

Assignment No :- 2

Page No. _____
Date _____

Q1]

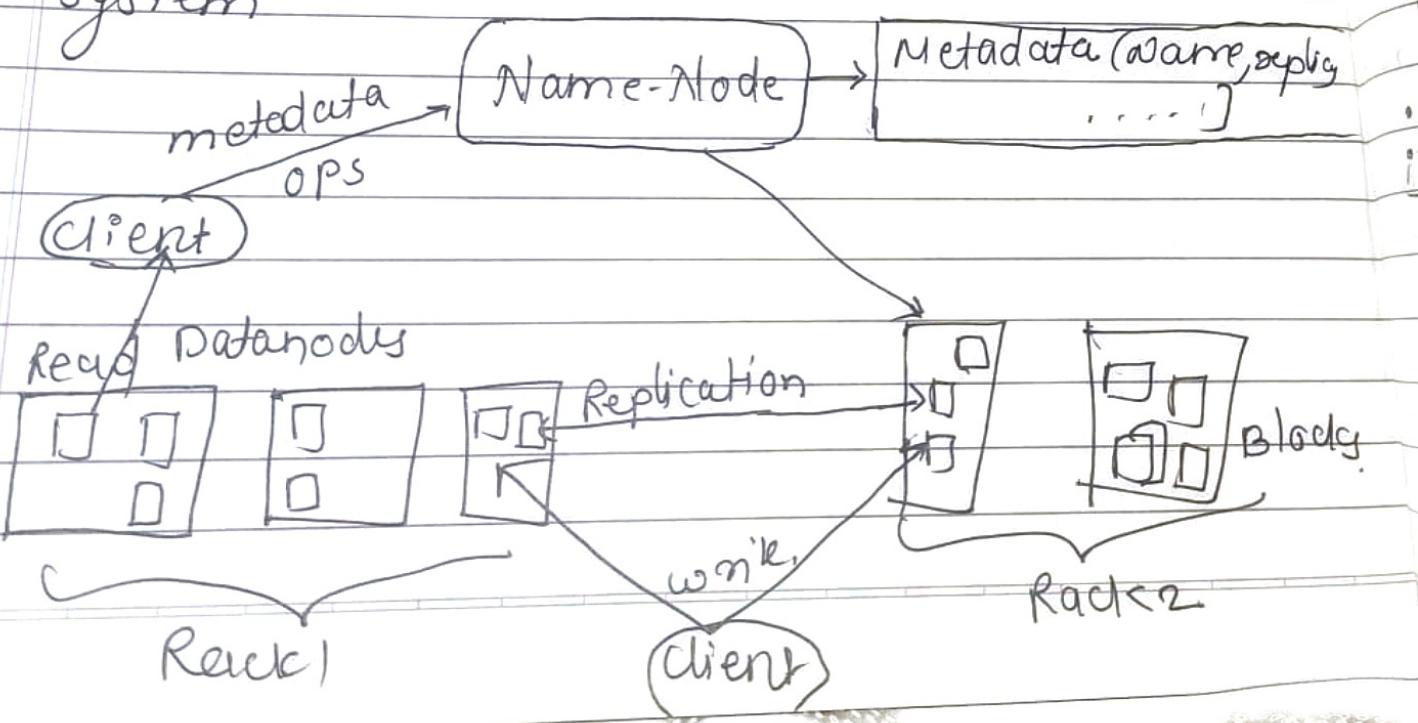
→ Write a note on cloud file system with Architecture.
 Cloud storage is the abstraction, pooling & sharing of store resources through the internet. File storing is the dominant technology used on NAS system & is responsible for organizing data & representing it to users. It is hierarchical structure allows us to navigation data from top to bottom easily, but increase processing time.

Cloud file storage is a storage service that is delivered over the internet billed on a pay-per-use basis & has an architecture based on common file level protocols such as Server Message Block (SMB), common Internet File System (CIFS) & Network File System (NFS).

i) Types of cloud file system.

i) GFS:- Amazon Web Services (AWS) use GFS cloud file system (GFS). GFS runs over Amazon's EC2, SimpleDB web service & S3.

ii) Hadoop File System:- is a distributed file system inspired by GFS that organizes files & store their data on a distributed computing system.



iii) Kosmos file system (KFS) :- is a open source project written in C++ by search startup kosmix. Three components are :- i) One or more chunk servers that store the data on their own hard disks, b) A metaserver that keeps an eye on the chunks server p c) An app in that quickly get rids of a single large file.

Q2] Discuss Hbase Data Model?

