

# **FINAL PROJECT ON DISEASE MODEL**

Anuj Kumar Shah

Niyanta Pandey

Dakshit Ashokbhai Golakiya

## **Executive Summary:**

Our project develops a sophisticated database system to manage and analyze disease-related data. Leveraging our knowledge in data modeling, database design, and analysis, we've crafted a tool that offers valuable insights into disease trends and the effectiveness of treatments.

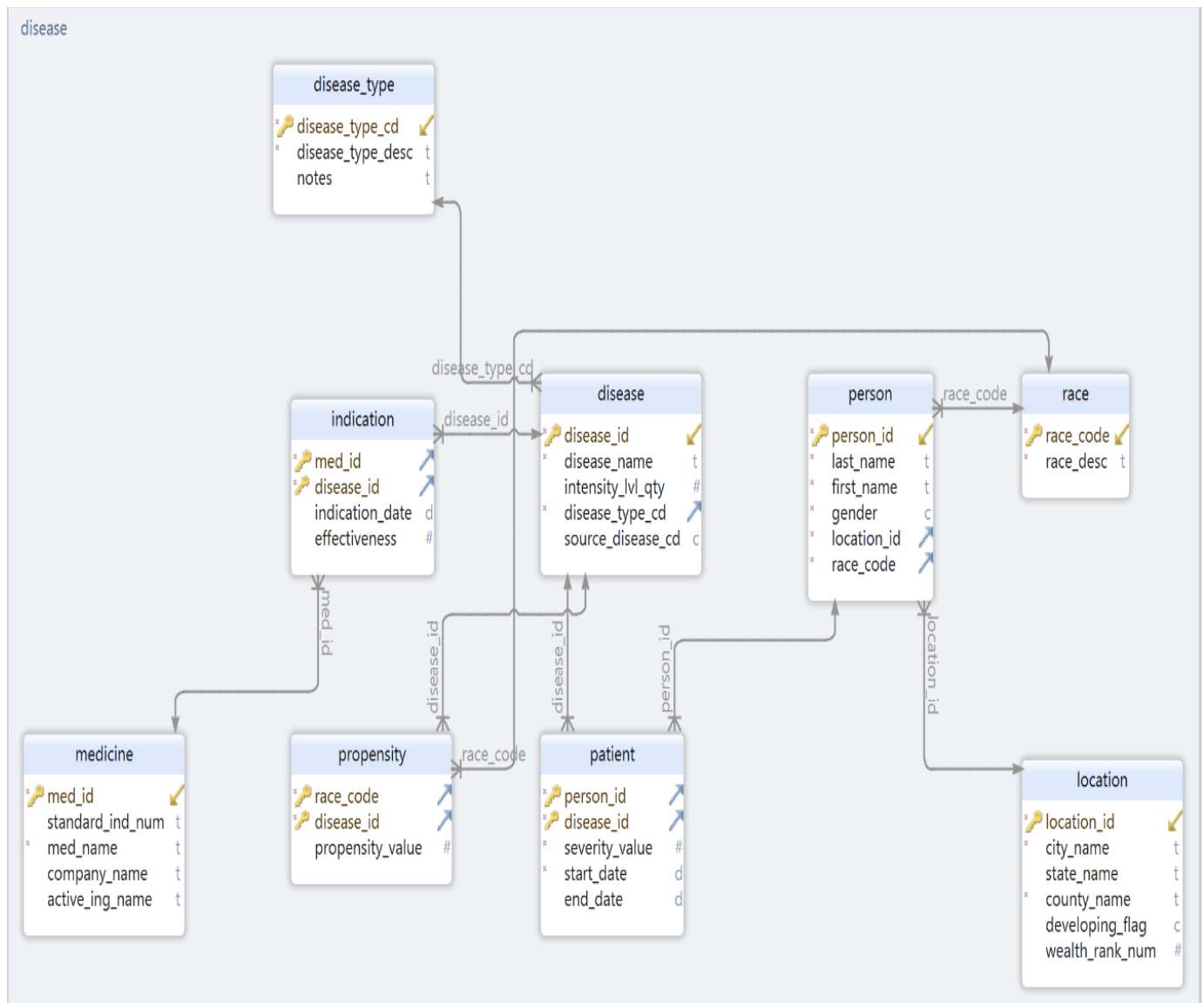
## **Introduction**

The goal of this project is to build a database system adept at handling disease data. We aim to utilize various data modeling and database management techniques learned during our course to derive meaningful insights into patterns in disease spread, treatment efficacy, and demographic distributions.

## **Background and Context**

In this project, we applied fundamental concepts such as Entity-Relationship (ER) modeling, OLTP, OLAP, and ETL processes. These concepts are vital in managing and interpreting extensive health-related data sets. We also explored transitioning to flexible, schema-less NoSQL databases like MongoDB or Neo4J, enhancing scalability and performance. Additionally, we researched using AWS components like RDS/Aurora, DynamoDB, Lambda, and Kinesis for robust, secure, and real-time data processing. For data warehousing, we delved into Snowflake's capabilities for scalability, automated performance tuning, and efficient data sharing.

## ER Diagram and Data Dictionary



Entity-Relationship Diagram (ERD) for a database related to healthcare information. It shows several tables with their columns, and the lines between them indicate relationships between the tables. Here's a detailed breakdown of the tables, their columns, and how they're connected:

## 1. **disease\_type:**

Disease\_type\_cd: (Primary Key) Code identifying the type of disease.

disease\_type\_desc: The description of the disease type.

notes: Any notes related to the disease type.

## 2. **disease:**

disease\_id: Unique identifier for the disease.

disease\_name: Name of the disease.

intensity\_lvl\_qty: The intensity level quantity of the disease.

disease\_type\_cd: A foreign key linking to the disease\_type\_cd in the disease\_type table.

source\_disease\_cd: The source code of the disease.

## 3. **indication:**

med\_id: The medication identifier.

disease\_id: A foreign key linking to the disease\_id in the disease table.

indication\_date: The date when the medication was indicated.

effectiveness: The effectiveness of the medication.

