UT-1

**Q.2)** - Hadoop File system was developed using distributed File system. It is run on commodity hardware. Unlike other distributed systems. HDFS is highly faulttolerant and designed using low-cost hardware.

Features of HDFS:
1) It is suitable for the distributed storage and processing.
2) Hot Hadoop provides a command interface to interact with HDFS.
3) The built-in servers of namenode and datanode helps user to easily check status of clusters.
4) streaming access to file system data.

HDFS Architecture:



@shavan

## NameNode:

The Namenode is a commodity hardware that contains the GNU/Linux operating system and the namenode software. It is a software that can be run on commodity hardware. The system having the namenode acts as the master server and it does following tasks:-

1) managing the file system namespace.
2) Regulates client's access to file.
3) It also executes file system operations such as renaming, closing and opening files.

## Data node:

The Datanode is a ~~community~~ commodity hardware using GNU/Linux operating system. and datanode software. Every node in a cluster, there will be datanode. These nodes manage the data storage of their system.

1) Datanodes perform read-write operations on file system as per client request.
2) They also perform operations such as block creation, deletion and replication according to instructions of namenode.

## Block:
Generally the user data is stored in files of HDFS. The file in a file system will be divided into one or more segments and / or stored in individual data nodes. These file segments are called as blocks.

UT-1

Q.3) Hadoop Ecosystem is a platform or a suite which provides various services to solve the big data problems. It includes Apache projects and various commercial tools and solutions. There are four major elements in Hadoop. i.e. Hadoop, MapReduce, YARN and Hadoop Common.

Hadoop Ecosystem

| Oozie workflow monitoring | Chukwa monitoring | Flume monitoring | Zookeeper Management |
|---|---|---|---|

Data management

| Hive (SQL) | Pig (Dataflow) | Mahout (ML) | (Avro) (RPC) | Sqoop (RDBMS) Connection |
|---|---|---|---|---|

Data Access

| Map Reduce (Cluster management) | YARN (Cluster and Resources management) |
|---|---|

Data processing

| HDFS (File system) | HBASE (Column DB storage) |
|---|---|

Data Storage

1) HDFS: HDFS is the primary or major component of hadoop ecosystem and is responsible for storing large data sets of structured or unstructured data across various nodes and thereby maintaining the metadata in form of log files.
HDFS consists of Name Node and Data node.

2) YARN
: Yet Another Resource Negotiator as the name replies, YARN is one who helps to manage the resources across the clusters. In short, it performs scheduling and resource allocation.
It consists of :
1) Resource manager
2) Nodes manager
3) Application manager.

3) MapReduce:
By making the use of distributed and parallel algorithms. MapReduce makes it possible to carry over the processing logic and helps to write applications which transform big data sets into manageble one.

4) PIG: Pig was basically developed by Yahoo which works on pig latin language, which is Query language similar to SQL.
It is platform for structuring the data flow , processing and analyzing huge data sets.

UT-1

4) HIVE! With the help of SQL methodology and interface, HIVE perform reading and writing large datasets.

5) mahout: mahout allows machine to learning learnability to a system or application. machine learning as name suggests helps the system to develop itself based on some patterns, user/enviromental interaction algorithms.