

Single point of failure (SPOF)

Reasons

- Hardware failure
- Overcapacity caused by metadata due to small data problem.

Cloudera (CDH - cloudera distributed hadoop)

(Horton data pattern) HDP

In 2019, Cloudera required +

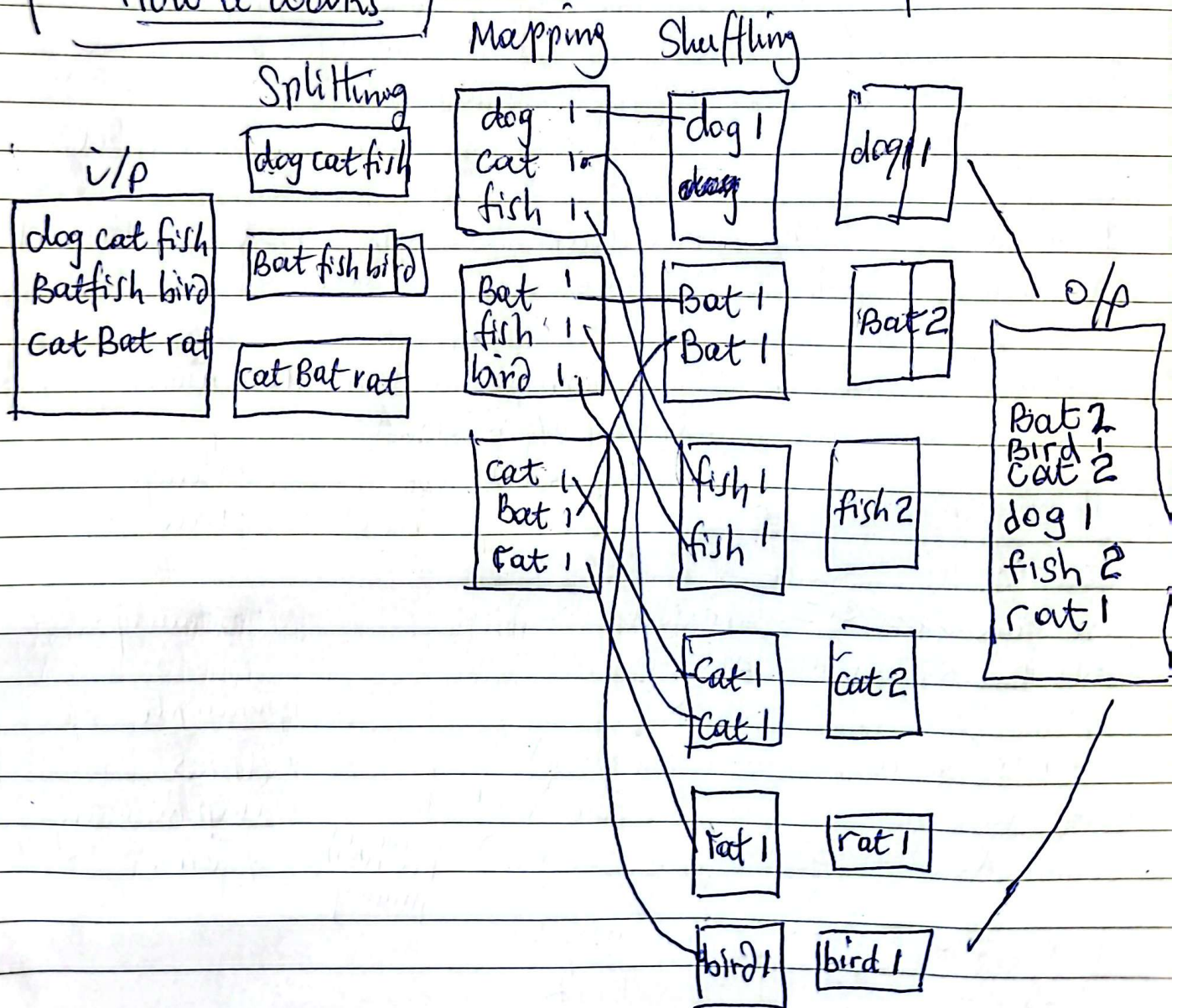
