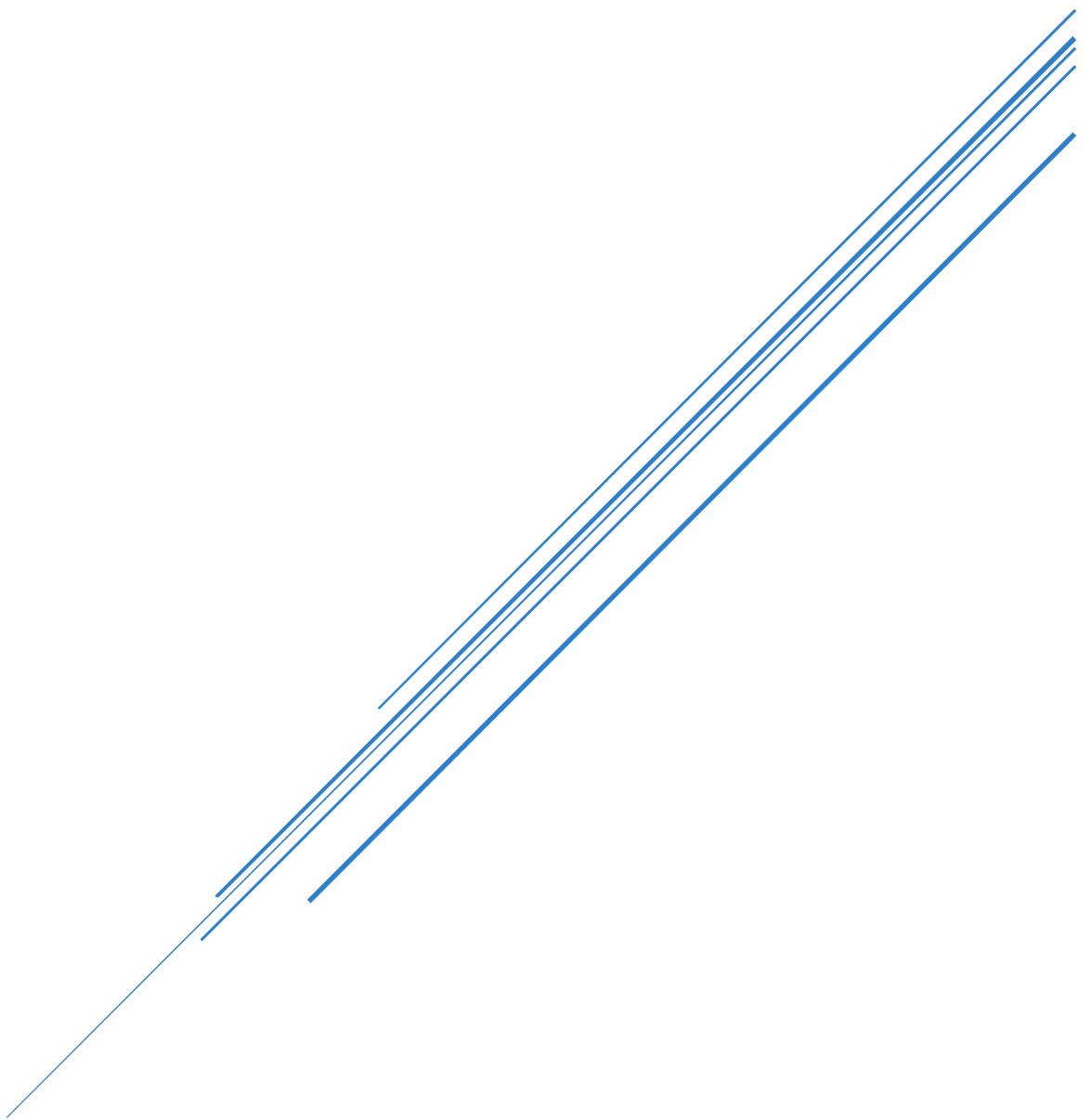


HealthFit Database Design

Design, Implementation, and Optimization for Analytics



Executive Summary – HealthFit Database Design

Author: Joanna Ronchi

Date: May 24, 2025

Business Problem

HealthFit, a health technology company, needed to upgrade its data infrastructure to keep pace with rapid business growth and the increasing complexity of wearable and clinical health data. The company's existing relational database model created data silos and scalability bottlenecks, which delayed critical health recommendations and reduced the value of its flagship HealthTrack platform.

Business Needs

- Real-Time Data Monitoring: continuous ingestion and analysis of health data from wearable devices.
- Scalability: support for millions of users with rapidly growing, diverse datasets.
- Integration of Heterogeneous Sources: combining biometric, clinical, and lifestyle data into unified records.
- Compliance: meeting stringent data privacy and security requirements under HIPAA and GDPR.

Proposed Solution

After evaluating options, HealthFit adopted a NoSQL document database (MongoDB) as the optimal solution. MongoDB's schema flexibility, sharding capability, and aggregation pipelines address the challenges of scalability, real-time analysis, and complex data structures.

Key Solution Elements

- Schema Flexibility: store nested, semi-structured data without rigid relational schemas.
- Horizontal Scaling: use sharding to distribute traffic across multiple servers.
- Real-Time Analysis: leverage aggregation pipelines to provide instant health alerts and insights.
- Data Security & Privacy: implement encryption, role-based access control, audit logging, and anonymization techniques.

Business Impact

- The database design directly supports HealthFit's mission to provide actionable, data-driven health recommendations to patients and providers.
- Improved Patient Care: faster generation of health alerts (e.g., abnormal heart rates, early warning signs).
- Operational Efficiency: reduced processing delays through scalable infrastructure.
- Compliance Assurance: lower legal and financial risk by aligning with healthcare privacy regulations.
- Business Growth Enablement: infrastructure capable of supporting millions of users and expanding product lines.

Conclusion

The HealthFit database design demonstrates the intersection of business analysis and data science. By aligning technical implementation with strategic business needs, the solution ensures HealthFit can deliver timely, accurate, and secure health insights at scale.

Table of Contents

Part 1: HealthFit Design Document.....	3
A1. Description of the business problem	3
Supporting the Business Mission	3
Top Business Needs	3
Top Solution Considerations	3
A2. Justification for NoSQL database solution for the business problem.....	3
Schema Flexibility.....	3
Horizontal Scalability	3
Supporting Real-Time Data Performance.....	3
Data Structure Compatibility.....	3
A3. Proposed NoSQL database type to solve the identified business problem.....	4
A4. How the business data will be used within the database solution.....	4
B. How the proposed database design addresses scalability concerns.....	4
NoSQL Database.....	4
Sharding and Horizontal Scaling.....	4
Real-Time Monitoring	4
Flexibility	4
Data Durability.....	5
Real-Time Analysis.....	5
C. Privacy and security measures	5
Compliance with HIPPA and GDPR.....	5
Security and Privacy of Health Data	5
Role-Based Access Control and Audit Logging	5
Backups	6
Specific MongoDB Features Strengthening Security.....	6
Part 2: Implementation	6
D1. Database instance screenshots (Python and JavaScript).....	6
Database Script Created in Visual Studio:	7
Screenshot of Database in MongoDB:	7
D2. Script to insert and map the data records into the database instance.	8
Successful Upload of Medical Data into D597_Task_2:	8
Successful Upload of Fitness Tracker Data into D597_Task_2:	8
MongoDB Screenshot – Successful Mapping of Medical Data into D597_Task_2.....	9
MongoDB Screenshot – Successful Mapping of Fitness Trackers into D597_Task_2.....	9
D3. Query Scripts.....	10
QUERY 1	10
QUERY 2	11
QUERY 3	13
D4. Optimization Techniques.....	16
QUERY 1 OPTIMIZATION.....	16

QUERY 2 OPTIMIZATION.....	18
QUERY 3 OPTIMIZATION.....	20
References.....	21

Part 1: HealthFit Design Document

A1. Description of the business problem

Supporting the Business Mission

HealthFit is a health-focused business selling a data solution that enables customers to track their health data. The company provides data-driven recommendations for customers and their healthcare providers to help them improve health.

Top Business Needs

The business needs to upgrade its current database solution to address scalability, improve performance, and integrate data from heterogeneous sources. The top business problems include data silos, resulting in delaying health recommendations, and reducing the value of Health Fit's HealthTrack platform. The business also suffers from scalability bottlenecks due to rapid growth. The bottlenecks slow down the processing of data and hinder business progress. Health Track's real-time data processing issues, caused by large volumes of incoming data, create issues during the processing stage of the data lifecycle, resulting in processing delays or incorrect outputs.

