Introduction to Databases Lecture 5 – Distributed databases and NoSQL

Gianluca Quercini

gianluca.quercini@centralesupelec.fr

Master DSBA 2020 - 2021



What you will learn

In this lecture you will learn:

- The limitations of the relational data model.
- What a **distributed database** is.
- How data is **distributed** across different machines.
- The availability-consistency trade-off (CAP theorem).
- The main characteristics of **NoSQL databases**.
- The families of NoSQL databases.

Relational data model limitations: impedance mismatch

Definition (Impedance mismatch)

Impedance mismatch refers to the challenges encountered when one needs to map objects used in an application to tables stored in a relational database.



Impedance mismatch: solutions

Object-oriented databases

- Data is stored as **objects**.
- Object-oriented applications save their objects as they are.
- Examples. ConceptBase, Db4o, Objectivity/DB.

Disadvantage

- Not as popular as relational database systems.
- Requires familiarity with object-oriented concepts.
- No standard query language.

Impedance mismatch: solutions

Object relational mappers (ORM)

- Use of libraries that map objects to relational tables.
- The application manipulates objects.
- The ORM library translates object operations into SQL queries.
- Examples. SQLAlchemy, Hibernate, Sequelize.

Disadvantage

- Abstraction. Weak control on how queries are translated.
- Portability. Each ORM has a different set of APIs.

Limitations of the relational model: graph data

Normalization

- In a relational databases, tables are normalized.
- Data on different entities are kept in different tables.
- This reduces redundancy and guarantees integrity.
- In a **normalized** relational database, links between entities are expressed with **foreign key constraints**.
- Need to join different tables (expensive operation).



Limitations of the relational model: data distribution

Objective of a relational database system

- Privilege data integrity and consistency.
- Different mechanisms to ensure integrity and consistency.
 - Primary and foreign key constraints.
 - Transactions.
- Mechanisms to enforce data integrity and consistency have a cost.
 - Manage transactions.
 - Check that new data complies with the given integrity constraints.
- Things get worse in distributed databases.
 - Data is distributed across several machines.
 - Join operations become very expensive.
 - Integrity mechanisms become very expensive.

Distributed database

Definition (Distributed database)

A **distributed database** is one where data is stored across several **machines**, a.k.a, **nodes**.

Shared-nothing architecture

- Each node has its own CPU, memory and storage.
- Nodes only share the network connection.

Pros/cons of a distributed database

- Allows storage and management of large volumes of data. ©
- Far more complex than a single-server database. ③

Distributed database



Distributing data: when?

Small-scale data

- Data distribution is not a good option when the data scale is small.
- With **small-scale data**, the performances of a distributed database are **worse** than a single-server database.
 - **Overhead.** We lose more time distributing and managing data than retrieving it.

Large-scale data

- If the data does not fit in a single machine, data distribution is the only option left.
- Distributed databases allow more concurrent database requests than single-server databases.

Distributing data: how?

Data distribution options

- **Replication.** Multiple copies of the same data stored on different nodes.
- Sharding. Data partitions stored on different nodes.
- Hybrid. Replication + Sharding.

Properties

- Location transparency: applications do not have to be aware of the location of the data.
- **Replication transparency**: applications do not need to be aware that the data is replicated.

- The same piece of data is replicated across different nodes.
 - Each copy is called a replica.
- **Replication factor.** The number of nodes on which the data is replicated.

A	
	•

	Departmen	t
codeD	nameD	budget
14	Administration	300,000
25	Education	150,000
62	Finance	600,000
45	Human Resources	150,000



	Departmen	t
codeD	nameD	budget
14	Administration	300,000
25	Education	150,000
62	Finance	600,000
45	Human Resources	150,000



	Departmen	t
codeD	nameD	budget
14	Administration	300,000
25	Education	150,000
62	Finance	600,000
45	Human Resources	150,000

Advantages

- Scalability. Multiple nodes can serve queries on the same data.
- Latency. Queries can be served by geographically proximate nodes.
- Fault tolerance. The database keeps serving queries even if some nodes fail.

Disadvantages

- **Storage cost.** Storage is used to keep multiple copies of the same data.
- Consistency. All replicas must be kept in sync.

Replica consistency

When a replica is updated, the other replicas must be updated as well.



