## 4.1   NoSQL (NOT ONLY SQL)

The term NoSQL was first coined by Carlo Strozzi in 1998 to name his lightweight, open-source, relational database that did not expose the standard SQL interface. Johan Oskarsson, who was then a developer at last. fm, in 2009 reintroduced the term NoSQL at an event called to discuss open-source distributed network. The #NoSQL was coined by Eric Evans and few other database people at the event found it suitable to describe these non-relational databases.

Few features of NoSQL databases are as follows:

1. They are open source.
2. They are non-relational.
3. They are distributed.
4. They are schema-less.
5. They are cluster friendly.
6. They are born out of 21<sup>st</sup> century web applications.

### 4.1.1   Where is it Used?

NoSQL databases are widely used in big data and other real-time web applications. Refer Figure 4.1. NoSQL databases is used to stock log data which can then be pulled for analysis. Likewise it is used to store social media data and all such data which cannot be stored and analyzed comfortably in RDBMS.

### 4.1.2   What is it?

**NoSQL** stands for Not Only SQL. These are non-relational, open source, distributed databases. They are hugely popular today owing to their ability to scale out or scale horizontally and the adeptness at dealing with a rich variety of data: structured, semi-structured and unstructured data. Refer Figure 4.2 for additional features of NoSQL. NoSQL databases.

1. **Are non-relational:** They do not adhere to relational data model. In fact, they are either key–value pairs or document-oriented or column-oriented or graph-based databases.
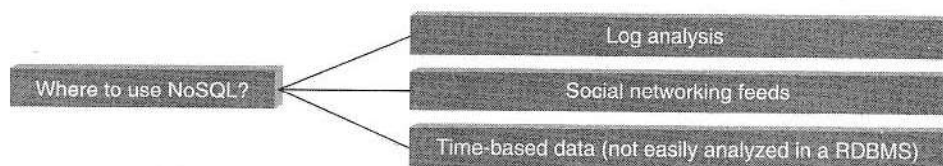


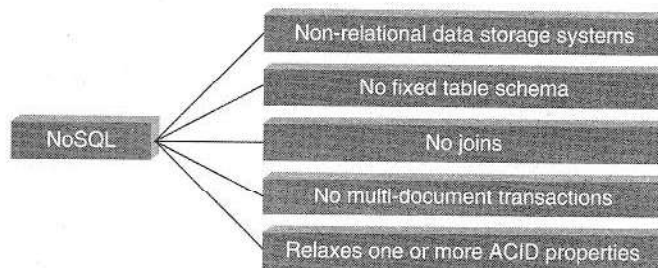**Figure 4.1**   Where to use NoSQL?



**Figure 4.2**   What is NoSQL?

2. **Are distributed:** They are distributed meaning the data is distributed across several nodes in a cluster constituted of low-cost commodity hardware.

3. **Offer no support for ACID properties (Atomicity, Consistency, Isolation, and Durability):** They do not offer support for ACID properties of transactions. On the contrary, they have adherence to Brewer's CAP (Consistency, Availability, and Partition tolerance) theorem and are often seen compromising on consistency in favor of availability and partition tolerance.

4. **Provide no fixed table schema:** NoSQL databases are becoming increasing popular owing to their support for flexibility to the schema. They do not mandate for the data to strictly adhere to any schema structure at the time of storage.

### 4.1.3 Types of NoSQL Databases

We have already stated that NoSQL databases are non-relational. They can be broadly classified into the following:

1. Key–value or the big hash table.
2. Schema-less.

Refer Figure 4.3. Let us take a closer look at key–value and few other types of schema-less databases:

1. **Key–value:** It maintains a big hash table of keys and values. For example, Dynamo, Redis, Riak, etc.
   ***Sample Key–Value Pair in Key–Value Database***

   | Key | Value |
   |------------|----------|
   | First Name | Simmonds |
   | Last Name | David |

2. **Document:** It maintains data in collections constituted of documents. For example, MongoDB, Apache CouchDB, Couchbase, MarkLogic, etc.
   ***Sample Document in Document Database***

   ```
   {
   "Book Name":        "Fundamentals of Business Analytics",
   "Publisher":        "Wiley India",
   "Year of Publication":     "2011"
   }
   ```

