



## Data Management and Warehousing Project

### Project Team:

names	Student number	Work percentage
Mohammed Khalid Altufayhi	444003827	20%
Anas Mohammed Alsubhi	444000066	20%
Faisal Hammad Alomari	444005427	20%
Nawaf Abdulrahman Alageel	44410754	20%
Albadar Ibrahim Almaymani	444000184	20%

## Introduction

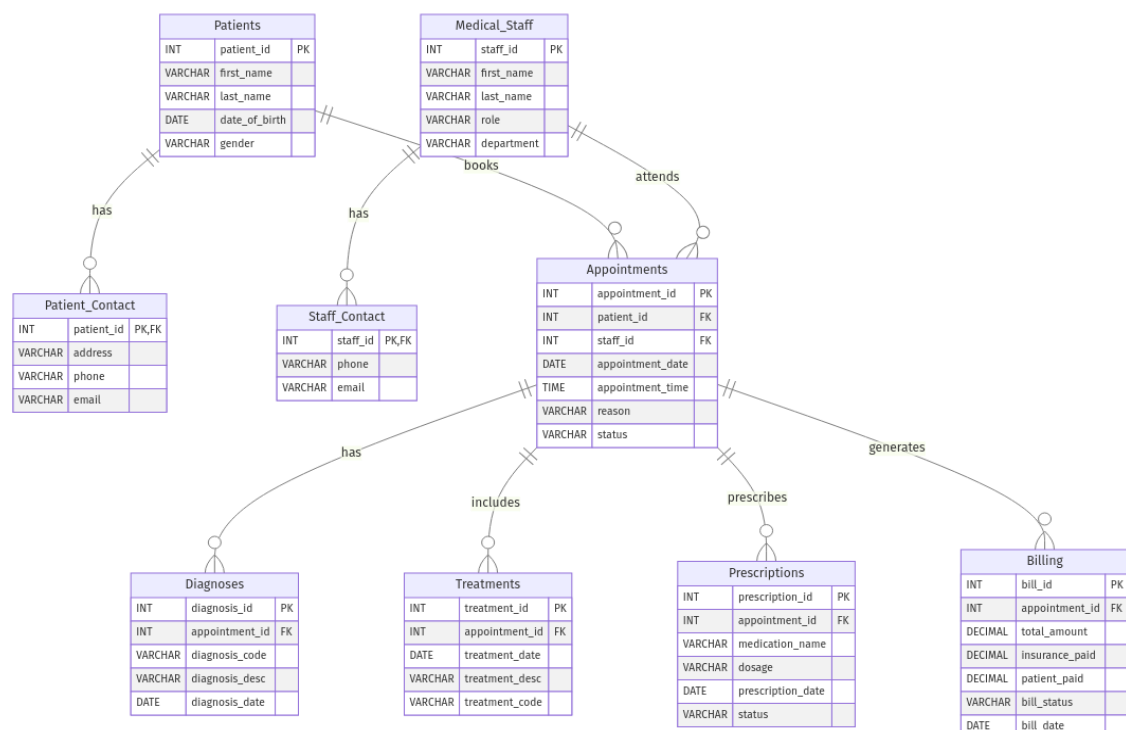
The healthcare domain often deals with massive datasets that contain sensitive and critical information, including patient records, diagnoses, treatments, prescriptions, and billing information. A well-designed data warehouse schema can help efficiently manage and analyze such data for better decision-making, operational insights, and compliance with healthcare standards. This report explores three different warehouse schemas designed for the given healthcare transactional database: the **Snowflake Schema**, the **Galaxy Schema**, and the **Star Schema**. Each schema is designed with a unique structure to meet various analytical needs, balancing performance, storage efficiency, and ease of querying.

