

# Data Engineering - CSL4030

## Lab Assignment - 01

---

Irwindeep Singh  
B22AI022

Sarthak Malviya  
B22AI034

## 1 Hotel Management System

We have designed a hotel management system incorporating MySQL for database management. Further, we use Python-Flask app for backend.

The codebase used for the project along-with a README.md file can be found at: [GitHub](#)

## 2 Entity-Relation (ER) Diagram

The following Entity-Relation (ER) Diagram represents the schema for the Hotel Management System. It outlines the main entities involved in the system, such as Room, Guest, Booking, Payment, and User, along with their respective attributes. The relationships between these entities, such as "Booked By" and "Made By," are also depicted to illustrate how the data is interconnected within the system.

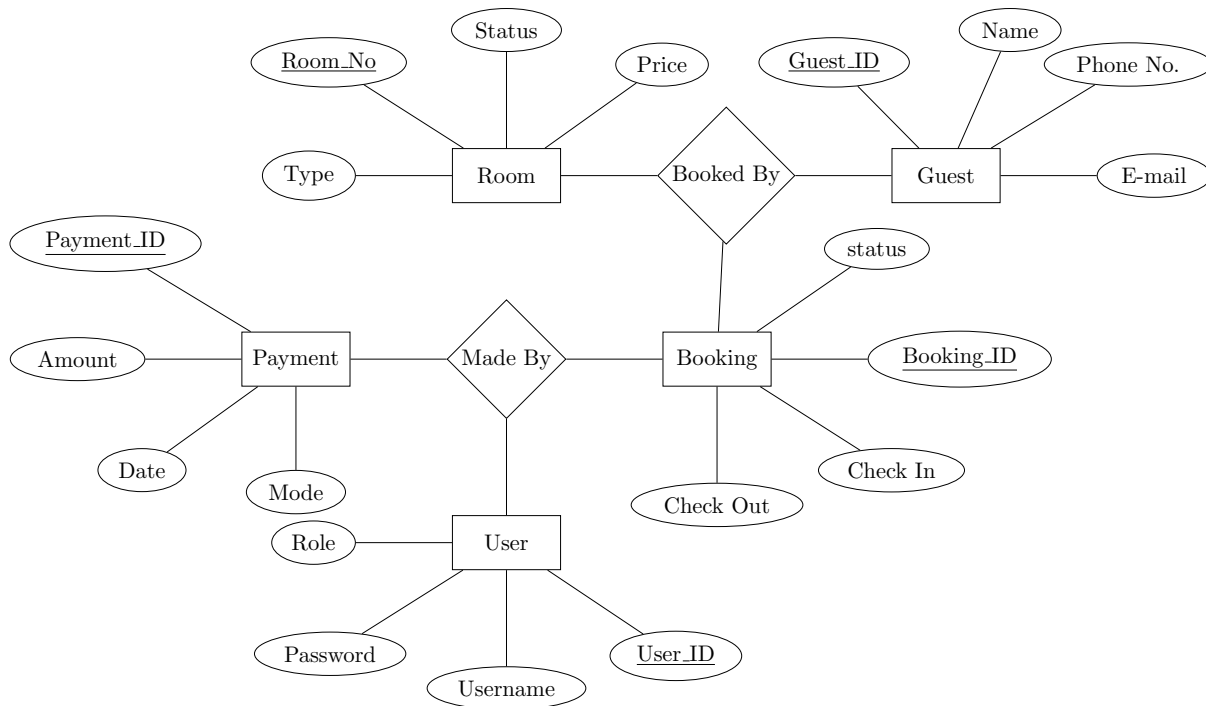


Figure 1: ER Diagram for Hotel Management System

### 3 Table Creation, Insertion of Dummy Records

We have used a dataset named `hotel_management` for this assignment, the same can be created with the following MySQL queries:

```
1 CREATE DATABASE hotel_management;
2 USE hotel_management;
```

