# DICE
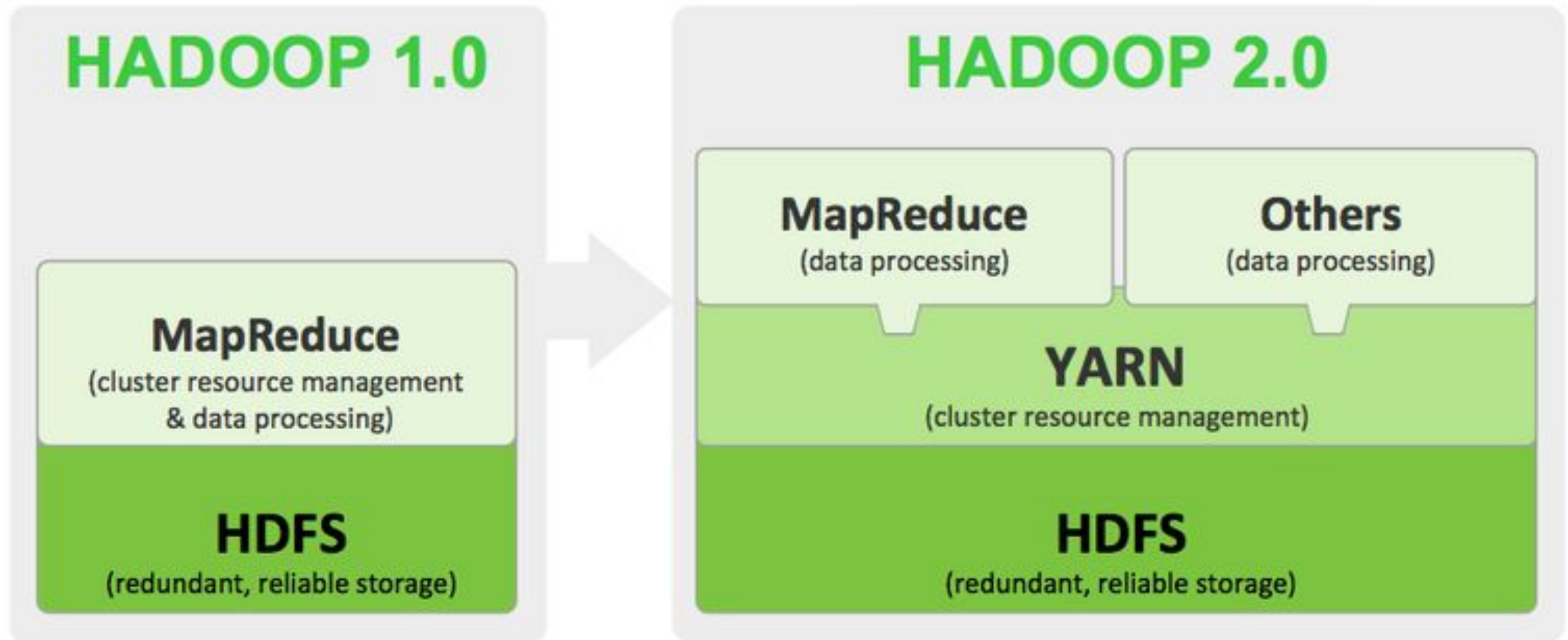## ANALYTICS

**Apache Pig**

DICE
ANALYTICS

# *YARN*
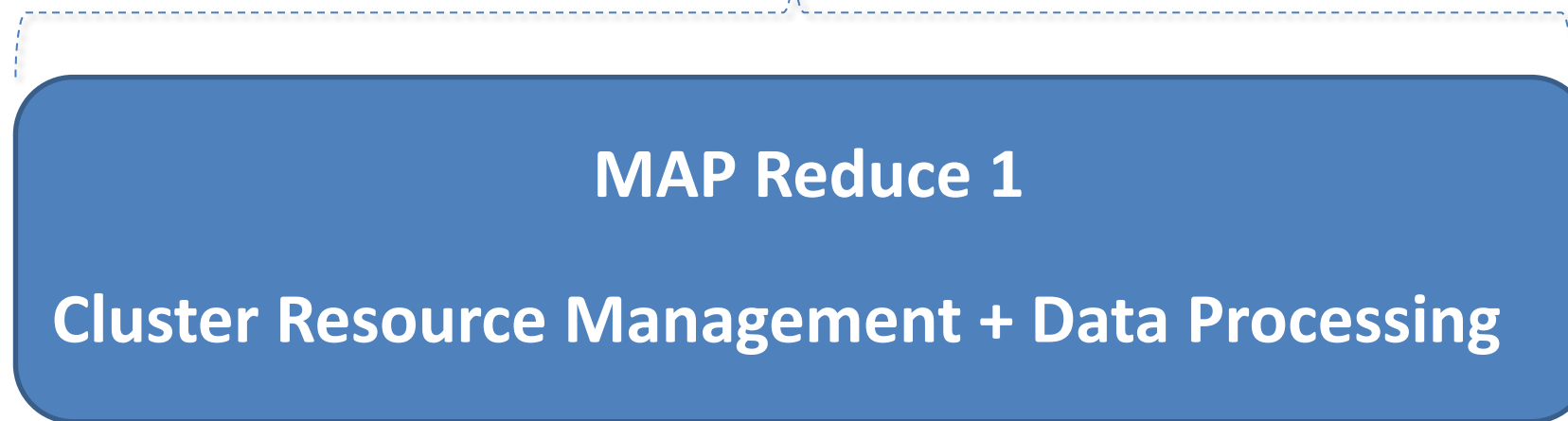
## *YET Another Resource Negotiator*
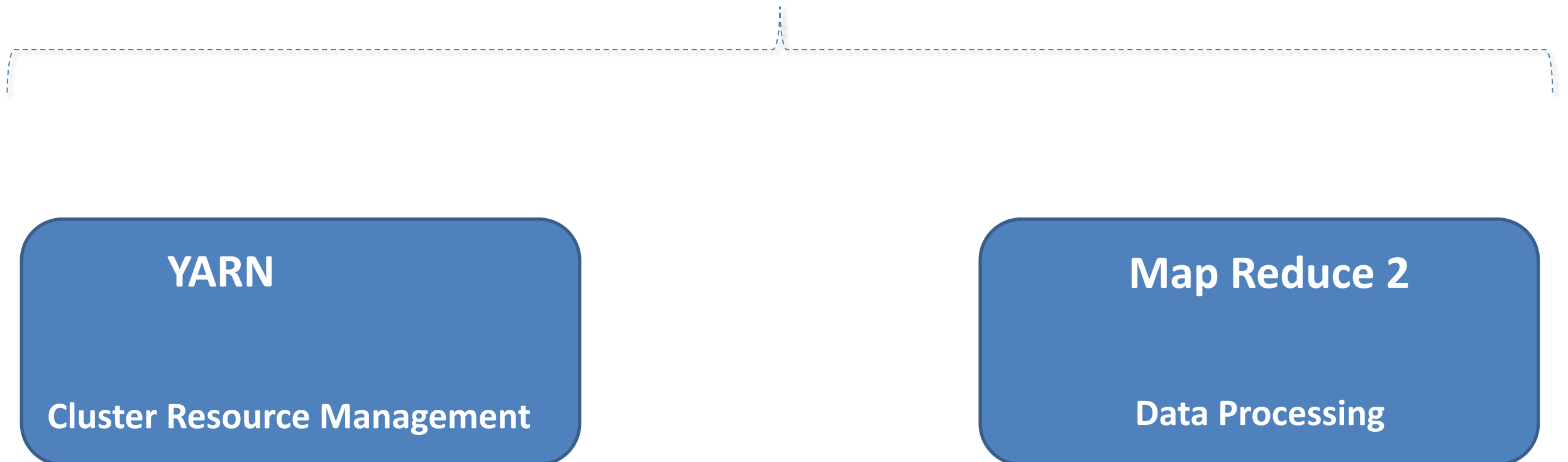
**https://www.youtube.com/watch?v=DMHf_xiSSgA**

- *Hadoop 2.0 Was Released in 2012*
- *YARN is part of the Hadoop 2.0 and onwards*
- *Previously cluster resource Management was part of the Map Reduce – Hadoop 1.0*

