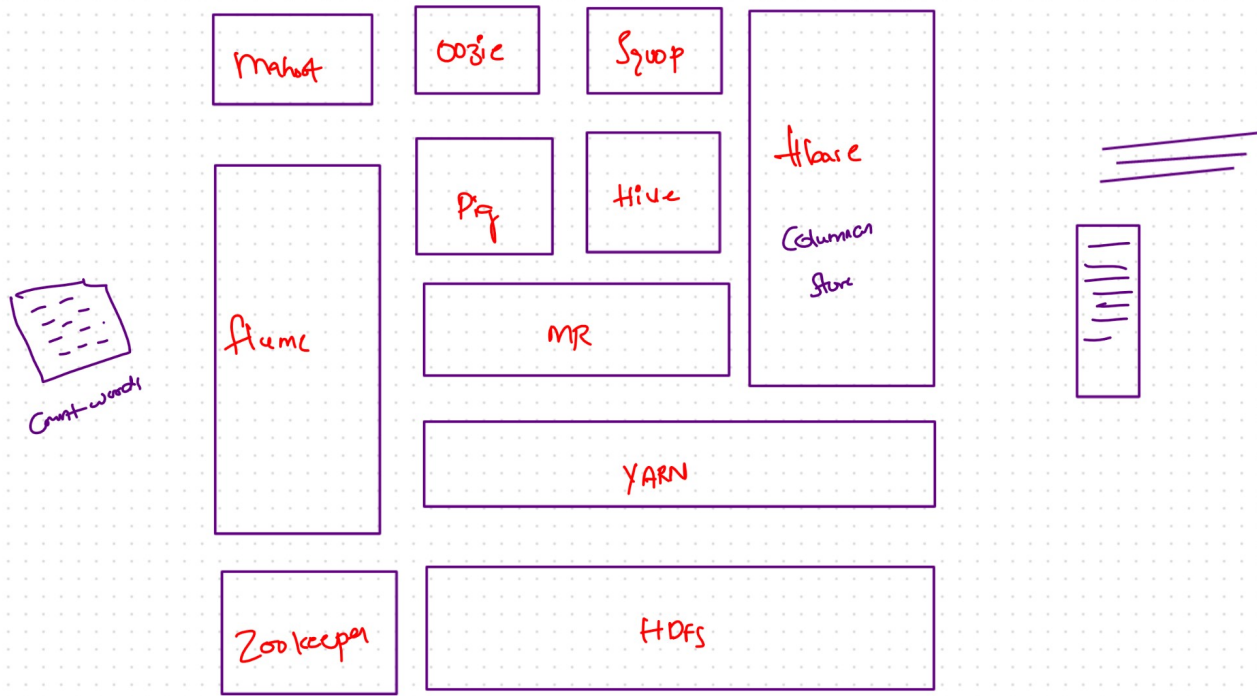


Hadoop Ecosystem Components



Extended components

4. Hive:

* developed by facebook in 2009

* Query engine (not a database)

Key feature: Abstracts MapReduce by translating SQL queries into MapReduce jobs

Solves: Avoids writing java code while retaining MapReduce performance

Use case: Data warehousing solutions

5. Pig:

* created by yahoo

Key feature: High-level scripting language (Pig Latin) for MapReduce

Use case: Data transformation without writing Java or SQL

Abstracts MapReduce in the backend

6. Sqoop:

Key feature: It helps to import & export data between Hadoop & relational databases

Use case: Data migration without writing MapReduce code

7. Oozie:

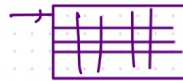
* Created by Cloudera

* key feature: XML-based workflow scheduling & automation

Use case: manage job dependencies using XML for scheduling

complex workflow

map {key:-}



8. HBase:

* NoSQL database modeled after Google's Big Table

* Column store (i.e. stores data by columns rather than rows)

Key feature: Allows real-time reads & writes on HDFS

Use cases: Time series data, messaging system

9. Mahout:

* Data science component

* key feature: provides machine learning libraries

Use case: implements algorithms for clustering, classification, collaborative filtering

10. flume:



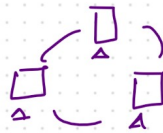
data ingestion & streaming

* messaging queue (similar to kafka)

