# Unit-5 — Map Reduce

## Map Reduce:

- programming model associated implementation for processing & generating large data sets
- Developed - Google - for large scale data processing in distributed computing environments
- Dividing large data set into smaller chunks, Processing them in a parallel Process across multiple nodes in a distributed system & combining into final output

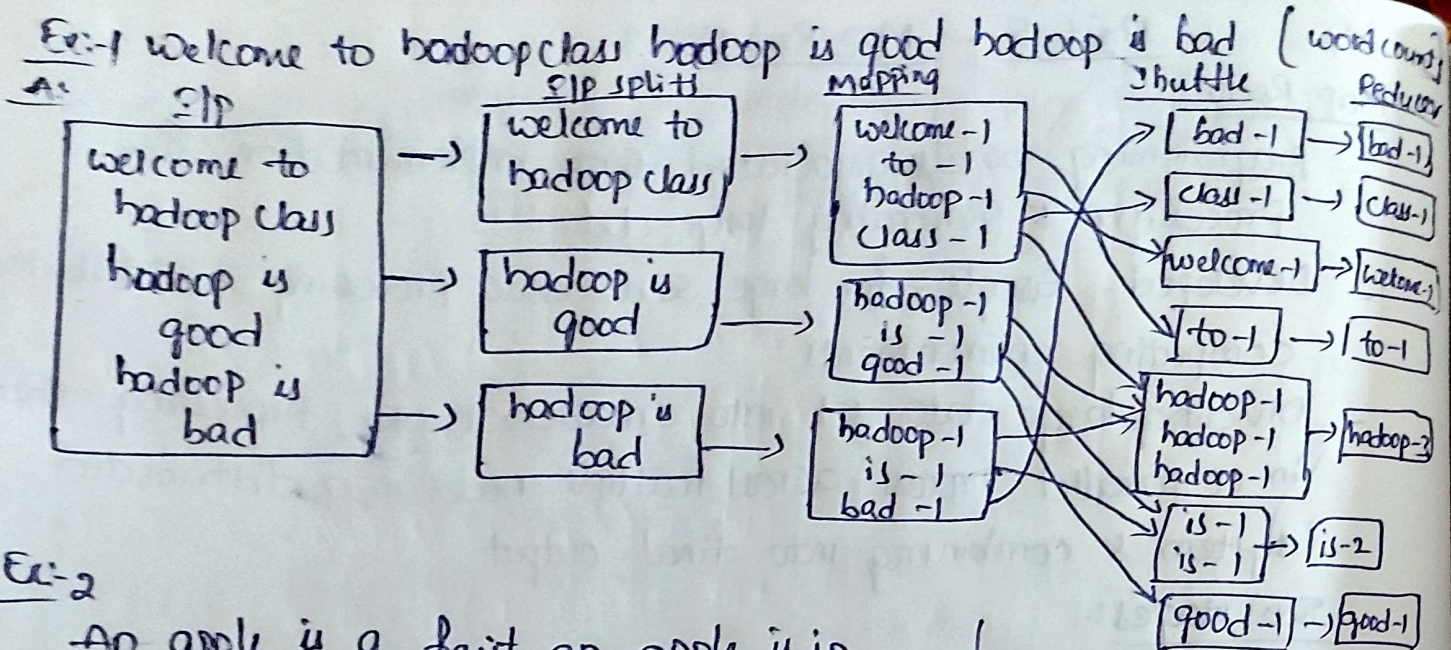### Two stages:-

- Map
- Reduce

### Map:- 
Each node in a system process a partition of input data & generates a set of intermediate key value pairs

- Typically specified by used in the form of map fn, which take a key - value pair & produces a set of intermediate key- value
- The intermediate key value shuffled & sorted based on their keys & send to reduce stage