namenode

Secondary namenode }

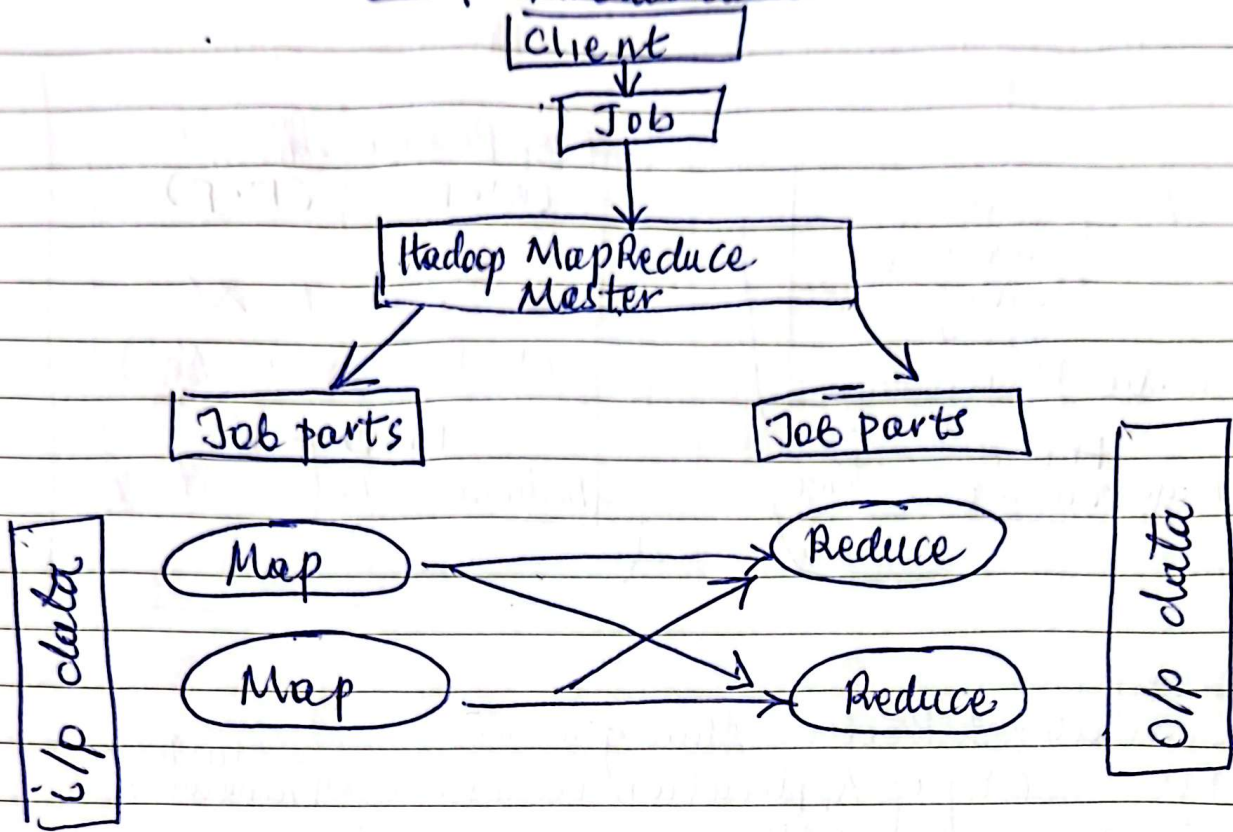
Datanode slave / compute / worker node ✓

Hadoop Distributed File System HDFS

How it works



Map Architecture



YARN

Yarn means "Yet Another Resource Negotiator".