Top Solution Considerations

To solve the problem, the business needs to increase the capacity of computing resources. HealthFit needs the new system to integrate complex facts – including biometric, clinical, and lifestyle data - for accurate analysis and insights. The business needs a solution to integrate and unify multiple source data, handle nested information, and unify patient information from various sources into a single document. HealthFit Innovations needs to connect the siloed and disconnected data to provide value for customers and ensure the accuracy of crucial health insights.

A2. Justification for NoSQL database solution for the business problem

Schema Flexibility

A NoSQL database will provide schema flexibility for storing varying data from multiple sources and categories. Since the company's health data varies dramatically between heart rate, blood glucose, step count, sleep and activity patterns, blood oxygen saturation, body temperature, or GPS data (for measuring distance traveled), the relational data model cannot effectively accommodate HealthFit Innovations' business needs.

Horizontal Scalability

No SQL solutions will allow the system to scale out using sharding, enabling scalability as the volume of data grows. NoSQL databases fit nicely into dynamically changing, semi-structured data from wearable devices because they do not sacrifice performance as they distribute the load between several nodes (Dey, 2020).

Supporting Real-Time Data Performance

The business collects real-time data from wearable devices, needs to support real-time data monitoring, and requires high write throughput.

Data Structure Compatibility

The NoSQL database will support document-based data more effectively than relational models by supporting nested data and arrays. The NoSQL data model directly represents hierarchical relationships and manages complex structures by including nested arrays and objects in the data

documents. The structure eliminates the need for complex tables and designs in relational databases, improving performance and easier access to data with simpler queries (Blefqih, 2024).

A3. Proposed NoSQL database type to solve the identified business problem

MongoDB will support storing and analyzing semi-structured customer health data, such as JSON documents from wearable devices and applications. The database will give the business flexibility and support inputs from various users, including healthcare providers and health-data reporting systems. MongoDB will support real-time data reporting and provide performance optimization capability, support indexing, aggregation functionality, and real-time analytics (Wolniak, 2023).

A4. How the business data will be used within the database solution

The data from HealthTrack will be used for real-time health-data monitoring and analysis to quickly identify and provide critical health alert data, such as heart rate measurements and alerts to customers. The database's functionality will enable the generation of health insights based on past data and identify potential health threats using predictive analytics. The solution will use data to analyze early warning signs of illness and allow users to make health-related decisions. The system will enable dashboard generation for various users, including patients and their healthcare providers, to manage and improve patient health, track improvement progress, and generate customized health tips.

B. How the proposed database design addresses scalability concerns

NoSQL Database

Rigid, relational schema designs are unsuitable for the HealthTrack business, primarily due to inflexibility. HealthFit Innovations solution needs a scalable database that supports millions of users uploading real-time, continuous, and diverse data related to varying health metrics. A document-oriented model offered by MongoDB will accommodate the scalability and performance and offer flexibility for the business to handle the massive amounts of data coming into the system.

Sharding and Horizontal Scaling

MongoDB allows for the distribution of data across servers, which is known as sharding. The nodes distribute data based on parameters such as `user_id` or `device_id`. The database will allow Heath Fit innovations to add more users and increase without sacrificing performance by adding servers to accommodate the growth. The additional servers can handle specified portions of read/write traffic, enabling throughput and preventing failure due to traffic increases. Health monitoring data from mobile devices is high-velocity and requires strong, scalable, and flexible systems. NoSQL databases support mobile device health data and efficiently scale (He, 2021).

Real-Time Monitoring

NoSQL databases allow for high-speed writes of health data from wearable devices and can effectively store time-series data. Healthcare IoT devices handle data better than traditional relational models (Dey, 2020).

Flexibility

The variability of health data requires flexible data schemas to monitor diverse and evolving data structures such as heart rate metrics, blood glucose levels, medication compliance logs, or sleep

and activity data. Each type of health data requires a customized model, and the NoSQL database will enable the business to structure this data for accurate analysis (Alshammari, 2019). MongoDB offers a schema-less design and adds data without typical relational model complexities. The flexibility of the MongoDB database will enable agile development and growth of the data from various users as the business expands its customer base.

Data Durability

NoSQL models provide data redundancy by storing data in multiple copies across servers and safeguarding the data against loss. In case of failure of a single node, another server assumes control over the data, minimizing downtime. The fault tolerance capability supports data availability, which is particularly important in environments that handle critical health data.

Real-Time Analysis

MongoDB offers a unique aggregation pipeline capability to process, transform, and analyze data. HealthFit data can use this feature to quickly and efficiently generate insights from real-time incoming patient/customer data and enables quick generation of trigger alerts such as abnormal heart rate (Zhang, 2019). Aggregation pipelines are efficient, flexible, and important, especially for health monitoring and IoT in healthcare settings. The business can provide critical health insights, offer value to customers, and advertise the feature's availability for effective marketing campaigns.

C. Privacy and security measures

Compliance with HIPPA and GDPR

The U.S. Health Insurance Portability and Accountability Act (HIPPA) and the European Union's General Data Protection Regulation (GDPR) Act require businesses dealing with health information to implement guardrails to protect the safety and privacy of individual health information. Pseudonymization replaces data that identifies individuals with fields such as patientID or userID. Anonymization is another technique HealthFit should use to perform the analysis while protecting individual privacy. The data masking methods will enable safe data sharing and meet the legal obligations of a company dealing with highly sensitive health data (Wang, Zhang, 2020).

Security and Privacy of Health Data

HealthFit Innovations has legal obligations to protect the privacy and security of its users. The database must provide field-level security to protect patient names, social security numbers, and other sensitive data (Fernandez, 2013). HealthFit needs to protect the data at rest and in transit. MongoDB supports the industry-standard AES-256 and Transport Level Security, offering additional client-level encryption (MongoDB, n.d., p. 1).

Role-Based Access Control and Audit Logging

The business must follow the principle of least privilege to protect patient and user personal information based on necessity, especially in cloud-based health data (Gayanajake, 2017). Users and patients should be able to access all their data and assign permissions to their healthcare providers. Data analysts and scientists should only have access to anonymous and pseudonymized data sets. The business will implement audit logging to track and monitor who and when accessed the data and made changes. Configuring the database with RBAC (Role-Based Access Controls) is important. An important MongoDB feature that will enable data confidentiality, integrity, and

access controls is the system built-in `security.authorization: enabled` setting, which configures MongoDB to require users' authentication based on their roles, for example admin, read-only, or root access.

Backups

MongoDB will provide automated backups to prevent data loss, which is critical in healthcare environments. The database will redundantly and securely store data in several locations using snapshots and encryption. MongoDB enables regular, scheduled backups at convenient times. The database offers replicas where one node writes, and a secondary node reads the data, preventing loss and failure. MongoDB documentation has extensive information on replication functionality, assuring the business that the solution will provide real-time configured backups and increase the security of data (MongoDB, n.d., p. 2).

