

Apache Oozie:-

* Oozie is a job scheduling system integrated with Hadoop.

* used to execute multiple jobs in parallel.

* Integrated Spark, Hive, scala

* Integrated scheduler for Hadoop.

* It is a workflow scheduler for Hadoop.

* It is a system which runs the workflow

of dependent jobs

* user permitted to create DAG (Directed

Acyclic Graph) of workflow which can run

in parallel and sequential

3 types of jobs

1. Oozie workflow jobs

- DAG is used to define jobs
edges specifies action.

2. Oozie coordinator jobs

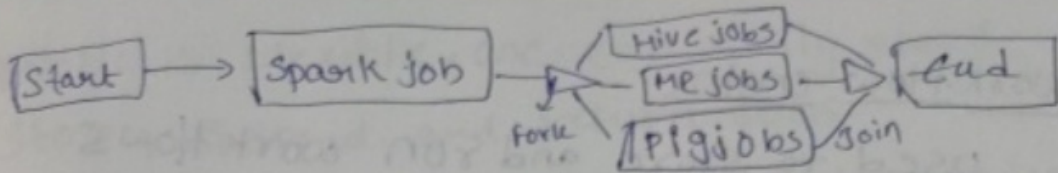
- workflow jobs triggered by time
and data availability.

3. Oozie bundle.

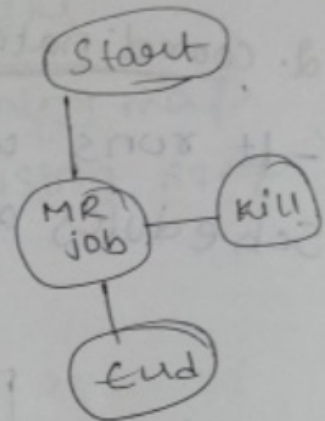
- Multiple coordinator.

How Oozie works:-

① Oozie workflow jobs



- output of the previous action is the input of the present.
- output of the present is the input of the future.



3/3/23

features of oozie:

- * Oozie has client API and command line interface which can be used to travel, control and monitor jobs from Java application.
- * using webservice API can control the jobs from anywhere.
- * provision to execute jobs which are schedule to run periodically.
- * provision to send email notification upon completion of jobs.