It is an important aspect of Hadoop ecosystem that is responsible for decoupling / separating the functionalities of "managing Resources" such as CPU, memories & Storage distribution from data processing as separate components. Consequently, Global Resource Mgr can now support diverse Applications

BEFORE 2012

AFTER 2012

such as Spark, MapReduce, Tez, Storm etc

Before 2012



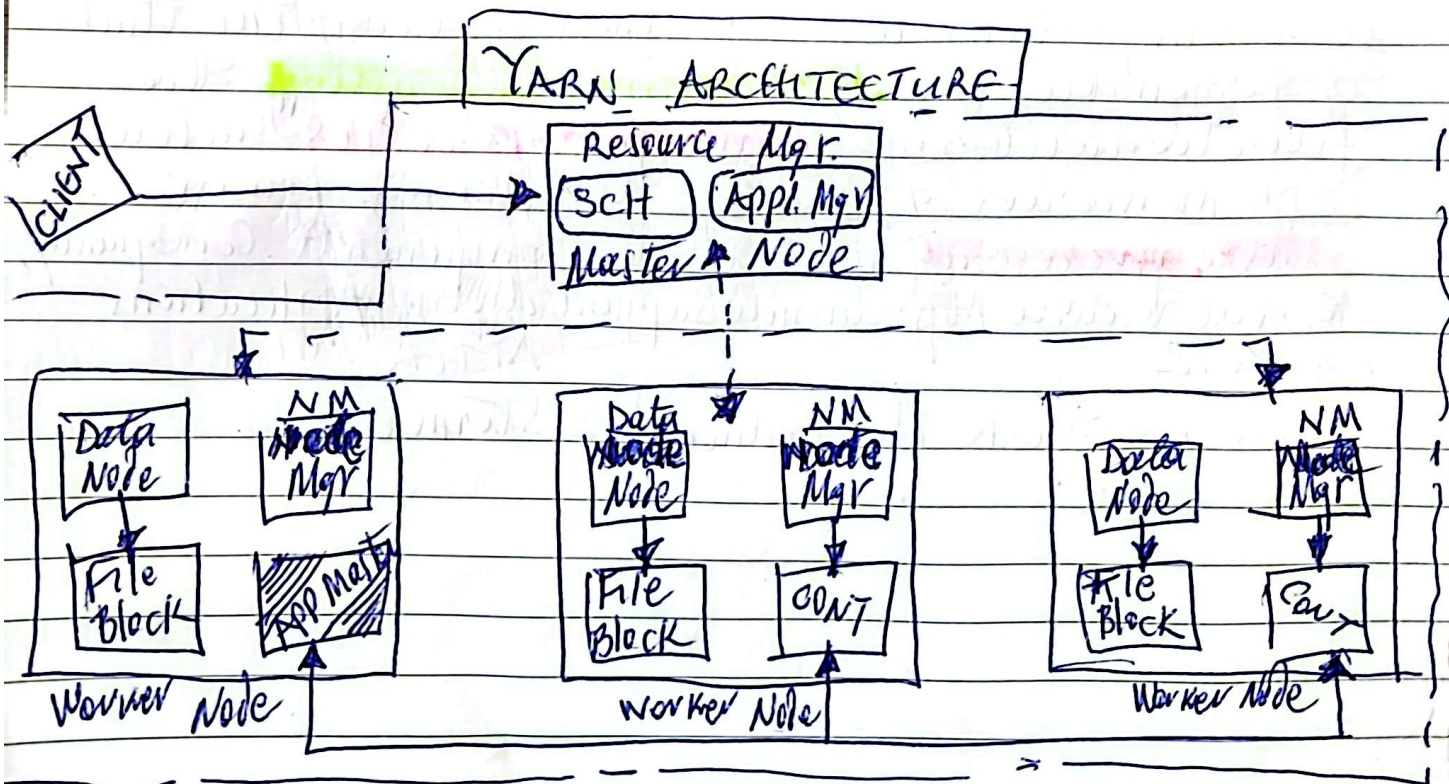
After 2012.



Benefit.

- 1) Centralized Resource Management
- 2) Flexibility of Application across the Framework.
- 3) Cost Effective: All in one Framework.
- 4) Better data utilization amongst Applications.
- 5) Less Data Movement within same framework/cluster Applications.

