Data, Data Storage, Data Collection Lecture 7: Data Storage

Romain Pascual

MICS, CentraleSupélec, Université Paris-Saclay

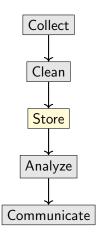
Recap

Recap: Data Storage

- Understand the difference between storage models (flat, tabular, hierarchical, columnar) and formats (CSV, JSON, Parquet)
- 2 Choose the right format based on your use case: Excel for business, CSV/JSON for sharing, Parquet for analytics
- 3 Consider the trade-offs between readability, size, speed, and interoperability when selecting a format
- For long-term storage, prioritize open standards and self-describing formats with proper metadata
- 5 Document your storage decisions and maintain data provenance for reproducibility
- Consider ethical implications including privacy, accessibility, and sustainability

Introduction

Within the lifecycle



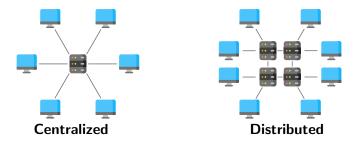
Session Objectives

At the end of this session, you should be able to:

- Explain the differences between SQL and NoSQL databases.
- Describe the four families of NoSQL databases and their use cases.
- Perform basic CRUD operations in MongoDB.

Why NoSQL databases?

NoSQL for Not only SQL



Goals

- Provides better scalability for distributed databases
- Provides better support for semi-structured databases
- (Addresses the object-relational impedance mismatch: store data as objects, like in code, instead of tables)

Aggregate

Aggregate: The Basic Unit of NoSQL

Unlike relational databases, we can control how data is distributed. The atomic blocks for distribution are called aggregates.

Definition (Aggregate)

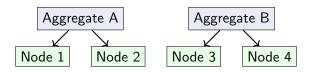
A collection of data that interacts as a single unit.

```
1 {
2    "user_id": 123,
3    "name": "Alice",
4    "balance": 1000,
5    "orders": [
6         {"order_id": 1, "amount": 99.99},
7         {"order_id": 2, "amount": 49.99}
8    ]
9 }
```

Aggregates for Distribution and Operations

Aggregates are the atomic blocks for:

- Replication: Copying data across nodes for redundancy.
- **Sharding**: Splitting data across nodes for scalability.
- Operations: CRUD operations access the complete unit.



Comparison with Relational Databases

- ullet \sim a row in relational databases: store a specific entity
- $\bullet \neq$ a row in relational databases: can nest other aggregates

```
1 {
2    "user_id": 123,
3    "name": "Alice",
4    "balance": 1000,
5    "orders": [
6         {"order_id": 1, "amount": 99.99},
7         {"order_id": 2, "amount": 49.99}
8    ]
9 }
```

Distributing a Relational Database

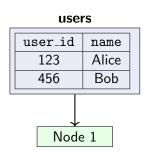
users

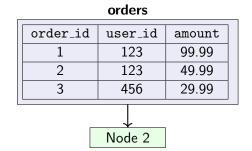
Alice
Alice
Bob

orders

order_id	user_id amount				
1	123	99.99			
2	123	49.99			
3	456	29.99			

Distributing a Relational Database

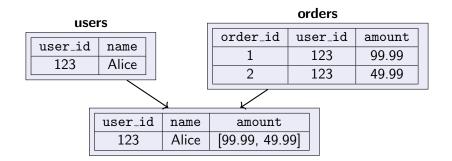




Issue

- Data for an entity is split across nodes, e.g., Alice + her orders
- Queries require JOINs, e.g., "Get all orders for Alice"

Denormalize the Table



Idea

 Simulate the JOINs such that all data for an entity is stored together (data co-location)

Example: Get All Orders for Alice

users

user_id	name
123	Alice
456	Bob

orders

```
        order_id
        user_id
        amount

        1
        123
        99.99

        2
        123
        49.99

        3
        456
        29.99
```

```
1 {
2    "user_id": 123,
3    "name": "Alice",
4    "balance": 1000,
5    "orders": [
6         {"order_id": 1, "amount": 99.99},
7         {"order_id": 2, "amount": 49.99}
8    ]
9 }
```

- SQL: Requires a JOIN between tables.
- NoSQL: Simply access the aggregate.