## 4. **medicine:**

med\_id: Unique identifier for the medicine.

standard\_id\_num: A standard identification number for the medicine.

med\_name: Name of the medicine.

company\_name: The name of the company that produces the medicine.

active\_ing\_name: The name of the active ingredient in the medicine.

## 5. **person:**

person\_id: Unique identifier for the person.

last\_name: The last name of the person.

first\_name: The first name of the person.

gender: The gender of the person.

location\_id: A foreign key linking to the location\_id in the location table.

race\_code: A foreign key linking to the race\_code in the race table.

## **6. race:**

race\_code: The code for the race.

race\_desc: The description of the race.

## **7. patient:**

person\_id: A foreign key linking to the person\_id in the person table.

disease\_id: A foreign key linking to the disease\_id in the disease table.

severity\_value: The value indicating the severity of the patient's condition.

start\_date: The start date of the disease or condition.

end\_date: The end date of the disease or condition.

## **8. location:**

location\_id: Unique identifier for the location.

city\_name: The name of the city.

state\_name: The name of the state.

county\_name: The name of the county.

developing\_flag: A flag indicating whether the location is in a developing area.

wealth\_rank\_num: A numerical rank of the location's wealth.

## **9. propensity:**

race\_code: A foreign key linking to the race\_code in the race table.

disease\_id: A foreign key linking to the disease\_id in the disease table.

propensity\_value: The value indicating the propensity of the race to a disease.

The connections between these tables are represented by the lines, which indicate the relationships based on foreign keys. For example, the `disease_id` in the `disease` table is linked to the `disease_id` in the `indication` table, meaning that the indications are related to specific diseases. Similarly, the `person_id` in the `person` table is linked to the `person_id` in the `patient` table, suggesting that the patient records are related to specific people in the `person` table. The above ERD shows how data related to diseases, patients, medications, and demographics are structured and interrelated in this database.

## Database and Table Creation in PostgreSQL

We developed a relational database in PostgreSQL, focusing on referential integrity and relational design. The schema was crafted to optimize query performance and data accuracy, reflecting the relationships established in the ER model.

```
Query    Query History
59  -- Creation of table for storing geographical locations
60  CREATE TABLE location(
61      location_id SERIAL PRIMARY KEY, -- Using SERIAL for auto-incrementing primary key
62      city_name VARCHAR(100) NOT NULL,
63      state_name VARCHAR(100),
64      county_name VARCHAR(100) NOT NULL,
65      developing_flag CHAR(1), -- Flag to indicate if the location is developing
66      wealth_rank_num INT CHECK (wealth_rank_num >= 1) -- Check constraint for wealth rank
67  );
68
69  -- Creation of table for storing personal details
70  CREATE TABLE person(
71      person_id SERIAL PRIMARY KEY, -- Using SERIAL for auto-incrementing primary key
72      last_name VARCHAR(50) NOT NULL,
73      first_name VARCHAR(50) NOT NULL,
74      gender CHAR(1) NOT NULL CHECK (gender IN ('M', 'F', 'U')), -- Check constraint for gender
75      location_id INT NOT NULL,
76      race_code CHAR(5) NOT NULL,
77      FOREIGN KEY(race_code) REFERENCES race(race_code)
78          ON UPDATE CASCADE ON DELETE CASCADE,
79      FOREIGN KEY(location_id) REFERENCES location(location_id)
80          ON UPDATE CASCADE ON DELETE CASCADE
81  );
82
83  -- Creation of table for storing patient disease records
84  CREATE TABLE patient(
85      person_id INT NOT NULL,
86      disease_id INT NOT NULL,
87      severity_value INT NOT NULL DEFAULT 1, -- Default severity
88      start_date DATE NOT NULL,
89      end_date DATE,
90      PRIMARY KEY(person_id, disease_id), -- Composite primary key
91      FOREIGN KEY(disease_id) REFERENCES disease(disease_id)
92          ON UPDATE CASCADE ON DELETE CASCADE,
93      FOREIGN KEY(person_id) REFERENCES person(person_id)
94          ON UPDATE CASCADE ON DELETE CASCADE
95  );
96
Data Output    Messages    Notifications
CREATE TABLE
Query returned successfully in 318 msec.
```

## Data Loading Process and Operational Queries on OLTP

We sourced data from credible healthcare datasets in CSV format and loaded it into our PostgreSQL database. This setup enables efficient querying and analysis, as demonstrated by the operational queries we executed:

```
-- Query to list all medicines and their indications for 'Tuberculosis'
SELECT
    m.med_name, -- Selecting the name of the medicine
    m.active_ing_name, -- Selecting the active ingredient of the medicine
    i.effectiveness, -- Selecting the effectiveness of the medicine for the disease
    i.indication_date -- Selecting the date when the medicine was indicated
FROM
    medicine m -- Starting from the medicine table
    JOIN indication i ON m.med_id = i.med_id -- Joining with the indication table to match medicines with their indications
WHERE
    i.disease_id = (SELECT disease_id FROM disease WHERE disease_name = 'Tuberculosis'); -- Filtering to include only those medicines indicated for T
```

Output Messages Notifications

Output Messages Notifications

med_name	active_ing_name	effectiveness	indication_date
character varying (25)	character varying (150)	numeric (5,2)	date
Dovprela	Pretomanid	90.00	1998-12-06

```
-- Query to summarize the number of disease incidences by location
SELECT
    l.city_name, -- Selecting the city name
    l.state_name, -- Selecting the state name
    COUNT(p.disease_id) AS disease_count -- Counting the number of diseases for each location
FROM
    patient p -- Starting from the patient table
    JOIN person per ON p.person_id = per.person_id -- Joining with the person table to get the location ID
    JOIN location l ON per.location_id = l.location_id -- Joining with the location table to get city and state names
GROUP BY
    l.city_name, l.state_name; -- Grouping the results by city and state names
```

