**Strictly for class use. Do NOT share outside the class now or in future**

9/14/24

1

# HDFS

All nodes in the cluster share the same namespace

Write-once-read-many - Clients can only append to existing files
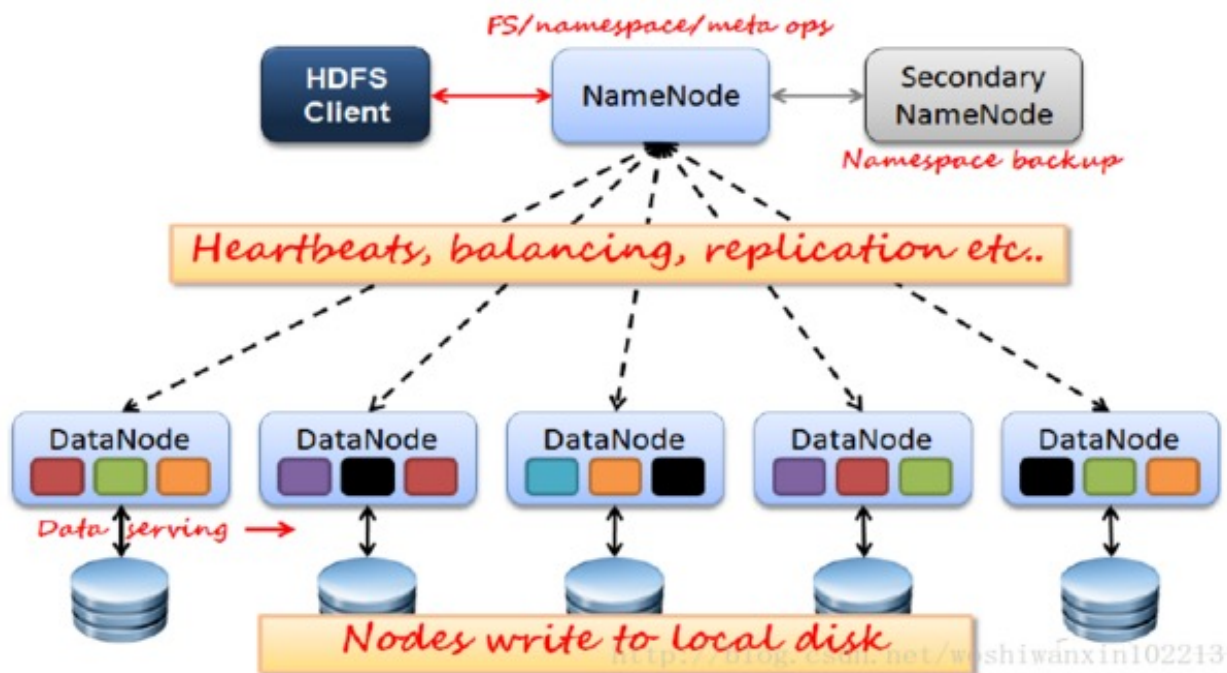
Clients can find location of blocks and directly accesses data from the DataNodes

Entire metadata is in the main memory of NameNode and is never swapped to disk

Metadata comprises of filenames, blocks for each file, replication factor, DataNodes for each block, and File attributes

DataNodes store the blocks in their local OS file system, can directly serve the data to the clients, and can send data to other DataNodes directly

