## Apache Oozie:-

* Oozie is a job scheduling system integration with Hadoop.

* used to execute multiple jobs in parallel.

* Integrated Spark, Hive, scala

* 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
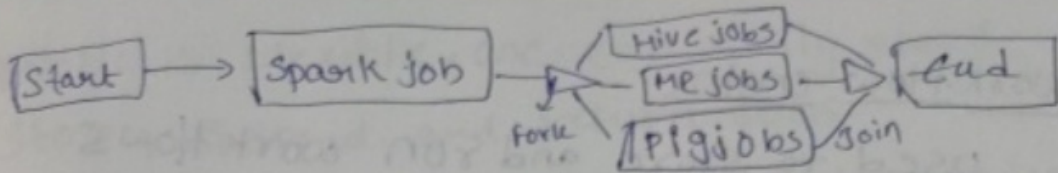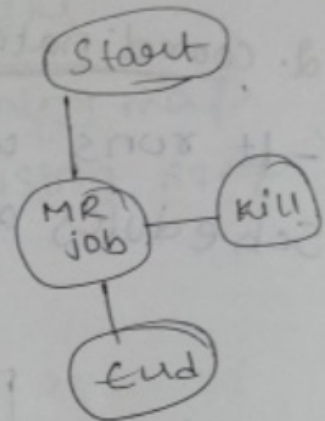   - workflow jobs triggered by time and data availability.

3. Oozie bundle.
   - Multiple coordinator.

## How Oozie works:-

① Oozie work flow jobs

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



31/3/23

## 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 webservices API can control the jobs from anywhere.

* provision to execute jobs which are schedule to run periodically.

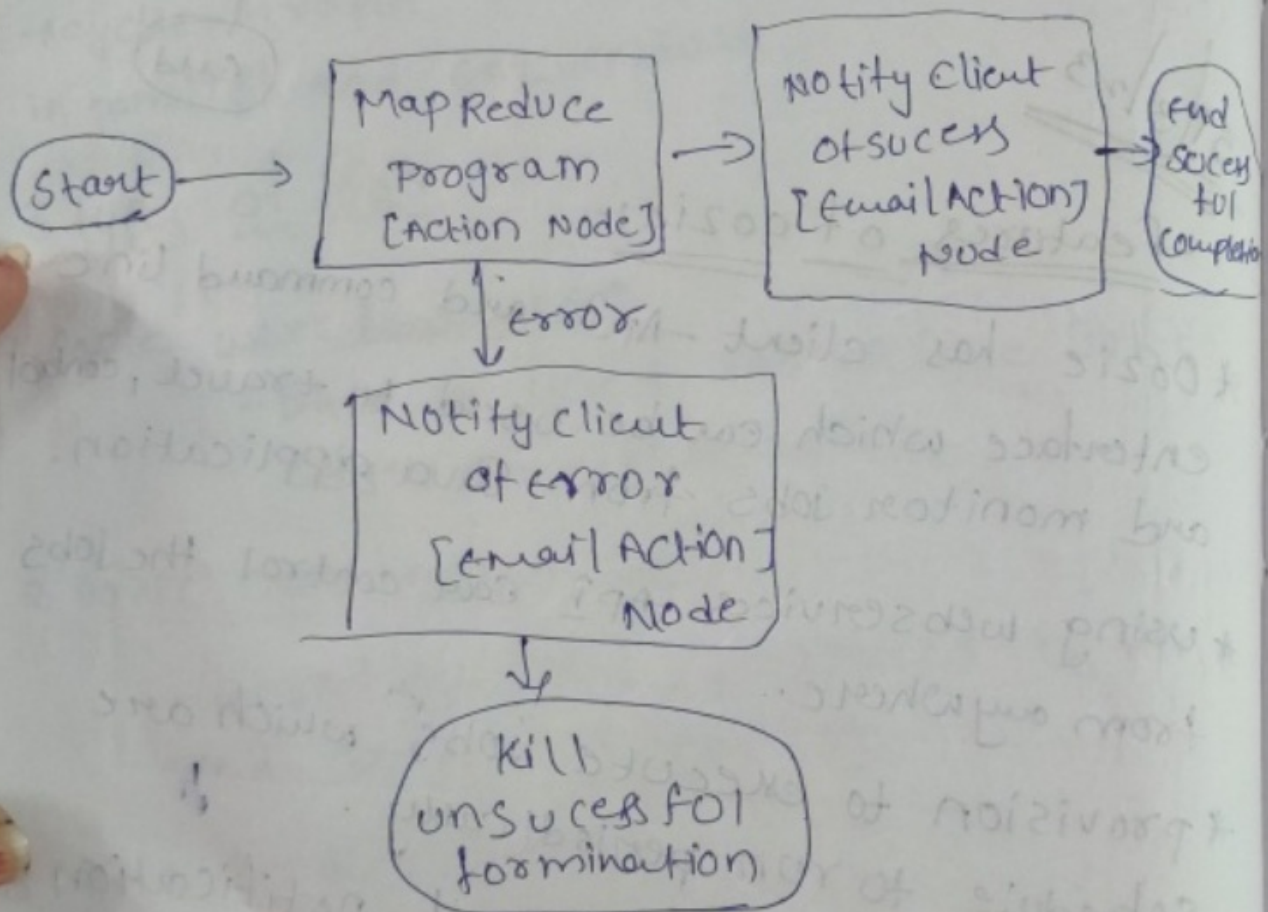* provision to send email notification upon completion of jobs.

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 predefine
schedules and availability of data.

```
(Start) ──→ ┌─────────────┐ ──→ ┌──────────────┐ ──→ (End
             │ Map Reduce  │     │ Notify Client │      Sucey
             │ Program     │     │ of sucess     │      ful
             │ [Action Node]│     │ [email Action]│     Completio
             └─────────────┘     │   Node        │
                    │            └──────────────┘
                    │ error
                    ↓
             ┌──────────────┐
             │ Notify client │
             │ of error      │
             │ [email Action]│
             │    Node       │
             └──────────────┘
                    ↓
              ( Kill
               unsucessful
               formination )
```

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

*Oozic 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 sroop or running sheu script.

*Control Node:-

control wook flow execution b/w actions

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

→ Start node:- it used to start the work flow 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

4