# Assignment11

## Task1:

Explain the below concepts with an example in brief.

● Nosql Databases

- NoSQL is an approach to database design that can accomodate a wide variety of data models, including key-value, document, columnar and graph formats.
- NoSQL, which stand for "not only SQL," is an alternative to traditional relational databases in which data is placed in tables and data schema is carefully designed before the database is built.
- NoSQL databases are especially useful for working with large sets of distributed data.
- NoSQL databases are very flexible in structure and can store all types of related data in one place
- User can retrieve the whole post with a single query avoiding joins thus increasing the performance
- Data on NoSQL databases scale out naturally and hence able to deal with the continuous streaming of posts
- Features:-
  - Generic Data Model :–
    Heterogeneous containers, including sets, maps, and arrays
  - Dynamic type discovery and conversion :–
    NoSQL analytics systems support runtime type identification and conversion so that custom business logic can be used to dictate analytic treatment of variation.
  - Non-relational and De-normalised :-
    Data is stored in single tables as compared to joining multiple tables.
  - Commodity hardware : –
    Adding more of the economical servers allows NoSQL databases to scale to handle more data.
  - Highly distributable :–
    Distributed databases can store and process a set of information on more than one device.
- Example: MongoDB, Cassandra

- Types of Nosql Databases

There are 4 basic types of NoSQL databases:
i)Key-Value Store :-
- Key value type basically, uses a hash table in which there exists a unique key and a
- pointer to a particular item of data.
- A bucket is a logical group of keys – but they don't physically group the data.
- There can be identical keys in different buckets. The data which is a collection of key value pairs is compressed as a document store quite similar to a key-value store, but the only difference is that the values stored (referred to as "documents") provide some structure and encoding of the managed data.
- There is no complexity around the Key Value Store database model as it can be
  implemented in a breeze
- Example- Riak, Amazon S3 (Dynamo)

ii) Document-based Store :-
- The data which is a collection of key value pairs is compressed as a document store quite similar to a key-value store but the only difference is that the values stored (referred to as "documents")provide some structure and encoding of the managed data.
- It stores documents made up of tagged elements.
- XML, JSON (Java Script Object Notation), BSON (which is a binary encoding of JSON objects) are some common standard encodings
- Example- CouchDB

iii)Column-based Store :-
- In column-oriented NoSQL database, data is stored in cells grouped in columns of data rather than as rows of data.
- Columns are logically grouped into column families.
- Column families can contain a virtually unlimited number of columns that can be created at runtime or the definition of the schema.
- Read and write is done using columns rather than rows.
- Each storage block contains data from only one column.
- Column oriented databases are faster in access compared to row oriented like RDBMS as data retrieval become easy when stored as seperate column.
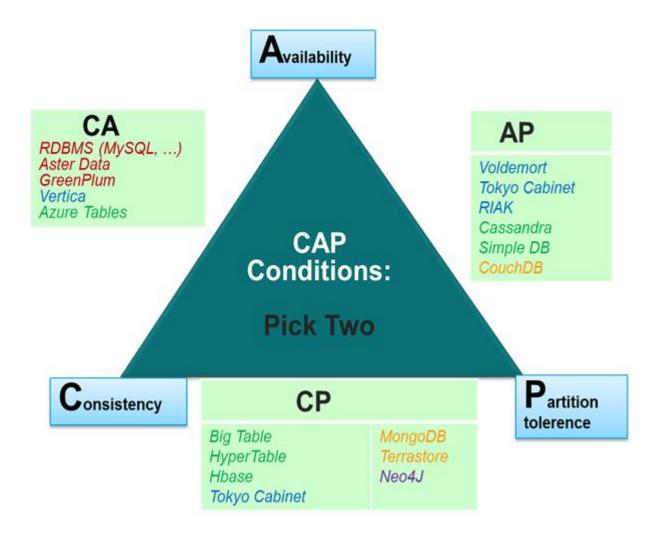
- HBASE uses column databases.
- Data is inserted in seperate column format
- Example- HBase, Cassandra

iv) Graph-based :-
- In a Graph Base NoSQL Database, you will not find the rigid format of SQL or the tables and columns representation, a flexible graphical representation is instead used which is perfect to address scalability concerns.
- Graph structures are used with edges,nodes and properties which provides index-free adjacency.
- A network database that uses edges and nodes to represent and store data.
- Data can be easily transformed from one model to the other using a Graph Base NoSQL database.
- Example- Neo4J

## ● CAP Theorem

- It stands for Consistency Availability Partition Tolerance
- Consistency-This means that the data in the database remains consistent after the
- execution of an operation. For example after an update operation, all clients see the
- same data.
- Availability-This means that the system is always on (service guarantee availability),
- no downtime.
- Partition Tolerance-This means that the system continues to function even if the
- communication among the servers is unreliable, i.e. the servers may be partitioned.