## Schema Designs

### 1. Snowflake Schema

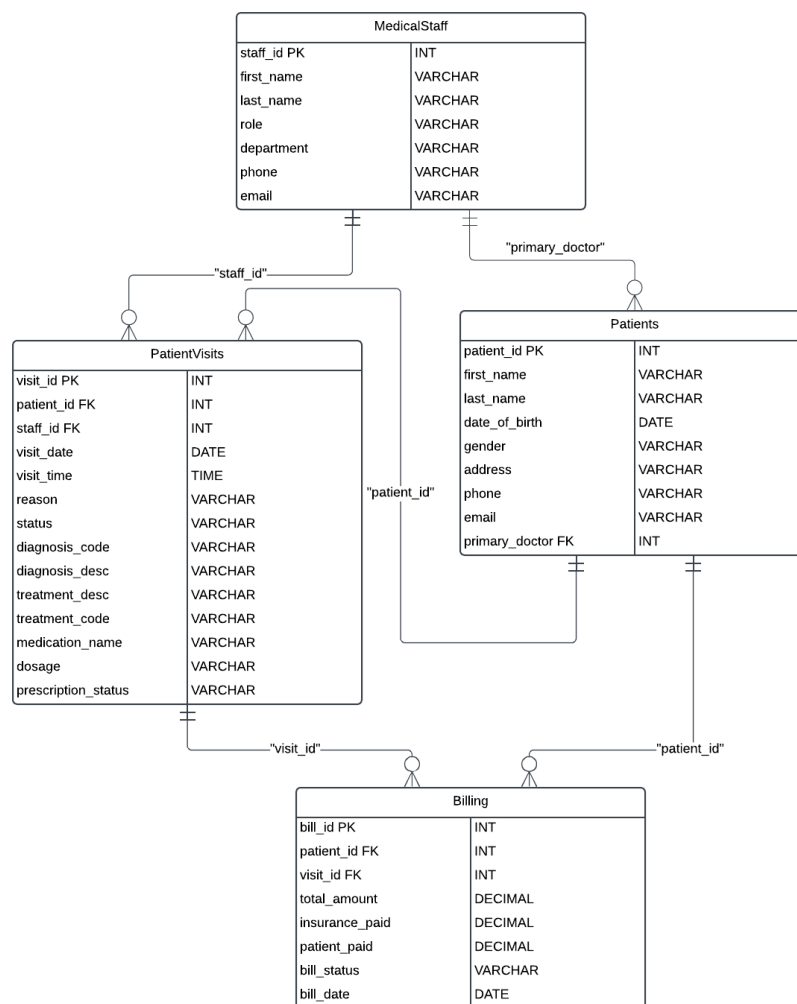
The **Snowflake Schema** is a normalized version that reduces redundancy by dividing the tables into multiple dimensions. This schema separates contact information for both patients and medical staff into dedicated tables to achieve a higher level of normalization.

- **Design:**
  - The Patients table is divided into Patients and Patient\_Contact, where Patient\_Contact stores contact information.
  - Similarly, the Medical\_Staff table is split into Medical\_Staff and Staff\_Contact.
  - Other tables such as Appointments, Diagnoses, Treatments, Prescriptions, and Billing remain as individual dimension tables linked to Patients and Medical\_Staff.



## 2. Galaxy Schema

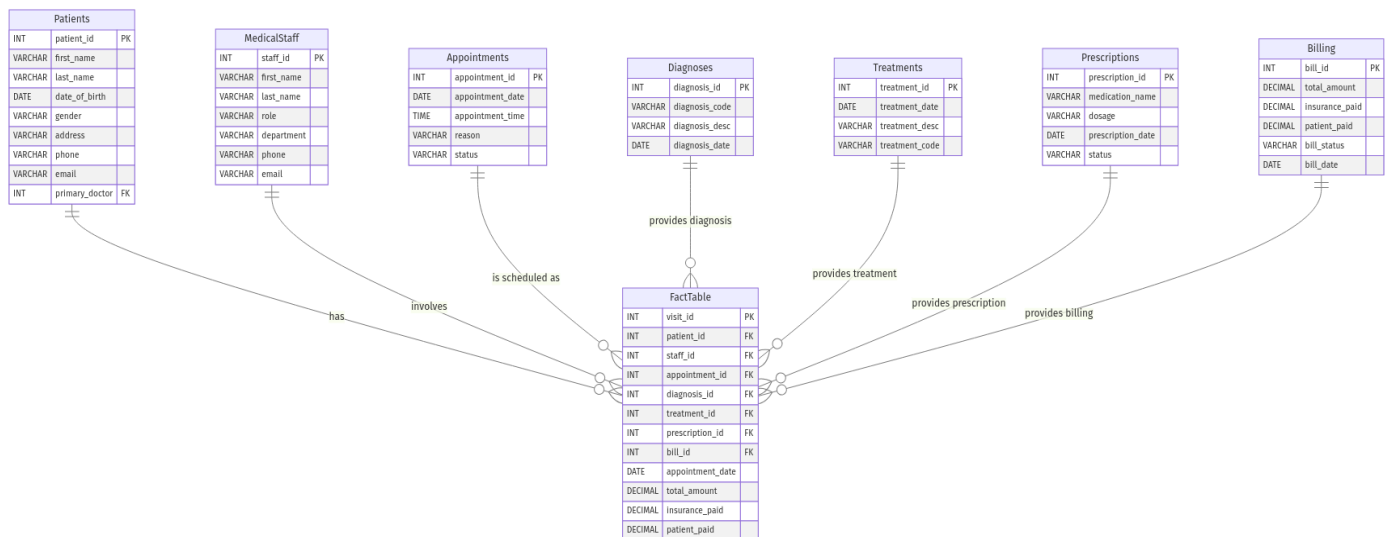
- The Galaxy Schema is a combination of the star and snowflake schemas, where key tables are slightly denormalized, reducing the number of tables without introducing significant redundancy.
- **Design:**
  - Patients: Contains patient info and contact details, with primary\_doctor linking to MedicalStaff.
  - PatientVisits serves as a consolidated table for Appointments, Diagnoses, Treatments, and Prescriptions, linking each patient visit to the appropriate staff.
  - Billing tracks financial details for each visit, connected to PatientVisits and Patients.