Denormalization in Aggregates

Why It Matters

- Performance: Avoids complex JOIN operations, which can be slow in distributed systems.
- **Simplicity**: Simplifies queries by keeping related data together.
- Flexibility: Allows for flexible schema design, accommodating changes more easily.

Trade-offs

- Redundancy: Denormalization can lead to data redundancy, which may increase storage requirements.
- Consistency: Ensuring data consistency can be more challenging with denormalized data.

Aggregates are Schemaless

Aggregates might not have the same attributes.

```
1 {
2    "user_id": 123,
3    "name": "Alice",
4    "balance": 1000,
5    "orders": [
6         {"order_id": 1, "amount": 99.99},
7         {"order_id": 2, "amount": 49.99}
8    ]
9 }
```

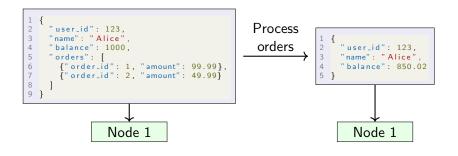
```
1 {
2    "user_id": 456,
3    "name": "Bob",
4    "balance": 0
5 }
```

- No need to fix a rigid schema.
- NULL values are avoided.

Atomic Operations on Aggregates

Operations on a single aggregate are atomic:

- The entire operation is completed as a single unit.
- No partial updates



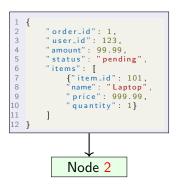
- Ensures data consistency within an aggregate.
- Simplifies error handling and recovery.

Cross-Aggregate Operations

```
1 {
2  "user_id": 123,
3  "name": "Alice",
4  "balance": 1000,
5  "orders": [
6  {"order_id": 1, "amount": 99.99},
7  {"order_id": 2, "amount": 49.99}
8  }
9 }

Node 1
```

Cross-Aggregate Operations



Cross-aggregate operations can involve several nodes:

- No guarantee that related aggregates are stored on the same node.
- Network latency and partition might cause inconsistencies.

Aggregate-Based Databases

■ Schemaless

- No rigid schema
- Flexible evolution

Aggregate Update

Atomicity (all or nothing)

≠ Denormalization

- Fast queries
- Redundancy (cheap storage)

✗ Cross-Aggregate

- No atomicity
- Availability vs. consistency

Distribution

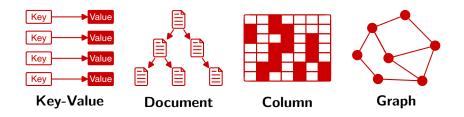
- Aggregate in a single node
- Control

Design

- Optimize for query patterns
- Eventual consistency

NoSQL Databases

Four families of NoSQL

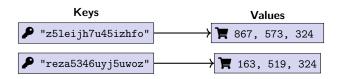


- ♠ NoSQL is not a single database or query language but a diverse family of technologies.
 - The first three NoSQL families use aggregates, while graph databases were designed ACID-compliant.

Key-Value Databases

Description

- Key: Usually (auto-generated) an alphanumeric string.
- Value: An aggregate.
- Distribution based on the key (little to no integrity constraint)
- Optimized for speed (read/write queries) and scalability.



Examples

Amazon DynamoDB, Redis.

Document-Oriented Databases

Description

- Data is modeled as nested key-value pairs.
- Supports searching aggregates based on their attribute values.
- Values are typically JSON, or XML documents.

```
1 {
2    "user_id": 123,
3    "name": "Alice",
4    "balance": 1000,
5    "orders": [
6         {"order_id": 1, "amount": 99.99},
7         {"order_id": 2, "amount": 49.99}
8    ]
9 }
```

Examples

MongoDB, CouchDB.

Column-Oriented Databases

Description

- Data is stored in columns rather than rows.
- Columns (or column families) as aggregates.
- Optimized for analytical queries on large datasets.

Name	Age	City
Alice	28	Paris
Bob	32	London
Charlie	25	Berlin

	✓		≯	
Row Storage			Colum	n Storag
Row 1	Alice			Alice
	28		Name	Bob
	Paris			Charlie
Row 2	Bob		Age	′ 28
	32			ı 32
	London			25
Row 3	Charlie			Paris
	25		City	London
	Berlin			Berlin

Examples

Apache Cassandra, HBase.