- It is a database that is open source platform & it is the implementation of storage architecture by Google's Big Table. HBase databases are column-oriented thus it makes it unique from other databases. One of the unique qualities of Hbase is it doesn't care about data types because we can store different data types or data for the same column in different rows.
- It contains different sets of tables that maintain the data in key-value format. Hbase is best suitable for sparse data sets which are very common in the case of big data. It can be used to manage structure of semi-structured data & it has many built-in features such as :- scalability, versioning, compression, garbage collection.
- Two types of data storage medium :-
- Row-oriented :- data is stored & retrieved one row at a time. This could lead to several problems. Suppose we want only some part of the data from the row but according to this approach you have to retrieve the complete row even if you don't need it.

- Column-Oriented :- data is stored & retrieved based on the columns. Thus the problem which were facing in the case of the row-oriented approach has been solved because in the column-oriented approach we can fill out the data which is required to us from the whole set of data with the help of corresponding columns. OA the read & write op's are slower than other but it can be efficient while performing operations on the entire database hence it permits very high compression rate.

column families

Row Key emp id	personal data			professional data	
	name	city	designation	salary	
1	A	AAA	managers	₹150,000	84
2	B	BBB	sr. engg	₹70,000	→
3	C	CCC	Jr. engg	₹50,000	

) Q3) Draw & Explain cloud file system GFS ?

→ A file system in cloud is a hierarchical storage system that provides shared access file data. Users can create, delete, modify read & write file & can organize them logically in directory trees for intuitive access

- cloud file storage is most appropriate for unstructured data or semi-structured data, such as documents, spreadsheets, presentation & other file-base data.
- CFS is a method for storing data in the cloud that provides servers & app in clouds the data through shared file system. This compatibility

- makes cloud file storage ideal for workloads that rely on shared file systems & provide simple integration without code changes.
- **Cloud File Systems**: It gives high redundant elastic mountable, cost effective & standard based file system. A fully featured scalable & stable cloud file system is provided by GFS.
 - GFS runs over Amazon's S3, EC2 & SimpleDB web services.
 - When using GFS, user can have complete control of the data & can be accessed as a standard n/w disk drive. It is highly secure. It can be mounted on server, client or access file via web page.

Q4) Explain Working of cloud - Data storage?

- Data store is a connection to a store of data, whether the data is stored in a database or in one or more files. The data store may be used as the source of data for a process.
- Datastore is a repository for storing, managing & distributing data sets on an enterprise level.
- Distributed Data store: Computer n/w where information is stored on more than one node by means of data replication it is termed as DDS. It is used to refer to distributed db where user store info on no. of nodes or a computer n/w in which store information on no. of peer n/w nodes.

- DDS is one in which file, scripts & images are stored in more than one server or volumes rather than a single server as in traditional system
- Eg:- Google's Big Table, Amazon's Dynamo
- These types of data store are non-relational databases that searches data quickly over a large multiple nodes.
- Data store Types
 - i] Big table :- is a distributed storage system that is used for managing & storing structured data at Google. Big table is designed to reliably scale to petabytes of data & thousands of machine. Big table has multiple goals like applicability, high availability, scalability, high performance. Big table is build on Google file system for storing the data for scheduling large scale data processing. It stored data in form of rows, columns & timestamp that means it maps with arbitrary string value like row key & column key as well as timestamp.
 - ii] Dynamo :- is a property key value structured storage system. It can act as database & also distributed hash table. Dynamo dynamically partitions a set of keys over a set of storage nodes. It is most powerful relational database available in world. Dynamo does not support replication. It is used to manage the state of service that have very high reliability requirements & need tight control over the trade b/w availability consistency, cost effective & performance.

Q5] How does cloud file system are different from Normal file system.

- cloud file storage Traditional File storage
- i) stored on remote server typically stored on local hard drives
- ii) Accessible from anywhere with internet limited to local network or device
- iii) Highly scalable, can easily scale up or down based on demand Limited scalability, capacity is tied to hardware
- iv) Build-in redundancy through data replication & backups Relies on manual backup or RAID for redundancy
- v) pay-as-you-go mode based on storage usage Upfront costs for hardware & maintenance
- vi) Managed & maintained by cloud service providers Requires manual management & updates
- vii) Build-in collaboration features such as shared folders & access controls Collaboration may require complex setups
- viii) often includes robust security measures provided by the cloud provider Security features depend on local configurations