### Reduce:- 
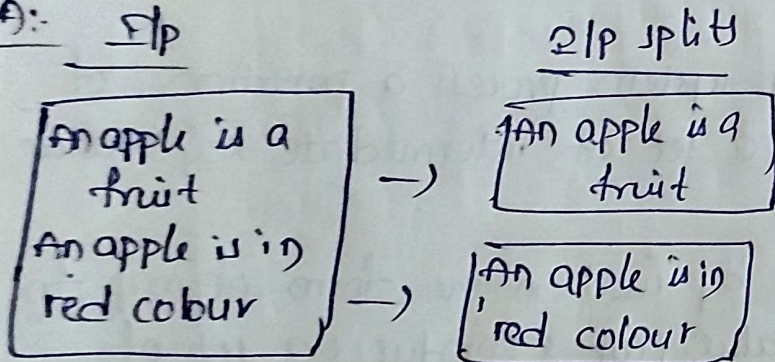process inter key value produced by map stage & generate final output.

### Algorithms using MapReduce:-

- Matrix vector multiplication
- word count map reduce Architecture
- MR has following phases
  - L I/p splits
  - L mapping
  - L shuffling
  - L sorting
  - L Reducing

Ex:-1 welcome to hadoop class hadoop is good hadoop is bad (word count)

A: I/p → I/p split → mapping → shuttle → Reducer

**I/p**
```
welcome to
hadoop class
hadoop is
good
hadoop is
bad
```

**I/p split**
```
welcome to
hadoop class
```
```
hadoop is
good
```
```
hadoop is
bad
```

**mapping**
```
welcome -1
to -1
hadoop -1
class -1
```
```
hadoop -1
is -1
good -1
```
```
hadoop -1
is -1
bad -1
```

**shuttle**
```
bad -1
```
```
class -1
```
```
welcome -1
```
```
to -1
```
```
hadoop -1
hadoop -1
hadoop -1
```
```
is -1
is -1
```
```
good -1
```

**Reducer**
```
bad -1
```
```
class -1
```
```
welcome -1
```
```
to -1
```
```
hadoop -3
```
```
is -2
```
```
good -1
```

Ex-2

An apple is a fruit an apple is in red colour.

A:-  **I/p**

```
An apple is a
fruit
An apple is in
red colour
```

**I/p split**

```
An apple is a
fruit
```
```
An apple is in
red colour
```

**final o/p**
```
bad -1
class -1
welcome -1
to -1
is -2
good -1
hadoop-3
```

**mapping**
```
An -1
apple -1
is -1
a -1
fruit -1
```
```
An -1
Apple -1
is -1
in -1
red -1
Colour -1
```

**shuttle**
```
An-1
An -1
```
```
Apple-1
Apple -1
```
```
is -1
is -1
```
```
a -1
```
```
fruit 1
```
```
in -1
```
```
red -1
```
```
colour -1
```

**Reducer**
```
An- 2
```
```
Apple -2
```
```
is- 2
```
```
a -1
```
```
fruit -1
```
```
in -1
```
```
red -1
```
```
colour -1
```

**final o/p**
```
An- 2
Apple-2
is- 2
a -1
fruit -1
in -1
red -1
colour -1
```

# YARN:-

- Yet Another Resource Negotiation
- Cluster management component
- Large scale distributed os for big data appl
- Yarn - resource manager created by separating the process engine & mgt far of MR
  - monitors & manages workload

## YARN Architecture:

↳ splitting job tracker responsibility mgt (management) & job scheduling / monitoring into separate daemons

## concept:

↳ Application: job submitted to the system
  eg. MR job

↳ container: Basic unit of allocation
  - Resource allocation across
  multiple resource type eg:- bcpu, container - 0: 2GB

## components

↳ client: submitting MR jobs

↳ Resource manager: manage use of resource across cluster
  ↳ 2 components ⟨ schedular
                  ⟨ Application manager