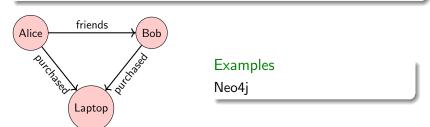
Row-based access reads entire row

Column-based access reads only needed columns

Graph Databases

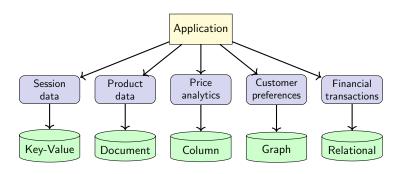
Description

- Relationships are first-class citizens.
- Nodes, edges, and properties are the core components.
- Optimized for traversing relationships.



Polyglot Persistence

- NoSQL databases complement rather than replace relational databases.
- Use different data storage technologies based on data type and usage.



MongoDB

MongoDB: Generalities

General Characteristics

- Document-oriented NoSQL database (JSON-like documents)
- Client-server architecture
- Open-source (Community Edition) and commercial (Enterprise Edition)
- Cloud-based solution: MongoDB Atlas

Key Advantages

- Flexible schema for evolving data requirements
- Horizontal scalability with sharding
- Rich query language and aggregation framework

Data Organization and Storage Format

Data Structure

- Document: nested dictionaries
- Collection: set of documents
- Database: set of collections

BSON (Binary JSON)

- Faster parsing and processing
- More space-efficient
- Supports Data types

Every document must have an _id field:

Auto-generated as ObjectId if not specified

```
1 {
2    "_id": ObjectId("507f1f77bcf86cd799439011"),
3    "title": "The Shawshank Redemption",
4    "year": 1994,
5    "directors": ["Frank Darabont"],
6    "actors": ["Tim Robbins", "Morgan Freeman"],
7    "rating": 9.3,
8    "awards": {"nominations": 7, "wins": 0}
9 }
```

Seting Up MongoDB

Local Installation: MongoDB Community Server

Server:

- Start server with: mongod
- Default port: 27017
- Can manage multiple databases

Client Options

- CLI: mongosh
- GUIs: Compass, Atlas UI
- Drivers for all major languages

Cloud Solution: MongoDB Atlas

- Fully-managed cloud database
- Connection through an API
- Web-based GUI for database management

Use the MongoDB shell to connect:

- mongosh: connect to the server (localhost:27017 by default)
- use <database>: switch to a specific database
- db.createCollection(<collection>): create a collection in the database

1 \$ mongosh

Use the MongoDB shell to connect:

- mongosh: connect to the server (localhost:27017 by default)
- use <database>: switch to a specific database
- db.createCollection(<collection>): create a collection in the database

```
1 $ mongosh
2 > use moviedb
3 switched to db moviedb
```

Use the MongoDB shell to connect:

- mongosh: connect to the server (localhost:27017 by default)
- use <database>: switch to a specific database
- db.createCollection(<collection>): create a collection in the database

```
1 $ mongosh
2 > use moviedb
3 switched to db moviedb
4 > db.createCollection("movies")
5 { ok: 1 }
```

Use the MongoDB shell to connect:

- mongosh: connect to the server (localhost:27017 by default)
- use <database>: switch to a specific database
- db.createCollection(<collection>): create a collection in the database

```
1 $ mongosh
2 > use moviedb
3 switched to db moviedb
4 > db.createCollection("movies")
5 { ok: 1 }
6 > show collections
7 movies
```

Connecting to a MongoDB Database

Use the MongoDB shell to connect:

- mongosh: connect to the server (localhost:27017 by default)
- use <database>: switch to a specific database
- db.createCollection(<collection>): create a collection in the database

```
$ mongosh
2 > use moviedb
3 switched to db moviedb
4 > db.createCollection("movies")
5 { ok: 1 }
6 > show collections
7 movies
```

▲ In the lab, we will manage these steps by accessing the API via Python

Manipulating Data with MongoDB

MongoDB provides a set of functions to apply CRUD operations on the data.

▲CRUD: Create, Update, Read, Delete

Format of a CRUD function in MongoDB

db.collection.function()

- db: current database
- collection: collection where the CRUD operation is applied.
- function: invoked function.