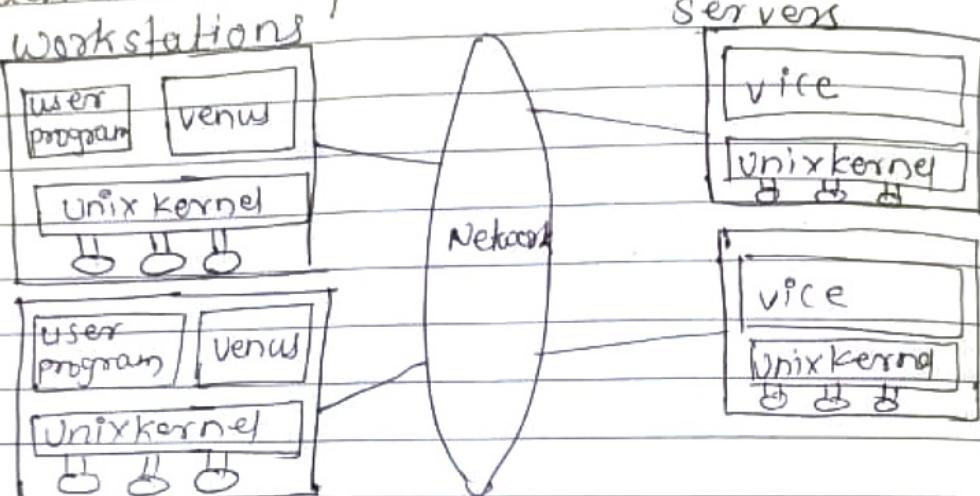
~~Q6) Explain AFS Architecture?~~

→ Andrew File System Architecture.

- i) Vice :- AFS provides a homogeneous, location-transparent file namespace to all client workstations by utilizing a group of trustworthy servers known as vice.
- ii) Venus :- caches files from vice & returns updated versions of those file to those file to the servers from which they originated.

Prog No.	
Date	

- File system architecture was largely inspired by the need for scalability. To increase the no. of clients a server can service, Venus performs as much work as possible rather than vice.



- AFS is a distributed file system. It uses the client server model, where all the files are stored on the server machines.

Q7) Difference b/w SAN & NAS.

- | | |
|---|--|
| • SAN | • NAS. |
| i) storage area n/w | • Network Attached Storage |
| ii) In SAN data is identified by disk block | • In NAS data is identified by file name as well as byte offset. |
| • In SAN file system is managed by servers | • File system is managed by Head Unit |
| • It is more costly & it is more complex | • It is less expensive. |
| • Protocol used in SAN are SCSI, SATA | • It is less complex than |
| • Suitable for that environment which has high speed traffic. | • Protocol :- File servers, CIFS (common Internet File System) |
| | • Not suitable for that environment which has high speed traffic |

- Has lower latency
- has higher latency
- It support virtualizatⁿ
- Do not support virtualizatⁿ

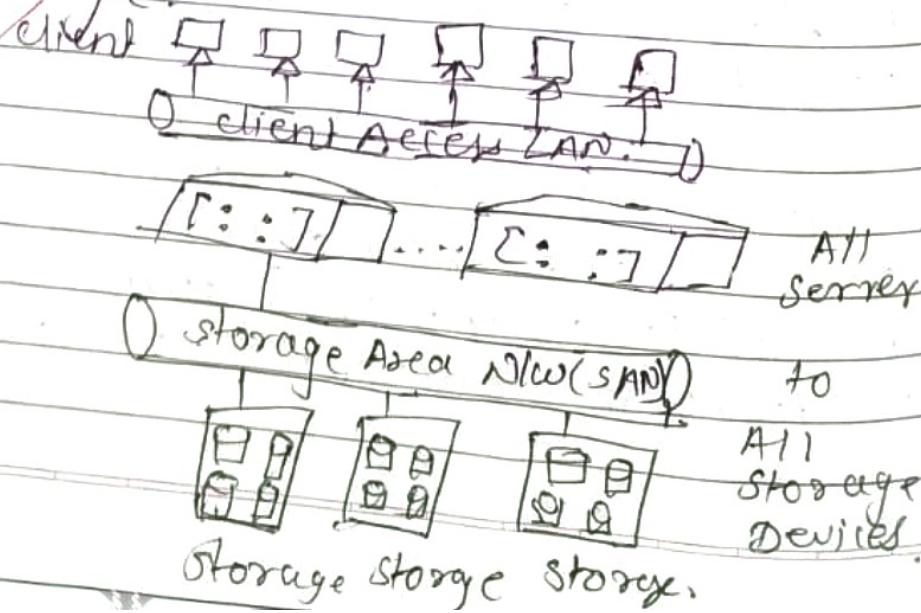
- Q8) Elaborates HDFS Architecture in details.
- Hadoop File System :- Hadoop is an open source software framework that supports data-intensive distributed appⁿ, licensed under the Apache v2 license. It provide software framework for distributed processing of large datasets in real time appⁿ.
 - Hadoop Distributed File System is a block-structured file system where each file is divided into blocks of pre-determined size. These blocks are stored across a cluster of one or several machines.
 - Blocks are the nothing but the smallest continuous location on your hard drive where data is stored. Similarly, HDFS store each file as blocks which are scattered throughout the Apache Hadoop cluster.
 - It is the primary storage system used by Hadoop application. HDFS operates as a distributed file system designed to run on commodity h/w.