Output Messages Notifications

Output Messages Notifications

city_name	state_name	disease_count
character varying (100)	character varying (100)	bigint
Charlottesville	Virginia	1
Danbury	Connecticut	1
Shoreline	Washington	1
Waterbury	Connecticut	1
Lynchburg	Virginia	1
Maricopa	Arizona	1
Youngstown	Ohio	1
Daly City	California	1
Eagan	Minnesota	1
Coral Gables	Florida	1
Burnsville	Minnesota	1
Rancho Palos Verdes	California	1
Lafayette	Louisiana	1

```
-- Query to analyze the propensity of races for different diseases
SELECT
  r.race_desc, -- Selecting the race description
  d.disease_name, -- Selecting the disease name
  pr.propensity_value -- Selecting the propensity value indicating how common the disease is among the race
FROM
  propensity pr -- Starting from the propensity table
  JOIN race r ON pr.race_code = r.race_code -- Joining with the race table to get race descriptions
  JOIN disease d ON pr.disease_id = d.disease_id -- Joining with the disease table to get disease names
ORDER BY
  pr.propensity_value DESC; -- Ordering by propensity value to see which combinations have the highest propensity
```

Output Messages Notifications



race_desc	disease_name	propensity_value
character varying (100)	character varying (100)	integer
Black or African American, not of Hispanic origin	Malaria	10
Unspecified	Dengue	10
Hispanic	Trachoma	10
Unspecified	Other intestinal infections	10
Unspecified	Lower respiratory infections	10
Asian	HIV/AIDS	10
White, not of Hispanic origin	Leprosy	10
White, not of Hispanic origin	Intestinal nematode infections	10
Hispanic	Diarrhoeal diseases	9
Unspecified	Hepatitis C (d)	9
Unspecified	Trachoma	9
Unspecified	HIV/AIDS	9
White, not of Hispanic origin	Diarrhoeal diseases	9

```
-- Query to count the number of diseases by their type
SELECT
  dt.disease_type_desc, -- Selecting the description of the disease type
  COUNT(d.disease_id) AS disease_count -- Counting the number of diseases for each type
FROM
  disease d -- Starting from the disease table
  JOIN disease_type dt ON d.disease_type_cd = dt.disease_type_cd -- Joining with the disease_type table to get type descriptions
GROUP BY
  dt.disease_type_desc; -- Grouping the results by disease type descriptions
```

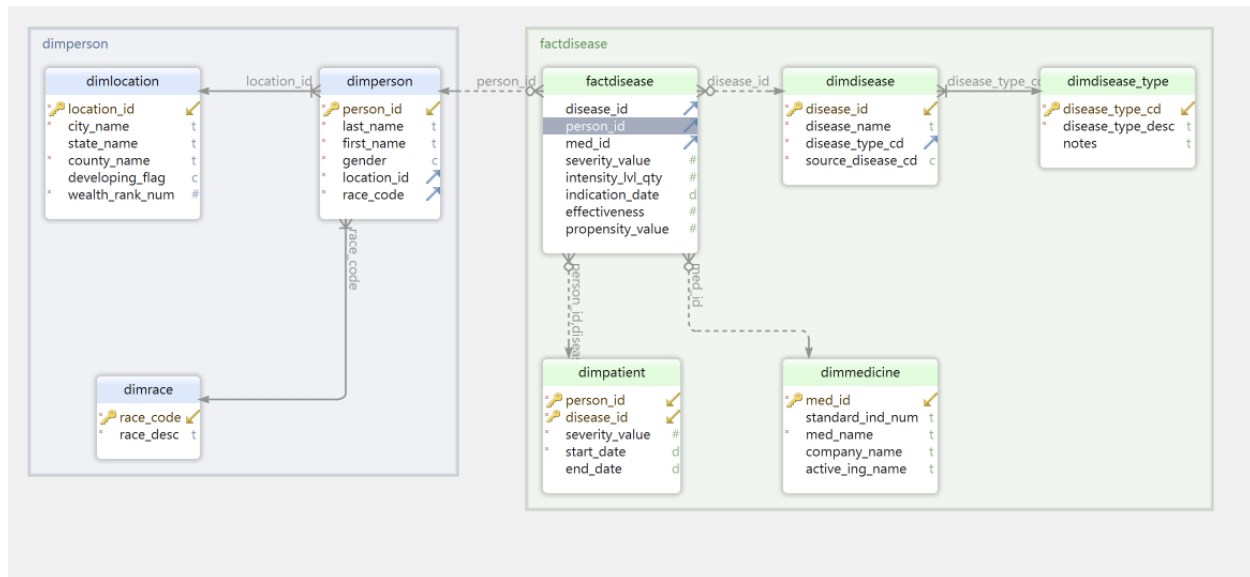
Output Messages Notifications

disease_type_desc character varying (1000)	disease_count bigint
Hepatitis B (d)	1
Meningitis	1
Trachoma	1
Tuberculosis	1
Childhood-cluster diseases	1
Obstructed labour	1
STDs excluding HIV	1
Hypertensive disorders	1
Maternal haemorrhage	1
Japanese encephalitis	1
Intestinal nematode infections	1
Diarrhoeal diseases	1
HIV/AIDS	1
Dengue	1
Maternal sepsis	1

## Dimensional Data Model and ELT in Data Warehouse:

In our data warehouse, we have adopted the Star Schema, comprising both fact and dimension tables. This structure enhances the efficiency of queries and analytical operations, enabling us to handle intricate analytical queries effortlessly.





The data workflow in our system follows an ETL (Extract, Transform, Load) approach. Initially, we extract data from the OLTP (Online Transaction Processing) system. Following this, the data undergoes transformation to align with the dimensional model of our data warehouse. Finally, it's loaded into the warehouse. This methodology guarantees that the data used for analytical queries is tailored for optimal analysis.