**☆schedular:** schedular of RM decides allocation of resource to running applications
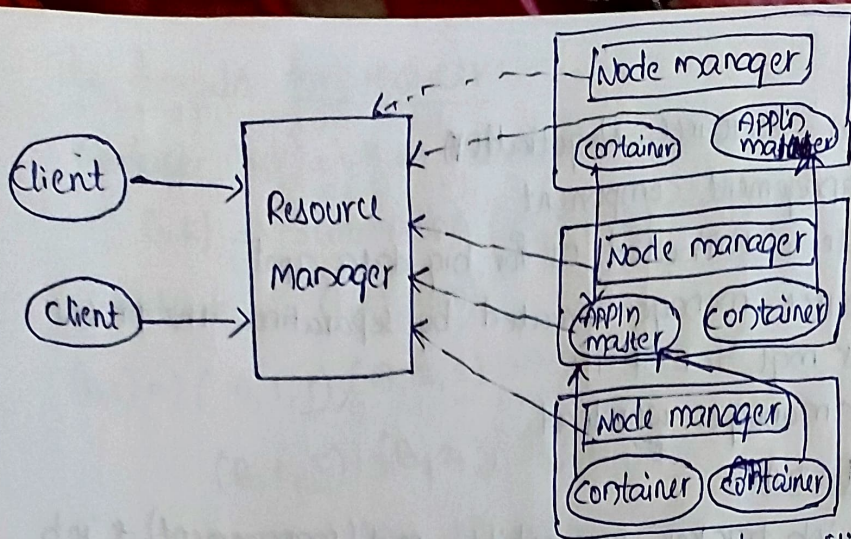
**Application manager:** → Accepting job submission
  → Negotiating - in appli master
  → Restarting

**Node manager:** launching / monitoring component containers on machines in cluster
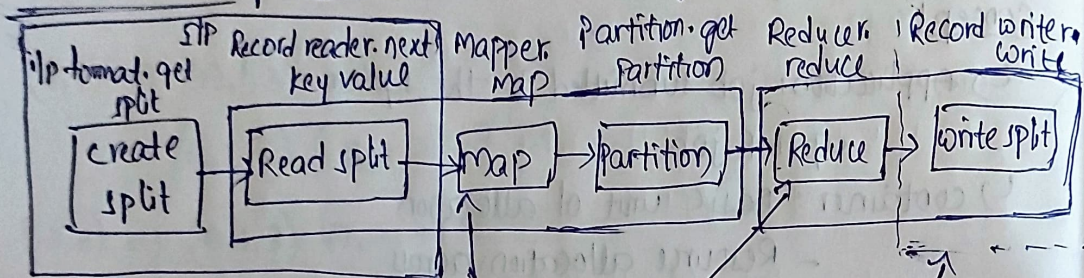  - Monitors resource usage

↳ MR application Cluster: checks task running the MR job

- Application master of & MR task run in container that are scheduled by resource manager & managed by node manager

Diagram: Client → Resource Manager; Node managers with Container, Appln master components.

## Understanding i/p's and o/p's in MR :-

the partitioner job is **logically** format map o/p to reducer & output

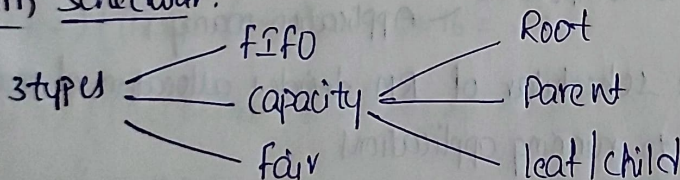| i/p format. get split | SIP Record reader. next key value | Mapper. map | Partition. get partition | Reducer. reduce | Record writer. write |
|---|---|---|---|---|---|
| create split → | Read split → | map → | partition → | Reduce → | write split |

↑ the i/p format & Record reader are responsible for determining what data to feed into map tar

↑ The map & reduce funs are typically written by the user to address the specific use case

↑ The record writer writes the reduce o/p to the destination data sink, which is the final resting place of the MR data flow

## Yarn scheduler :

3 types
- FIFO
- capacity
- fair

- Root
- Parent
- leaf / child

## capacity scheduler :



```
              Root
              100%
     ┌──────────┼──────────┐
     ↓          ↓          ↓
   Ad hoc    work flow   Preference
   20%-40%   60-80%      20%-80%
```

leaf of root { Ad hoc 20%-40% }

Min capacity of root leaf add upto 100%.
20+60+20 = 100

leaf of work flow {
- Ingest 35-50%
- ETL 65%-80%
- low 20-50%
- High 80%-100%
} leaf of preference