YARN ARCHITECTURE



Components of YARN

Client, Resource Mgr, Scheduler, Application Mgr,
Data Node, File Block, Node Mgr, Container, Application Master,
(- Master Worker & - Node Worker.)

□ Master Node & Worker Node



Apache Hive is a Data Warehousing Infrastructure built on top of Hadoop to provide: (HQL)

- * SQL-like Query language called Hive Query lang.
- * Structured data analysis on large Datasets stored in HDFS.

- * Batch processing rather than Real-time transactions
This batch processing is also referred to as Write-Once-Read-Many (WORM).

~ Why use Hive.

- Familiarity to SQL-like Interface (HQL) makes it easy
- Scalability & fault tolerance by running Queries using MapReduce or Tez.
- Supports multiple formats e.g. text, ORC, parquet, Avro
- Efficient query ^{Execution} using partition & Indexing.

Main components of Hive.

1) Hive Client: Interfaces like CLI, JDBC, ODBC drivers & Beeline

2) Hive Services:

- Drivers: Parses Queries, execution plans
- Metastore: Stores Metadata (Schema, Partition, locations)
- Execution Engine: Converts Queries into MapReduce, Tez or Spark jobs.

3) Hive Storage (HDFS) stores actual Data files.

Execution Flow of a Query.

1. Query submission via CLI, Beeline or JDBC
2. Query parsing and compilation
3. Execution Plan generation (MapReduce/Tez)
4. Query Execution & result retrieval.

Primitive Data Types.

Int, BigInt, Float, Double, Boolean, String & Date/Timestamp

Complex Data type.

Array, Map, Struct.

Array = Array(int)

Map = Map(String, int)

Struct = Struct(name: String, age: Int)

Syntaxes / Statements Common.

- (i) Create database mydb;
- (ii) Use mydb;
- (iii) Create table employees (

id int,

name String,

salary float

) Row format delimited
field terminated by ',';

For EXTERNAL TABLE
Create External Table

location: ' ' ;

- (iv) Insert into employees
Values (1, 'John', 5000)
(2, 'Kelvin', 6000)
(3, 'Curry', 7000);

- (v) Select * from employees;

Working
Managed Table

Create table
(id INT, name)

Loading

hdfs dfs

LOAD DATA

INTO TABLE

Partition

partition

Create

i

p

Bucket

Create

id

) cu

Hive F

String

Date

Aggregate

Perfor

Use

10 e

Exc

⇒ S

Enab

⇒ se

Working with Tables:

<u>Managed Table by Hive</u>	<u>External table by 3rd party</u>
Create table tablename (id int, name String);	Create External table tablename (id int, name String) LOCATION 'hdfs dfs/path/';

Loading Data

hdfs dfs -put employees.text /user/hive/warehouse/
LOAD DATA PATH '/user/hive/warehouse/employees.text'
INTO TABLE employees;

Partitioning & Bucketing

Partitioning (Improves Query Performance)

Create table sales (
id int, amount float)
partitioned by (year int);

Bucketing (Distributes data into files).

Create table users (
id int, names String)
) clustered by (id) into 4 buckets;

Hive Functions

String function: Select Upper(name), Lower(name) from Emp;

Date function: Select Current date, Year('2025-03-18') from Emp;

Aggregate fn: Select Count(*), Sum(salary) from Emp;

Performance Optimization

Use ORC, or Parquet formats instead of txt.

To enable Partitioning & bucketing, use Vector Query Execution.

⇒ Set hive.vectorized.execution.enabled = true;

Enable Cost Based Optimization (CBO)

⇒ Set hive.cbo.enabled = true;

Hive Integration with other tools include:

- Apache Spark (Spark SQL Can run on Hive)
- Apache HBase (Hive Can Query NoSQL data)
- Apache Airflow (for job Scheduling).

Hive Applications with Real-World Use Cases:

- log Analysis (e.g. Web Server logs)
- ETL pipelines (Extract-Transform-load processes)
- Data Warehousing (Aggregations, Reporting).