

100 Marks

Date: 16/4/25

* Theory Assignment - 1 *

Explain GFS along with diagram.

GFS stands for Google File System, developed by Google to manage huge amount of data to manage efficiently across thousands of inexpensive servers.

It is scalable, fault-tolerant, & distributed file system designed specifically for large distributed data-intensive applications.

Features of GFS -

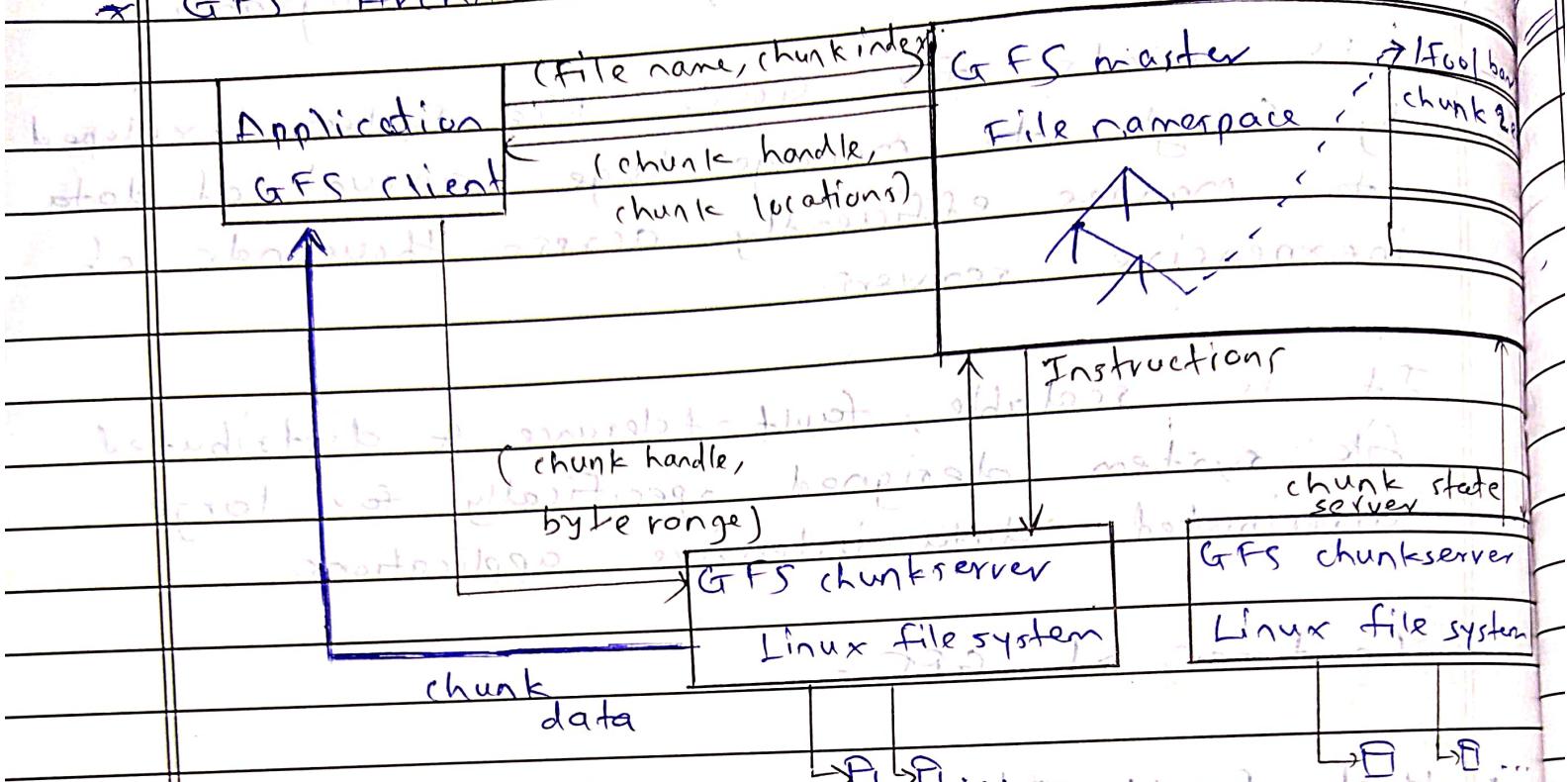
- High Fault Tolerance - continues working even if hardware fails frequently.
- High Throughput - optimized for large data reads/writes rather than low-latency small file access.
- chunk-based storage - files are divided into fixed-size chunks (typically 64MB), each identified by unique chunk handle.
- master-slave Architecture - a single master server manages metadata, & chunk servers store actual data.
- Replication - Each chunk is replicated (default 3 copies) across different chunk servers.