3. **Column:** Each storage block has data from only one column. For example: Cassandra, HBase, etc.
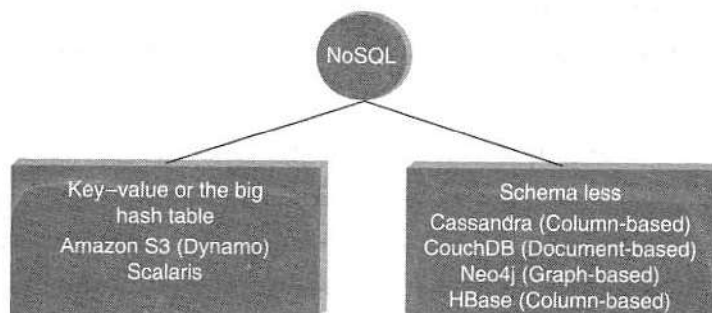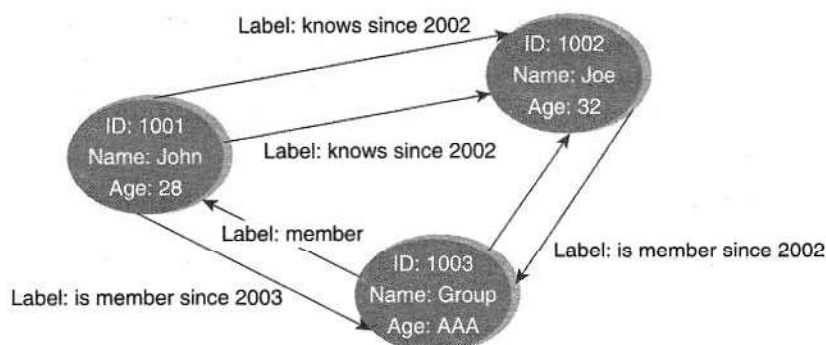


**Figure 4.3** Types of NoSQL databases.

4. **Graph:** They are also called network database. A graph stores data in nodes. For example, Neo4j, HyperGraphDB, etc.

   ***Sample Graph in Graph Database***



Refer Table 4.1 for popular schema-less databases.

### 4.1.4   Why NoSQL?

1. It has scale out architecture instead of the monolithic architecture of relational databases.
2. It can house large volumes of structured, semi-structured, and unstructured data.
3. **Dynamic schema:** NoSQL database allows insertion of data without a pre-defined schema. In other words, it facilitates application changes in real time, which thus supports faster development, easy code integration, and requires less database administration.
4. **Auto-sharding:** It automatically spreads data across an arbitrary number of servers. The application in question is more often not even aware of the composition of the server pool. It balances the load of data and query on the available servers; and if and when a server goes down, it is quickly replaced without any major activity disruptions.
5. **Replication:** It offers good support for replication which in turn guarantees high availability, fault tolerance, and disaster recovery.

### 4.1.5   Advantages of NoSQL

Let us enumerate the advantages of NoSQL. Refer Figure 4.4.

1. **Can easily scale up and down:** NoSQL database supports scaling rapidly and elastically and even allows to scale to the cloud.

**Table 4.1**   Popular schema-less databases

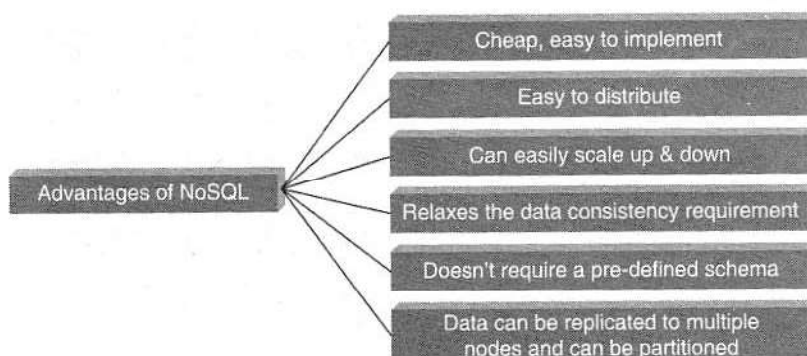| Key–Value Data Store | Column-Oriented Data Store | Document Data Store | Graph Data Store |
|---|---|---|---|
| • Riak | • Cassandra | • MongoDB | • InfiniteGraph |
| • Redis | • HBase | • CouchDB | • Neo4j |
| • Membase | • HyperTable | • RavenDB | • AllegroGraph |

**Figure 4.4** Advantages of NoSQL.

(a) *Cluster scale:* It allows distribution of database across 100+ nodes often in multiple data centers.
(b) *Performance scale:* It sustains over 100,000+ database reads and writes per second.
(c) *Data scale:* It supports housing of 1 billion+ documents in the database.

2. **Doesn't require a pre-defined schema:** NoSQL does not require any adherence to pre-defined schema. It is pretty flexible. For example, if we look at MongoDB, the documents (equivalent of records in RDBMS) in a collection (equivalent of table in RDBMS) can have different sets of key–value pairs.