Specific MongoDB Features Strengthening Security

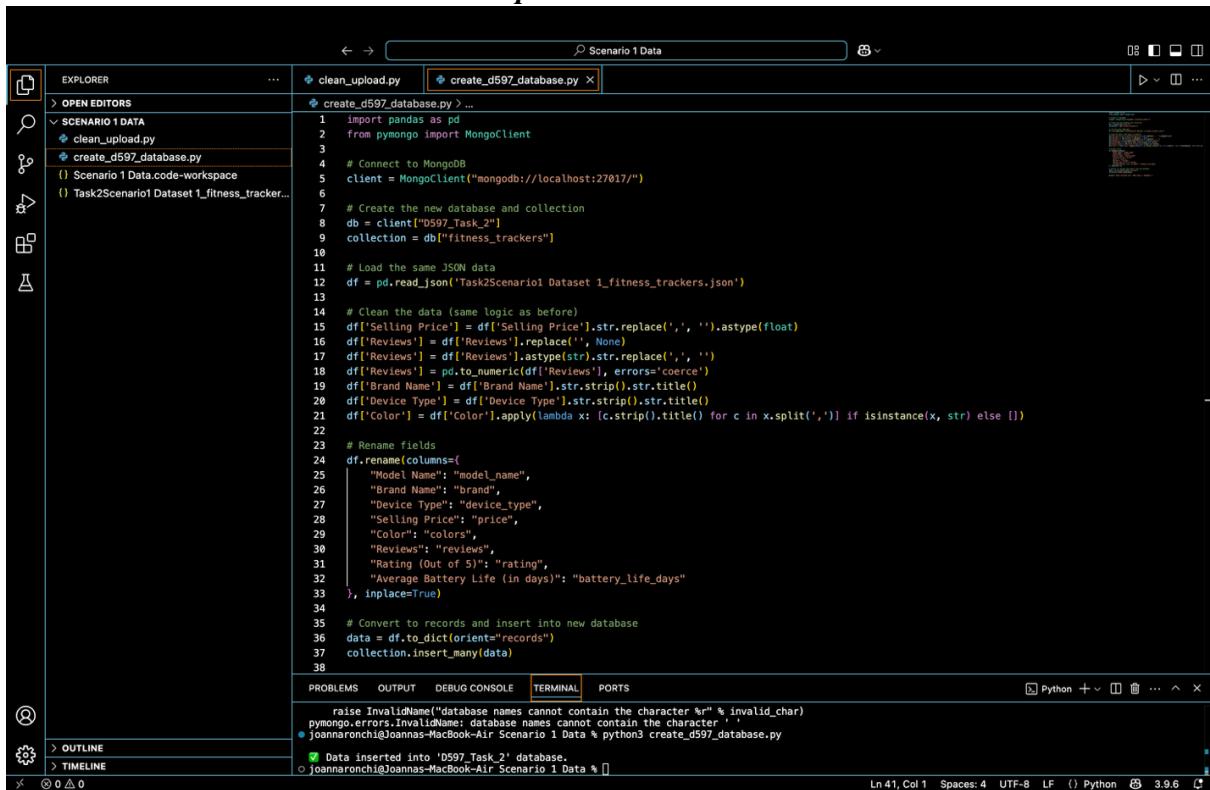
MongoDB comes with built-in security features. It is important to configure settings to prevent attacks and protect sensitive data. The recommended settings include:

- `net.port`—The feature changes the default listens to on position, reducing cyberattack vulnerabilities that scan MongoDB ports.
- `net.bindIp`—This feature limits access to the database to only specified IP addresses and prevents unauthorized access and misappropriated data use.
- `javascript-enabled` – set to disable or prevent Java script injection attacks. MongoDB will reduce the attack surface from untrusted queries or that are not sanitized inputs (MongoDB, n.d., p. 3)

Part 2: Implementation

D1. Database instance screenshots (Python and JavaScript)

Database Script Created in Visual Studio:



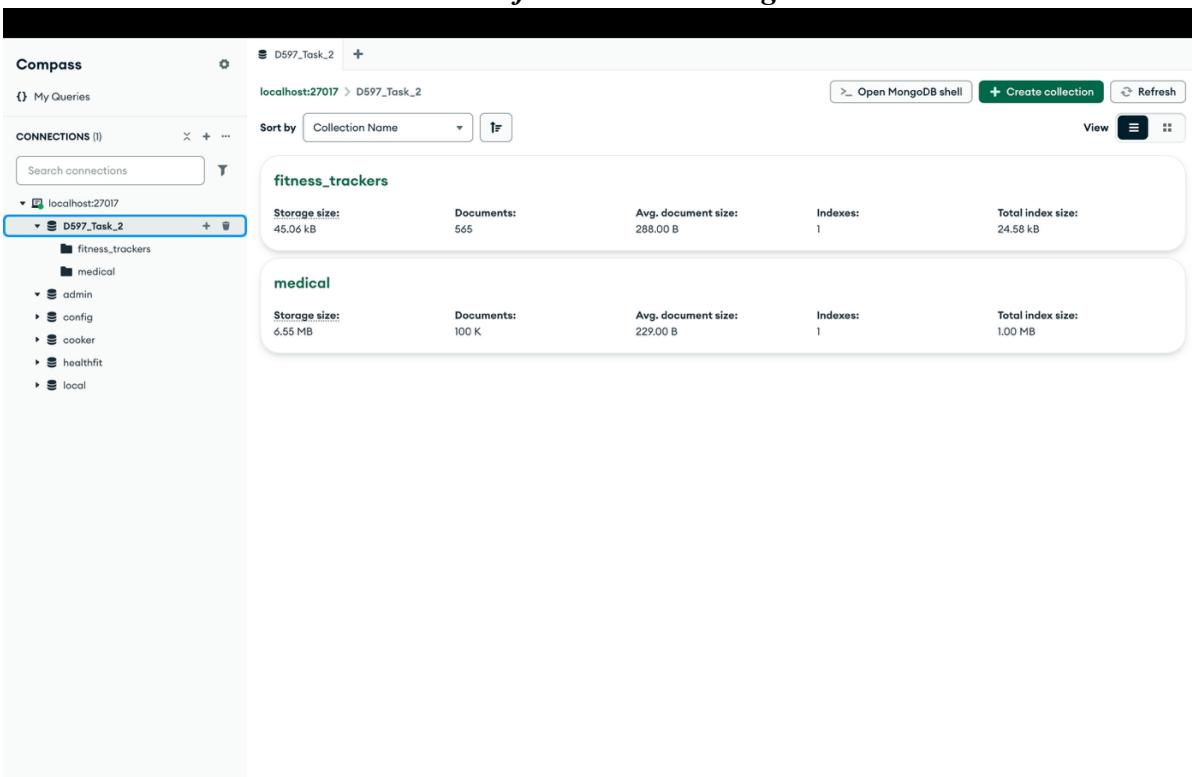
```

Scenario 1 Data> Scenario 1 Data> clean_upload.py > create_d597_database.py > ...
1 import pandas as pd
2 from pymongo import MongoClient
3
4 # Connect to MongoDB
5 client = MongoClient("mongodb://localhost:27017/")
6
7 # Create the new database and collection
8 db = client["D597_Task_2"]
9 collection = db["fitness_trackers"]
10
11 # Load the same JSON data
12 df = pd.read_json('Task2Scenario1 Dataset 1_fitness_trackers.json')
13
14 # Clean the data (same logic as before)
15 df['Selling Price'] = df['Selling Price'].str.replace(',', '').astype(float)
16 df['Reviews'] = df['Reviews'].replace('', None)
17 df['Reviews'] = df['Reviews'].astype(str).str.replace(',', '')
18 df['Reviews'] = pd.to_numeric(df['Reviews'], errors='coerce')
19 df['Brand Name'] = df['Brand Name'].str.strip().str.title()
20 df['Device Type'] = df['Device Type'].str.strip().str.title()
21 df['Color'] = df['Color'].apply(lambda x: [c.strip().title() for c in x.split(',')]) if isinstance(x, str) else []
22
23 # Rename fields
24 df.rename(columns={
25     "Model Name": "model_name",
26     "Brand Name": "brand",
27     "Device Type": "device_type",
28     "Selling Price": "price",
29     "Color": "colors",
30     "Reviews": "reviews",
31     "Rating (Out of 5)": "rating",
32     "Average Battery Life (in days)": "battery_life_days"
33 }, inplace=True)
34
35 # Convert to records and insert into new database
36 data = df.to_dict(orient="records")
37 collection.insert_many(data)
38
raise InvalidName("database names cannot contain the character %s" % invalid_char)
pymongo.errors.InvalidName: database names cannot contain the character ' '
joannarochi@Joannas-MacBook-Air: Scenario 1 Data % python3 create_d597_database.py
Data inserted into 'D597_Task_2' database.
joannarochi@Joannas-MacBook-Air: Scenario 1 Data %

```

The screenshot shows a Visual Studio Code interface. The left sidebar has sections for EXPLORER, OPEN EDITORS, SCENARIO 1 DATA, and OUTLINE. The SCENARIO 1 DATA section contains files: clean_upload.py and create_d597_database.py. The right pane displays the content of create_d597_database.py. The code uses pandas to read a JSON file, clean the data, and then insert it into a MongoDB database named 'D597_Task_2'. A terminal tab at the bottom shows the command `python3 create_d597_database.py` was run and completed successfully.

Screenshot of Database in MongoDB:



The screenshot shows the Compass MongoDB interface. On the left, the connections sidebar lists 'localhost:27017' and 'D597_Task_2'. The main area shows two databases: 'fitness_trackers' and 'medical'. The 'fitness_trackers' database has a storage size of 45.06 kB, 565 documents, an average document size of 288.00 B, 1 index, and a total index size of 24.58 kB. The 'medical' database has a storage size of 6.55 MB, 100 K documents, an average document size of 229.00 B, 1 index, and a total index size of 1.00 MB.

D2. Script to insert and map the data records into the database instance.