### 3. Star Schema

The **Star Schema** is a denormalized structure with a central fact table, FactTable, that links to all dimension tables, creating a "star" shape. This schema is optimized for fast querying and is suited for large-scale analytical processing.

- **Design:**
  - A central FactTable holds references to dimension tables: PATIENTS, MEDICALSTAFF, APPOINTMENTS, DIAGNOSES, TREATMENTS, PRESCRIPTIONS, and BILLING.
  - This schema is designed to allow rapid aggregation of data across various dimensions with minimal joins, making it ideal for data warehousing and reporting in high-performance scenarios.



## Pros and Cons of Each Schema

### Snowflake Schema

- **Pros:**
  - Reduces data redundancy, leading to potentially smaller storage requirements.
  - Normalized structure makes updates easier and minimizes duplication.
- **Cons:**
  - Complex queries due to multiple joins across tables, which may slow down performance.
  - Less suitable for large-scale data analysis as normalization requires more relational navigation.

## Galaxy Schema

- **Pros:**
  - Flexibility: Allows for analysis of multiple processes or events with shared data.
  - Provides a balanced approach to data redundancy and storage efficiency.
- **Cons:**
  - Increased Redundancy, the consolidated **PatientVisits** table may have some redundancy, as each visit stores multiple details.
  - Combining tables may limit flexibility if future use cases require further breakdown or normalization of data.

## Star Schema

- **Pros:**
  - Fast query performance due to minimal joins with a central fact table.
  - Ideal for analytical and reporting purposes with high volumes of data.
- **Cons:**
  - Increased data redundancy, leading to larger storage requirements.
  - Updating data is more complex due to the denormalized structure.

---

## Implementing Schema

We chose the **Star Schema** because it is the best fit for large-scale analytical queries in the healthcare domain. The Star Schema's denormalized structure (less splitting of tables) allows for faster query performance, which is essential for analyzing healthcare data quickly and efficiently. In this setup, all essential details are linked to a single **FactTable**, simplifying analysis across different tables.



The database is created and selected as HealthCare to organize all the healthcare-related tables.

```
1  -- Dimension Table: PATIENTS
2  CREATE TABLE PATIENTS (
3      patient_id INT PRIMARY KEY,
4      first_name VARCHAR(100),
5      last_name VARCHAR(100),
6      date_of_birth DATE,
7      gender VARCHAR(10),
8      address VARCHAR(255),
9      phone VARCHAR(20),
10     email VARCHAR(100),
11     primary_doctor INT
12 );
```

- **Description:** This table stores basic information about each patient, such as their name, birth date, contact details, and gender. It also includes a primary\_doctor field, linking each patient to a primary care doctor.
- **Purpose:** Serves as a dimension to provide patient-specific information when analyzing healthcare events.

```
1  -- Dimension Table: MEDICALSTAFF
2  CREATE TABLE MEDICALSTAFF (
3      staff_id INT PRIMARY KEY,
4      first_name VARCHAR(100),
5      last_name VARCHAR(100),
6      role VARCHAR(50),
7      department VARCHAR(50),
8      phone VARCHAR(20),
9      email VARCHAR(100)
10 );
```

- **Description:** Contains data about the medical staff, including their role (doctor, nurse), department, and contact details.
- **Purpose:** Acts as a dimension table for information about healthcare providers, helping to track which staff member interacted with which patient.

```
1  -- Dimension Table: APPOINTMENTS
2  CREATE TABLE APPOINTMENTS (
3      appointment_id INT PRIMARY KEY,
4      appointment_date DATE,
5      appointment_time TIME,
6      reason VARCHAR(255),
7      status VARCHAR(20)
8  );
```