#### 3.1 Table Creation

We have Used 6 tables named - rooms, payments, cancellations, guests, bookings, staff which were created with the following queries:

```
1 CREATE TABLE rooms(room_no INT PRIMARY KEY, type VARCHAR(20), status VARCHAR(20), price FLOAT);
2 CREATE TABLE payments(payment_id INT PRIMARY KEY, mode VARCHAR(20), amount FLOAT, payment_time DATETIME);
3 CREATE TABLE cancellations(cancel_id INT AUTO_INCREMENT PRIMARY KEY, room_no INT, guest_id INT,
4     cancel_time DATETIME, note TEXT);
5 CREATE TABLE guests(guest_id INT PRIMARY KEY, fname VARCHAR(20), lname VARCHAR(20), phone_no VARCHAR(20),
6     email VARCHAR(255));
```

```
1 CREATE TABLE bookings(
2     booking_id INT PRIMARY KEY,
3     room_no INT,
4     guest_id INT,
5     booking_time DATETIME,
6     total_occupants INT,
7     travel_mode VARCHAR(20),
8     check_in DATE,
9     check_out DATE,
10    payment_id INT
11 );
12
13 CREATE TABLE staff(
14     id INT PRIMARY KEY,
15     email_id VARCHAR(255),
16     password VARCHAR(255),
17     fname VARCHAR(20),
18     lname VARCHAR(20),
19     role VARCHAR(20)
20 );
```

#### 3.2 Insertion of Dummy Records

To insert dummy records into the dataset, we must first connect to the dataset (via `localhost`). We have utilised `MySQLdb` (a python library to view and manipulate `MySQL` databases. The following program assumes that we have the password for `MySQL` user stored in the operating system environment. We are first populating the rooms, then we select 5 random rooms and populate bookings with those 5 rooms. With bookings known, we are updating room status and populating payments.

```

1 import MySQLdb
2 from MySQLdb import Error
3 import os
4
5 try:
6     connection = MySQLdb.connect(
7         host="localhost",
8         user="root",
9         password=os.getenv('MySQL'),
10        database="hotel_management"
11    )
12    print("Connected to Database Successfully!")
13 except Error as e:
14    print(f"Error Occurred: {e}")

```

Now, we can insert dummy records in all tables by obtaining a `cursor` object from the established connection, and with that `cursor`, we can execute queries and populate tables as follows:

```

1 cursor = connection.cursor()
2
3 # Creating empty rooms with different types and prices
4 for room_no in range(101, 201):
5     if room_no <= 125: type, price = "Standard", 2500.0
6     elif room_no <= 150: type, price = "Deluxe", 3200.0
7     elif room_no <= 175: type, price = "Executive", 4000.0
8     else: type, price = "Premium", 6500.0
9
10    cursor.execute(f"INSERT INTO rooms VALUES({room_no}, '{type}', 'Empty', {price})")
11
12 connection.commit()

```

```

1 import numpy as np
2
3 # Assuming we have lists containing all the variables used
4 for id in range(1, 6):
5     room_no = np.random.randint(low=101, high=200)
6     query = f"{id}, {room_no}, '{id}', '{book_times[id-1]}', {occupants[id-1]}, '{travel_modes[id-1]}', "
7     query += f"'{in_dates[id-1]}', '{out_dates[id-1]}', {id}"
8     cursor.execute(f"INSERT INTO bookings VALUES({query})")
9     cursor.execute(f"UPDATE rooms SET status='Booked' WHERE room_no={room_no}")
10
11    query = f"{id}, '{fnames[id-1]}', '{lnames[id-1]}', '{phone_nos[id-1]}', '{emails[id-1]}'"
12    cursor.execute(f"INSERT INTO guests VALUES({query})")
13
14    cursor.execute(f"SELECT price FROM rooms WHERE room_no={room_no}")
15    amount = cursor.fetchone()[0]
16    cursor.execute(f"SELECT DATEDIFF(booking_till, booking_date) AS number_of_days WHERE room_no={room_no}")
17    amount *= cursor.fetchone()[0]
18    cursor.execute(f"INSERT INTO payments VALUES({id}, 'Cash', {amount}, {book_times[id-1]})")
19
20 connection.commit()

```

## 4 Normalisation

- **1NF (First Normal Form):**
  - Each table has a primary key.
  - All columns contain atomic values (e.g., VARCHAR for strings, INT for integers, FLOAT for prices).
  - **Conclusion:** The database is in 1NF.
- **2NF (Second Normal Form):**
  - All non-key attributes are fully functionally dependent on the entire primary key in all the tables.
  - **Conclusion:** The database is in 2NF.
- **3NF (Third Normal Form):**
  - There are no transitive dependencies in any of the table.
  - The database is in 2NF.
  - **Conclusion:** The database is in 3NF.
- **BCNF (Boyce-Codd Normal Form):**
  - For every functional dependency  $X \rightarrow Y$ ,  $X$  is the superkey in all the tables.
  - The database is in 3NF.
  - **Conclusion:** The database is in BCNF

## 5 Hashing and Indexing Schemes

Hashing and indexing are critical mechanisms for optimizing query performance in MySQL databases. Both techniques are designed to enhance data retrieval efficiency, but they operate in fundamentally different ways.

- **Hashing:** Hashing is a technique used to distribute data across a hash table to speed up access. In MySQL, hashing is employed primarily in the context of hash indexes, which are a type of in-memory index used by the MEMORY storage engine. When a hash index is used, the hash function converts the key into a hash value, which determines the index location. This method is effective for equality searches but inefficient for range queries since hash functions do not maintain any order of data.
- **Indexing:** Indexing, on the other hand, involves creating data structures that improve the speed of data retrieval operations. MySQL supports various types of indexes, with B-tree and R-tree indexes being the most common:
  - **B-Trees:** B-tree (Balanced Tree) indexes are organized as a balanced tree structure, allowing efficient searching, insertion, and deletion operations. *They are the default index type for most MySQL storage engines.*
  - **R-Trees:** R-trees are used for indexing multi-dimensional information, such as spatial data.

## 6 Hash Function Design

To design a hash function (for rooms table) incorporating common characters of "B22AI022" and "B22AI034", we have utilized the following python program:

```

1 def hash_function(key: int, table_size: int) -> int:
2     key_str, PRIME_1, PRIME_2 = str(key), 31, 7
3
4     hash_value = 0 # Initialize hash value
5     for char in key_str:
6         hash_value = (hash_value * PRIME_1 + int(char)) % table_size
7
8     common_chars = 'B22AI0' # Common characters in "B22AI022" and "B22AI034"
9     for char in common_chars:
10        hash_value = (hash_value * PRIME_2 + ord(char)) % table_size
11
12    return hash_value % table_size
13
14 # Example usage:
15 key1, key2 = 22022022, 22023034
16 table_size = 97
17
18 print(f"Hash value for key '{key1}': {hash_function(key1, table_size)}")
19 print(f"Hash value for key '{key2}': {hash_function(key2, table_size)}")

```

```

1 Output:
2     Hash value for key '22022022': 65
3     Hash value for key '22023034': 10

```

## 7 Clustering Indexing

Since, clustering indexing is already applied `rooms` due to the presence of the primary key (`room_no`). This means that the rows of the table are physically ordered on disk according to the `room_no` key.

## 8 Secondary Indexing

With the connection to MySQL client made above, we can apply secondary indexing (i.e., indexing on attribute other than primary attribute) on the `rooms` table with the following python program:

```

1 cursor.execute("CREATE INDEX idx_room_type ON rooms (type)")
2 connection.commit()
3
4 # Verifying the indexing
5 cursor.execute("SHOW INDEX FROM rooms")
6 indices = cursor.fetchall()
7 for index in indices:
8     print(index)

```

```

1 Output:
2 ('rooms', 0, 'PRIMARY', 1, 'room_no', 'A', 100, None, None, '', 'BTREE', '', '', 'NO')
3 ('rooms', 1, 'idx_room_type', 1, 'type', 'A', 4, None, None, 'YES', 'BTREE', '', '', 'NO')