Successful Upload of Medical Data into D597_Task_2:

```

upload_medical.py
1 #!/usr/bin/python3
2 from pymongo import MongoClient
3
4 # Step 1: Load the JSON file
5 df = pd.read_json('medical.json')
6
7 # DEBUG: Show column names before cleaning
8 print("Original columns:", df.columns.tolist())
9
10 # Step 2: Clean column names
11 df.columns = df.columns.str.strip().str.lower().str.replace(' ', '_')
12
13 # DEBUG: Show cleaned column names
14 print("Cleaned columns:", df.columns.tolist())
15
16 # Step 3: Now try cleaning date fields using standardized names
17 if 'birth_date' in df.columns:
18     df['birth_date'] = pd.to_datetime(df['birth_date'], errors='coerce')
19
20 if 'last_appointment_date' in df.columns:
21     df['last_appointment_date'] = pd.to_datetime(df['last_appointment_date'], errors='coerce')
22
23 # Step 4: Clean gender field (if it exists)
24 if 'gender' in df.columns:
25     df['gender'] = df['gender'].astype(str).str.strip().str.title()
26
27 # Step 5: Print preview and missing data
28 print("\nCleaned data preview:\n", df.head())
29 print("\nMissing values:\n", df.isnull().sum())
30
31 # Step 6: Upload to MongoDB
32 client = MongoClient("mongodb://localhost:27017/")
33 db = client["D597_Task_2"]
34 collection = db["medical"]
35 collection.delete_many({})
36 collection.insert_many(df.to_dict(orient='records'))
37

```

PROBLEMS OUTPUT DEBUG CONSOLE TERMINAL PORTS

```

medications      0
allergies        0
last_appointment_date  0
tracker          0
dtype: int64

```

Cleaned medical.json uploaded to 'medical' collection in D597_Task_2.

Successful Upload of Fitness Tracker Data into D597_Task_2:

```

upload_fitness_trackers.py
1 import pandas as pd
2 from pymongo import MongoClient
3
4 # 1. Load the JSON file
5 df = pd.read_json('Task2Scenario1 Dataset 1_fitness_trackers.json')
6
7 # 2. Clean 'Selling Price' column
8 df['Selling Price'] = df['Selling Price'].str.replace(',', '').astype(float)
9
10 # 3. Clean 'Reviews' column
11 df['Reviews'] = df['Reviews'].replace('', None)
12 df['Reviews'] = df['Reviews'].astype(str).str.replace(' ', '')
13 df['Reviews'] = pd.to_numeric(df['Reviews'], errors='coerce')
14
15 # 4. Standardize text formatting
16 df['Brand Name'] = df['Brand Name'].str.strip().str.title()
17 df['Device Type'] = df['Device Type'].str.strip().str.title()
18
19 # 5. Convert 'Color' field to list
20 df['Color'] = df['Color'].apply(lambda x: [c.strip().title() for c in x.split(',') if isinstance(x, str) else []])
21
22 # 6. Rename fields for MongoDB (no spaces, lowercase with underscores)
23 df.rename(columns={
24     "Model Name": "model_name",
25     "Brand Name": "brand",
26     "Device Type": "device_type",
27     "Selling Price": "price",
28     "Color": "colors",
29     "Reviews": "reviews",
30     "Rating (Out of 5)": "rating",
31     "Average Battery Life (in days)": "battery_life_days"
32 }, inplace=True)
33
34 # 7. Remove any duplicates
35 df.drop_duplicates(inplace=True)
36

```

PROBLEMS OUTPUT DEBUG CONSOLE TERMINAL PORTS

```

● joannarochi@joannas-MacBook-Air Scenario 1 Data % /u
● sbin/python3 "/Users/joannarochi/Desktop/NGU MDA
ata Management D597/Task 2/Scenario 1/Scenario 1 Data
/Task2Scenario1 Dataset 1_fitness_trackers.json"
● Data cleaned and uploaded successfully.
● joannarochi@joannas-MacBook-Air Scenario 1 Data % python3 upload_fitness_trackers.py
● Data cleaned and uploaded successfully.
○ joannarochi@joannas-MacBook-Air Scenario 1 Data %

```

MongoDB Screenshot – Successful Mapping of Medical Data into D597_Task_2

The screenshot shows the Compass MongoDB interface with the 'medical' collection selected. The interface includes a sidebar for connections and a main area for document viewing. The documents list shows four entries, each representing a patient record with their details.

```

_id: ObjectId('68325b3154bdde7171856e91')
patient_id: 1
name: "Scott Webb"
date_of_birth: "4/28/1967"
gender: "M"
medical_conditions: "None"
medications: "None"
allergies: "None"
last_appointment_date: 2022-07-26T00:00:00.000+00:00
tracker: "Band 4"

_id: ObjectId('68325b3154bdde7171856e92')
patient_id: 2
name: "Rachel Frederick"
date_of_birth: "4/4/1977"
gender: "M"
medical_conditions: "None"
medications: "None"
allergies: "None"
last_appointment_date: 2023-02-14T00:00:00.000+00:00
tracker: "Band 3"

_id: ObjectId('68325b3154bdde7171856e93')
patient_id: 3
name: "Eric Kline"
date_of_birth: "5/18/1926"
gender: "M"
medical_conditions: "Watch"
medications: "Yes"
allergies: "None"
last_appointment_date: 2021-04-24T00:00:00.000+00:00
tracker: "Band 51"

_id: ObjectId('68325b3154bdde7171856e94')

```

MongoDB Screenshot – Successful Mapping of Fitness Trackers into D597_Task_2

The screenshot shows the Compass MongoDB interface with the 'fitness_trackers' collection selected. The interface includes a sidebar for connections and a main area for document viewing. The documents list shows three entries, each representing a fitness tracker with its specifications.

```

_id: ObjectId('683240c08aa9d37ac574de53')
brand: "Xiaomi"
device_type: "Fitnessband"
model_name: "Smart Band 5"
colors: Array (1)
price: 2499
Original Price: "2,999"
Display: "AMOLED Display"
rating: 4.1
Strap Material: "Thermoplastic polyurethane"
battery_life_days: 14
reviews: NaN

_id: ObjectId('683240c08aa9d37ac574de54')
brand: "Xiaomi"
device_type: "Fitnessband"
model_name: "Smart Band 4"
colors: Array (1)
price: 2099
Original Price: "2,499"
Display: "AMOLED Display"
rating: 4.1
Strap Material: "Thermoplastic polyurethane"
battery_life_days: 14
reviews: NaN

_id: ObjectId('683240c08aa9d37ac574de55')
brand: "Xiaomi"
device_type: "Fitnessband"
model_name: "HMSH01GE"
colors: Array (1)
price: 1722
Original Price: "2,099"
Display: "LCD Display"
rating: 3.5
Strap Material: "Leather"

```

D3. Query Scripts

QUERY 1

Query1_HealthFit_Patient_Demographics - Stage 1&2 (out of 4)

The screenshot shows the Compass MongoDB interface with the following details:

- Left Sidebar (Connections):** Shows connections to localhost:27017, D597_Task_2 (with fitness_trackers and medical collections selected), admin, config, cooker, examples, recipes, todos, users, healthfit, and local.
- Top Bar:** Shows the connection as "localhost:27017 > D597_Task_2 > medical".
- Aggregations Tab:** Selected.
- Stages:**
 - Stage 1 \$addFields:** Contains the following code:

```

1 // Convert date_of_birth from
2 // string to date
3 {
4   date_of_birth: {
5     $dateFromString: {
6       dateString: "$date_of_birth",
7       format: "%m/%d/%Y"
8     }
9   }
10 }

```
 - Output after \$addFields stage (Sample of 10 documents):** Displays two sample documents. The first document is for "Scott Webb" with a date of birth of 1967-04-28T00:00:00+00:00. The second document is for "Rachel Frederick" with a date of birth of 1977-04-04T00:00:00+00:00.
 - Stage 2 \$addFields:** Contains the following code:

```

1 // calculate age
2 {
3   age: {
4     $dateDiff: {
5       startDate: "$date_of_birth",
6       endDate: "$$NOW",
7       unit: "year"
8     }
9   }
10 }

```
 - Output after \$addFields stage (Sample of 10 documents):** Displays two sample documents. The first document for "Scott Webb" has an age of 50. The second document for "Rachel Frederick" has an age of 48.
- Bottom Buttons:** PREVIEW, STAGES, TEXT, WIZARD.