Create

Insert Operations in MongoDB

- insertOne: adds a new document to a collection
- insertMany: adds multiple documents to a collection

```
db.movies.insertOne({
    "title": "Inception",
    "year": 2010,
    "genres": ["Action", "Sci-Fi", "Thriller"],
    "director": "Christopher Nolan",
    "actors": ["Leonardo DiCaprio", "Ken Watanabe"],
    "rating": 8.8
})
```

▲ If the collection does not exist, any insert operation will create it.

Read

Query Operations in MongoDB

- find: retrieve documents matching a query
- findOne: retrieve a single document
- Query operators: \$gt, \$lt, \$in, \$regex, etc.

```
1 SELECT title FROM movies WHERE genres = 'Sci-Fi'
```

Update

Update Operations in MongoDB

- updateOne: update the first document matching the filter
- updateMany: update all documents matching the filter
- Update operators:
 - \$set: Set field value
 - \$unset: Remove field
 - \$push: Add to array
 - \$inc: Increment numeric value

There is also replaceOne to replace the entire document.

Delete

Delete Operations in MongoDB

- deleteOne: delete the first matching document
- deleteMany: delete all matching documents

```
1 // Delete a specific movie
db.movies.deleteOne({"title": "Inception"})
 // Delete all movies from a year
5 db. movies. deleteMany({"year": 2010})
```

Delete operations cannot be undone: always verify filters

Aggregation

Aggregation operations process multiple documents and return computed results.

Simple Aggregation Methods

- countDocuments: count documents in a collection
- distinct: find distinct values for a field

```
1 // Count movies by genre
2 db.movies.countDocuments({"genres": "Sci-Fi"})
 // Find all unique directors
5 db.movies.distinct("director")
```



A For complex aggregations, use aggregation pipelines.

Aggregation Pipelines

Definition (Aggregation Pipelines)

A sequence of data processing stages where each stage transforms the documents as they flow through sequentially.

Common Pipeline Stages

- \$match: Filter documents
- \$group: Group documents
- \$sort: Sort documents

- \$unwind: Deconstruct array
- \$addFields: Add fields

Aggregation Example: What Does it Compute?

```
1 db.movies.aggregate(
2 // Pipeline
   {"$match": {"rating": {"$gte": 8.0}}}, // Stage 1
    {"$group": {
                                               // Stage 2
          "_id": "$director".
5
          "count": {"$sum": 1},
6
          "avg_rating": {"$avg": "$rating"}
7
8
    \{" \$ sort" : \{" avg_rating" : -1\} \},
                                         // Stage 3
    {" $limit": 3}
                                               // Stage 4
11 ])
```

Aggregation Example: What Does it Compute?

```
1 db.movies.aggregate(
2 // Pipeline
    {"$match": {"rating": {"$gte": 8.0}}}, // Stage 1
   {"$group": {
                                              // Stage 2
         "_id": "$director".
5
         "count": {"$sum": 1},
6
         "avg_rating": {"$avg": "$rating"}
7
8
     \{"\$sort": \{"avg_rating": -1\}\},
                                        // Stage 3
    {" $limit": 3}
                                              // Stage 4
```

This pipeline finds the top 3 directors with the highest average ratings among movies rated 8.0 or higher, showing how many qualifying movies each director has and their average rating.

Data Modeling in MongoDB

Embedding

- Store related data in same document
- Better for read performance
- Use when data is frequently accessed together
- Use when data has one-to-few relationships

Referencing

- Store references between documents
- Better for write performance
- Use when data changes frequently
- Use for many-to-many relationships

```
1 {
2  "title": "The Dark Knight",
3  "actor_ids": [
4    ObjectId("..."),
5    ObjectId("...")
6  ]
7 }
```

Takeaways about MongoDB

- Document-oriented database (JSON-like documents stored in collections)
- Offers both local and cloud-based solutions

SQL	MongoDB
Tables	Collections
Rows	Documents
Columns	Fields
Joins	Embedded documents or references
GROUP BY	\$group aggregation stage

Conclusion

Takeaways: NoSQL

- NoSQL databases address challenges of distributed systems and semi-structured data
- 2 The aggregate is the fundamental unit of organization in NoSQL databases
- 3 Data modeling is crucial for performance in NoSQL databases
- 4 Four main families of NoSQL databases: Key-value, Document-oriented, Column-oriented, Graph databases
- Solution NoSQL databases complement rather than replace relational databases (polyglot persistence)

Forget SQL vs NOSQL!

