

Big Data

- Volume
- Variety
- Velocity
- Veracity
- Volume

Need for a new technology stack ?

Monolithic vs Distributed

One powerful system

Cluster of systems

(Commodity machines)

Compute, Memory, Storage } $\begin{matrix} 2\times \text{resources} \\ \neq \\ 2\times \text{performance} \end{matrix}$

Vertical Scaling

Horizontal Scaling
(True Scaling)

Design a good Big Data System? Consider

- Storage (distributed storage)
- Process (distributed processing)
- Scalability (\uparrow number of nodes)

Hadoop → first framework designed to solve Big Data problems

2006 2012 Now
1.0 → 2.0 → 3.0

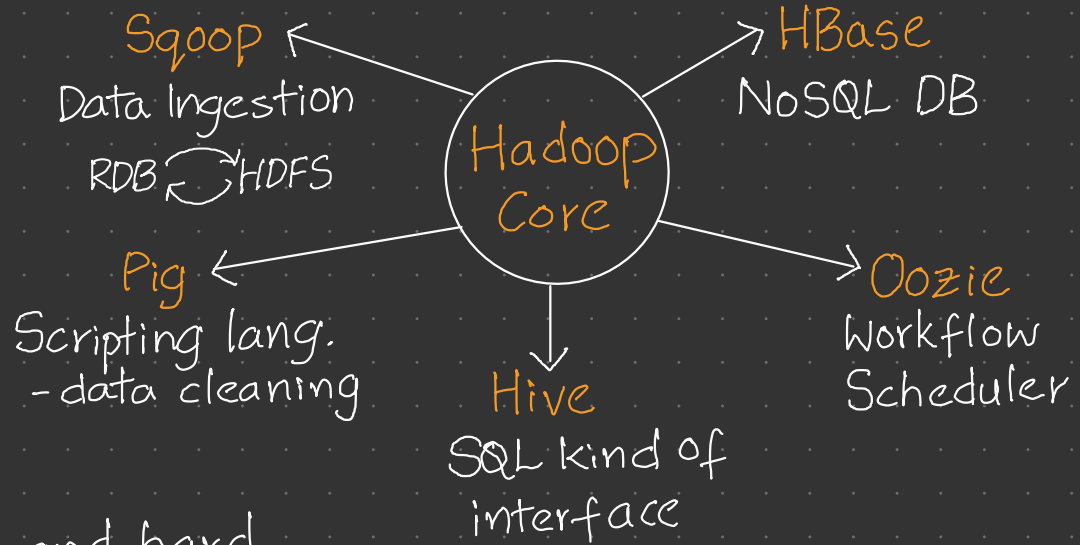
Hadoop Ecosystem

Core Components

- HDFS
- Map Reduce
- YARN (added in 2.0)

Challenges

- MR is very slow and hard
- Need to learn different components for different tasks



Advantages of cloud

- Scalable
- Minimum OpEx
- Agility
- Geo Distribution
- Disaster Recovery

Cloud Types

- Public : AWS, Azure, GCP
- Private : when the data is very confidential
- Hybrid

On-Premise

Vs

Cloud

- | | |
|--------------------------------------------------------|----------------------------------------------|
| - Buy the needed Infra, like space to hold the servers | - Infra is taken care by cloud providers |
| - Buy servers | - No need to buy servers |
| - Setup a cooling system | - Setup the cluster with a click of a button |
| - Hire a technical team to install and maintain s/w | - Low maintenance cost |
| - Huge upfront cost / Capex & Opex | - No upfront cost / Capex |
| - Not very scalable | - Highly scalable |

Spark → general - purpose in-memory compute engine

- Spark can act as a replacement for MapReduce (not for Hadoop)
- It is a plug and play compute engine
- Needs storage and resource manager to function.

not bound to
HDFS, and YARN, only.