Query	Query History
247	-- Creating dimension table for Person in the data warehouse
248	<b>CREATE TABLE</b> disease_dw.dimPerson( 249     Person_Id INT <b>PRIMARY KEY</b> , -- Primary Key for Person 250     Last_Name VARCHAR(50) <b>NOT NULL</b> , -- Last name of the person 251     First_Name VARCHAR(50) <b>NOT NULL</b> , -- First name of the person 252     Gender CHAR(1) <b>NOT NULL</b> , -- Gender of the person 253     Location_Id INT <b>NOT NULL</b> , -- Foreign key to Location 254     Race_Code CHAR(5) <b>NOT NULL</b> , -- Foreign key to Race 255 <b>FOREIGN KEY</b> (Race_Code) <b>REFERENCES</b> disease_dw.dimRace(Race_Code) 256 <b>ON UPDATE CASCADE ON DELETE CASCADE</b> , -- Cascade updates/deletes to Race 257 <b>FOREIGN KEY</b> (Location_Id) <b>REFERENCES</b> disease_dw.dimLocation(Location_Id) 258 <b>ON UPDATE CASCADE ON DELETE CASCADE</b> -- Cascade updates/deletes to Location 259     ); 260
261	-- Populating the dimPerson table from the existing person table
262	<b>INSERT INTO</b> disease_dw.dimPerson
263	<b>SELECT</b> Person_Id, Last_Name, First_Name, Gender, Location_Id, Race_Code <b>FROM</b> person;
264	
267	-- Creating dimension table for Medicine in the data warehouse
268	<b>CREATE TABLE</b> disease_dw.dimMedicine( 269     Med_Id INT <b>PRIMARY KEY</b> , -- Primary Key for Medicine 270     Standard_Ind_Num VARCHAR(250), -- Standard Industry Number 271     Med_Name VARCHAR(25) <b>NOT NULL</b> , -- Name of the Medicine 272     Company_Name VARCHAR(150), -- Name of the Company that manufactures the Medicine 273     Active_Ing_Name VARCHAR(150) -- Active Ingredient in the Medicine 274     ); 275
276	-- Populating the dimMedicine table from the existing medicine table
277	<b>INSERT INTO</b> disease_dw.dimMedicine
278	<b>SELECT</b> Med_Id, Standard_Ind_Num, Med_Name, Company_Name, Active_Ing_Name <b>FROM</b> medicine;
279	
280	
281	-- Drop the existing dimPatient table if it exists
282	<b>DROP TABLE IF EXISTS</b> disease_dw.dimPatient;
283	
284	
285	-- Creating dimension table for Patient in the data warehouse

Servers (1)  
PostgreSQL 15  
Databases (14)  
Assignment3  
Disease  
Casts  
Catalogs  
Event Triggers  
Extensions  
Foreign Data Wrappers  
Languages  
Publications  
Schemas (1)  
public  
Aggregates  
Collations  
Domains  
FTS Configurations  
FTS Dictionaries  
FTS Parsers  
FTS Templates  
Foreign Tables  
Functions  
Materialized Views  
Operators  
Procedures  
Sequences  
Tables (9)  
Trigger Functions  
Types  
Views  
Subscriptions  
GoTData  
aboutme  
afceast  
books  
college\_grouping  
flights  
library  
mydb  
postgres

Disease/postgres@PostgreSQL 15

Query History

```
320 c.med_id AS med_id,  
321 a.severity_value AS severity_value,  
322 b.intensity_lvl_qty AS intensity_lvl_qty,  
323 c.indication_date AS indication_date,  
324 d.propriensity_value AS propriensity_value  
325 FROM  
326 patient a  
327 INNER JOIN disease b ON a.disease_id = b.disease_id  
328 INNER JOIN indication c ON a.disease_id = c.disease_id  
329 INNER JOIN medicine e ON c.med_id = e.med_id  
330 INNER JOIN propensity d ON a.disease_id = d.disease_id;  
331  
332 -- Sample queries to view data from the dimensional model and fact table  
333 SELECT * FROM disease_dw.dimDisease;  
334 SELECT * FROM disease_dw.factDisease;  
335  
336
```

Data Output Messages Notifications

disease_id [PK] integer	disease_name character varying (100)	disease_type_cd character	source_disease_cd character
1	1 Tuberculosis	U003	TB
2	2 STDs excluding HIV	U004	STD
3	3 HIV/AIDS	U009	HIV
4	4 Diarrhoeal diseases	U010	DHD
5	5 Meningitis	U011	MGTS
6	6 Hepatitis B (d)	U017	HPB
7	7 Hepatitis C (d)	U018	HPC
8	8 Malaria	U029	MLR
9	9 Leprosy	U030	LP
10	10 Dengue	U031	DG
11	11 Trachoma	U032	TRC
12	12 Intestinal nematode infections	U043	INI
13	13 Other intestinal infections	U044	OII
14	14 Lower respiratory infections	U045	LRI
15	15 Upper respiratory infections	U046	URI

```
320 c.med_id AS med_id,  
321 a.severity_value AS severity_value,  
322 b.intensity_lvl_qty AS intensity_lvl_qty,  
323 c.indication_date AS indication_date,  
324 d.propriensity_value AS propriensity_value  
325 FROM  
326 patient a  
327 INNER JOIN disease b ON a.disease_id = b.disease_id  
328 INNER JOIN indication c ON a.disease_id = c.disease_id  
329 INNER JOIN medicine e ON c.med_id = e.med_id  
330 INNER JOIN propensity d ON a.disease_id = d.disease_id;  
331  
332 -- Sample queries to view data from the dimensional model and fact table  
333 SELECT * FROM disease_dw.dimDisease;  
334 SELECT * FROM disease_dw.factDisease;  
335  
336
```