1. Workflow engine:

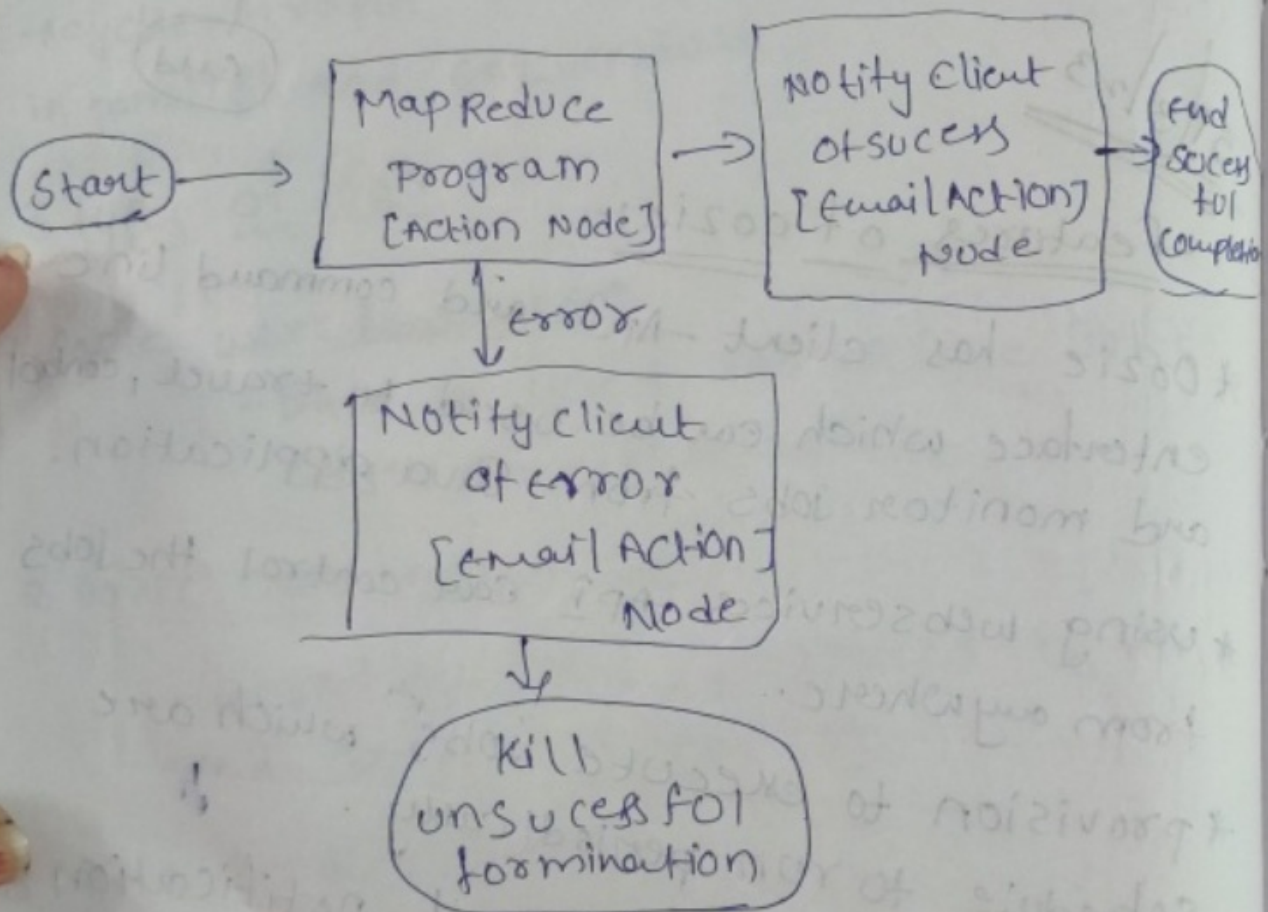
used to store and run workflows

composed of hadoop jobs.

- eg: map reduce, Pig, Hive.

2. Coordinator engines:

- It runs workflow jobs based on predefined schedules and availability of data.



*Oozie is scalable and can manage the final execution of thousands of workflow (each consist of dozens of jobs) in a hadoop cluster.

* Oozie is flexible, one can easily start, stop, suspend and return jobs.

* Oozie workflow consist of Action nodes and control flow nodes.

Action node:-

- Represents work flows tasks.

(eg.) moving files to HDFS, running map reduce jobs import data using sqoop or running shell script.

Control Node:-

control workflow execution b/w actions

→ allowing conditional logic where different branches may be followed depend up on the result of the earlier node.

→ Start node:- it used to start the workflow job.

→ end node: it signals the end of the job

→ error node:- it designated to the occurrence of the error and notify the error to notify print

15/3/23

Hadoop

- It is an open source software framework used for storing data and running applications on a group of commodity hardware.
- open source.
- Distributed data processing
- map - Reduce programming.

Why hadoop

- handle any data type
 - structured / unstructured data.
- schema, schema less
 - ↳ transaction → job.

- high volume / low volume of data.
- All kind of analytic application
- * proven with petabyte scale.
- * Grows with business
- * capacity and performance grows.

features

- No vendor lock in
- Rich eco system and community development
- More affordable
- Scalable

- cost effective
- flexible
- fast
- Resilience to failure.