- Q9] Explain the Data storage (EDS,DAS,SAN,NAS)
- i] EDS (Enterprise Data storage) :- Enterprise storage is a centralized repository for business-critical information that provides data sharing, data management & data protection across multiple computer system.
 - Enterprise storage may include a SAN, NAS,DAS.

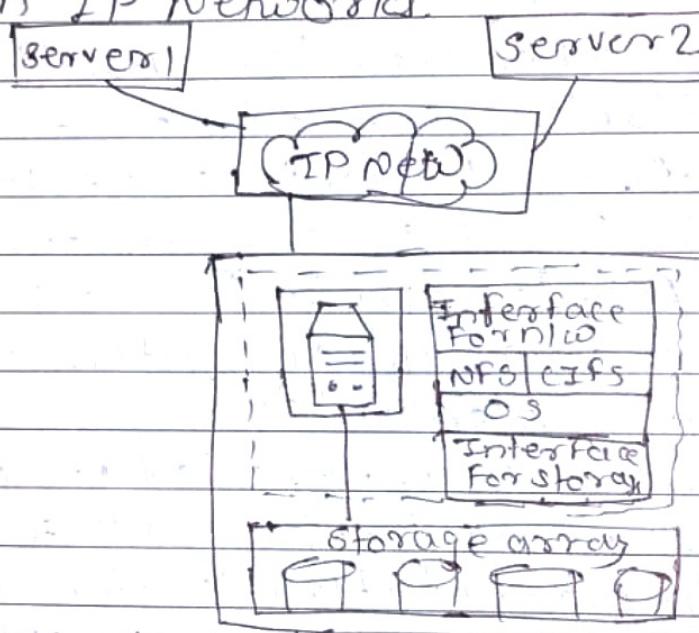
- Direct Attached Storage is hard disk drives, solid-state drives connected directly inside or outside to a single computer or server they cannot be accessed by other computer or server. DAS is not networked through Ethernet or FC switches.
- Storage device may includes one or more drives built into a server with an appropriate host bus adapter, may be configured as RAID array. Eg:- hard drives, disc drives & storage on external



- Storage Area Network :- is a dedicate high-performance n/w or subnetwork dedicated to storage that is independent of an organization's common user Network.
- It interconnects pools of disk or solid-state storage & shares it to multiple servers so each one can access data as if it was directly attached. A SAN is a dedicated n/w that provides access to consolidated, block level data storage.



- Network-attached storage :- is a file level computer data storage server connected to a computer n/w providing data access to a heterogeneous group of clients. It is specialization for serving file either by its hardware, s/w or configuration.
- NAS devices are storage arrays or gateways that support file-based storage protocols such as NFS & CIFS & are typically connected via an IP network.



- NAS device is a storage device connected to a n/w that allow storage & retrieval of data from central location for authorised n/w user & named clients.

Q10) Explain the Data Intensive Technologies for cloud computing.

→ DI system encompass terabytes to petabytes of data. Such systems require massive storage & intensive computational power in order to execute complex queries & generate timely results.

- Data Intensive computing is "a class of parallel computing applications which use a data parallel approach to processing large volumes of data"
- 1) Processing approach :- current data-intensive computing platforms use a "divide & conquer" parallel processing approach combining multiple processors & disks in large computing clusters connected using high-speed comm' n/w.
 - It allows the data to be partitioned among the available computing resources & process independently to achieve performance & scalability based on the amount of data.
- 2) System Architecture :- for data-intensive computing an array of system architectures have been implemented.
 - a) MapReduce :- is a parallel programming model proposed by Google & available open source implementation known as Hadoop. It aims at supporting distributed computation on large datasets by using a large no. of computers with scalability & fault tolerance guarantees.
 - Map & reduces are two primitives in functional programming languages, such as Lisp, Haskell, etc. A map function processes a fragment of a key-value pairs list to generate a list of intermediate key-value pairs.
 - Reduce function merges all intermediate values associated with a same key & produces a list of key-value pairs as output.

b) High Performance computing cluster (HPCC) also known as DAS (Data Analytics Supercomputer) is an open source, data-intensive computing system platform developed by LexisNexis Risk Solut. DAS is platform designed to refine, link & fuse large amount of data from disparate sources for complex analysis of queries.

- Architecturally, DAS is a HPCC based on commodity servers ~~blue~~, which can be scaled up to thousands of processors to handle any amount of data and runs on the Linux OS.
- Custom system software & middleware parts were created & layered to provide H execution environment & distributed file system support that is essential for data intensive computing on the base of Linux O.S.

QY