Data Output Messages Notifications

disease_id integer	person_id integer	med_id integer	severity_value integer	intensity_lvl_qty integer	indication_date date	effectiveness double precision	propriensity_value integer
1	1	61	3	5	6 1998-12-06	90	7
2	1	61	3	5	6 1998-12-06	90	4
3	1	61	3	5	6 1998-12-06	90	8
4	1	61	3	5	6 1998-12-06	90	5
5	1	61	3	5	6 1998-12-06	90	2
6	1	46	3	1	6 1998-12-06	90	7
7	1	46	3	1	6 1998-12-06	90	4
8	1	46	3	1	6 1998-12-06	90	8
9	1	46	3	1	6 1998-12-06	90	5
10	1	46	3	1	6 1998-12-06	90	2
11	1	31	3	1	6 1998-12-06	90	7
12	1	31	3	1	6 1998-12-06	90	4
13	1	31	3	1	6 1998-12-06	90	8
14	1	31	3	1	6 1998-12-06	90	5
15	1	31	3	1	6 1998-12-06	90	2
16	1	16	3	1	6 1998-12-06	90	7
17	1	16	3	1	6 1998-12-06	90	4
18	1	16	3	1	6 1998-12-06	90	8
19	1	16	3	1	6 1998-12-06	90	5
20	1	16	3	1	6 1998-12-06	90	2

## Analytical Queries:

Our analytical queries delve into diverse facets of disease data. For instance, one such query focuses on examining the occurrence of diseases over time, uncovering trends and patterns. The valuable insights derived from these queries are crucial for shaping healthcare policies and strategies.

### 1. Trend Analysis of Disease Incidence Over Time:

Our analysis involves tracking the annual trends of various diseases, based on the recorded number of cases.

338  
339  
340  
341  
342  
343  
344  
345  
346  
347  
348  
349  
350  
351  
352

```
-- Analytical Queries
-- Trend Analysis of Disease Incidence Over Time:
-- This query analyzes the yearly trend of different diseases based on the number of cases.

SELECT d.disease_name, COUNT(pt.disease_id) AS incidence_count, EXTRACT(YEAR FROM pt.start_date) AS year
FROM patient pt
JOIN disease d ON pt.disease_id = d.disease_id
GROUP BY d.disease_name, EXTRACT(YEAR FROM pt.start_date)
ORDER BY year, incidence_count DESC;
```

Data Output Messages Notifications

	disease_name character varying (100)	incidence_count bigint	year numeric
1	Tuberculosis	5	2023
2	STDs excluding HIV	5	2023
3	Hepatitis C (d)	4	2023
4	Leprosy	4	2023
5	Meningitis	4	2023
6	Lower respiratory infections	4	2023
7	HIV/AIDS	4	2023
8	Trachoma	4	2023
9	Hepatitis B (d)	4	2023
10	Intestinal nematode infections	4	2023
11	Malaria	4	2023
12	Other intestinal infections	4	2023
13	Upper respiratory infections	4	2023
14	Diarrhoeal diseases	4	2023
15	Dengue	4	2023

### 2. Geographical Distribution of Specific Diseases

We implemented analytical queries to identify the distribution of diseases across cities and states in the United States.

Query
Query History

```

354 -- Geographical Distribution of Specific Diseases:
355 -- This query identifies the distribution of diseases across cities and states.
356
357 SELECT l.city_name, l.state_name, d.disease_name, COUNT(p.disease_id) AS cases
358 FROM patient p
359 JOIN person pr ON p.person_id = pr.person_id
360 JOIN location l ON pr.location_id = l.location_id
361 JOIN disease d ON p.disease_id = d.disease_id
362 GROUP BY l.city_name, l.state_name, d.disease_name
363 ORDER BY cases DESC;
364
365
366
367
368

```

Data Output
Messages
Notifications

	city_name character varying (100)	state_name character varying (100)	disease_name character varying (100)	cases bigint
1	Elkhart	Indiana	Hepatitis C (d)	1
2	Santa Rosa	California	Tuberculosis	1
3	Rancho Santa Margarita	California	Intestinal nematode infections	1
4	Worcester	Massachusetts	Leprosy	1
5	Cleveland	Ohio	Malaria	1
6	Nampa	Idaho	Diarrhoeal diseases	1
7	Perris	California	Hepatitis C (d)	1
8	Fargo	North Dakota	Leprosy	1
9	Montgomery	Alabama	Tuberculosis	1
10	Roswell	New Mexico	Upper respiratory infections	1
11	Charlottesville	Virginia	Lower respiratory infections	1
12	Hayward	California	Dengue	1
13	Burnsville	Minnesota	Other intestinal infections	1
14	Omaha	Nebraska	Trachoma	1
15	Waterbury	Connecticut	STDs excluding HIV	1
16	Marietta	Georgia	Leprosy	1
17	Daly City	California	Leprosy	1
18	Taylor	Michigan	Dengue	1
19	Fort Wayne	Indiana	Diarrhoeal diseases	1
20	Woburn	Massachusetts	HIV/AIDS	1
21	Palm Springs	California	Lower respiratory infections	1
22	Shoreline	Washington	Meningitis	1

### 3. Medication Effectiveness Analysis

This analytical query evaluates the effectiveness of medications for different diseases.

```

367
368 -- Medication Effectiveness Analysis:
369 -- This query evaluates the effectiveness of medications for different diseases.
370
371 SELECT d.disease_name, m.med_name, AVG(i.effectiveness) AS average_effectiveness
372 FROM indication i
373 JOIN disease d ON i.disease_id = d.disease_id
374 JOIN medicine m ON i.med_id = m.med_id
375 GROUP BY d.disease_name, m.med_name
376 ORDER BY d.disease_name, average_effectiveness DESC;
377
378

```

Data Output Messages Notifications

	disease_name character varying (100)	med_name character varying (25)	average_effectiveness numeric
1	Dengue	Dengvaxia	94.0000000000000000
2	Diarrhoeal diseases	Truberzi	93.0000000000000000
3	Hepatitis B (d)	Procomvax	88.0000000000000000
4	Hepatitis C (d)	Hexyon	85.0000000000000000
5	HIV/AIDS	Biktarvy	70.0000000000000000
6	Intestinal nematode infections	Ayvakyt	87.0000000000000000
7	Leprosy	Synagis	75.0000000000000000
8	Lower respiratory infections	Esbriet	97.0000000000000000
9	Malaria	Artesunate Amivas	93.0000000000000000
10	Meningitis	Trumenba	85.0000000000000000
11	Other intestinal infections	Sutent	89.0000000000000000
12	STDs excluding HIV	penicillin g benzathine	89.0000000000000000
13	Trachoma	Evotaz	83.0000000000000000
14	Tuberculosis	Dovprela	90.0000000000000000
15	Upper respiratory infections	Cayston	90.0000000000000000