Date

File location

* GFS Architecture - (Multi-Host)



* Components

① **Client** - sender for reading/writing requests.
communicates with master to take chunk

Server

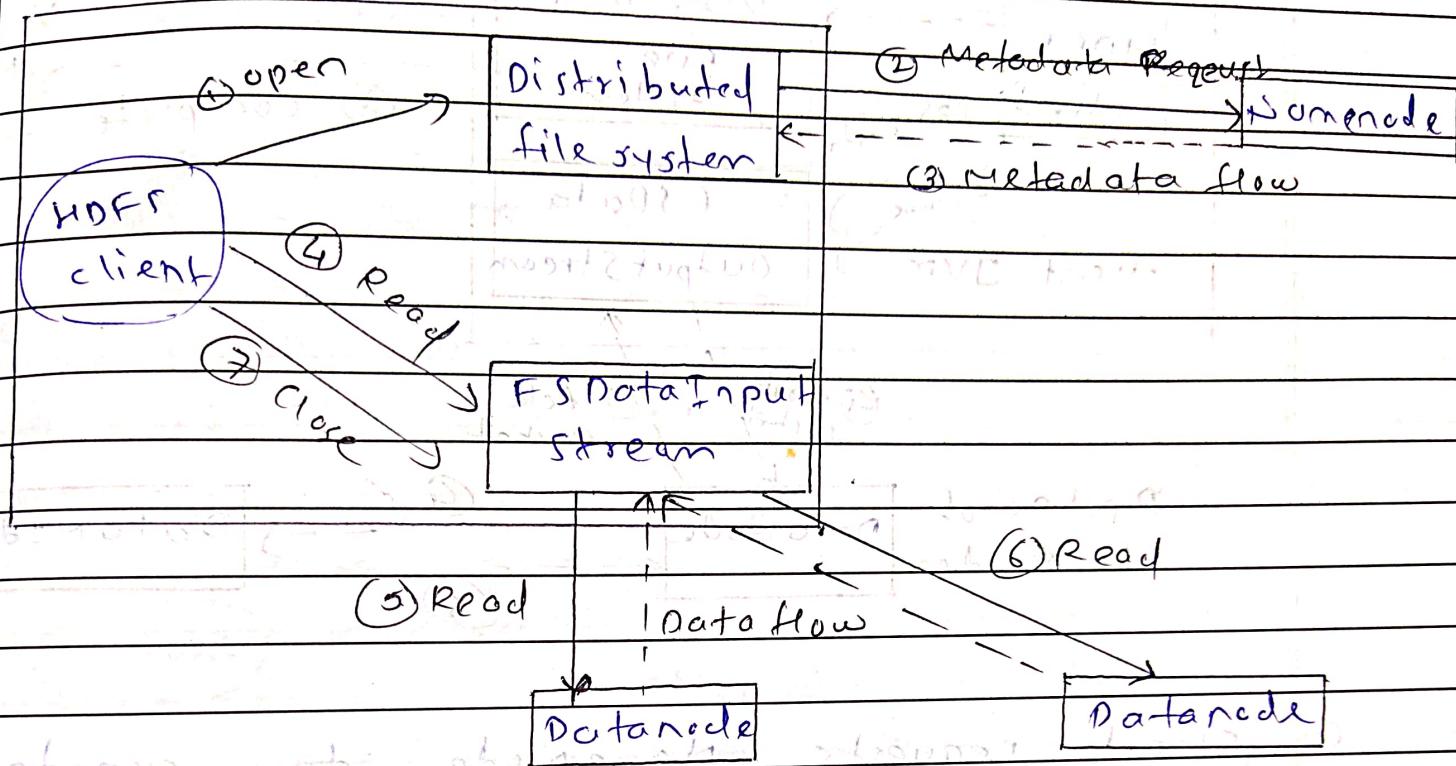
② **Master Server** - maintains metadata (name, size, mapping of files to chunks, locations, chunk size, etc.)

③ **Chunkserver** - store actual file data in chunks. Each chunk is stored on multiple servers. It also contains a location table to locate data.

File location in cloud and on authorized servers (using S3).