**A**vailability

**CA**
RDBMS (MySQL, ...)
Aster Data
GreenPlum
Vertica
Azure Tables

**AP**
Voldemort
Tokyo Cabinet
RIAK
Cassandra
Simple DB
CouchDB

CAP Conditions:

Pick Two

**C**onsistency

**CP**
Big Table
HyperTable
Hbase
Tokyo Cabinet

MongoDB
Terrastore
Neo4J

**P**artition tolerence

● HBase Architecture

HBASE is a distributed column oriented database built on top of hadoop to provide
real time access to Big Data.
HBase is composed of three types of servers in a master slave type of architecture.
Region servers serve data for reads and writes.
HBase Master process handles the Region assignment, DDL (create, delete tables) operations
Zookeeper maintains a live cluster state.
The Hadoop Data Node stores the data that the Region Server is managing.
All HBase data is stored in HDFS files.
The Name Node maintains metadata information for all the physical data blocks that comprise the files.

● HBase vs RDBMS

| HBase | RDBMS |
|---|---|
| HBase is a distributed, column-oriented data Storage system | RDBMS is row-oriented databases |
| HBase table have Flexible schema | RDBMS tables have fixed-schema |
| Transaction are done in Single row only | Transactions are done in multiple row |
| Max data size 1PB | Max data size TBs |
| HBase uses Java client API and Jruby | RDBMS uses SQL (Structured query Language) to query the data |
| Read/write throughput limits -> Millions of queries/second | Read/write throughput limits -> 1000s queries/second |

## Task2

Execute blog present in below link

https://acadgild.com/blog/importtsv-data-from-hdfs-into-hbase/

## Step1:

Inside Hbase shell give the following command to create table along with 2 column family.

**Create 'bulktable', 'cf1', 'cf2'**

Step2:

Come out of HBase shell to the terminal and also make a directory for Hbase in the local drive; So since you have your own path you can use it.

**mkdir hbase**

**cd hbase**

```
[acadgild@localhost ~]$ mkdir hbase
[acadgild@localhost ~]$ cd hbase
[acadgild@localhost hbase]$
```

Step3:

Create a file inside the HBase directory named bulk_data.tsv with tab

separated data inside using below command in terminal.

**Cat&gt;bulk_data.tsv**

**1 Amit 4**

**2 Girija 3**

**3 Jatin 5**

**4 Swati 3**

```
acadgild@localhost:~/hbase
[acadgild@localhost ~]$ mkdir hbase
[acadgild@localhost ~]$ cd hbase
[acadgild@localhost hbase]$ cat>bulk_data.tsv
1 Amit 4
2 Girija 3
3 Jatin 5
4 Swati 3
^C
You have new mail in /var/spool/mail/acadgild
[acadgild@localhost hbase]$
```

Step4:

Our data should be present in HDFS while performing the import task to Hbase.

In real time projects, the data will already be present inside HDFS.

Here for our learning purpose, we copy the data inside HDFS using below

commands in terminal.

Commands:

**hadoop fs -mkdir /hbase**

**hadoop fs -put bulk_data.tsv /hbase/**

**hadoop fs -cat /hbase/bulk_data.tsv**

Step5:

After the data is present now in HDFS.In terminal, we give the following

command along with arguments&lt;tablename&gt; and &lt;path of data in HDFS&gt;

Command:

**hbase org.apache.hadoop.hbase.mapreduce.ImportTsv**

**–Dimporttsv.columns=HBASE_ROW_KEY,cf1:name,cf2:exp**

**bulktable /hbase/bulk_data.tsv**

Observe that the map is done 100% although we get an error afterward.

For now, ignore the error message due to our task is to map data in HBase table.

Now,also let us check whether we actually got the data inside HBase by using the below command.

**Scan 'bulktable'**

```
acadgild@localhost:~/install/hbase/hbase-1.2.6/bin                    —    □    ×

hbase(main):015:0> scan 'bulktable'
ROW                 COLUMN+CELL
 1                  column=cf1:name, timestamp=1519742797906, value=Amit
 1                  column=cf2:exp, timestamp=1519743112316, value=4
 2                  column=cf1:name, timestamp=1519743112366, value=girja
 2                  column=cf2:exp, timestamp=1519743112421, value=3
 3                  column=cf1:name, timestamp=1519743112484, value=jatin
 3                  column=cf2:exp, timestamp=1519743112530, value=5
 4                  column=cf1:name, timestamp=1519743112584, value=swati
 4                  column=cf2:exp, timestamp=1519743117188, value=3
4 row(s) in 0.0640 seconds

hbase(main):016:0>
```