

Big Data Analytics

6. NoSQL Databases

Lars Schmidt-Thieme

Information Systems and Machine Learning Lab (ISMLL)
Institute of Computer Science
University of Hildesheim, Germany

original slides by Lucas Rego Drumond, ISMLL

Outline

1. Introduction
2. Key-Value Stores
3. Document Databases
4. Graph Databases
5. Column Databases and Object Databases

Outline

1. Introduction
2. Key-Value Stores
3. Document Databases
4. Graph Databases
5. Column Databases and Object Databases

Relational vs Big-Data Technologies

- ▶ **Structure:** In relational data bases, the data is structured but in big data technologies the data is raw.
- ▶ **Process:** Relational data bases were initially created for transactional processing, but big data technologies are for data analysis.
- ▶ **Entities:** A relational data base is big if it has many entities and many rows per each entity. But in big data technologies, the number of entities is not necessary high. The amount of information that exist for entities is big.
- ▶ **horizontal scaling:** Adding new nodes in big data technologies can be done with low cost while in relational data base, either it fixed or it is expensive.

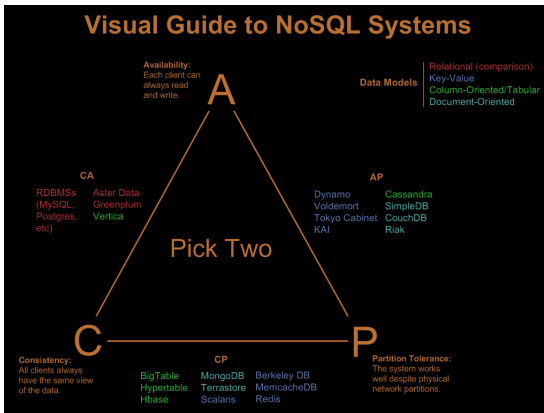
NoSQL Databases

- ▶ A NoSQL or Not Only SQL database provides a mechanism for storage and retrieval of data that is modeled in means **other than** the tabular relations used in relational databases.
- ▶ Motivations for this approach include simplicity of design and horizontal scaling.
- ▶ The data structure differs from the RDBMS, and therefore some operations are faster in NoSQL and some in RDBMS.
- ▶ Most NoSQL stores lack true ACID transactions.

NoSQL Databases / Types and Implementations

- ▶ **Key-value:** Dynamo, FoundationDB, MemcacheDB, Redis, Riak
- ▶ **Document:** Clusterpoint, Couchbase, MarkLogic, MongoDB
- ▶ **Graph:** Allegro, Neo4J, OrientDB, Virtuoso
- ▶ **Column:** Accumulo, Cassandra, HBase
- ▶ **Object:** Db4o, ZODB

...



<http://blog.scottlogic.com/2014/08/04/mongodb-vs-couchdb.html>

Outline

1. Introduction
2. Key-Value Stores
3. Document Databases
4. Graph Databases
5. Column Databases and Object Databases

Key-Value stores

- ▶ Key-Value stores use the associative array as their fundamental data model.
- ▶ In this model, data is represented as a collection of key-value pairs.
- ▶ The key-value model is one of the simplest non-trivial data models.

Example:

```
{  
  "Great Expectations": "John",  
  "Pride and Prejudice": "Alice",  
  "Wuthering Heights": "Alice"  
}
```

Outline

1. Introduction
2. Key-Value Stores
- 3. Document Databases**
4. Graph Databases
5. Column Databases and Object Databases

Document Databases

- ▶ A document-oriented database is a computer program designed for storing, retrieving, and managing document-oriented information.
- ▶ In contrast to relational databases and their notions of "Relations" (or "Tables"), these systems are designed around an abstract notion of a "Document".
- ▶ Documents inside a document-oriented database are not required to have all the same sections, slots, parts, or keys.

Example 1 / Schema-less

```
1 {
2   _id: ObjectId(7df78ad8902c),
3   FirstName: "Bob",
4   Age: 35,
5   Address: "5_Oak_St.",
6   Hobby: "sailing"
7 }
```

```
1 {
2   _id: ObjectId(5df78ad8902c),
3   FirstName: "Jonathan",
4   Age: 37,
5   Address: "15_Wanamassa_Point_Road",
6   Languages: [ 'English', 'German' ]
7 }
```

Example 2 / Foreign Keys

```
1 {
2   _id: ObjectId(5df78ad8902c),
3   FirstName: "Jonathan",
4   Age: 37,
5   Address: "15_Wanamassa_Point_Road",
6   Children: [ ObjectId(5df78ad89020), ObjectId(5df78ad89021),
7             ObjectId(5df78ad89022), ObjectId(5df78ad89023) ]
8 },{
9   _id: ObjectId(5df78ad89020),
10  FirstName: "Michael",
11  Age: 10,
12 },{
13  _id: ObjectId(5df78ad89021),
14  FirstName: "Jennifer",
15  Age: 8,
16 },{
17  _id: ObjectId(5df78ad89022),
18  FirstName: "Samantha",
19  Age: 5,
20 },{
21  _id: ObjectId(5df78ad89023),
22  FirstName: "Elena",
23  Age: 2,
24 }
```

Example 3 / Embedded Documents

```
1 {
2   _id: ObjectId(5df78ad8902c),
3   FirstName: "Jonathan",
4   Age: 37,
5   Address: "15_Wanamassa_Point_Road",
6   Children: {
7     { FirstName: "Michael", Age: 10 },
8     { FirstName: "Jennifer", Age: 8 },
9     { FirstName: "Samantha", Age: 5 },
10    { FirstName: "Elena", Age: 2}
11  }
12 }
```

Document Databases

- ▶ Documents are addressed in the database via a unique key that represents that document.
- ▶ Documents can be retrieved by their
 - ▶ key
 - ▶ content
- ▶ Documents are organized through
 - ▶ Collections
 - ▶ Tags
 - ▶ Non-visible Metadata
 - ▶ Directory hierarchies
 - ▶ Buckets

RDBMS vs. Document Database Terminology

relational database	document database
Database	Database
Table	Collection
Tuple/Row	Document
Column	Field
Primary key	Primary key (e.g., default key <code>_id</code> in mongodb)
Foreign key	Foreign key

- ▶ In document databases, join operations often are replaced by **embedding documents**.
- ▶ Document databases are **schema-less**
 - ▶ one can add a new field at any time
 - ▶ a document basically is a dictionary (like Python objects)

Big Data Document Databases

Document databases are useful for big data for two main reasons:

1. Document databases can be **horizontally partitioned / sharded**.
 - ▶ Documents are distributed over different nodes.
 - ▶ Using a **partition/sharding function**
 - ▶ e.g., as hash function
 - ▶ works exactly the same ways as for RDBMS.

Big Data Document Databases

Document databases are useful for big data for two main reasons:

1. Document databases can be **horizontally partitioned / sharded**.
 - ▶ Documents are distributed over different nodes.
 - ▶ Using a **partition/sharding function**
 - ▶ e.g., as hash function
 - ▶ works exactly the same ways as for RDBMS.
2. Documents are **sparse representations**,
 - ▶ only some fields/keys have values, while tuples/rows are **dense**.
 - ▶ all columns have values stored explicitly
 - ▶ also NULL is explicitly stored.

Outline

1. Introduction
2. Key-Value Stores
3. Document Databases
- 4. Graph Databases**
5. Column Databases and Object Databases

Graph databases

A graph database is a database that uses graph structures with nodes, edges, and properties to represent and store data.

Graph databases

- ▶ Nodes represent entities such as people, businesses, accounts, or any other item you might want to keep track of.
- ▶ Properties are relevant information that relate to nodes.
- ▶ Edges are the lines that connect nodes to nodes
- ▶ Most of the important information is really stored in the edges.
- ▶ Meaningful patterns emerge when one examines the connections and interconnections of nodes, properties, and edges.

Graph databases

- ▶ Compared with relational databases, graph databases are often faster for associative data sets
- ▶ They map more directly to the structure of object-oriented applications.
- ▶ As they depend less on a rigid schema, they are more suitable to manage ad hoc and changing data with evolving schemas.
- ▶ Graph databases are a powerful tool for graph-like queries.

Graph queries

- ▶ Reachability queries
- ▶ shortest path queries
- ▶ Pattern queries

Outline

1. Introduction
2. Key-Value Stores
3. Document Databases
4. Graph Databases
5. Column Databases and Object Databases

Column Databases

- ▶ Column databases are a tuple (a key-value pair) consisting of three elements:
 - ▶ Unique name: Used to reference the column
 - ▶ Value: The content of the column.
 - ▶ Timestamp: The system timestamp used to determine the valid content.
- ▶ Differences to a relational database

Example

```
{  
  street: name: "street", value: "1234 x street", timestamp: 123456789,  
  city: name: "city", value: "san francisco", timestamp: 123456789,  
  zip: name: "zip", value: "94107", timestamp: 123456789,  
}
```

Object Database

- ▶ An object database is a database management system in which information is represented in the form of objects as used in object-oriented programming.
- ▶ Most object databases also offer some kind of query language, allowing objects to be found using a declarative programming approach (OQL)
- ▶ Access to data can be faster because joins are often not needed.
- ▶ Many object databases offer support for versioning.
- ▶ They are specially suitable in applications with complex data.