Query1_HealthFit_Patient_Demographics - Stage 3&4

The screenshot shows the Compass MongoDB interface with the following details:

- Left Sidebar (Connections):** Same as the previous screenshot.
- Top Bar:** Shows the connection as "localhost:27017 > D597_Task_2 > medical".
- Aggregations Tab:** Selected.
- Stages:**
 - Stage 3 \$addFields:** Contains the following code:

```

1 // create age group
2 {
3   age_group: {
4     $cond: [
5       {
6         $lt: ["$age", 30]
7       },
8       {
9         $lt: ["Under 30",
10           {
11             $cond: [
12               {
13                 $lt: ["$age", 60]
14               },
15               "30-59",
16               "60 and over"
17             ]
18           }
19         ]
20       }
21     ]
22   }
23 }

```
 - Output after \$addFields stage (Sample of 10 documents):** Displays two sample documents. The first document for "Scott Webb" is in the "Under 30" group. The second document for "Rachel Frederick" is in the "30-59" group.
 - Stage 4 \$group:** Contains the following code:

```

1 // group by age and gender
2 {
3   _id: {
4     gender: "$gender",
5     age_group: "$age_group"
6   },
7   count: {
8     $sum: 1
9   }
10 }

```
 - Output after \$group stage (Sample of 6 documents):** Displays two sample documents. The first document has an _id of {"gender": "M", "age_group": "Under 30"} and a count of 21984. The second document has an _id of {"gender": "F", "age_group": "30-59"} and a count of 13181.
- Bottom Buttons:** PREVIEW, STAGES, TEXT, WIZARD.

Query 1 Results

The screenshot shows the Compass MongoDB interface with the following details:

- Connections:** localhost:27017, D597_Task_2, medical.
- Aggregations:** \$addFields, \$addFields, \$addFields, \$group.
- Results:** ALL RESULTS (Showing 1 ~ 6 count results)
- Data:**

 - gender: "M", age_group: "60 and over", count: 21993
 - gender: "F", age_group: "60 and over", count: 21984
 - gender: "F", age_group: "30-59", count: 14914
 - gender: "M", age_group: "30-59", count: 15034
 - gender: "M", age_group: "Under 30", count: 13181
 - gender: "F", age_group: "Under 30", count: 12894

QUERY 2

Query2_HealthFit_Patient_Demographics - Stage 1&2 (out of 3)

The screenshot shows the Compass MongoDB interface with the following details:

- Connections:** localhost:27017, D597_Task_2, medical.
- Aggregations:** \$group, \$sort.
- Stages:**
 - Stage 1 \$group:**

```

1 // _id defines how to group the documents:
2 //by tracker, age group, and gender.
3
4 //count calculates how many documents (use
5 //are in each group.
6
7 {
8   _id: {
9     tracker: "$tracker",
10    age_group: "$age_group",
11    gender: "$gender"
12  },
13  count: {
14    $sum: 1
15  }
16 }

```
 - Stage 2 \$sort:**

```

1 //Sorts the grouped results in descending
2 //showing the most common tracker-demogra...
3 //combinations first.
4
5 {
6   count: -1
7 }

```
- Output:**
 - Output after \$group stage (Sample of 10 documents):
 - _id: Object, count: 1
 - Output after \$sort stage (Sample of 10 documents):
 - _id: Object, count: 1578
 - _id: Object, count: 1513

Query2_HealthFit_Patient_Demographics - Stage 3 (out of 3)

The screenshot shows the Compass interface for a MongoDB database named 'localhost:27017 > D597_Task_2 > medical'. The 'Aggregations' tab is selected. The aggregation pipeline consists of three stages:

```

1 //showing the most common tracker-demographic combinations first.
2
3
4
5 {
6   count: -1
7 }
  
```

Stage 3 (\$project):

```

1 //Removes the _id object and flattens
2 //the fields for easy reading or export
3
4 {
5   _id: 0,
6   tracker: "$_id.tracker",
7   age_group: "$_id.age_group",
8   gender: "$_id.gender",
9   count: 1
10 }
  
```

The output after the \$project stage (sample of 10 documents) shows the flattened fields:

count	tracker	age_group	gender
1578	"Band 5"	"60 and over"	"F"
1513	"Band 5"	"60 and over"	"M"

Query 2 Results

The screenshot shows the Compass interface for the same database and collection. The 'Aggregations' tab is selected, and the results page displays 20 count results:

count	tracker	age_group	gender
1578	"Band 5"	"60 and over"	"F"
1513	"Band 5"	"60 and over"	"M"
1300	"Amazfit Bip"	"60 and over"	"M"
1245	"Amazfit Bip"	"60 and over"	"F"
1081	"Band 3"	"60 and over"	"M"
1057	"Amazfit Bip S"	"60 and over"	"M"

QUERY 3**Query3_BatteryLife_by_Tracker_Brand – Stage 1 & 2 (out of 7)**

The screenshot shows the Compass MongoDB interface with the aggregation pipeline for Query3_BatteryLife_by_Tracker_Brand. The pipeline consists of two stages:

```

Stage 1: $addFields
1 // Convert DOB to a date
2
3 {
4   date_of_birth: {
5     $dateFromString: {
6       dateString: "$date_of_birth",
7       format: "%m/%d/%Y"
8     }
9   }
10 }

```

Output after \$addFields stage (Sample of 10 documents):

```

_id: ObjectId('68325b3154bdde7171856e91')
patient_id: 1
name: "Scott Webb"
date_of_birth: 1967-04-28T00:00:00.000+00:00
gender: "M"
medical_conditions: "None"
medications: "No"
allergies: "None"
last_appointment_date: 2022-07-26T00:00:00.000+00:00

```

```

_id: ObjectId('68325b3154bdde7171856e92')
patient_id: 2
name: "Rachel Frederick"
date_of_birth: 1977-04-04T00:00:00.000+00:00
gender: "M"
medical_conditions: "None"
medications: "No"
allergies: "None"
last_appointment_date: 2023-02-14T00:00:00.000+00:00

```

Stage 2: \$addFields

```

1 // calculate age
2
3 {
4   age: {
5     $dateDiff: {
6       startDate: "$date_of_birth",
7       endDate: "$$NOW",
8       unit: "year"
9     }
10 }

```

Output after \$addFields stage (Sample of 10 documents):

```

_id: ObjectId('68325b3154bdde7171856e91')
patient_id: 1
name: "Scott Webb"
date_of_birth: 1967-04-28T00:00:00.000+00:00
gender: "M"
medical_conditions: "None"
medications: "No"
allergies: "None"
last_appointment_date: 2022-07-26T00:00:00.000+00:00

```

```

_id: ObjectId('68325b3154bdde7171856e92')
patient_id: 2
name: "Rachel Frederick"
date_of_birth: 1977-04-04T00:00:00.000+00:00
gender: "M"
medical_conditions: "None"
medications: "No"
allergies: "None"
last_appointment_date: 2023-02-14T00:00:00.000+00:00

```

Query3_BatteryLife_by_Tracker_Brand – Stage 3 & 4 (out of 7)

The screenshot shows the Compass MongoDB interface with the aggregation pipeline for Query3_BatteryLife_by_Tracker_Brand. The pipeline consists of three stages:

Stage 3: \$match

```

1 // Filter only 60+ patients
2
3 {
4   age: {
5     $gte: 60
6   }
7 }

```

Output after \$match stage (Sample of 10 documents):

```

_id: ObjectId('68325b3154bdde7171856e93')
patient_id: 3
name: "Eric Kline"
date_of_birth: 1926-05-18T00:00:00.000+00:00
gender: "F"
medical_conditions: "Watch"
medications: "Yes"
allergies: "None"
last_appointment_date: 2021-04-24T00:00:00.000+00:00

```

```

_id: ObjectId('68325b3154bdde7171856e94')
patient_id: 4
name: "James Rodriguez"
date_of_birth: 1954-07-20T00:00:00.000+00:00
gender: "M"
medical_conditions: "None"
medications: "No"
allergies: "None"
last_appointment_date: 2022-05-26T00:00:00.000+00:00

```

Stage 4: \$lookup

```

1 // Join fitness_trackers info
2
3 {
4   from: "fitness_trackers",
5   localField: "tracker",
6   // in 'medical'
7   foreignField: "model_name",
8   // in 'fitness_trackers'
9   as: "tracker_info"
10 }

```

Output after \$lookup stage (Sample of 10 documents):

```

_id: ObjectId('68325b3154bdde7171856e93')
patient_id: 3
name: "Eric Kline"
date_of_birth: 1926-05-18T00:00:00.000+00:00
gender: "F"
medical_conditions: "Watch"
medications: "Yes"
allergies: "None"
last_appointment_date: 2021-04-24T00:00:00.000+00:00

```

```

_id: ObjectId('68325b3154bdde7171856e94')
patient_id: 4
name: "James Rodriguez"
date_of_birth: 1954-07-20T00:00:00.000+00:00
gender: "M"
medical_conditions: "None"
medications: "No"
allergies: "None"
last_appointment_date: 2022-05-26T00:00:00.000+00:00

```

Stage 5: \$unwind

Query3_BatteryLife_by_Tracker_Brand – Stage 5 & 7 (out of 7)

The screenshot shows the Compass interface for MongoDB, displaying the aggregation pipeline for the 'medical' collection. The pipeline consists of several stages:

- Stage 5 (\$unwind):** Unwinds the 'tracker_info' array to work with each item individually.
- Stage 6 (\$group):** Groups by tracker brand to calculate average battery life.