Explain in detail Read & write operations in HDFS.

→ RRead



client contacts NameNode to request location of file blocks (metadata only).

NameNode responds with list of DataNodes holding replicas of each block of file.

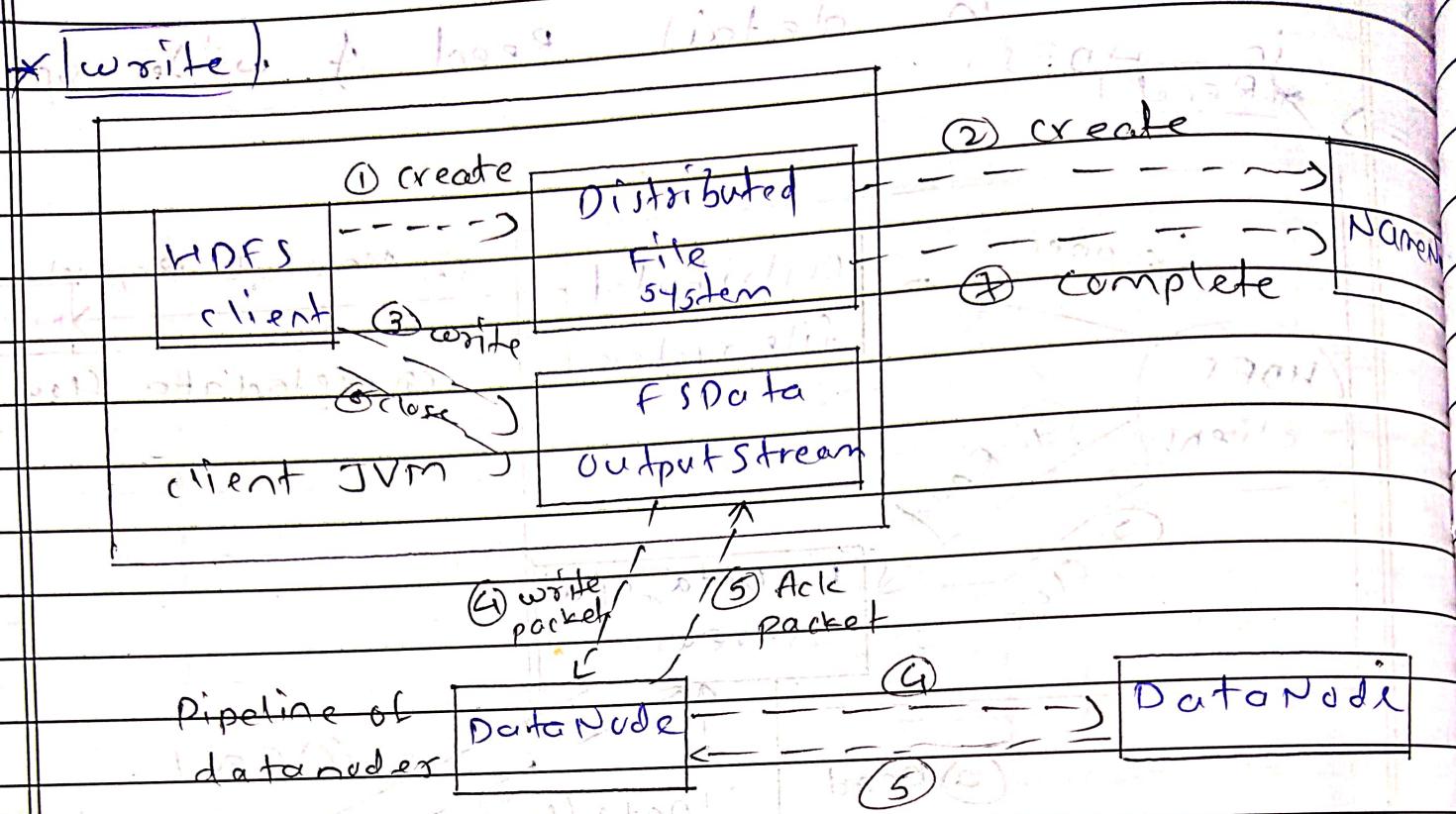
client directly contacts nearest DataNode

DataNode reads block by block from DataNodes & assembled into complete file

on client side.

If DataNode - unavailable - client

automatically switches to another replica provided by NameNode.



