theoretical scenario but we treat this Task as a real world problem.*

Moreover, in the context of Task 4, where we aim to optimize a slow-running query, gathering statistics is a key step. It should be one of the first actions taken when addressing performance issues, especially if statistics haven't been computed recently or if the data has undergone significant changes.

In conclusion, incorporating the **COMPUTE STATS** command is not only useful but crucial for ensuring that Impala has the necessary information to optimize query performance efficiently. To sum up, these commands are very useful for Task 4 and should be included as part of the optimization process.

From our observations, we identified that a significant delay was attributed to the multiple joins in the query. A proposed solution to mitigate this inefficiency is to restructure our database. By adding a column in the **Attended** table to store student names, we can potentially eliminate the need for additional joins:

```
1. ALTER TABLE Attended ADD COLUMNS (sname STRING);
```

To further streamline our data retrieval, we consider removing the Student table - if it is allowed by our client:

```
1. drop table Student;
```

Post these alterations, our modified query looks like:

```
1. SELECT Attended.sname as Student_Name
2. FROM Attended JOIN Course
3. WHERE (Attended.cid = Course.cid AND
4. Course.title='Artificial Intelligence' AND Attended.ac_year='2021-2022');
```

In scenarios where **table removal is not feasible from the client's policy** for example, we can resort to the **STRAIGHT\_JOIN** command to control the order of join operations, ensuring that *Impala does not reorder the join clauses*.

To be more precise, the **STRAIGHT\_JOIN** command can be invoked to prevent Impala from internally reordering the join clauses. Instead, it depends on the join clauses being sequentially ordered in the query. An illustration of this command is:

```
1. SELECT STRAIGHT_JOIN Student.name as Student_Name
2. FROM Student JOIN Attended JOIN Course
3. WHERE (Student.sid = Attended.sid AND Attended.cid = Course.cid AND
4. Course.title='Artificial Intelligence' AND Attended.ac_year='2021-2022');
```

In furtherance of our optimization efforts, we advocate for a change in the data type of the **ac\_year** column in the **Attended** table. Transitioning from CHAR(9) to either VARCHAR(9) or STRING is prudent. The rationale behind this recommendation stems from the fact that the CHAR type, in *its current version, lacks the support of Impala Codegen*. Both VARCHAR and STRING types, on the other hand, enjoy this support, and the performance enhancements from Codegen unequivocally surpass the advantages of a fixed-width CHAR. The command to execute this modification is:

```
1. -- Modifying the data type of the `ac_year` column in the "Attended" table from
2. CHAR(9) to VARCHAR(9) or STRING for optimal performance.
3. ALTER TABLE Attended CHANGE ac_year ac_year STRING;
```

We choose the **STRING** variable type, due to the observed performance benefits of **STRING** over CHAR in *our version of Impala*.

In conclusion, by adopting this systematic approach of diagnostic commands followed by potential solutions, we aim to ensure the efficient execution of our queries, thereby optimizing the overall performance of our Impala database. A sample image of a part of the above-described procedure follows below:

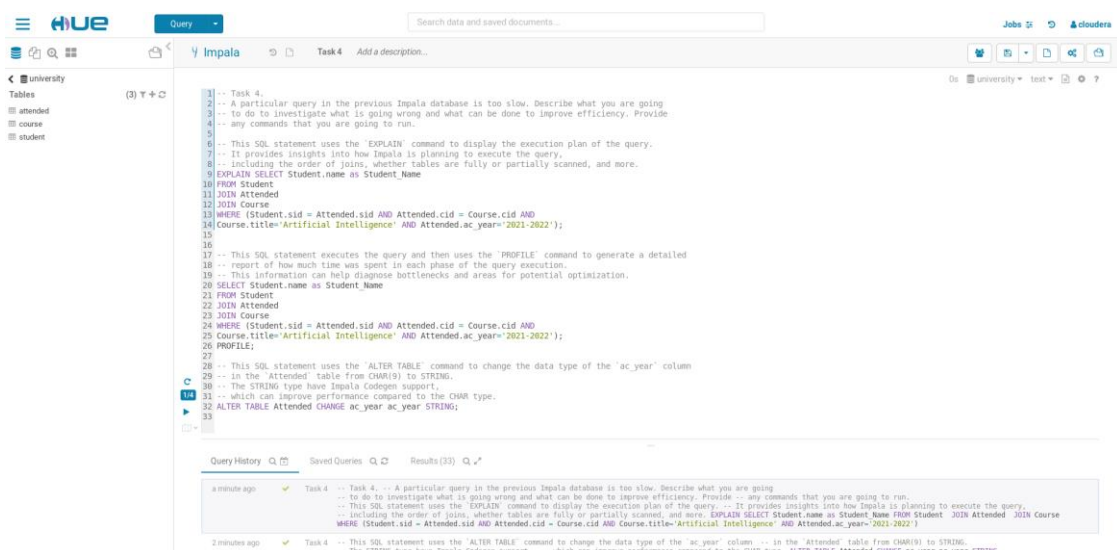


Figure 11: A sample picture demonstrating the procedure followed for the Task 4.

## Conclusion - Task 4

To sum up, in Task 4, our central objective was to pinpoint and remedy inefficiencies within a specific query in our Impala database that exhibited subpar performance. To address this concern, we meticulously formulated a systematic methodology, harmoniously melding diagnostic techniques with actionable solutions. **It's worth noting that this meticulous approach was not exclusive to a singular query; we consistently applied it across all the queries we tackled in Task 3, ensuring a holistic optimization strategy.**

At the outset, we harnessed the diagnostic power of the **EXPLAIN** command. This command offered an intricate glimpse into the intended execution strategy for the query, enabling us to discern potential bottlenecks such as unoptimized join sequences or exhaustive table scans. Such revelations were invaluable in our endeavor to streamline the query's performance.

To delve deeper into the actual runtime dynamics, we turned to the **PROFILE** command. This tool, with its granular breakdown, illuminated the exact duration each phase of the query took, furnishing us with a clear picture of where inefficiencies might lurk.

For a more synthesized perspective, we also leveraged the **SUMMARY** command. By providing a condensed overview of the various execution phases, it allowed us to juxtapose and contrast this summary against the in-depth profile, focusing our optimization efforts even further.

Through our analysis, it became evident that a primary source of delay was the multiple joins in our query. To mitigate this, we contemplated a structural adjustment to our database. By introducing a new column to store student names directly in the **Attended** table, we envisioned a reduction in the need for extraneous joins, potentially hastening our query processes.

Moreover, when table alterations weren't a viable option, we proposed the deployment of the **STRAIGHT JOIN command**. This command would ensure the sequence of join operations remained as scripted in the query, preventing any inadvertent internal reordering by Impala.

In addition, we discerned an opportunity to optimize the data type selection for certain columns. Notably, the transition from CHAR(9) to STRING for the ac\_year column in the Attended table was considered, based on the benefits the latter type offered in our version of Impala.

In conclusion, through Task 4, we have embarked on a holistic journey of query optimization, blending diagnostic insights with strategic solutions. By adhering to this comprehensive strategy **across all queries from Task 3**, we are poised to elevate the performance of our Impala database, ensuring it operates at its optimal potential.

For a comprehensive examination of the query associated with **Task 4**, we kindly direct you to the [Abstract](#) section of the current document.

**Note: For this Task all the evaluation metrics were applied to every query of the Task 3 in order to identify the worst/slowest query and optimize it. We showcase only one query for clarity and proof of work, but we did it for each query in order to find the slowest (the one we showcase) during our research phase.**

