



Name :- Aditya S. Khandagale

Div :- ATEC - 2

Roll No :- 03

Sub :- Data Engineering

Assignment - I

Q 1

- a) A healthcare organization wants to digitize patient records. Explain the Structured, Semi Structured data type they would encounter and how the data lifecycle supports this transformation.

Ans When a healthcare organization digitizes patient record, it will face three main types of data

i) Structured Data :-

a) organized and stored in relational databases

b) Examples : Patient ID, Name, Age, Contact, Blood Test Results

c) Easy to query using SQL (rows and columns)

ii) Semi-Structured Data :-

a) Data that has some structure but not strictly tabular

b) Example :- Electronic prescription, lab reports in XML / JSON format, wearable device logs

c) flexible and can evolve without fixed schema

iii] Unstructured Data :-

- a] No predefined structure, often large in size
- b] Example : MRI Scans, X-ray images, doctor's notes, audio / video consultation recording
- c] Need Big Data / NoSQL Systems for storage and analysis

Data lifecycle Support

- Collection :- Patient data gathered from hospital system, labs and IoT devices
- Storage :- Structured data → relation DB, Semi structured → JSON / XML, unstructured → image repositories
- Processing :- Transform into usable formats
- Analysis :- Data mining for disease patterns, predictive analytics
- Sharing :- Secure sharing between doctor and insurance agencies
- Archival :- old records stored or deleted as per compliance

- b] A logistics company wants to design a system to analyze vehicle telemetry data. Describe the phases of the data lifecycle applicable to their system.

Ans Telemetry data includes GPS, speed, fuel consumption, engine temperature, etc. Data lifecycle phases applicable

- i] Data collection :- Sensors in vehicles capture telemetry continuously
- ii] Data storage :- Data stored in real-time database or cloud (time series DB like influxDB)
- iii] Data Processing :- Streaming engines (Kafka, Spark) process live data
- iv] Data Analysis :- Analytics dashboards track fuel usage, routes, and driver behavior
- v] Data Sharing :- Reports shared with managers for route optimization
- vi] Data Archival :- Historical telemetry stored for compliance or deleted after retention period

Q2

- a] A university plans to deploy a research database. Compare and analyze key-value and document-based NoSQL database for this use case.

Ans • Key - Value Store

- i] Structure :- Data stored as key-value pairs (like a dictionary)
- ii] Pros :- Extremely fast lookups, good for caching search IDs.
- iii] Cons :- No support for complex queries, poor handling of relationships
- iv] Example : Redis, Amazon DynamoDB



- Document - Based Store
- i] Structure :- JSON / BSON documents with flexible Schema
- ii] Pros : Can store complex research data (metadata, experiment results, references) in one document
- iii] Good for queries, indexing and semi-structured data
- iv] Example :- MongoDB, CouchDB

Analysis :-

for a research databases where data varies (papers, metadata, authors, experiment datasets) document - based NoSQL is better because of flexible schema and query support. key value is faster but too simple for research use cases

- b) An e-commerce platform receives high volume clickstream data. Explain How Big data characteristic (volume, variety, velocity) influence the data engineering process

Ans • Volume :- Massive clickstream logs generated daily by millions of users. Require distributed storage (HDFS, NoSQL)

• Variety : clickstream data include page visits device type, location, session logs (structured + unstructured). Needs schema - flexible system like MongoDB



- **Velocity** :- Data generated at high speed in real time. Needs stream processing frameworks (Kafka, SparkStreaming)

Impact on Data Engineering

- Data pipelines must be Scalable
- ETL processes need to handle multi format data
- Real-time analytics required for personalization and fraud detection

Big Data characteristic drive design of data ingestion, storage and processing pipelines

Q3

- a) A Startup is evaluating NoSQL option for a recommendation engine. Distinguish between Document and Graph database in term of performance and structure

i] Document Database

- Store data as JSON document
- Best for user profiles, purchase history
- Performance : High for reading / writing document
- Limitation : Not efficient for complex relationships

ii] Graph Database

- Store entities as node and relationships as edges
- Best for recommendation systems

Q4 a) A food delivery app aggregates reviews and ratings. Explain how MongoDB's JSON document structure help store heterogeneous data

Ans MongoDB uses a flexible JSON-like format allowing heterogeneous data in a single collection

Example

```
{
  "restaurant": "Pizza Hut",
  "review": "Great taste!",
  "rating": 4.5,
  "user": {
    "name": "Johny",
    "location": "Mum"
  },
  "Images": ["img1.jpg", "img2.jpg"]
}
```

- Some reviews may have images, others may not
- No fixed Schema - easy to handle varied user data
- Supports nested fields, arrays and flexible attributes

b] A medical records system must be queried efficiently. Analyze how mongoDB index and aggregation can support analytical dashboards

Ans

- Indexes :-

Improve query performance (e.g. Searching patients by ID, disease or date)

- Example

```
db.records.createIndex({patientId: 1})
```



- Aggregation:-
- used to group and summarize data for dashboards
- Example :- Average patient age , disease - wise case counts
- Query

db.records.aggregate([

```
{ $group : { _id : "$disease" , total :  
{ $sum : 1 } } }  
])
```

with indexes , queries become fast and with aggregation , insights can be visualized in real time dashboards for medical decision-making

- Q5 a) An IoT data hub captures reading every second discuss the advantages of MongoDB in handling high-ingestion use cases

Ans when an IoT data hub captures sensor reading every second the database must handle high ingestion rate , scalability and flexible schema MongoDB provides several advantages

i] High write Throughput

- MongoDB is optimized for rapid insert operation allowing millions of sensor readings per second
- uses memory-mapped storage engine for fast ingestion

ii] Horizontal Scalability

- IoT hubs generate huge amounts of data continuously



- MongoDB Supports sharding, distributing data across multiple servers for balanced load
- iii] Flexible Schema
- IoT data may vary (Temperature, pressure, humidity)
- MongoDB's JSON / BSON documents allow inserting new sensor fields without altering schema

b] A developer uses MongoDB compass to model data Analyze its usefulness in query visualization and performance tuning

Ans MongoDB Compass is a GUI tool that helps developer's interact with MongoDB visually instead of command-line. Its usefulness includes

- i] Data Modeling and Exploration
 - Developers can view JSON documents in tree / table view
 - Helpful for understanding schema in semi-structured IoT or e-commerce datasets
- ii] Query visualization
 - Provide a visual query builder
 - users can filter, sort and project fields without writing full mongo queries
 - Queries show real time result with execution stats
- iii] Performance Tuning
 - Show query execution slow plans with detail

like index usage and Scan performance

- Helps identify slow queries and optimize them using indexes

iv) Index Management

- Allows creation and monitoring of indexes from GUI
- Developers can test how indexes affect query response time

v) Aggregation pipeline Builder

- Provides drag and drop interface to build multi-stage pipelines for analytics
- useful for dashboards and reports