Output after \$unwind stage (Sample of 10 documents):

```

1 // Unwind tracker_info array
2 // to work with each item individually
3 // joining or matching on array items
4 +
5   path: "$tracker_info",
6   preserveNullAndEmptyArrays: true
7
  
```

Output after \$group stage (Sample of 9 documents):

```

1 // Group by tracker brand to calculate
2 // average battery life
3 +
4   _id: "$tracker_info.brand",
5   average_battery_life: {
6     $avg: "$tracker_info.battery_life_days"
7   },
8   user_count: {
9     $sum: 1
10 }
11
  
```

Query3_BatteryLife_by_Tracker_Brand – Stage 7 (out of 7)

The screenshot shows the Compass interface for MongoDB, displaying the final stage of the aggregation pipeline:

- Stage 7 (\$sort):** Sorts the documents by average battery life.

Output after \$sort stage (Sample of 9 documents):

```

1 // Sort by battery life
2 +
3   average_battery_life: -1
4
  
```

Query 3 Results

Compass

localhost:27017 > D597_Task_2 > medical

Documents 100.0K Aggregations Schema Indexes 1 Validation

ALL RESULTS

Showing 1–9 count results

`$addFields` `$addFields` `$match` `$lookup` `$unwind` `$group` `$sort` `Edit` `Explain` `Export` `Run` `Options`

`_id: "Huami"`
`average_battery_life: 18.642364847846494`
`user_count: 87346`

`_id: "Oppo"`
`average_battery_life: 14`
`user_count: 2142`

`_id: "Realme"`
`average_battery_life: 12.347676419965577`
`user_count: 8134`

`_id: "Honor"`
`average_battery_life: 11.24505415655858`
`user_count: 58534`

`_id: "Boat"`
`average_battery_life: 7.983078424991865`
`user_count: 6146`

`_id: "Xiaomi"`
`average_battery_life: 7`
`user_count: 4172`

`_id: "Huawei"`
`average_battery_life: 7`
`user_count: 6062`

`user_count: 2142`

`_id: "Realme"`
`average_battery_life: 12.347676419965577`
`user_count: 8134`

`_id: "Honor"`
`average_battery_life: 11.24505415655858`
`user_count: 58534`

`_id: "Boat"`
`average_battery_life: 7.983078424991865`
`user_count: 6146`

`_id: "Xiaomi"`
`average_battery_life: 7`
`user_count: 4172`

`_id: "Huawei"`
`average_battery_life: 7`
`user_count: 6062`

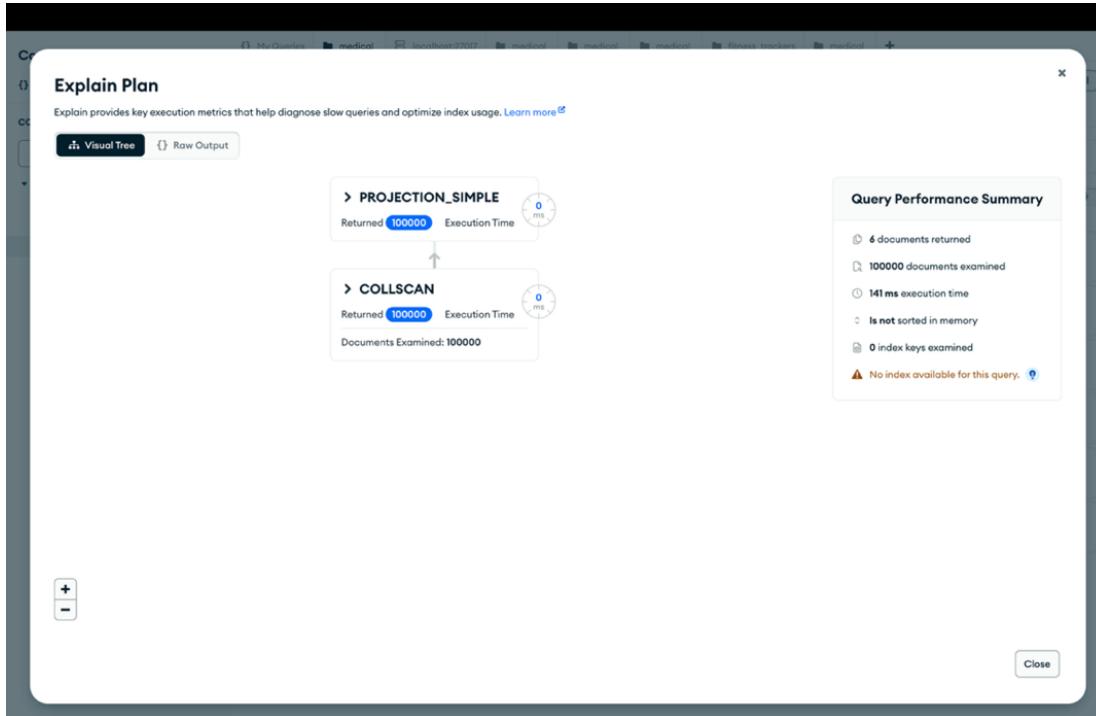
`_id: "Infinix"`
`average_battery_life: 4`
`user_count: 6182`

`_id: null`
`average_battery_life: null`
`user_count: 7969`

D4. Optimization Techniques

QUERY 1 OPTIMIZATION

Query 1 Execution Without Optimization (141ms)



Query 1 Optimization Techniques

The screenshot shows the Compass interface for MongoDB, displaying the aggregation pipeline for 'Query_1'. The pipeline consists of two stages: Stage 1 (\$match) and Stage 2 (\$group). Stage 1 filters documents based on gender ('M' or 'F') and age group ('Under 30', '30-59', '60 and over'). Stage 2 groups by gender and age group, counts the documents, and sums a value. The output after each stage is shown as a sample of 10 and 6 documents respectively.

```

Stage 1 ($match)
1: {
  gender: { $in: ["M", "F"] },
  age_group: { $in: ["Under 30", "30-59", "60 and over"] }
}

Stage 2 ($group)
1: {
  _id: { gender: "$gender", age_group: "$age_group" },
  count: { $sum: 1 }
}

```

The screenshot shows the Compass MongoDB interface. The left sidebar displays connections and databases, including 'D597_task_2' and 'D597_optimization'. The main area shows an aggregation pipeline for the 'medical_optimization' collection. The pipeline consists of three stages:

```

1. $match: {
  _id: {
    gender: "$gender",
    age_group: "$age_group"
  }
}
2. $group: {
  count: {
    $sum: 1
  }
}
3. $sort: {
  "_id.gender": 1,
  "_id.age_group": 1
}

```

The results show two documents from the output stage:

- Document 1: `_id: Object count: 15034`
- Document 2: `_id: Object count: 13181`

Below the pipeline, a preview of the output after the `$sort` stage is shown, containing a sample of 6 documents.

Query 1 After Optimization (92ms)

The screenshot shows the 'Explain Plan' dialog in Compass. It displays the execution plan for the query, which includes three stages:

- GROUP**: Returned 6 documents, Execution Time 0 ms.
- FETCH**: Returned 100000 documents, Execution Time 15 ms.
- IXSCAN**: Returned 100000 documents, Execution Time 34 ms. Index Name: `gender_1`, Multi Key Index: no.

The 'Query Performance Summary' section provides the following metrics:

- 6 documents returned
- 100000 documents examined
- 92 ms execution time
- Is not sorted in memory
- 100000 index keys examined

It also notes that the query used the `gender` index.