Hadoop Infrastructure

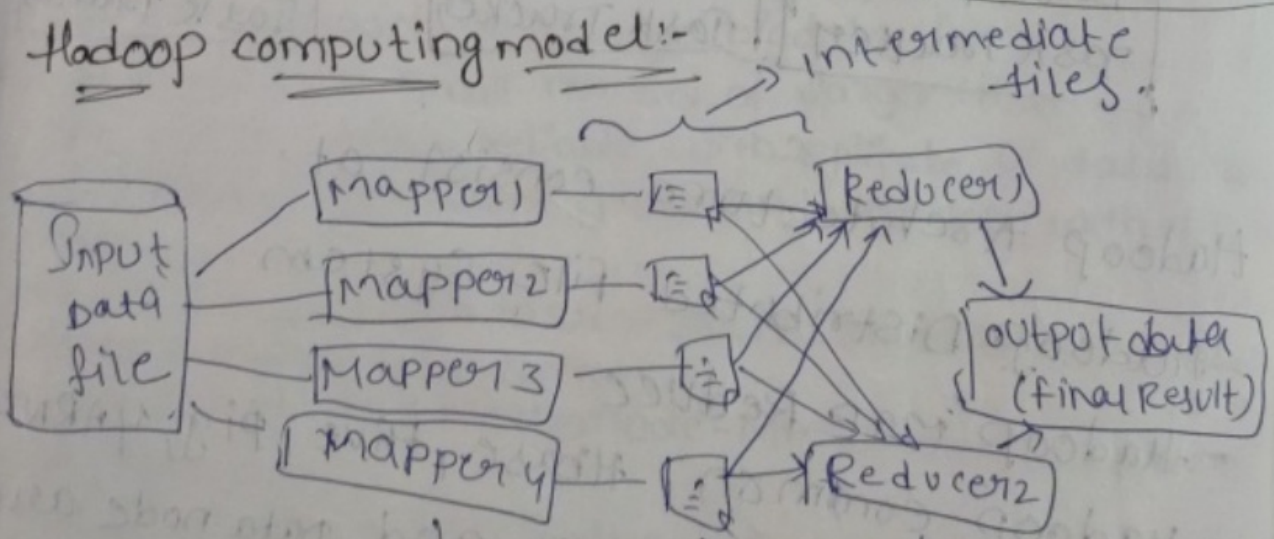
Two Infrastructure model of Hadoop

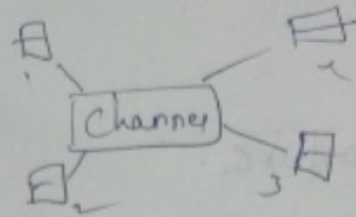
1. Data model
2. computing model

Distributed data model vs Hadoop data model

- | | |
|---------------------------------------|------------------------------------------------|
| - deals with tables and relation | - deals with flat files in any format |
| - must have schema for data | - No schema |
| - data fragmentation and Partitioning | - files are automatically divided into blocks. |