- **Description:** Holds information on appointments, including the date, time, reason, and status (completed, canceled).
- **Purpose:** This dimension allows for analyzing patient visit patterns, including the timing and purpose of appointments.

```
1  -- Dimension Table: DIAGNOSES
2  CREATE TABLE DIAGNOSES (
3      diagnosis_id INT PRIMARY KEY,
4      diagnosis_code VARCHAR(20),
5      diagnosis_desc VARCHAR(255),
6      diagnosis_date DATE
7  );
```

- **Description:** Stores diagnosis information, such as a code (ICD code), description, and the date the diagnosis was made.
- **Purpose:** Helps in tracking the conditions diagnosed for each patient, which can be useful for both reporting and longitudinal studies.

```
1  -- Dimension Table: TREATMENTS
2  CREATE TABLE TREATMENTS (
3      treatment_id INT PRIMARY KEY,
4      treatment_date DATE,
5      treatment_desc VARCHAR(255),
6      treatment_code VARCHAR(20)
7  );
```

- **Description:** Contains information about treatments provided to patients, including a description, date, and a procedure code.
- **Purpose:** This table enables the analysis of treatments given, allowing insights into healthcare practices and patient outcomes.

```
1  -- Dimension Table: PRESCRIPTIONS
2  CREATE TABLE PRESCRIPTIONS (
3      prescription_id INT PRIMARY KEY,
4      medication_name VARCHAR(100),
5      dosage VARCHAR(50),
6      prescription_date DATE,
7      status VARCHAR(20)
8  );
```

- **Description:** Tracks prescriptions issued to patients, with details about the medication, dosage, date, and current status (active, expired).
- **Purpose:** This table allows analysis of medication usage and patterns, useful for tracking prescribed medications over time.

```

1  -- Dimension Table: BILLING
2  CREATE TABLE BILLING (
3      bill_id INT PRIMARY KEY,
4      total_amount DECIMAL(10,2),
5      insurance_paid DECIMAL(10,2),
6      patient_paid DECIMAL(10,2),
7      bill_status VARCHAR(20),
8      bill_date DATE
9  );

```

- **Description:** Stores billing and payment details for healthcare services, including the total amount, amount covered by insurance, and any remaining balance.
- **Purpose:** Facilitates financial reporting and analysis of patient payments, insurance contributions, and outstanding balances.

```

1  -- Fact Table: FACT
2  CREATE TABLE FactTable (
3      visit_id INT PRIMARY KEY,
4      patient_id INT,
5      staff_id INT,
6      appointment_id INT,
7      diagnosis_id INT,
8      treatment_id INT,
9      prescription_id INT,
10     bill_id INT,
11     appointment_date DATE,
12     total_amount DECIMAL(10,2),
13     insurance_paid DECIMAL(10,2),
14     patient_paid DECIMAL(10,2),
15     FOREIGN KEY (patient_id) REFERENCES PATIENTS(patient_id),
16     FOREIGN KEY (staff_id) REFERENCES MEDICALSTAFF(staff_id),
17     FOREIGN KEY (appointment_id) REFERENCES APPOINTMENTS(appointment_id),
18     FOREIGN KEY (diagnosis_id) REFERENCES DIAGNOSES(diagnosis_id),
19     FOREIGN KEY (treatment_id) REFERENCES TREATMENTS(treatment_id),
20     FOREIGN KEY (prescription_id) REFERENCES PRESCRIPTIONS(prescription_id),
21     FOREIGN KEY (bill_id) REFERENCES BILLING(bill_id)
22 );

```

- **Description:** The central FactTable links all dimension tables (Patients, Medical Staff, Appointments, Diagnoses, Treatments, Prescriptions, and Billing). It stores the primary key (visit\_id) for each visit or interaction, linking to each relevant dimension.
- **Purpose:** This table is essential for fast data retrieval and analysis across multiple dimensions. By centralizing the links to all patient interactions (appointments, diagnoses, treatments, prescriptions, and billing), this table supports quick aggregation and reporting, making it ideal for analytical queries.



## Conclusion

In this **Star Schema** design, each dimension table represents an entity in the healthcare domain, while the **FactTable** serves as the central table linking all patient interactions. This structure allows efficient querying across multiple dimensions, such as analyzing patient care, tracking medical staff interactions, and monitoring billing details. This schema is well-suited for analytical purposes, providing a balanced approach to data organization and query performance.