QUERY 2 OPTIMIZATION

Query 2 Optimization Techniques

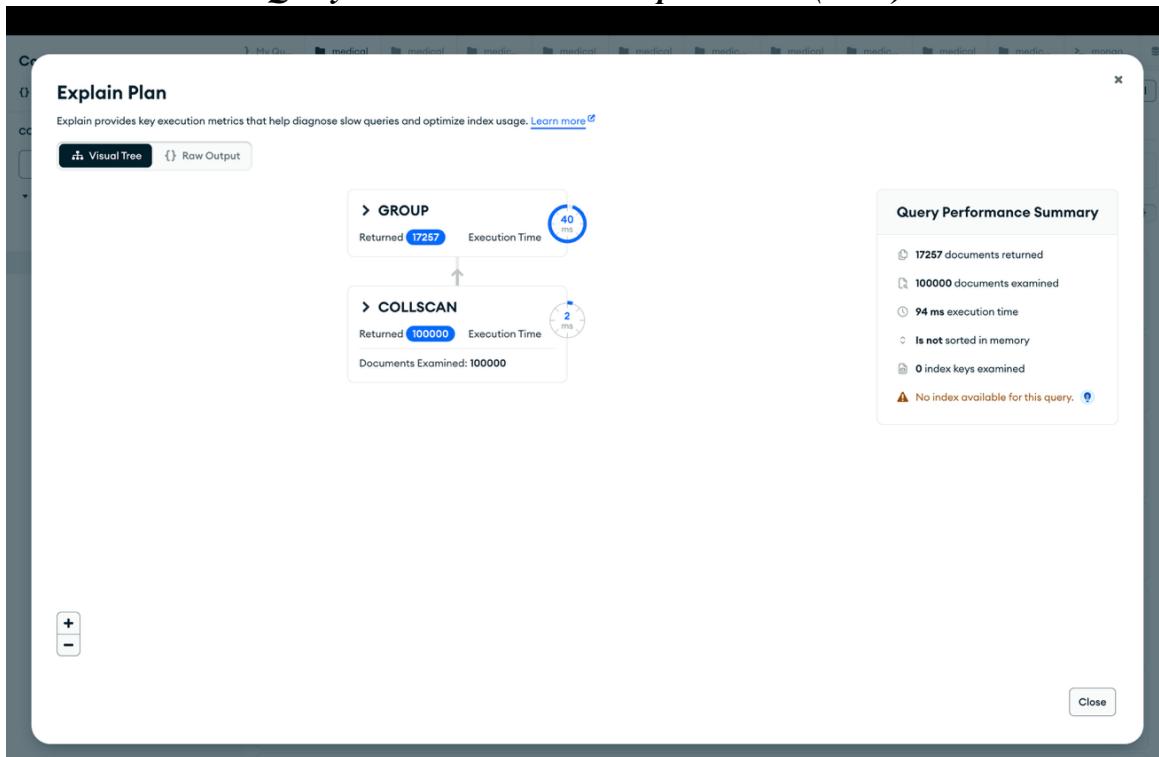
The screenshot shows the Compass MongoDB interface with the following details:

- Left Sidebar (Connections):** Shows various databases and collections, including `localhost:27017`, `D597_Task_2`, and `D597_optimization`.
- Top Bar:** Shows the current connection as `localhost:27017 > D597_optimization > medical_optimization` and includes tabs for `Documents`, `Aggregations` (selected), `Schema`, `Indexes`, and `Validation`.
- Aggregation Pipeline:**
 - Stage 1 (\$sort):** Sorts by `tracker`, `age_group`, and `gender` in ascending order. The output shows 10 sample documents.
 - Stage 2 (\$group):** Groups documents by `tracker`, `age_group`, and `gender`. The output shows 10 sample documents with a count for each group.
- Bottom Buttons:** Includes `SAVE`, `CREATE NEW`, `EXPORT TO LANGUAGE`, `PREVIEW` (selected), `STAGES`, `TEXT`, `WIZARD`, and `Run`.

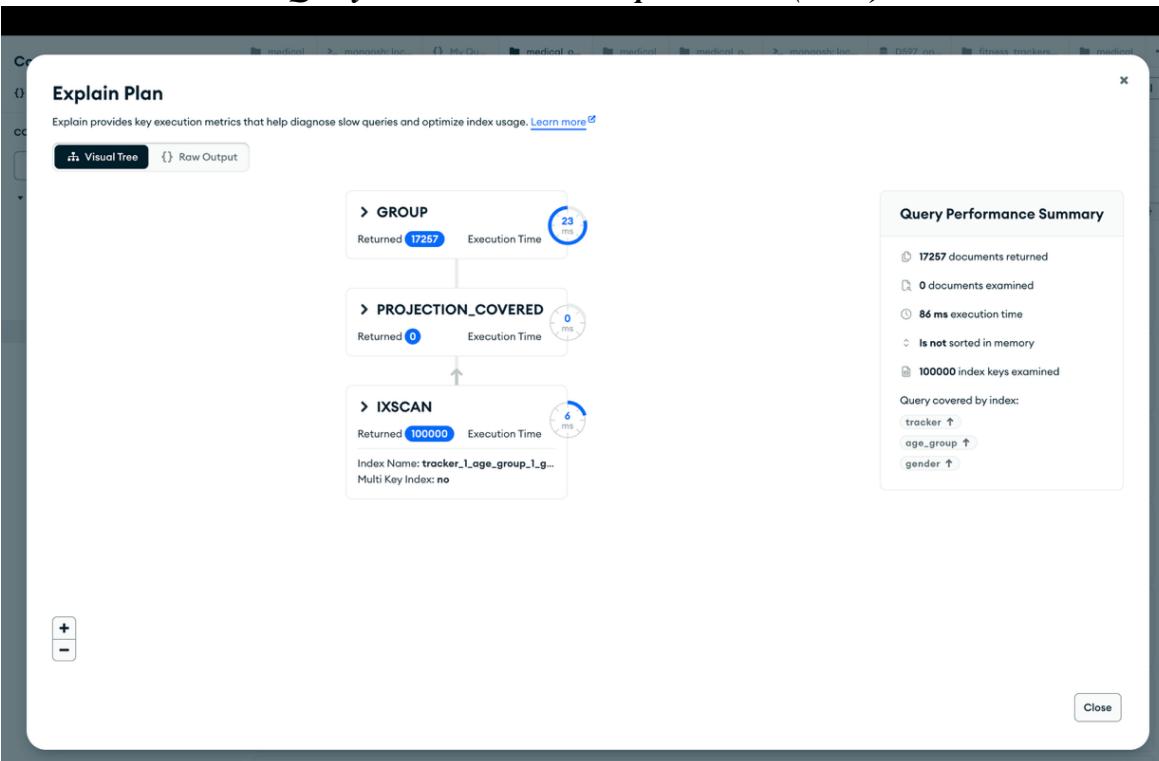
The screenshot shows the Compass MongoDB interface with the following details:

- Left Sidebar (Connections):** Shows various databases and collections, including `localhost:27017`, `D597_Task_2`, and `D597_optimization`.
- Top Bar:** Shows the current connection as `localhost:27017 > D597_optimization > medical_optimization` and includes tabs for `Documents`, `Aggregations` (selected), `Schema`, `Indexes`, and `Validation`.
- Aggregation Pipeline:**
 - Stage 3 (\$project):** Reshapes the output to promote fields from `_id` to top-level fields. The output shows 10 sample documents.
 - Stage 4 (\$sort):** Sorts the grouped results by count in descending order. The output shows 10 sample documents.
- Bottom Buttons:** Includes `SAVE`, `CREATE NEW`, `EXPORT TO LANGUAGE`, `PREVIEW` (selected), `STAGES`, `TEXT`, `WIZARD`, and `Run`.