#### 4. Correlation between Race and Disease Propensity.

We tried to examine if certain races have higher propensities for specific diseases.

```

382 --Correlation Between Race and Disease Propensity:
383 -- This query examines if certain races have higher propensities for specific diseases.
384
385 SELECT r.race_desc, d.disease_name, AVG(p.propensity_value) AS average_propensity
386 FROM propensity p
387 JOIN race r ON p.race_code = r.race_code
388 JOIN disease d ON p.disease_id = d.disease_id
389 GROUP BY r.race_desc, d.disease_name
390 ORDER BY average_propensity DESC;
391
392
393

```

Data Output Messages Notifications

	race_desc character varying (100)	disease_name character varying (100)	average_propensity numeric
1	White, not of Hispanic origin	Intestinal nematode infections	10.0000000000000000
2	Black or African American, not of Hispanic origin	Malaria	10.0000000000000000
3	Unspecified	Lower respiratory infections	10.0000000000000000
4	Unspecified	Dengue	10.0000000000000000
5	Asian	HIV/AIDS	10.0000000000000000
6	Unspecified	Other intestinal infections	10.0000000000000000
7	Hispanic	Trachoma	10.0000000000000000
8	White, not of Hispanic origin	Leprosy	10.0000000000000000
9	Unspecified	HIV/AIDS	9.0000000000000000
10	White, not of Hispanic origin	Diarrhoeal diseases	9.0000000000000000
11	Unspecified	Hepatitis C (d)	9.0000000000000000
12	Hispanic	Diarrhoeal diseases	9.0000000000000000
13	Unspecified	Trachoma	9.0000000000000000
14	Black or African American, not of Hispanic origin	Meningitis	8.0000000000000000
15	Black or African American, not of Hispanic origin	Other intestinal infections	8.0000000000000000
16	Black or African American, not of Hispanic origin	Diarrhoeal diseases	8.0000000000000000
17	Black or African American, not of Hispanic origin	STDs excluding HIV	8.0000000000000000
18	Black or African American, not of Hispanic origin	Tuberculosis	8.0000000000000000
19	Hispanic	Malaria	8.0000000000000000
20	White, not of Hispanic origin	Malaria	8.0000000000000000
21	Unspecified	Tuberculosis	7.0000000000000000
22	Asian	Malaria	7.0000000000000000

## 5. Patient Demographics for Specific Disease

The below analytical query provides demographical insights such as gender and race for patients with specific diseases.

```

393
394 -- Patient Demographics for Specific Diseases:
395 -- This query provides demographic insights (gender, race) for patients with specific diseases.
396 SELECT d.disease_name, pr.gender, r.race_desc, COUNT(*) AS patient_count
397 FROM patient p
398 JOIN person pr ON p.person_id = pr.person_id
399 JOIN race r ON pr.race_code = r.race_code
400 JOIN disease d ON p.disease_id = d.disease_id
401 GROUP BY d.disease_name, pr.gender, r.race_desc;
402
403
404

```

Data Output Messages Notifications

	disease_name character varying (100)	gender character	race_desc character varying (100)	patient_count bigint
1	HIV/AIDS	M	White, not of Hispanic origin	1
2	Lower respiratory infections	F	Black or African American, not of Hispanic origin	1
3	STDs excluding HIV	F	Unspecified	1
4	Malaria	F	Black or African American, not of Hispanic origin	1
5	STDs excluding HIV	F	Black or African American, not of Hispanic origin	1
6	Dengue	M	Asian	1
7	Other intestinal infections	F	Hispanic	1
8	Dengue	F	Unspecified	1
9	Leprosy	M	Hispanic	1
10	Leprosy	M	Black or African American, not of Hispanic origin	1
11	Trachoma	M	Black or African American, not of Hispanic origin	2
12	Tuberculosis	M	White, not of Hispanic origin	1
13	Dengue	M	Black or African American, not of Hispanic origin	1
14	Meningitis	M	White, not of Hispanic origin	1
15	Lower respiratory infections	F	Hispanic	1
16	Malaria	M	White, not of Hispanic origin	1
17	STDs excluding HIV	F	White, not of Hispanic origin	1
18	Hepatitis C (d)	F	White, not of Hispanic origin	1
19	Intestinal nematode infections	M	Black or African American, not of Hispanic origin	1
20	Hepatitis C (d)	M	Black or African American, not of Hispanic origin	2
21	Hepatitis B (d)	M	Black or African American, not of Hispanic origin	2
22	Diarrhoeal diseases	M	Black or African American, not of Hispanic origin	2

## NoSQL Database Structure for Disease Model

### MongoDB (Document-Oriented NoSQL Database):

In MongoDB, instead of tables, we have collections of documents. Each document can have a different structure. A patient document in MongoDB can look like the following:

```
{
  "patient_id": "P123",
  "name": "John Doe",
  "age": 45,
  "diseases": [
    {
      "disease_id": "D456",
      "name": "Diabetes",
      "treatments": [
        {
          "treatment_id": "T789",
          "name": "Insulin Therapy",
          "effectiveness": 85
        }
      ]
    }
  ]
}
```

- For instance, a Patient document could include embedded documents for Diseases and Treatments, allowing for a more holistic view of the patient's health record in a single document. This structure allows adding new diseases or treatments to a patient's record without altering the schema for other patients.
- This model is flexible; as new diseases or treatment modalities are introduced, they can easily be added to the patient's record without altering the structure of other documents or collections.



Advantages of MongoDB (Document-Oriented NoSQL Database):

**Documents and Collections:** Each document can have a varied structure, containing all the data related to a single entity, which allows embedding arrays and other documents directly.