## Abstract

In the current Abstract section, we present a concise collection of all the SQL code utilized in **Tasks 3 and 4** of our assignment, demonstrating our hands-on interaction with the Apache Impala database system. This section serves as a reference point for the actual code implemented, stripped of extensive contextual or explanatory information. Here, you'll find the commands used to create the database and tables, insert data, retrieve specific information, and diagnose and optimize query performance. While a brief description accompanies each code snippet, detailed discussions and explanations are located in the main body of the document. The Abstract, thus, serves as a 'code snapshot', providing a quick overview of the command structure used in our Impala interactions.

```

27. -- Task 3:
28. -- 3a) Create the Impala database & the required tables.
29.
30. -- To replicate the desired scheme - will create issues
31. -- with 3b and varchar needs to be casted to String.
32. -- CREATE DATABASE IF NOT EXISTS university;
33. -- USE university;
34.
35. -- CREATE TABLE Student(
36. -- sid INT,
37. -- name VARCHAR(50)
38. -- );
39.
40. -- CREATE TABLE Course (
41. -- cid INT,
42. -- title VARCHAR(50),
43. -- description STRING
44. -- );
45.
46. -- CREATE TABLE Attended (
47. -- sid INT,
48. -- cid INT,
49. -- ac_year CHAR(9),
50. -- grade DECIMAL(3,1)
51. -- );
52.
53. -- If a database named "university" exists. If it does not exist,
54. -- it creates a new database with that name.
55. CREATE DATABASE IF NOT EXISTS university;
56.
57. -- This statement is used to select the "university"
58. -- database for use in subsequent operations.
59. USE university;
60.
61. -- This statement creates a new table named "Student" with columns "sid"
62. -- (an integer representing the student id) and
63. -- "name" (a string representing the student's name).
64. CREATE TABLE Student(
65. sid INT,
66. name STRING
67. );
68.
69. -- This statement creates a new table named "Course" with columns "cid"
70. -- (an integer representing the course id),
71. -- "title" (a string representing the course title),
72. -- and "description" (a string representing the course description).
73. CREATE TABLE Course (
74. cid INT,
75. title STRING,
76. description STRING
77. );
78.
79. -- This statement creates a new table named "Attended" with columns "sid"