```

## 9 Clustering vs Secondary Indexing

### 9.1 Clustering Index (Primary Key Index)

- **Storage:**
  - **Physical Order:** The table data is physically ordered on disk according to the clustering index. This means that the rows are stored in the same order as the clustering index.
  - **Space Efficiency:** Because the data is stored in a sorted order, it can be more space-efficient when accessing data sequentially. However, if there are frequent updates, inserts, or deletes, maintaining this order can result in fragmentation and potential space overhead.
- **Execution Time:**
  - **Data Retrieval:** Accessing data using the clustering index is very efficient, especially for range queries. For example, retrieving all records within a specific range is fast because the data is stored in a sorted order.
  - **Search Performance:** Searching for a specific row using the clustering index is quick because it directly accesses the row using the indexed key.
  - **Insert/Update/Delete Performance:** These operations can be slower due to the need to maintain the sorted order of the table, which might involve reordering and shifting of rows.

### 9.2 Secondary Index (Non-Clustered Index)

- **Storage:**
  - **Separate Storage:** The secondary index is stored separately from the table data. It contains a copy of the indexed column(s) and pointers to the corresponding rows in the table.
  - **Space Overhead:** There is additional space overhead because the secondary index stores duplicate information. The more indexes you have, the more storage space is required.
- **Execution Time:**
  - **Data Retrieval:** Retrieving data using a secondary index can be efficient for specific queries where the indexed column is used. However, it involves an additional lookup step to fetch the actual row from the table.
  - **Search Performance:** Searching for a specific value in a secondary index is typically fast because secondary indexes use data structures like B-trees or hash tables, which allow for efficient lookups.
  - **Insert/Update/Delete Performance:** These operations can be slower with secondary indexes compared to tables with only a clustering index. Every time a row is inserted, updated, or deleted, all secondary indexes need to be updated, which can add overhead.

### 9.3 Summary

- **Clustering Index:**
  - **Pros:** Efficient for range queries, good for data retrieval when the data is frequently accessed in a sorted order.
  - **Cons:** Can cause performance issues with frequent inserts/updates due to the need to maintain data order; less flexible as there is only one clustering index per table.
- **Secondary Index:**
  - **Pros:** Flexible and allows for multiple secondary indexes; improves query performance for specific column searches; good for queries where the indexed columns are frequently used.
  - **Cons:** Additional storage space is required for indexes; can result in slower insert/update/delete operations due to the need to maintain multiple indexes.

## 10 Additional Queries

1. Add information about the inclusion of a new wing comprising 5 deluxe suites:

```
1  INSERT INTO rooms
2  VALUES
3      (201, 'Deluxe', 'Empty', 3200.0),
4      (202, 'Deluxe', 'Empty', 3200.0),
5      (203, 'Deluxe', 'Empty', 3200.0),
6      (204, 'Deluxe', 'Empty', 3200.0),
7      (205, 'Deluxe', 'Empty', 3200.0);
```

2. Prepare a report on all guest bookings and cancellations in the month of August, 2024:

```
1  SELECT
2      'Booking' AS type,
3      booking_id AS id,
4      room_no,
5      booking_time AS time,
6  FROM bookings
7  WHERE booking_date >= '2024-08-01 00:00:00' AND booking_date < '2024-09-01 00:00:00'
8
9  UNION ALL
10
11 SELECT
12     'Cancellation' AS type,
13     cancellation_id AS id,
14     room_no,
15     cancel_time AS time,
16  FROM cancellations
17  WHERE cancel_time >= '2024-08-01 00:00:00' AND cancellation_date < '2024-09-01 00:00:00'
18  ORDER BY date;
```

3. Cancel all guest bookings made after 7PM for August 15, 2024:

```
1  INSERT INTO cancellations (room_no, guest_id, cancel_time, note)
2  SELECT room_no, guest_id, '2024-08-31 00:00:00', 'Additional Query #3'
3  FROM bookings
4  WHERE booking_time >= '2024-08-15 19:00:00';
5
6  DELETE FROM bookings
7  WHERE booking_time >= '2024-08-15 19:00:00';
```