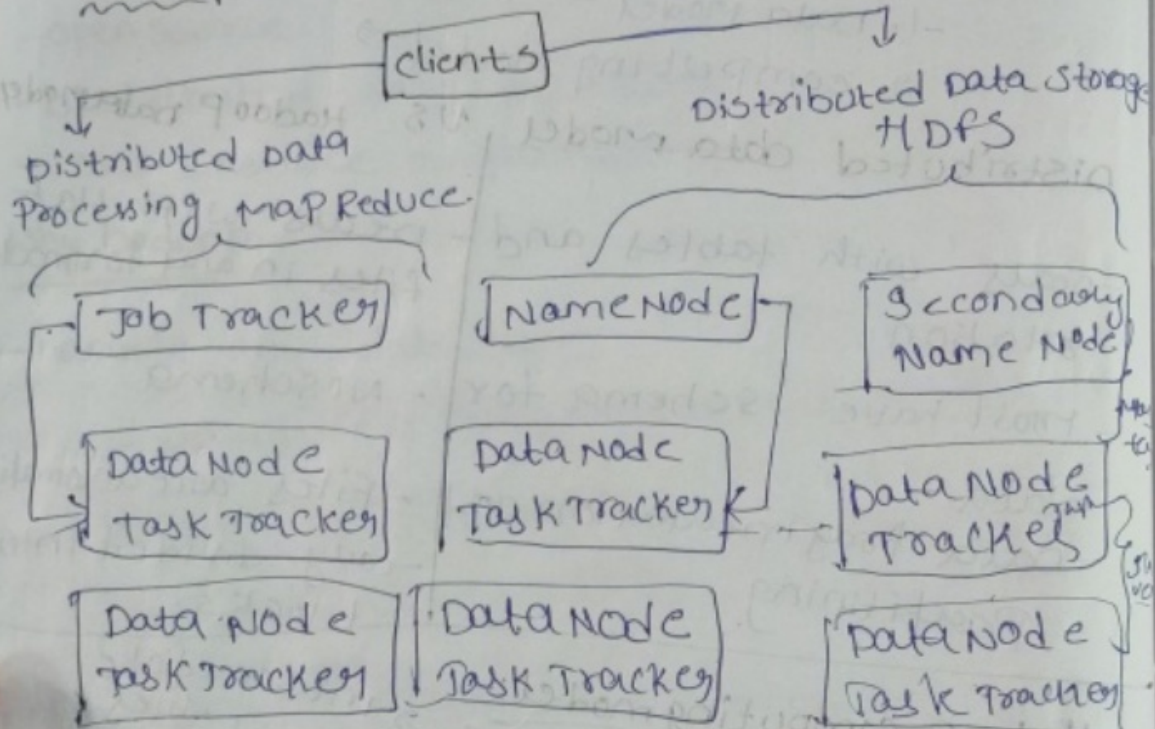
Hadoop computing model:-





16/3/23

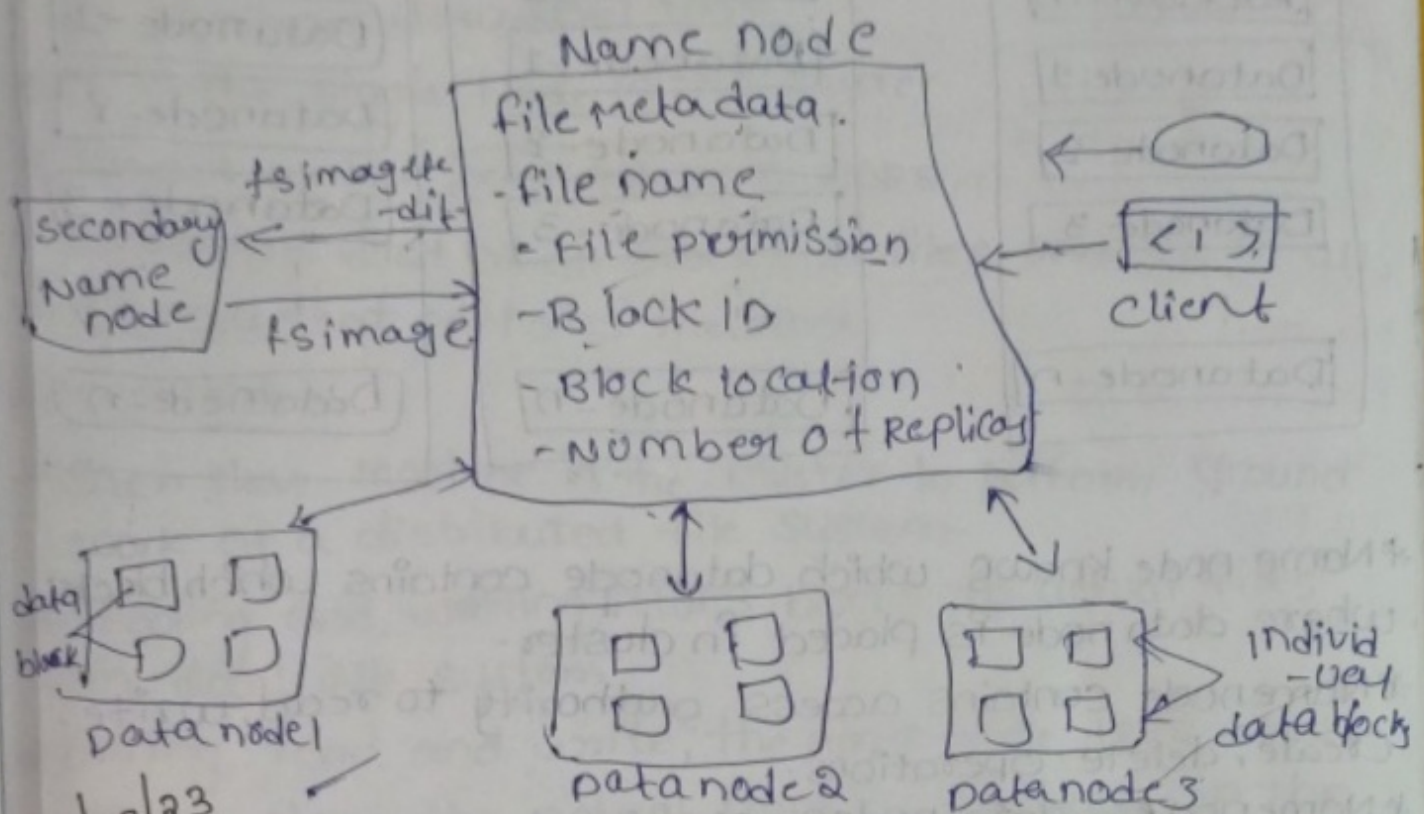
Hadoop Architecture:-



Hadoop Architecture consist of

- Hadoop Distributed file System
- Hadoop Map Reduce
- Hadoop common, Hbase, Hive, Pig, YARN etc,
- * Name Node is master and Data node are Slaves
- * Job Tracker is the master and task tracker is the slave

* every slave node come with task track daemon and data node synchronize with the job tracker and namenode respectively.



17/03/23

* HDFS Working model: