

## UNIT 3 - NOSQL Databases.

### Introduction

- ↳ NOSQL databases are designed to manage large volumes of unstructured & semi-structured data
- ↳ offer more flexibility
- ↳ "not only SQL"
- ↳ Supports consistency rather than strict ACID properties of RDBMS
- ↳ generally avoids join operations.

### History

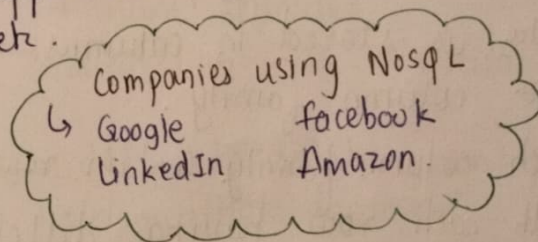
- ↳ first used by Carlo Strozzi in 1998
- ↳ he mention this name for his open source Database System, where no provision of SQL Query interface.
- ↳ 2009 NOSQL becomes in practice

### Advantages

- ↳ good resource scalability, low cost, support semi-structured data, no static schema, fast processing, etc.

### Disadvantages

- ↳ not a defined standard
- ↳ limited query capability.



### Four Types of NoSQL Database.

- ① Key-value store database
- ② Document Database
- ③ column-store database
- ④ Graph database

### Key-Value store database

- ↳ Data is stored as pairs of keys & values.
- ↳ Keys are unique and their corresponding values can be any type of data (string, JSON, binary, etc.)
- ↳ specially designed for storing data as schema-free data

### Advantages :-

- ↳ Simple & efficient
- ↳ Highly scalable & fault tolerant

### Disadvantages :-

- ↳ limited query capability  
ie complex query is challenge
- ↳ for many-to-many relationship, it show poor performance

### Examples :-

Redis  
DynamoDB  
Cassandra.  
Azure Table Storage



Document stores database.

- ↳ Data is stored in documents that contain complex structures, including nested objects and arrays.
- ↳ formats like JSON or XML.
- ↳ each document contain nested structures & various data types, allowing for flexibility

Advantages:-

- ↳ excellent for managing semi-structured data
- ↳ Easy retrieval & indexing of documents.

Disadvantages:-

- ↳ Handling multiple docs is complex
- ↳ Aggregate operations may not work accurately

Examples:-

MongoDB  
CouchDB

Column store database / column-family stores

- ↳ data is stored in columns rather than rows, organized into column family.
- ↳ Each column family contain rows identified by unique key, with each row having different set of columns.

Advantages:-

- ↳ Highly scalable & efficient
- ↳ good support for large query

Disadvantages:-

- ↳ limited support for ad-hoc & transactions queries
- ↳ complex in data modelling

Examples:-

- ↳ Apache Cassandra
- ↳ HBase
- ↳ BigTable
- ↳ HyperTable

Graph Database:-

- ↳ Data is represented as nodes (entity) and edges (relationships), making it ideal for complex interconnected data

OR

Data is stored as graph and their relationships are stored as link between them whereas entity acts as a node.

Advantages:-

- ↳ flexible model
- ↳ effective for applications like social networks

Disadvantages:-

- ↳ may struggle with heavy write workloads
- ↳ limited support for ad-hoc query
- ↳ does not offer better choice over others.

Applications

- ↳ Neo4j
- ↳ OrientDB



## MongoDB

- ↳ open source document database management system.
- ↳ one of the popular NoSQL database, written in c++.
- ↳ use Binary JSON (BSON) to store documents, which allows complex data types & structures.
- ↳ Imp. features include High performance, High availability.
- ↳ Has its own ad-hoc query language
- ↳ As a cache memory, it automatically uses all free memory available on the machine.

### Features

- High performance
- Good support for embedded data models.
- Replication & High availability
- Document oriented model
- Schema less Design
- Support horizontal scalability
- Aggregation framework
- flexible Querying
- Integration with Big Data Technologies

### Architecture

**Database & Collection** - MongoDB organizes data into database containing multiple collections. Each collection hold multiple documents.

**Documents** - Each document is a self-contained unit with its own structure, represent as key-value pairs

**BSON Format** - Documents are stored in BSON format, which supports various datatypes like arrays & nested objects.

### Use Cases

- CMS (Content Management Systems)
- Real Time Analytics
- IOT
- E-commerce Applications
- Social Networks.

### Limitations

- High memory usage
- Limited Nesting
- Joins not supported
- Limited Data Size.

### Column-Oriented Database :-

#### Apache Cassandra

- ↳ Wide column store, which combines features of both key-value and tabular databases (ie, column oriented)

↳ designed to handle large amount of data across many servers with high availability and no single point of failure.

### Features

- High Availability
- Scalability
- Partitioned Row Store
- Consistency
- CQL (Cassandra QL)
- Flexible Schema



## Use Cases

- ↳ Real time Analytics
- ↳ Social Media
- ↳ IOT Appl<sup>n</sup>
- ↳ E-commerce
- ↳ Content management System.

## RDBMS & NoSQL database Comparison.

### RDBMS

- ① Structured data in tables (rows & columns)
- ② Pre defined schema; rigid structure.
- ③ Vertically scalable (add more power to existing hardware)
- ④ SQL is standard language
- ⑤ Support complex transactions with ACID properties
- ⑥ Optimized for complex queries & joins.
- ⑦ Strong data integrity
- ⑧ Use case - banking system
- ⑨ Eg. MySQL, Oracle, MSSQL Server

### NoSQL

- ① Unstructured / semi-structured data
- ② Dynamic ~~static~~ schema; flexible structure
- ③ Horizontally scalable (add more power to servers)
- ④ Varies by database type.
- ⑤ Simpler transactions; may not fully support ACID.
- ⑥ High performance for read/write operations.
- ⑦ Weak data integrity.
- ⑧ Use case - social media, IOT
- ⑨ Eg. MongoDB, Cassandra, Redis, Couchbase

## NoSQL Database Development Tools.

### ① Map Reduce

↳ A programming model used for processing & generating large datasets with a distributed algorithm on cluster.

Components :-

Map function - processes input data & produce key-value pairs.

Reduce function - aggregates results from Map function to produce a final output.

↳ used for batch processing of large datasets.

↳ used commonly in conjunction with Hadoop to analyze big data



## Advantages -

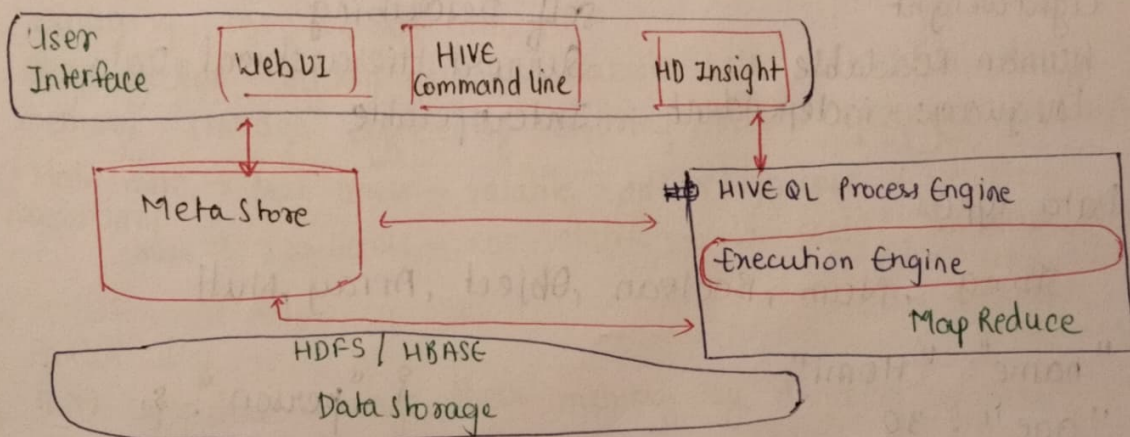
- ↳ Scalable to handle vast amount of data across distributed systems.
- ↳ Fault-tolerant.

## HIVE

- ↳ data warehouse infrastructure tool.
- ↳ built on top of Hadoop.
- ↳ designed for querying & managing large datasets using sql like interface called HiveQL.
- ↳ used to summarize big data, analysis of big data.
- ↳ supports ad-hoc queries.
- ↳ SQL type scripts can be created for MapReduce operations using HIVE

### Features :-

- ↳ SQL-like query language.
- ↳ Batch Processing.
- ↳ Schema management.
- ↳ scalable
- ↳ fast



Working from Textbook Pg 5-16.

## xml - Extensible Markup Language

store & share data in structured format.

looks like HTML but is more flexible

Features → Tags : opening & closing Tags

Hierarchy : data is stored in tree like structure

Attributes : extra info can be added to tags

Custom Tags : create your own tags

Eg: 

```
<person>
  <name> John </name>
  <age> 30 </age>
</person>
```

## Types of XML schemas

DTD - Document Type Defn  
XSD - XML Type Defn.

	DTD	XSD
Syntax	Non-XML	XML based
datatype support	Basic	Extensive (numbers, dates)
Namespace support	No	Yes
	less flexible	more flexible

## JSON (Java Script Object Notation)

stores & exchanges lightweight data  
key value pairs

features :-

Lightweight

Human Readable

Language-independent

self Describing

Support Hierarchical Data

Interoperable

## Data types

String, Num, Boolean, Object, Array, Null

"name": "Homi"

"age": 30

"isStudent": true

"address": {"city": "NY"}

"marks": [85, 90, 95]

"middle name": null.

{ "person": {

"name": "Homi",

"age": 30 }

}

<name>John</name>

<age>30</age>

</person>