Rack awareness for fault tolerance

**Strictly for class use. Do NOT share outside the class now or in future**

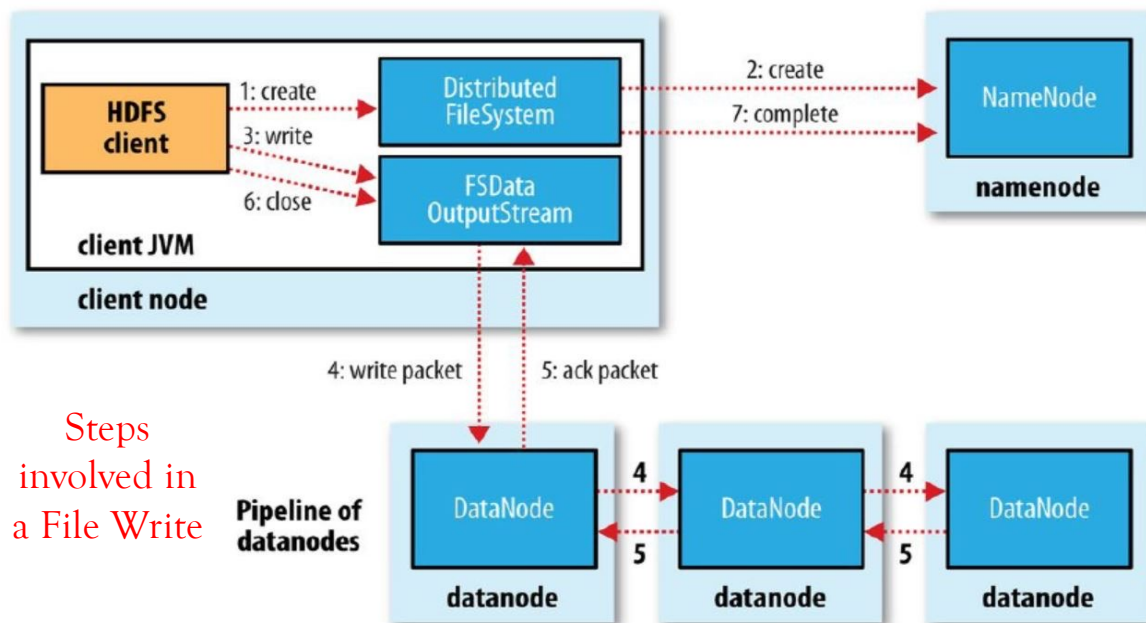9/14/24

# SOME HDFS COMMANDS

hdfs dfs -ls /path/to/directory

hdfs dfs -mkdir /path/to/new_directory

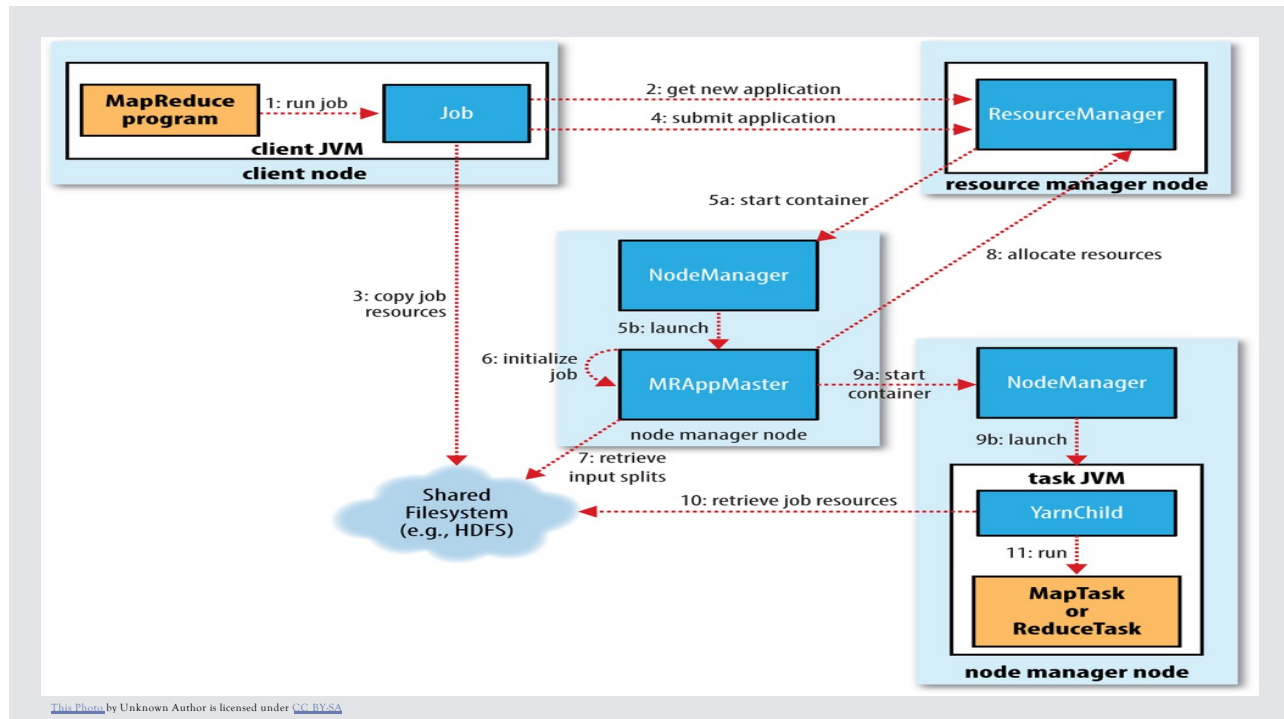hdfs dfs -copyFromLocal /local/path/to/file /hdfs/path/to/destination

