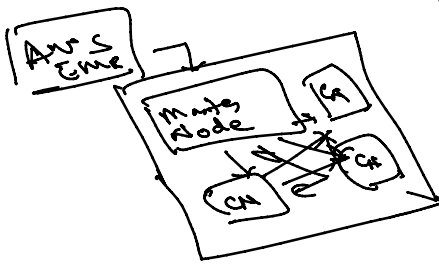


## ① Amazon EMR (Elastic Map Reduce) :-

① Managed cluster platform that simplifies running of Big data framework.



① Master/Primary Node :-  
Coordinating b/w core nodes. Give a lot of task or store

② Core Node :- Store data (HDFS) & run the task.

③ Task Node :- Run the task. It is not store any data.

Give PS  
Spark Hadoop Distributed File system

## YARN (Yet Another Resource Negotiator) :-

- ① Help the AWS EMR to Allocate the Resources
- ② Coordinate & schedule task among the node of the cluster.

Cluster Resource Manager - Manage the cluster & scheduling job.

Hadoop :- Distributed processing. Process Huge Amount of Data.  
HDFS

Spark :- Distributed processing Run on RAM & Cache memory. It bit costlier with Hadoop. Faster compared to Hadoop.

Cost → Spot Instance.

Adv. of EMR :-

- ① Low Cost  $\rightarrow$  Spot Instance. (50-80%)
- ② Scalability  $\rightarrow$  Scale in & scale out
- ③ Integration with other AWS services:-  
① S3

④ Security :- Amazon EMR integrate with AWS IAM for access control

Spark :-

- ① Distributed Processing
- ② Run on the RAM/Cache memory

Architecture

- ① Driver :- Central Component of Spark  
App Run main & high level for
- ② Executor :- Execute task
- ③ Cluster Managers Manage Resource of the cluster
- ④ Distributed Storage Amazon S3,  
Apache Cassandra,  
HDFS