Synchronous updates

- Updates are propagated immediately to the other replicas.
- Small inconsistency window. The replicas will be inconsistent for a short interval of time. ⁽ⁱ⁾
- If updates are frequent, the database might be too busy propagating updates than serving queries. ③

Asynchronous updates

- Updates are propagated at regular intervals.
- More efficient when updates are frequent. ©
- Long inconsistency window. 🔅

Master-slave replication

- Write operations are only possible on the master node.
- The master node propagates the updates to the slave nodes.
- **Read** operations are served by both the master and the slave nodes.



Master-slave replication

- Prevents write conflicts. ©
 - Only one replica is written at any given time.
- Single point of failure. 🙂
 - If the master fails, write operations are unavailable.
 - Algorithms exist to elect a new master.

• Read conflicts are possible. ③



Master-slave replication read conflict

Two read operations on the same data might return different values.



Peer-to-peer replication

• Read and write operations are possible on any node.



Peer-to-peer replication

- No single point of failure. ©
- Write and read conflicts are possible. ③



Sharding

Sharding

- Data is partitioned into balanced, non-overlapping shards.
- Shards are distributed across the nodes.





Sharding

Advantages

- Load balance. Data can be uniformly distributed across nodes.
- Inconsistencies cannot arise (non-overlapping shards).

Disadvantages

- When a node fails, all its partitions are lost.
- Join operations might need to be performed across nodes.
- When data is added, shards might need to be rebalanced.

Combining replication and sharding



Replication consistency

Keeping in sync all replicas of the same data.

Cross-record consistency

Ensuring the coherence of data in related records. Related records might be on different nodes.



Definition (Distributed transactions)

A **distributed transaction** is a sequence of read/write operations that are applied on data that reside on multiple nodes and are executed as an atomic operation.



Distributed transaction

- The nodes need to **coordinate** before committing the transaction operations on their data.
- The coordination requires an exchange of messages between the **transaction managers** on different nodes.



Distributed transaction

- Data being manipulated by a transaction is **locked**.
 - Locked data is unavailable for both read and write operations.
- Locking guarantees the **consistency** of the database.
- Locking reduces the **availability** of the database.



The CAP theorem

Consistency (C), Availability (A), Partition tolerance (P)

- Consistency. Replicas are in sync and related records are coherent across all nodes.
- Availability. A database can still execute read/write operations when some nodes fail.
- Partition tolerance. The database can still operate when a network partition occurs.



The CAP theorem

Theorem (CAP, Brewer 1999)

Given the three properties of **consistency**, **availability** and **partition tolerance**, a networked shared-data system can have at most two of these properties.

Proof

Suppose that the system is **partition tolerant (P)**. When a network partition occurs, we have two options.

- Allow write operations. This makes the database available (A), but not consistent (C).
 - Some of the replicas might not be synced due to the network partition.

Disable write operations. This makes the database consistent (C) but not available (A).

The CAP theorem

Theorem (CAP, Brewer 1999)

Given the three properties of **consistency**, **availability** and **partition tolerance**, a networked shared-data system can have at most two of these properties.

Proof

- The only way that we can have a **consistent (C)** and **available (A)** database is when network partitions do not occur.
- But if we assume that network partitions never occur, the system is not partition tolerant (P).

Consistency vs Availability

- Relational databases favor consistency over availability.
 - They take a transactional approach to data consistency.
- NoSQL databases favor availability over consistency.
 - In many contexts strong consistency is not necessary.



Alice does not see Bob's post between *t*1 and *t*2. **Is it really an issue?**

ACID vs BASE

ACID (strong consistency)

- Atomicity (A). "All or nothing".
- Consistency (C). From a consistent state to a consistent state.
- Isolation (I). Serializability of transactions.
- Durability (D). Upon commit, all the updates are permanent.

BASE (availability)

- Basic Availability (BA). The database appears to work most of the time.
- Soft state (S). Write and read inconsistencies can occur.
- Eventually consistent (E). The database will be consistent at some point.

NoSQL databases

NoSQL: interpretations of the acronym

- Non SQL: strong opposition to SQL.
- Not only SQL: NoSQL and SQL coexistence.

Goals

- Address the object-relational impedance mismatch.
- Provide better scalability for **distributed databases**.
- Provide a better modeling of **semi-structured data**.

NoSQL databases

Families

- Key-value databases.
- Document-oriented databases.
- Column-oriented databases.
- Graph databases.
- The first three families use the notion of **aggregate** to model the data.
 - They differ in how the aggregates are organized.
- Graph databases are somewhat outliers.
 - They were not conceived for data distribution in mind.
 - They were born ACID-compliant.

There is not a single NoSQL database and there is not a "NoSQL" query language.

Gianluca Quercini

Introduction to Databases



- An aggregate is a data structure used to store the data of a specific entity.
 - In that, it is similar to a row in a relational table.
- We can **nest** an aggregate into another aggregate.
 - This is a huge difference from a row in a relational table.
- An aggregate is a **unit of data** for **replication** and **sharding**.
 - All data in an aggregate will never be split across two shards.
 - All data in an aggregate will always be available on one node.
 - Unlike a relational database, we can control how data is distributed.

Denormalized table

- In a relational database, the following table would not be in **first normal form**.
- The column *categories* contains a list of values.
 - Searching for all products in category kitchen would be hard with SQL.

article_id	name	producer	categories
234543	Bamboo utensil spoon	KitchenMaster	home, kitchen, spatulas

In a relational database, we can address this problem by normalizing the table.

		0		
(alan	LUC2	•••	Herch	n i
Gian	iucu	्र	ucicii	

First normal form

- The following table is in first normal form.
- But we introduced redundancy.
 - What if we update the producer name of the article 234543?
 - In a distributed database, the rows corresponding to this article might be on **different nodes**.

article_id	name	producer	categories
234543	Bamboo utensil spoon	KitchenMaster	home
234543	Bamboo utensil spoon	KitchenMaster	kitchen
234543	Bamboo utensil spoon	KitchenMaster	spatulas

Gianluca Quercini

Introduction to Databases

Master DSBA 2020 - 2021 36 / 58

Second normal form

- To avoid redundancy, we split the table into three tables in **second normal form**.
- In a distributed database, the rows in these tables might be on different nodes.
 - We might need cross-node join operations, which are very expensive.

article		article	_category	cate	gory	
article_id	name	producer	article_id	category_id	category_id	name
234543	Bamboo utensil spoon	KitchenMaster	234543	1	1	kitchen
			234543	2	2	home
			234543	3	3	spatulas

Aggregate

- In an aggregate, list of values are allowed.
- Searching for all products in category kitchen is supported.

```
{
    "article_id": 234543,
    "name": "Bamboo utensil spoon",
    "producer": "KitchenMaster",
    categories: ["home", "kitchen", "spatulas"]
}
```

Mata in an aggregate is never split across different nodes.

- **Denormalization** is allowed in the aggregate.
- Data that are queried together are stored in the same node.

```
"code_employee": 12353,
  "first_name": "John",
  "last_name": "Smith",
  "salary": 50000,
  "position": "Assistant director",
  department: {
      "dept_code": 12,
      "dept_name": "Accounting",
      budget: 120000
    }
}
```

- Aggregates are schemaless.
- Aggregates might not have the same attributes.

```
{
    "code_employee": 12353,
    "first_name": "John",
    "last_name": "Smith",
    "salary": 50000,
    "position": "Assistant director",
    department: {
        "dept_code": 12,
        "dept_name": "Accounting",
        budget: 120000
    }
}
```

```
{
    "code": 345321,
    "first_name": "Jennifer",
    "last_name": "Green",
}
```

We don't need to fix a rigid the schema. NULL values are avoided.

Gianluca Quercini

```
"code_employee": 12353,
  "first_name": "John",
  "last_name": "Smith",
  "salary": 50000,
  "position": "Assistant director",
  departments: [
      "dept_code": 12,
      "dept_name": "Accounting",
      budget: 120000
    },
      "dept_code": 145,
      "dept_name": "HR",
      budget: 250000
    3
}
```

```
"code_employee": 12353,
 "first_name": "John",
 "last_name": "Smith",
 "salary": 50000,
 "position": "Assistant director",
 departments: [
      "dept_code": 12,
      "dept_name": "Accounting",
      budget: 120000
   },
                               We can update atomically the
      "dept_code": 145,
      "dept_name": "HR",
                                   salary of an employee. How would
      budget: 250000
                                   we represent the same in a rela-
    3
                                   tional database?
}
```

- We use a **denormalized table** (same as aggregate).
- **However**, we have no guarantees that the rows relative to the employee John Smith will be stored in the same node.

code_emp	first_name	last_name	salary	position	dept_code	dept_name	budget
234543	John	Smith	50000	Assistant director	12	Accounting	120000
234543	John	Smith	50000	Assistant director	145	HR	250000