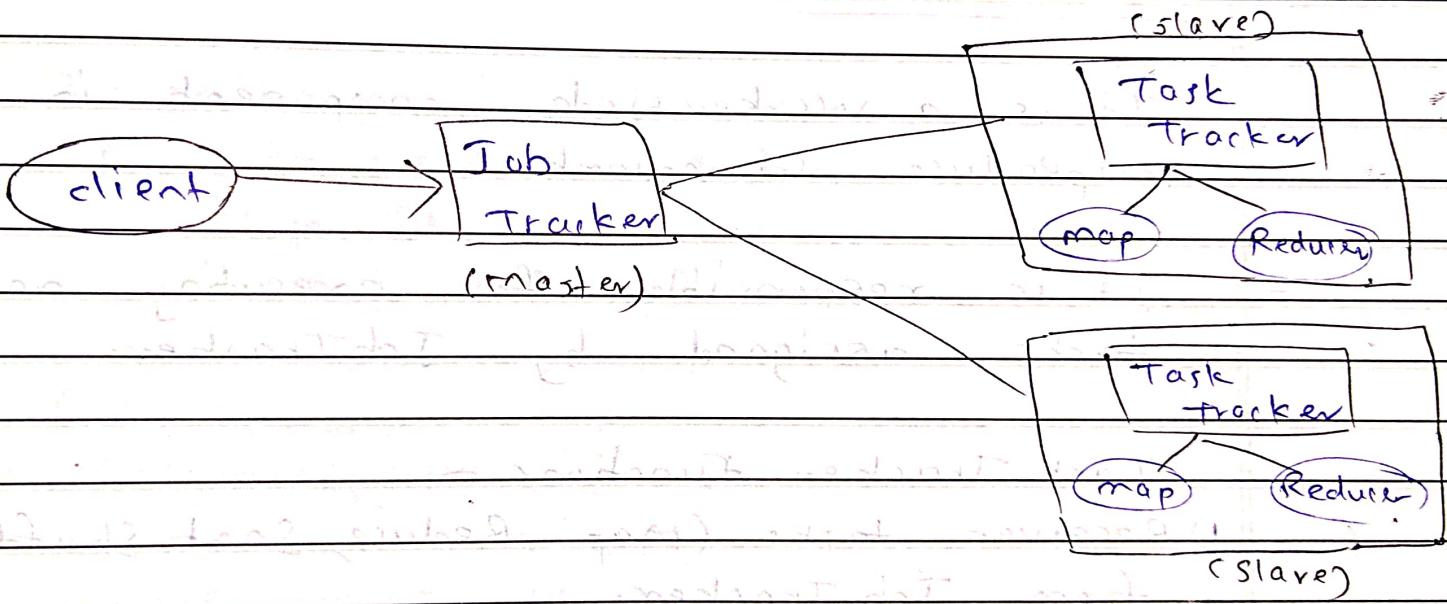
- ① Client requests NameNode to create a new file entry.
- ② NameNode checks if file already exists.
- ③ Client divides file into blocks (default: 128 MB).
- ④ NameNode provides a list of DataNodes for each block where replication will be stored (typically 3 replicas).
- ⑤ Client writes first block to first DataNode (replication begins).
- ⑥ DataNode forwards block to second DataNode.
- ⑦ second DataNode forwards it to 3rd DataNode (pipelining).

- (6) After successful writing & replication of each block, acknowledgment is sent back through pipeline to client.
- (7) Once all blocks are written & acknowledged, NameNode updates its metadata.

Note - Write once, Read many times.

g) Explain with Diagram -

① Job Tracker ② Task Tracker



①

Job Tracker

- master node component in MapReduce framework.
- It is responsible for managing resources & coordinating execution of MapReduce jobs across cluster.

- b) (3) NoSQL databases.
- ⇒ NoSQL - "Not Only SQL"
- It provides a way to store & retrieve data that is different from traditional RDBMS.
 - They are designed for large volumes of data, unstructured data, scalability.

- Characteristics -
- Schema-less
 - High-Performance
 - Variety of Data Models.
 - Horizontal Scalability
 - Flexible Data types

- Types -
- i) Document-based
- Stores data in JSON-like documents.
 - Best for nested data.
 - e.g. - MongoDB, couchDB
 - e.g. - document data stored -
- ```

{
 "id": 101,
 "name": "Shubham",
 "courses": ["ML", "DSA"]
}

```
- No strict schema.