Query 2 Execution Without Optimization (94ms)

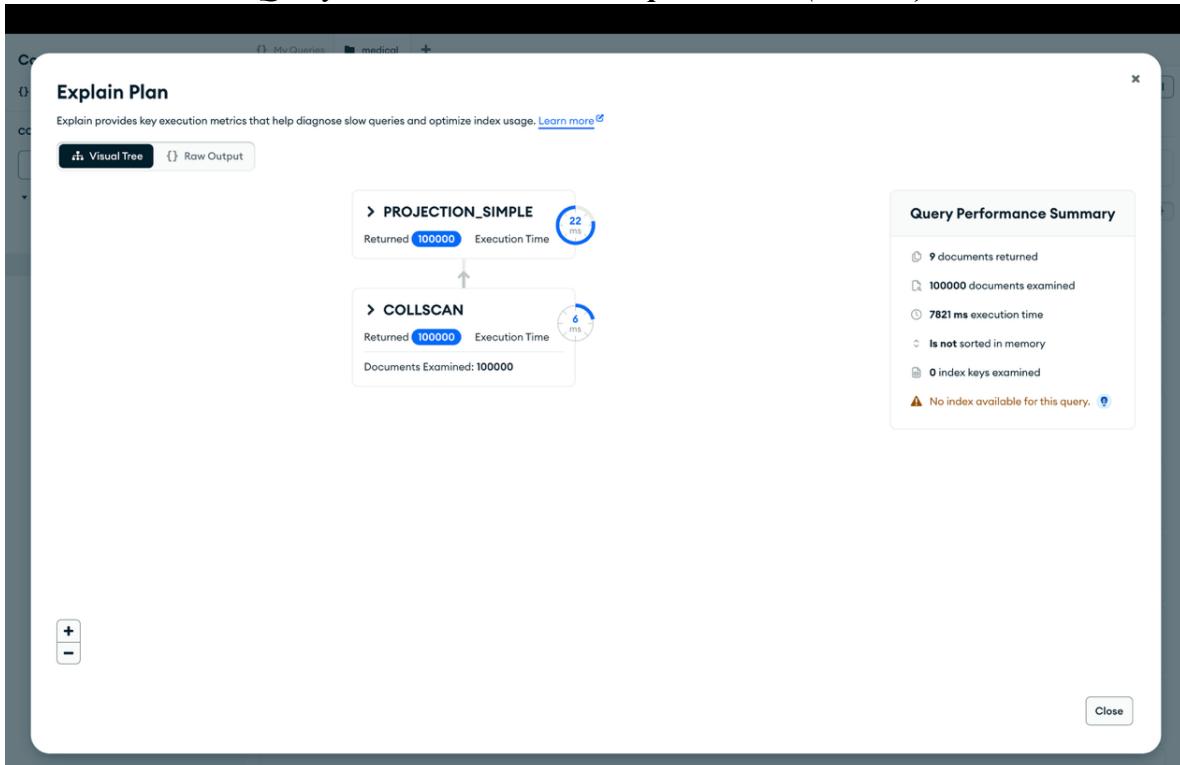


Query 2 Execution With Optimization (86ms)

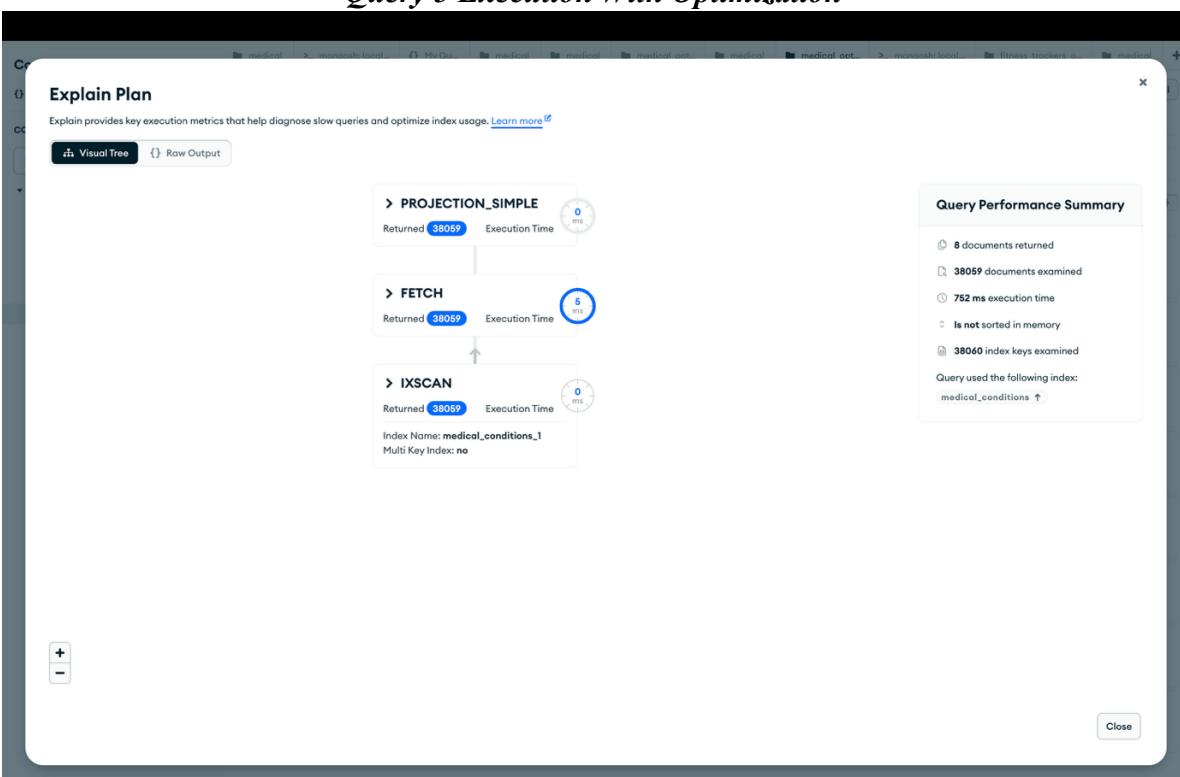


QUERY 3 OPTIMIZATION

Query 3 Execution Without Optimization (7821ms)



Query 3 Execution With Optimization



References

- Y. He, S. Lee, R. Zhang, & M. S. Kim. (2021). *Scalable storage of real-time personal health data using NoSQL databases in smart health applications*. *Sensors*, 21(2), 535. <https://doi.org/10.3390/s21020535>
- M. Dey, M. Kabir, & A. Ghosh (2020). *Evaluating NoSQL databases for scalable real-time health monitoring systems*. *Journal of Big Data*, 7(1). <https://doi.org/10.1186/s40537-020-00320-6>
- G. Alshammari, A. Alhaidari, & A. Almogren. (2019). *NoSQL for healthcare big data: A performance evaluation of MongoDB and Cassandra*. *IEEE Access*, 7, 149848–149855. <https://doi.org/10.1109/ACCESS.2019.2946975>
- Y. Zhang, J. Dang, & L. Zhang. (2019). *Real-time health monitoring and analytics with big data*. *IEEE Transactions on Industrial Informatics*, 15(1), 384–393. <https://doi.org/10.1109/TII.2018.2856580>
- Fernández-Alemán, J. L., Señor, I. C., Lozoya, P. Á. O., & Toval, A. (2013). *Security and privacy in electronic health records: A systematic literature review*. *Journal of Biomedical Informatics*, 46(3), 541–562. <https://doi.org/10.1016/j.jbi.2012.12.003>
- MongoDB, Inc. (n.d. 1). *Client-side field level encryption*. MongoDB Documentation. <https://www.mongodb.com/docs/manual/core/security-client-side-encryption/>
- Gajanayake, R., Iannella, R., & Sahama, T. (2017). Security and privacy in the use of cloud computing for health data. *Health Information Science and Systems*, 5(1). <https://doi.org/10.1007/s13755-017-0020-5>
- MongoDB, Inc. (n.d. 2). *Replication*. MongoDB Documentation. <https://www.mongodb.com/docs/manual/replication/>
- Xie, H., Wang, F., & Zhang, L. (2020). Data masking approaches for privacy-preserving electronic health record sharing. *BMC Medical Informatics and Decision Making*, 20(Suppl 3), 132. <https://doi.org/10.1186/s12911-020-01125-4>
- MongoDB, Inc. (n.d. 3). *Security checklist*. MongoDB Documentation. <https://www.mongodb.com/docs/manual/administration/security-checklist/>
- Dey, M., Kabir, M. A., & Ghosh, A. (2020). Evaluating NoSQL databases for scalable real-time health monitoring systems. *Journal of Big Data*, 7(1), 1–25. <https://doi.org/10.1186/s40537-020-00320-6>
- Wolniak, R. (2023). Functioning of real-time analytics in business. *Scientific Papers of Silesian University of Technology: Organization and Management Series*, (172), 659–672. <https://doi.org/10.29119/1641-3466.2023.172.40>

Belefqih, S., Zellou, A., & Berquedich, M. (2024). Semantic schema extraction in NoSQL databases using BERT embeddings. *Data Science Journal*, 23(1), 1–15.
<https://doi.org/10.5334/dsj-2024-057>