The update of the salary of a single employee might be a **cross-node operation**.

```
"code_employee": 12353,
"first_name": "John",
"last_name": "Smith",
"salary": 50000,
"position": "Assistant director",
departments: [
    "dept_code": 12,
    "dept_name": "Accounting",
    budget: 120000
  },
    "dept_code": 145,
                           Updating the information on a
    "dept_name": "HR",
                               department is a non-atomic op-
    budget: 250000
                               eration
```

```
"code_employee": 12353,
"first_name": "John",
"last_name": "Smith",
"salary": 50000,
"position": "Assistant director",
departments: [
    "dept_code": 12,
    "dept_name": "Accounting",
    budget: 120000
  },
    "dept_code": 145,
                            <sup>™</sup> We'll see how to alleviate this
    "dept_name": "HR",
                                problem when we introduce Mon-
    budget: 250000
                                goDB.
```

Aggregate-based NoSQL databases

- Aggregates are schemaless.
 - No need to adhere to a rigid schema.
 - Flexible evolution of the database.
- Normalization is not required.
 - We accept some redundancies in exchange of faster queries.
 - Remember: storage hardware is **cheap** today.
- All data in an aggregate is stored in a **single node**.
 - With aggregates, we are in control of how the data is distributed.
- In general, updates on an aggregate are atomic operations.
 - If an update entails many write operations, either all are executed or none.
- Cross-aggregate updates are **not guaranteed** to be atomic.
 - Multi-aggregate transactions might be supported and used if necessary.

Idea

Data are modeled as key-value pairs.

- Key: alphanumeric string, usually auto-generated by the database.
- Value: an aggregate.
- Query: get an aggregate given its key.



Idea

- Data is partitioned based on the key.
- Partitions are distributed across different nodes.
- Little to no checks on integrity constraints.
- Goal. High scalability and fast read/write queries.



Application scenarios

Scenario 1. Session store.

- A Web application starts a session when a user logs in.
- The application stores session data in the database.
 - User profile information, messages, personalized themes...
- Each session is assigned a **unique identifier** (the key).
- Session data is only queried by the identifier.



Application scenarios

Scenario 2. Shopping cart.

- An e-commerce website may receive billions of orders in seconds.
- Each shopping cart has a **unique identifier** (the key).
- Shopping cart data is only queried by the identifier.
- Shopping cart data can be easily replicated to handle node failures.



Existing key-value databases

- Amazon DynamoDB. One of the first NoSQL databases.
- Riak.
- Redis. Possibility of tuning data persistence.
- Voldemort.



Document-oriented databases

Idea

• Data is modeled as **key-value pairs**, and searching aggregates based on their **attribute values** is supported.

Database Collection Document

product_id: 12234345 name: "Bamboo utensil spoon" categories: ["home", "kitchen", "spatulas"]

Document

product_id: 98761 name: "Mini round cocotte" categories: ["home", "kitchen", "dining"] It is possible to search for all products in category *kitchen*.

Document-oriented databases

Existing document-oriented databases

• MongoDB, CouchDB, OrientDB.

Database	
Collection	
Document	
t id: 12234345	
"Bamboo utensil spoon"	
ries: ["home", "kitchen", "spatulas"]	
Document	
t id: 98761	
"Mini round cocotte"	
ries: ["home". "kitchen". "dining"]	
	Collection Document Ict_id: 12234345 : "Bamboo utensil spoon" pories: ["home", "kitchen", "spatulas"] Document Ict_id: 98761 : "Mini round cocotte" terriori "kitchere", "kitchere" 1

Idea

• Similar to document-oriented database but. an aggregate can be broken into smaller data units called **columns**.



Idea

- Columns can be organized into column families.
- Columns in the same family are stored on the same node.



Idea

• The value of a column can be an aggregate (wide column).



Existing column-oriented databases

• Cassandra, HBase, BigTable (Google).



Graph databases

Idea

- Their data model is optimized for storing and retrieving graph data.
- Relationships are first-class citizens.
 - In relational databases they are implicit in foreign key constraints.
 - In aggregate-based NoSQL stores, they are represented with nested aggregates or references.
- Existing graph databases: Neo4j, InfiniteGraph, AllegroGraph.

NoSQL databases: conclusions

Polyglot persistence

- NoSQL databases are not going to replace relational databases.
- Use of different data storage technologies based on the data type.
- This is called polyglot persistence.



References

• Hoffer, Jeffrey A. *Modern Database Management.* 10/e. Pearson Education India, 2011. • Click here