## 2) key-value store

- Data stored as key-value pairs.
- Simple & fast for basic lookup op.
- e.g. - Redis, DynamoDB
- e.g. - Data is stored like this -

| key        | value                     |
|------------|---------------------------|
| user:1     | {"name": "Om", "age": 25} |
| session:23 | {"token": "abcd"}         |

- ideal for caching, session management

## 3) column-based

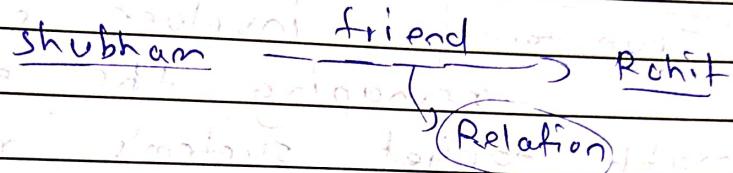
- Data stored as key-value in columns rather than rows.
- Efficient for large analytical queries
- columns grouped into families
- e.g. - Apache Cassandra, HBase
- e.g. Data is stored like -

| Student | Name | Course |
|---------|------|--------|
| 1       | Shub | ML     |
| 2       | Jay  | DS     |
| 3       | Om   | DS     |

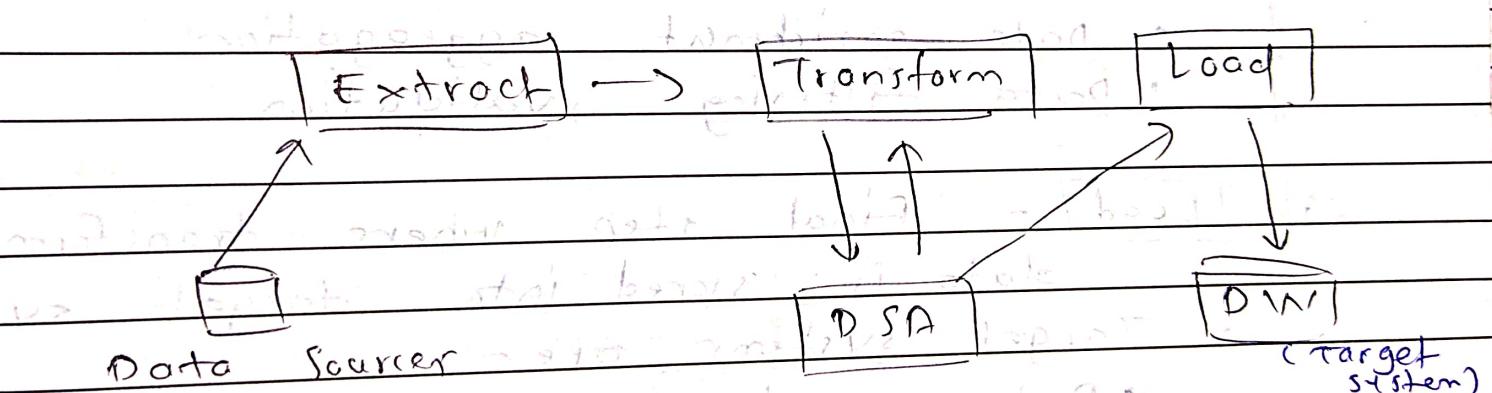
- Good for fast retrieval applications & massive data sets.

#### 4) Graph-based

- Data stored as nodes & edges
- Designed to store relationships between entities
- Neo 4J, Amazon Neptune
- e.g. -



#### 4. Explain ETL using details



- ETL (Extract, Transform, Load) is critical step in building a consistent data pipeline for analytical & BI.

#### ① Extract

- Process of gathering data from various source systems
- Sources can be -
  - Databases (SQL/NoSQL)

- APIs
  - Spreadsheets
  - Cloud Storage
  - Logs & sensor data
- Goal is to collect raw data without losing information.

- (2) Transform - involves cleaning, modifying, reshaping the extracted data to match target system's requirement
- tasks performed -
    - Data Cleaning (removing errors, duplicates)
    - Data Standardization
    - Data enrichment, aggregation
    - Data filtering, validation

- (3) Load - Final step where transformed data is stored into target system
- Target systems are -
    - Data Warehouses (Snowflake, Amazon Redshift)
    - Databases (PostgreSQL, MySQL)
    - Data Lakes (AWS S3, Azure Data Lake)
  - Some ETL Tools examples -
    - ① Apache Nifi, Talend Open Studio
    - ② AWS Glue, Google Cloud DataFlow
    - ③ IBM DataStage, Informatica