**Denormalization:** Related data might be stored together in the same document for quicker read access, reducing the need for joins

**Schema-less Design:** The schema can evolve over time without requiring migrations, making it easier to adapt to changing data needs.

### Neo4J (Graph-Based NoSQL Database):

In Neo4J, data is stored as nodes and edges. Nodes represent entities and edges represent relationships. For the disease model, patients, diseases, and treatments would be nodes. Relationships like "HAS\_DISEASE" or "UNDERGOES\_TREATMENT" would be edges connecting these nodes. This structure is particularly powerful for querying complex relationships, like tracing the spread of a disease through a network of patients or analyzing the efficacy of treatments across various demographics.

**Patient Node:** Labeled as "Patient" with properties such as name: "John Doe".

**Disease Node:** Labeled as "Disease" with properties like name: "Diabetes".

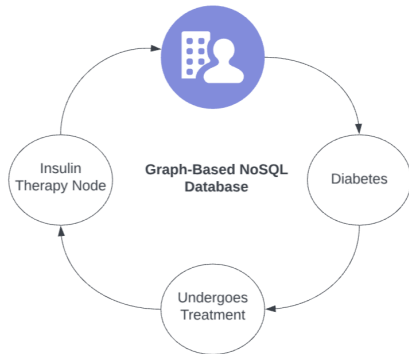
**Treatment Node:** Labeled as "Treatment" with properties such as name: "Insulin Therapy".

### Relationships:

HAS\_DISEASE Relationship connects "John Doe" (Patient Node) to "Diabetes" (Disease Node). This relationship indicates that John Doe has been diagnosed with Diabetes.

UNDERGOES\_TREATMENT Relationship connects "Diabetes" (Disease Node) to "Insulin Therapy" (Treatment Node). This relationship shows that the standard treatment for Diabetes in this case is Insulin Therapy.

The overall structure forms a chain from the patient, through their disease, to their treatment, clearly and intuitively representing the relationships between these entities.



### Neo4J (Graph-Based NoSQL Database):

**Nodes and Relationships:** Data is structured as nodes (entities) and edges (relationships), which is ideal for datasets where relationships are as important as the data itself.

**Property Graph:** Both nodes and relationships can have properties, allowing them to store data.

**Schema-free:** There's no fixed schema, enabling you to add different properties to nodes and relationships as needed.

### AWS Architecture for Disease Model

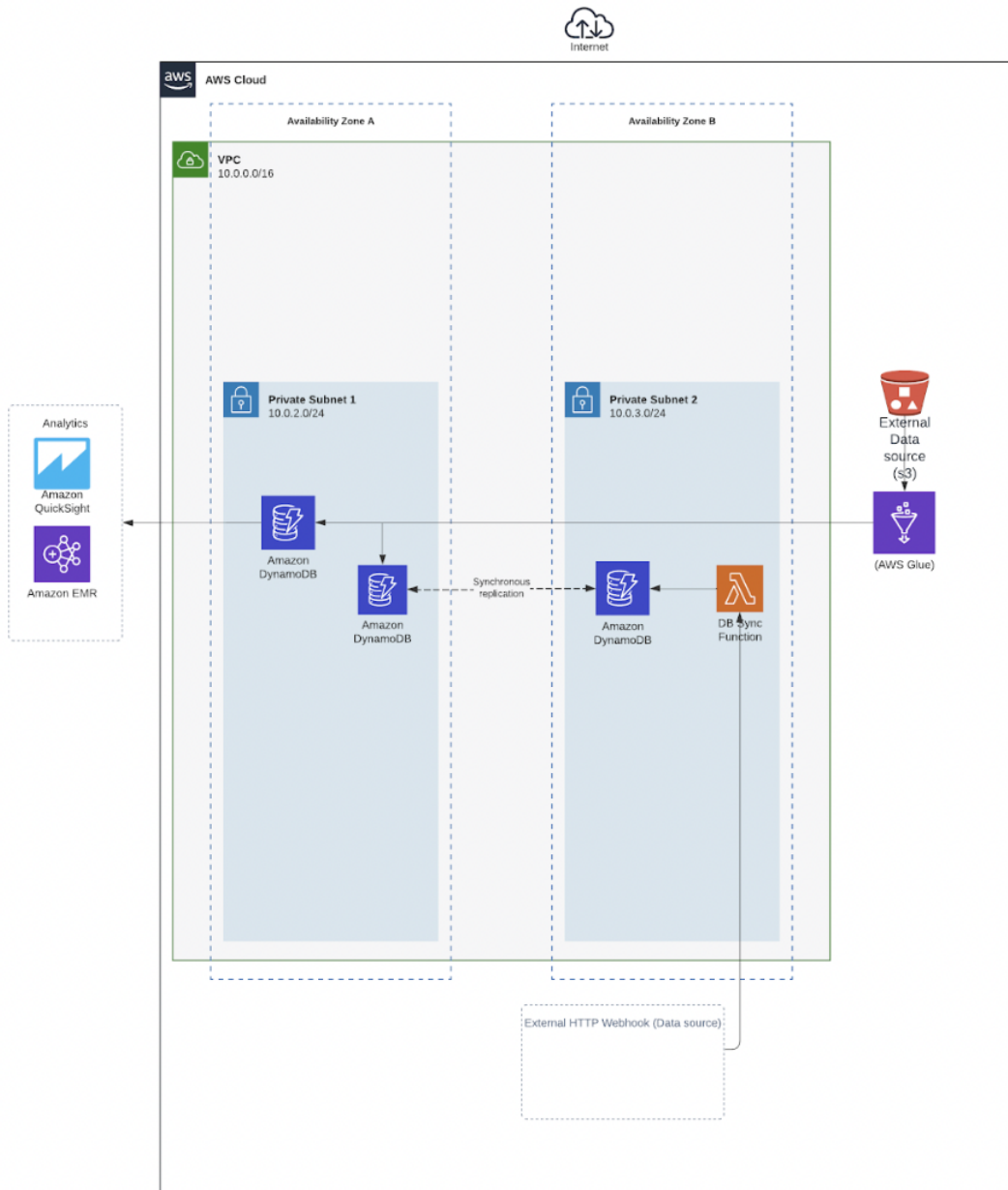
#### Components of AWS Architecture

1. Amazon EC2: Host application in the cloud.
2. Amazon S3: Store backups, logs, and other files.
3. AWS Glue: For ETL process

4. AWS DynamoDB: As NoSQL databases for different types of data storage needs.

5. AWS Lambda: To run code in response to triggers (such as changes in data or system states).

6. Amazon EMR on AWS: To handle data-processing jobs



## Data Processing and Loading:

1. **Batch Processing (AWS Glue):** Regularly scheduled ETL jobs that process historical patient data for analysis, for example, aggregating monthly treatment success rates.
2. **Real-Time Processing (AWS Lambda and Amazon Kinesis):** Real-time monitoring of patient vitals streamed through Kinesis, with Lambda functions triggering alerts for abnormal readings.
3. **Ensuring Resilience, Performance, and Security:**
  - a) **Resilience:** Use of Amazon RDS Multiple availability zones (AZ) deployments ensures that if one AZ goes down, the database can failover to another, ensuring continuity of critical healthcare applications.
  - b) **Performance:** Utilize Amazon ElastiCache to store frequently accessed data like common disease-treatment mappings, reducing the load on the main database.
  - c) **Security:** Implement strict IAM roles to control access to patient data, ensuring that only authorized personnel can view or modify sensitive information.

## Scalability

Snowflake	PostgreSQL
<b>Scenario:</b> Suppose our disease model database needs to rapidly scale up for a project analyzing extensive data from thousands of patients to identify genetic predispositions to certain diseases.	<b>Scenario:</b> The same genomic data analysis in a PostgreSQL-based warehouse.
<b>Snowflake's Approach:</b> In Snowflake, we can instantly scale up computational resources (warehouses) to handle this increased demand without impacting your storage costs. This means we can run complex genomic data analyses and queries concurrently without performance degradation.	<b>PostgreSQL's Approach:</b> To manage this, we might need to provision additional hardware or implement sharding (distributing the database load across multiple machines or partitions), which can be complex and time-consuming.

## Maintenance and Performance

Snowflake	PostgreSQL
<b>Scenario:</b> Running a complex query to correlate genetic markers with the efficacy of cancer treatments.	<b>Scenario:</b> The same complex query in a PostgreSQL environment.
<b>Snowflake's Approach:</b> Snowflake automatically optimizes query performance. Its architecture separates compute and storage, enabling efficient query execution without manual tuning.	<b>PostgreSQL's Approach:</b> This would likely require manual performance tuning, such as creating and maintaining indexes, which can be labor-intensive and require deep database expertise.
<b>Benefit:</b> This is particularly advantageous for dynamically changing query patterns often seen in medical research.	<b>Challenge:</b> For complex and evolving queries, this manual tuning can become a bottleneck in data analysis.

## Data Sharing and Security

Snowflake	PostgreSQL
<b>Scenario:</b> <i>Sharing anonymized patient treatment outcomes with external research institutions for collaborative studies.</i>	<b>Scenario:</b> <i>Attempting the same data sharing using PostgreSQL.</i>
<b>Snowflake's Approach:</b> Snowflake allows us to securely share a portion of our data warehouse, like specific datasets, without replicating the data. We can grant access to external researchers to these datasets with fine-grained control.	<b>PostgreSQL's Approach:</b> Typically, this involves creating database replicas or exporting data, which can be complex to manage and may introduce security vulnerabilities.
<b>Benefit:</b> This facilitates secure and efficient collaboration without increasing the risk of data breaches or duplicating data.	<b>Challenge:</b> This method increases the complexity of data management and can pose challenges in ensuring data security and compliance, especially crucial in healthcare data handling.

In summary, while PostgreSQL is a powerful and widely-used relational database system, Snowflake offers specific advantages in a data warehousing context, particularly for large-scale, complex, and collaborative disease model analyses. Snowflake's architecture provides scalability, ease of maintenance, and secure data sharing capabilities, which are highly beneficial for dynamic and data-intensive fields like healthcare research and genomics.

## Conclusion:

In conclusion, our project successfully designed and implemented a comprehensive database system focused on managing and analyzing disease-related data. By integrating principles of data modeling, database design, and analytical techniques, we developed a robust tool capable of providing critical insights into disease patterns, treatment effectiveness, and demographic distributions.

Utilizing a relational database setup in PostgreSQL allowed us to maintain referential integrity and optimize query performance. Our structured approach, from the conceptual ER model to the physical data dictionary, ensured a clear understanding of the database's organization, enabling efficient data handling and accurate insights.

The implementation of ETL processes facilitated a seamless flow of data from the OLTP system to our data warehouse, adhering to a dimensional model that supports complex analytical queries. This structured data workflow proved to be essential for optimizing the data for analysis, allowing for sophisticated trend analysis and demographic studies.

Our exploration into NoSQL databases like MongoDB and Neo4J demonstrated the flexibility and scalability of schema-less and graph-based models, respectively. These systems provided alternative methods for managing and querying data, showing promise for handling complex, unstructured datasets.

Incorporating AWS cloud services into our architecture offered robust, secure, and real-time data processing capabilities. Services like AWS Glue, Lambda, Kinesis, RDS, and ElastiCache contributed to a highly resilient, performant, and secure infrastructure.

Furthermore, the comparison between Snowflake and PostgreSQL highlighted the distinct advantages of Snowflake in terms of scalability, maintenance, and data sharing, particularly in the context of large-scale, collaborative research efforts.

Overall, our project stands as a testament to the power of modern database technologies in transforming healthcare data into actionable insights, thus supporting the advancement of healthcare decisions and strategies. The methodologies and technologies we've employed demonstrate the potential to significantly impact the



healthcare industry by enhancing the ability to conduct in-depth analyses of disease trends and treatment outcomes.