Amazon S3,
ADLS Gen 2,
GIC storage,
Local storage

Mesos,
Kubernetes

- Spark is 10x - 100x faster than MR

Polyglot

- Provides high-level APIs in Java, Scala, Python and R
- Spark is written in Scala

	Database	Data Warehouse	Data Lake
Workloads	- OLTP	- OLAP	- OLAP
Data type	- Structured or Semi-Structured	- Structured and/or Semi-Structured	- Structured, Semi-Structured and/or unstructured
Schema	- Schema-on-write	- Schema-on-write	- Schema-on-read
Freshness	- Operational current data	- Current and historical data	- Current and historical data
Prior Processing	- ETL	- ETL	- ELT
Cost	- Higher	- High	- Low
Examples	- PostgreSQL, Oracle, MySQL	- Teradata, AWS Redshift	- Amazon S3, ADLS Gen 2
Users	- Application developers	- Business Analysts	- Data Scientists

Data Engineering Flow



Data Pipeline Hadoop



Data Pipeline Azure



Data Pipeline AWS



Serverless Computing

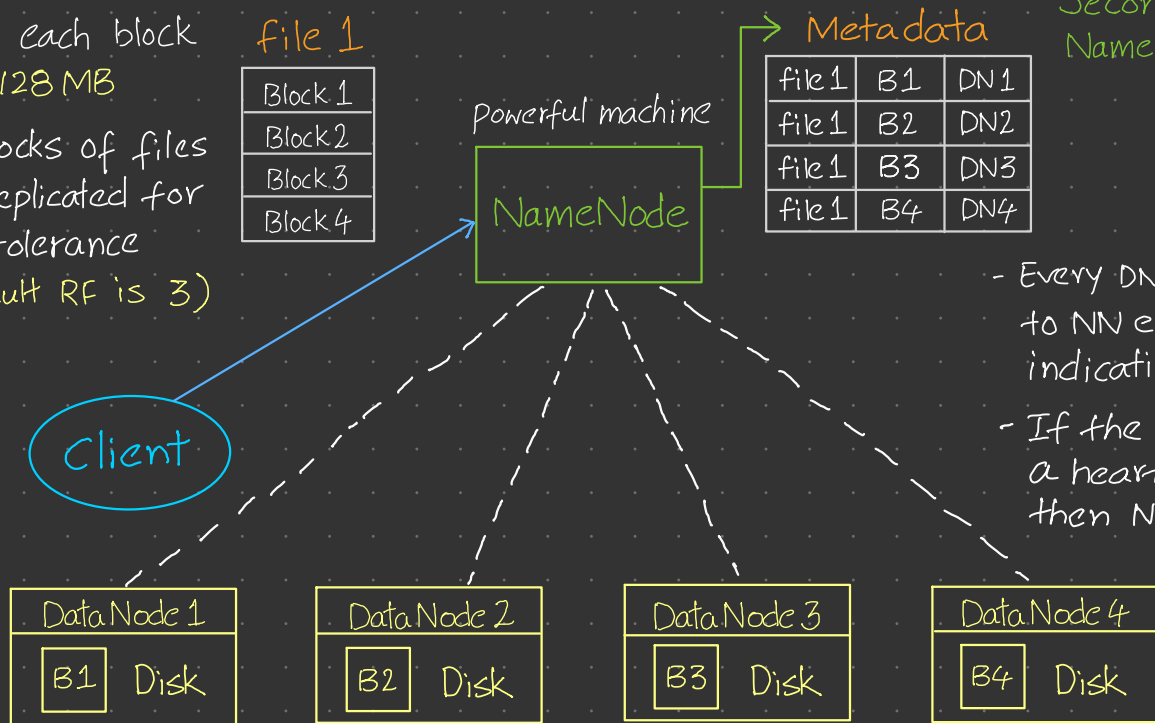
- Resources are not dedicated
↓
No guaranteed performance
- Less expensive
- Good for scheduled jobs
↓
that are not particular about the time in which they get finished
- Eg: Athena, Synapse serverless

Serverful Computing

- Resources are dedicated
↓
guaranteed performance
- More expensive
- Good for ad-hoc jobs
↓
that require immediate execution
- Eg: - Redshift, Synapse dedicated pool

HDFS Architecture

- Usually each block is of 128 MB
- the blocks of files gets replicated for fault tolerance (default RF is 3)



Secondary NameNode → backup for NameNode

Heartbeat

- Every DN sends a heartbeat to NN every 3 seconds indicating that it's alive
- If the NN does not receive a heartbeat for 30 seconds, then NN assumes DN is dead

Data locality

↓
Move the computation close to where the actual data resides on the node

Commodity machines