   {_id: 101,"BookName": "Fundamentals of Business Analytics", "AuthorName": "Seema Acharya", "Publisher": "Wiley India"}
   {_id:102, "BookName":"Big Data and Analytics"}

3. **Cheap, easy to implement:** Deploying NoSQL properly allows for all of the benefits of scale, high availability, fault tolerance, etc. while also lowering operational costs.

4. **Relaxes the data consistency requirement:** NoSQL databases have adherence to CAP theorem (Consistency, Availability, and Partition tolerance). Most of the NoSQL databases compromise on consistency in favor of availability and partition tolerance. However, they do go for eventual consistency.

5. **Data can be replicated to multiple nodes and can be partitioned:** There are two terms that we will discuss here:

   (a) *Sharding:* Sharding is when different pieces of data are distributed across multiple servers. NoSQL databases support auto-sharding; this means that they can natively and automatically spread data across an arbitrary number of servers, without requiring the application to even be aware of the composition of the server pool. Servers can be added or removed from the data layer without application downtime. This would mean that data and query load are automatically balanced across servers, and when a server goes down, it can be quickly and transparently replaced with no application disruption.

   (b) *Replication:* Replication is when multiple copies of data are stored across the cluster and even across data centers. This promises high availability and fault tolerance.

## 4.1.6 What We Miss With NoSQL?

With NoSQL around, we have been able to counter the problem of scale (NoSQL scales out). There is also the flexibility with respect to schema design. However there are few features of conventional RDBMS that are greatly missed. Refer Figure 4.5.
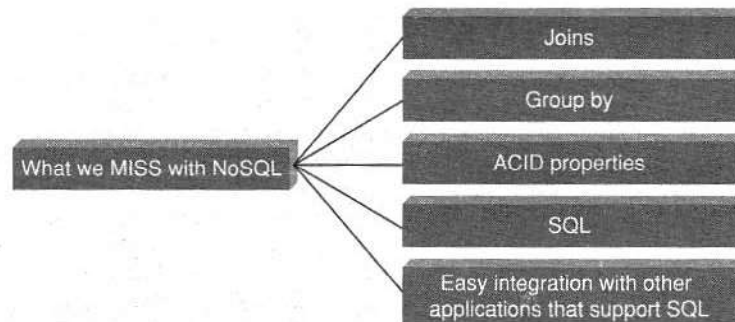
**Figure 4.5** What we miss with NoSQL?

NoSQL does not support joins. However, it compensates for it by allowing embedded documents as in MongoDB. It does not have provision for ACID properties of transactions. However, it obeys the Eric Brewer's CAP theorem. NoSQL does not have a standard SQL interface but NoSQL databases such as MongoDB and Cassandra have their own rich query language [MongoDB query language and Cassandra query language (CQL)] to compensate for the lack of it. One thing which is dearly missed is the easy integration with other applications that support SQL.

## 4.1.7 Use of NoSQL in Industry

NoSQL is being put to use in varied industries. They are used to support analysis for applications such as web user data analysis, log analysis, sensor feed analysis, making recommendations for upsell and cross-sell, etc. Refer Figure 4.6.
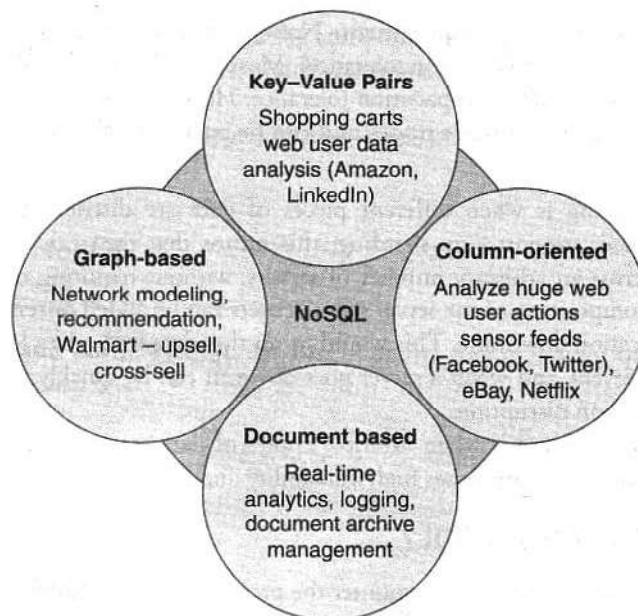


**Figure 4.6** Use of NoSQL in industry.

### 4.1.8 NoSQL Vendors

Refer Table 4.2 for few popular NoSQL vendors.

**Table 4.2** Few popular NoSQL vendors

| Company | Product | Most Widely Used by |
|---------|---------|---------------------|
| Amazon | DynamoDB | LinkedIn, Mozilla |
| Facebook | Cassandra | Netflix, Twitter, eBay |
| Google | BigTable | Adobe Photoshop |

### 4.1.9 SQL versus NoSQL

Refer Table 4.3 for few salient differences between SQL and NoSQL.

**Table 4.3** SQL versus NoSQL

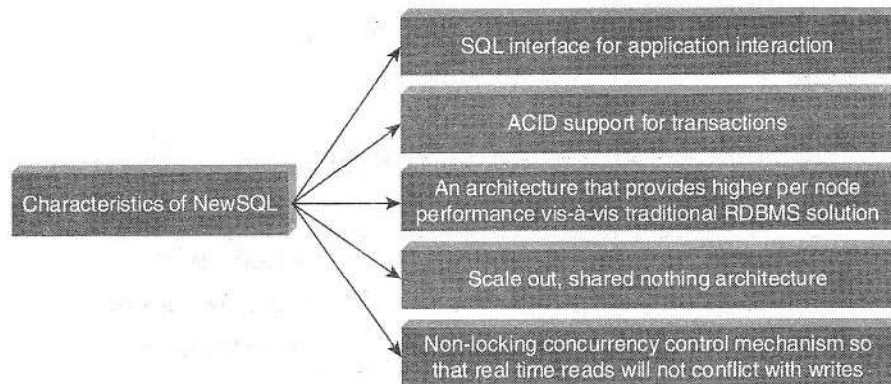| SQL | NoSQL |
|-----|-------|
| Relational database | Non-relational, distributed database |
| Relational model | Model-less approach |
| Pre-defined schema | Dynamic schema for unstructured data |
| Table based databases | Document-based or graph-based or wide column store or key–value pairs databases |
| Vertically scalable (by increasing system resources) | Horizontally scalable (by creating a cluster of commodity machines) |
| Uses SQL | Uses UnQL (Unstructured Query Language) |
| Not preferred for large datasets | Largely preferred for large datasets |
| Not a best fit for hierarchical data | Best fit for hierarchical storage as it follows the key–value pair of storing data similar to JSON (Java Script Object Notation) |
| Emphasis on ACID properties | Follows Brewer's CAP theorem |
| Excellent support from vendors | Relies heavily on community support |
| Supports complex querying and data keeping needs | Does not have good support for complex querying |
| Can be configured for strong consistency | Few support strong consistency (e.g., MongoDB), some others can be configured for eventual consistency (e.g., Cassandra) |
| Examples: Oracle, DB2, MySQL, MS SQL, PostgreSQL, etc. | Examples: MongoDB, HBase, Cassandra, Redis, Neo4j, CouchDB, Couchbase, Riak, etc. |

**Figure 4.7**   Characteristics of NewSQL.

## 4.1.10   NewSQL

There is yet another new term doing the rounds – "NewSQL". So, what is NewSQL and how is it different from SQL and NoSQL?

What is that we love about NoSQL and is not there with our traditional RDBMS and what is that we love about SQL that NoSQL does not have support for? You guessed it right!!! We need a database that has the same scalable performance of NoSQL systems for On Line Transaction Processing (OLTP) while still maintaining the ACID guarantees of a traditional database. This new modern RDBMS is called NewSQL. It supports relational data model and uses SQL as their primary interface.

### 4.1.10.1   Characteristics of NewSQL

Refer Figure 4.7 to learn about the characteristics of NewSQL. NewSQL is based on the shared nothing architecture with a SQL interface for application interaction.

## 4.1.11   Comparison of SQL, NoSQL, and NewSQL

Refer Table 4.4 for a comparative study of SQL, NoSQL and NewSQL.

**Table 4.4**   Comparative study of SQL, NoSQL and NewSQL

|  | SQL | NoSQL | NewSQL |
|---|---|---|---|
| Adherence to ACID properties | Yes | No | Yes |
| OLTP/OLAP | Yes | No | Yes |
| Schema rigidity Adherence to data model | Yes Adherence to relational model | No | Maybe |
| Data Format Flexibility | No | Yes | Maybe |
| Scalability | Scale up Vertical Scaling | Scale out Horizontal Scaling | Scale out |
| Distributed Computing | Yes | Yes | Yes |
| Community Support | Huge | Growing | Slowly growing |