```



```
80. -- (an integer representing the student id), "cid" (an integer representing the course id),
81. --
    "ac_year" (a character string with a fixed length of 9 characters representing the academic
    year),
82. -- and "grade" (a decimal number with a maximum of 3 digits,
83. -- 1 of which can be after the decimal point, representing the grade).
84. -- The FOREIGN KEY constraints "foreign1" and "foreign2" are also created
85. -- to link "sid" and "cid" with the "Student" and "Course" tables respectively.
86. CREATE TABLE Attended (
87. sid INT,
88. cid INT,
89. ac_year CHAR(9),
90. grade DECIMAL(3,1)
91. );
92.
93. -- Drop the database.
94. -- DROP DATABASE university CASCADE;
95. -- 3b) Give an example command that inserts an entry to the Student table (use your
96. -- own details for that entry) - 2822212 Dimitris Matsanganis.
97.
98. INSERT INTO Student (sid, name)
99. VALUES (2822212, "Dimitris Matsanganis");
100.
101.
102. --
    3c) Write a statement that retrieves all the names of the students that have attended the
103. -- course having title "Artificial Intelligence" during the academic year "2021-2022".
104.
105. -- This query retrieves the names of all students who attended the
106. -- course titled 'Artificial Intelligence' during the academic year 2021-2022.
107. -- First, joins the Student, Attended, and Course tables
108.
109. SELECT Student.name as Student_Name
110. FROM Student JOIN Attended JOIN Course
111. WHERE
112. --
    Ensures we're looking at records where the student ID in the Student and Attended tables match
113. Student.sid = Attended.sid AND
114. --
    Ensures we're looking at records where the course ID in the Attended and Course tables match
115. Attended.cid = Course.cid AND
116. -- Filters for the specific course title
117. Course.title='Artificial Intelligence' AND
118. -- Filters for the specific academic year
119. Attended.ac_year='2021-2022';
120. --
    3d) Write a statement that retrieves the titles and the average grades of all the courses for
121. -- which the average grade of the students that attended them is lower than 6.
122.
123. -- This query retrieves the titles and average grades of all courses for which
124. -- the average grade of the attending students is lower than 6.
125.
126. SELECT
127. -- Gets the course title
128. Course.title as Course_Title,
129. -- Calculates the average grade
130. avg(Attended.grade) as Average_Grade
131. FROM
132. -- Joins the Course and Attended tables
133. Course JOIN Attended
134. WHERE
135. --
    Ensures we're looking at records where the course ID in the Course and Attended tables match
136. Course.cid = Attended.cid
137. WHERE
```

```

138.      --
      Ensures we're looking at records where the course ID in the Course and Attended tables match

139.      Course.cid = Attended.cid
140.  GROUP BY
141.      -- Groups the results by course title
142.      Course.title
143.  HAVING
144.      -- Filters for courses where the average grade is less than 6
145.      avg(Attended.grade)<6;

1.  -- Task 4:
2.  --
      A particular query in the previous Impala database is too slow. Describe what you are going
3.  --
      to do to investigate what is going wrong and what can be done to improve efficiency. Provide
4.  -- any commands that you are going to run.
5.
6.  --
      This SQL statement uses the `EXPLAIN` command to display the execution plan of the query.
7.  -- It provides insights into how Impala is planning to execute the query,
8.  -- including the order of joins, whether tables are fully or partially scanned, and more.
9.  EXPLAIN SELECT Student.name as Student_Name
10. FROM Student
11. JOIN Attended
12. JOIN Course
13. WHERE (Student.sid = Attended.sid AND Attended.cid = Course.cid AND
14. Course.title='Artificial Intelligence' AND Attended.ac_year='2021-2022');
15.
16. --
      This SQL statement executes the query and then uses the `PROFILE` command to generate a deta
      iled
17. -- report of how much time was spent in each phase of the query execution.
18. -- This information can help diagnose bottlenecks and areas for potential optimization.
19. SELECT Student.name as Student_Name
20. FROM Student
21. JOIN Attended
22. JOIN Course
23. WHERE (Student.sid = Attended.sid AND Attended.cid = Course.cid AND
24. Course.title='Artificial Intelligence' AND Attended.ac_year='2021-2022');
25. PROFILE;
26.
27. -- This SQL command provides a summarized view of the execution phases for the query,
28. -- making it easier to pinpoint areas that might be causing delays.
29. SELECT Student.name as Student_Name
30. FROM Student
31. JOIN Attended
32. JOIN Course
33. WHERE (Student.sid = Attended.sid AND Attended.cid = Course.cid AND
34. Course.title='Artificial Intelligence' AND Attended.ac_year='2021-2022');
35. SUMMARY;
36.
37. --
      This SQL statement uses the `ALTER TABLE` command to change the data type of the `ac_year` c
      olumn
38. -- in the `Attended` table from CHAR(9) to STRING.
39. -- The STRING type have Impala Codegen support,
40. -- which can improve performance compared to the CHAR type.
41. ALTER TABLE Attended CHANGE ac_year ac_year STRING;
42.
43. -- To further optimize, statistics can be collected for the tables involved,
44. -- aiding Impala in making more informed decisions during query planning.
45. -- This can potentially lead to more efficient execution plans.
46. COMPUTE STATS Student;
47. COMPUTE STATS Attended;
48. COMPUTE STATS Course;
49.

```



```
50. -- Proposing structural changes to enhance the performance of the query.
51. -- Introducing a student name (sname) column to the Attended table.
52. ALTER TABLE Attended ADD COLUMNS (sname STRING);
53.
54. -- Deleting the Student table to streamline data retrieval.
55. drop table Student;
56.
57. -- Modified query post the structural alterations.
58. SELECT Attended.sname as Student_Name
59. FROM Attended JOIN Course
60. WHERE (Attended.cid = Course.cid AND
61. Course.title='Artificial Intelligence' AND Attended.ac_year='2021-2022');
62.
63. --
    Deploying the `STRAIGHT_JOIN` command to prevent Impala from reordering the join clauses.
64. SELECT STRAIGHT_JOIN Student.name as Student_Name
65. FROM Student JOIN Attended JOIN Course
66. WHERE (Student.sid = Attended.sid AND Attended.cid = Course.cid AND
67. Course.title='Artificial Intelligence' AND Attended.ac_year='2021-2022');
68.
69. -- Modifying the data type of the `ac_year` column in the "Attended" table from CHAR(9)
70. to VARCHAR(9) or STRING for optimal performance.
71. ALTER TABLE Attended CHANGE ac_year ac_year STRING;
```