hdfs dfs -copyToLocal /hdfs/path/to/file /local/path/to/destination

hdfs dfs -cp /hdfs/src/path /hdfs/destination/path

hdfs dfs -chown <user>:<group> /hdfs/path/to/file_or_directory

Steps involved in a File Write

**Strictly for class use. Do NOT share outside the class now or in future**

9/14/24

This Photo by Unknown Author is licensed under CC BY-SA
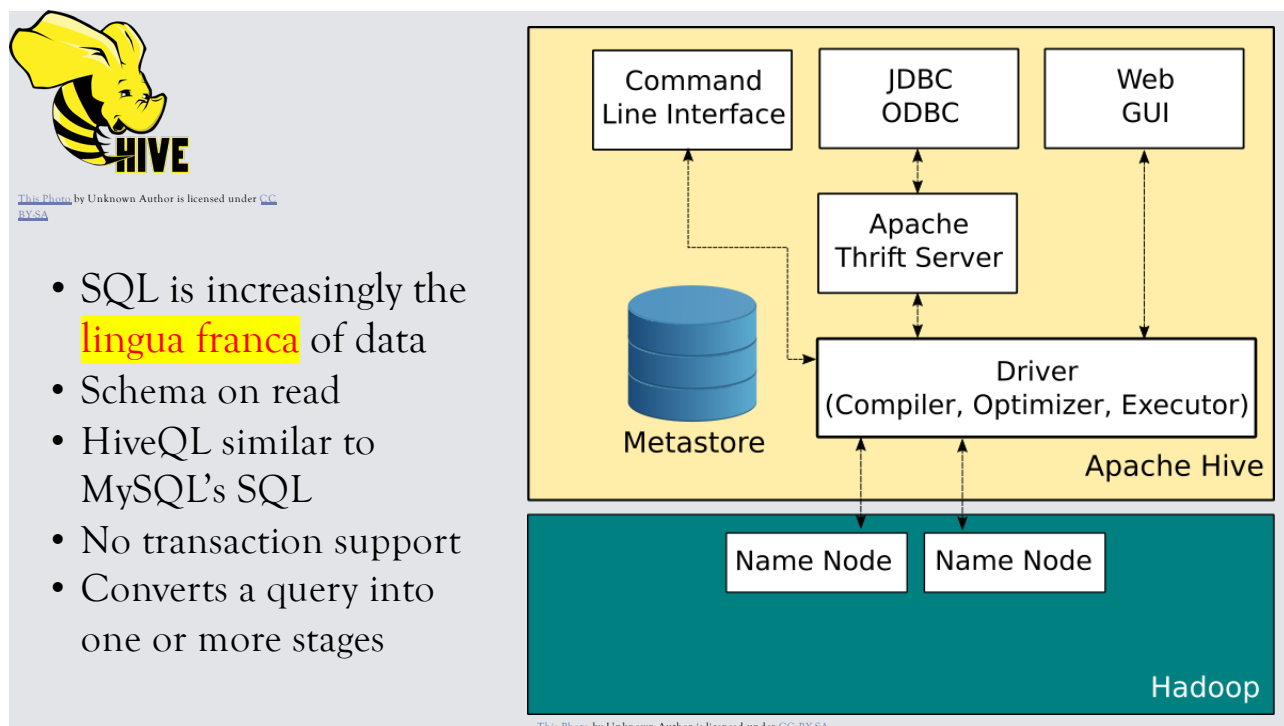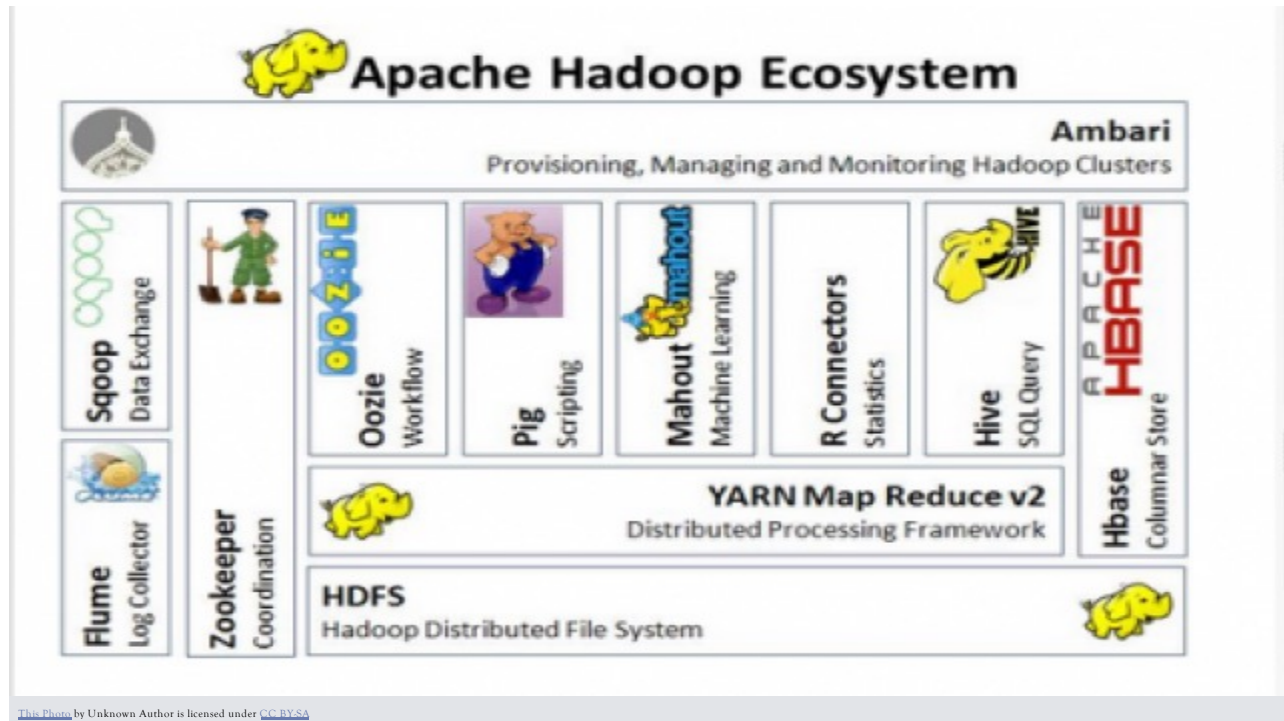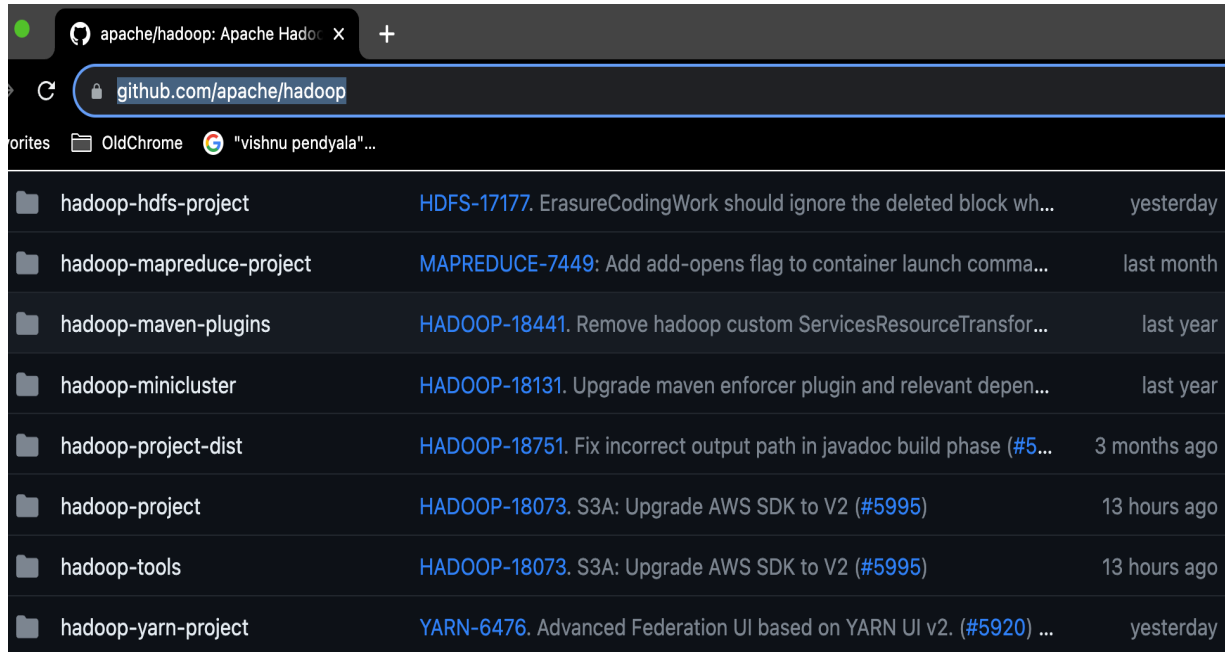
# Running Hadoop

hadoop jar $HADOOP_HOME/share/hadoop/tools/lib/hadoop-streaming-*.jar \
  -mapper <MapperScript>  -reducer <ReducerScript> \
  -input <InputDirectory>    -output <OutputDirectory>

• Hadoop's default mapper is the Identity Mapper

**Strictly for class use. Do NOT share outside the class now or in future**

9/14/24

Apache Hadoop Ecosystem

- SQL is increasingly the lingua franca of data
- Schema on read
- HiveQL similar to MySQL's SQL
- No transaction support
- Converts a query into one or more stages

# What are the issues with Hadoop?

Users must think in "map" and "reduce" terms and provide functions that align with the paradigm

Persisting intermediate results helps with fault tolerance, but HDD access between iterations slows it down substantially

Waiting for the map and reduce phases can also slow it down

**Strictly for class use. Do NOT share outside the class now or in future**

9/14/24

## HDD access hasn't kept up with time

| Storage/Access Type | Approximate Speed |
|---|---|
| RAM (Random Access Memory) | 20 GBps - 100 GBps (or higher) |
| SSD (Solid-State Drive) | 100 MBps - 6 GBps (or higher) |
| HDD (Hard Disk Drive) | 50 MBps - 200 MBps (or higher) |
| Broadband Internet (Typical) | 1 Mbps - 1 Gbps (or higher) |
| Low-Latency Internet (e.g., Fiber) | 1 Gbps - 10 Gbps (or higher) |
| Internet Latency (Round-trip) | 10 ms - 100 ms (or higher) |

SUMMARY?