* key feature: collects log or event data from various sources

* Use case: Real-time analysis, monitoring

11 Zookeeper:



→ Developed by Yahoo

→ Key features: co-ordinates distributed system to maintain consistency

Use Case: Critical for ensuring reliability in Hadoop Cluster

Note:

Hadoop framework Properties

(1) loosely coupled framework:

→ components can be removed or replaced while the system working

eg: MapReduce being replaced by Spark for processing

(2) Easy Integrations:

→ Can connect with other big data technologies (e.g. Spark)

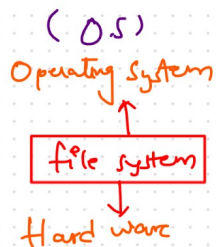
→ Can integrate with non-big data technologies

eg: (MySQL, other relational databases)

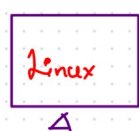
Key Terminology

HDFS = Hadoop Distributed file System

① File System



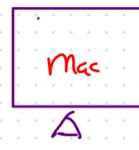
Servers ^{speed} _{secured}



ext3/ext4



NTFS/FAT32



APFS



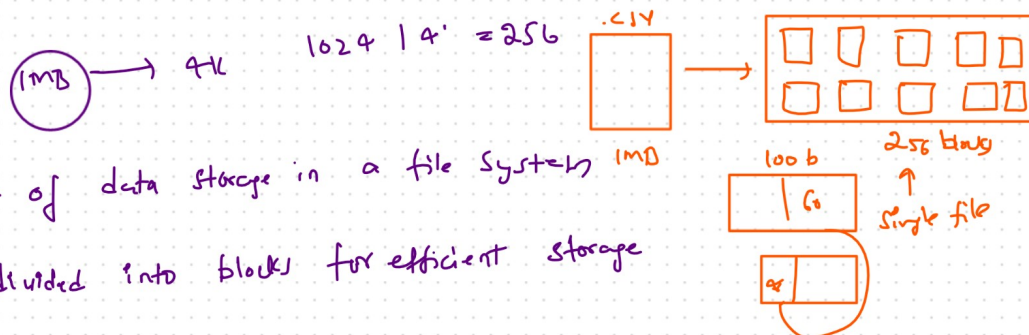
HDFS (distributed file system)

(Storage device)

- It is a method of data structure that operating systems use to manage files on storage devices
- Acts as a layer b/w software (O.S) & hardware (Storage devices)
- determines how data is stored, retrieved and organized

2. Blocks

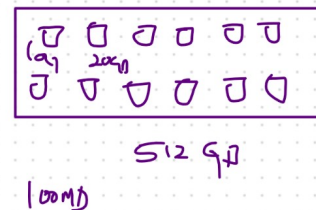
- + Smallest unit of data storage in a file system
- + files are divided into blocks for efficient storage



ex:

Windows (NTFS) : 4KB block size

HDFS : 128 MB default block size



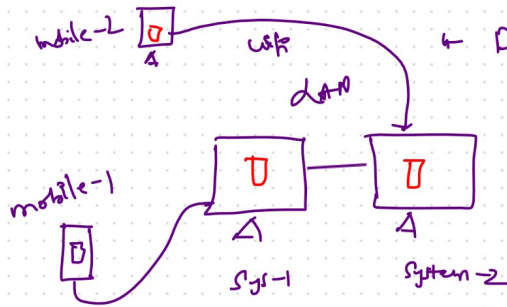
Benefits :

- + Quick & efficient for data storage & retrieval
- + efficient use of storage space

3. Types of File Systems

→ Standalone File Systems

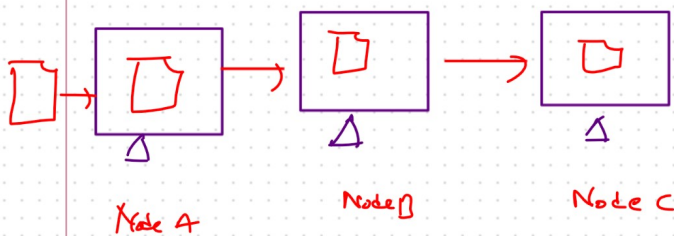
→ Distributed File Systems



→ files are stored & managed on a single machine
→ even in connected environment, files stay on one machine

Standalone File System

Distributed File System (HDFS)



⇒ files are stored across multiple machines in a cluster

more scalable — can add hundreds of machine to increase storage

3 Node cluster