*Hadoop 1.0*

**MAP Reduce 1**

**Cluster Resource Management + Data Processing**

*Hadoop 2.0*

**YARN**

**Cluster Resource Management**

**Map Reduce 2**

**Data Processing**

YARN to MapReduce: You only do Data Processing because you are best at it. Let me do cluster resource management
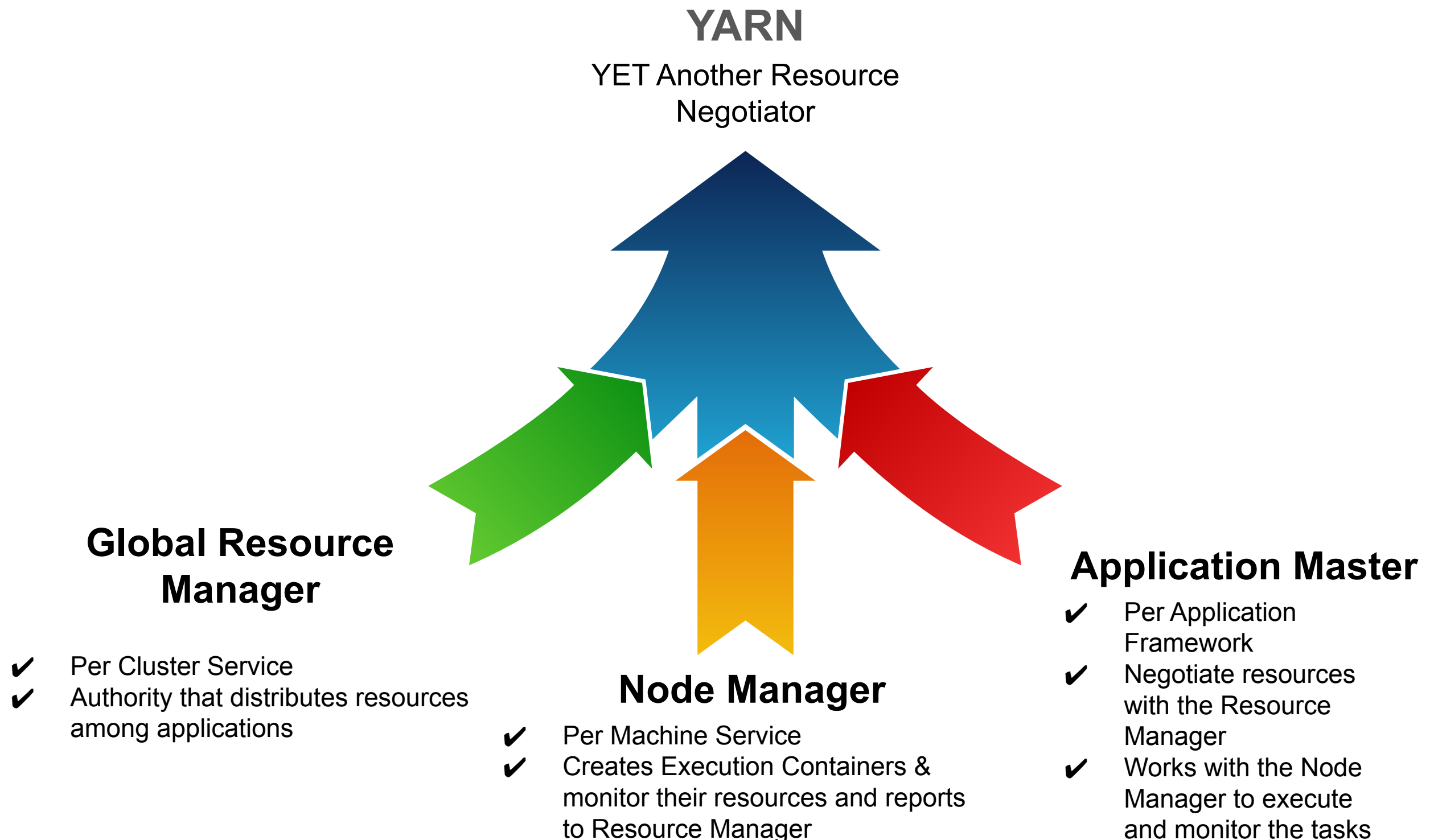
DICE
ANALYTICS

# Hadoop 2.0

Said Good Bye to

# NOW
its

# MAP Reduce 2

Which Only Does Data Processing

- YARN Divide Node Resources into *Containers*

  Container = Fixed **CPU** + Fixed **Memory**

**YARN**

YET Another Resource Negotiator



**Global Resource Manager**

✔ Per Cluster Service
✔ Authority that distributes resources among applications

**Node Manager**

✔ Per Machine Service
✔ Creates Execution Containers & monitor their resources and reports to Resource Manager

**Application Master**

✔ Per Application Framework
✔ Negotiate resources with the Resource Manager
✔ Works with the Node Manager to execute and monitor the tasks

# DATA ACCESS

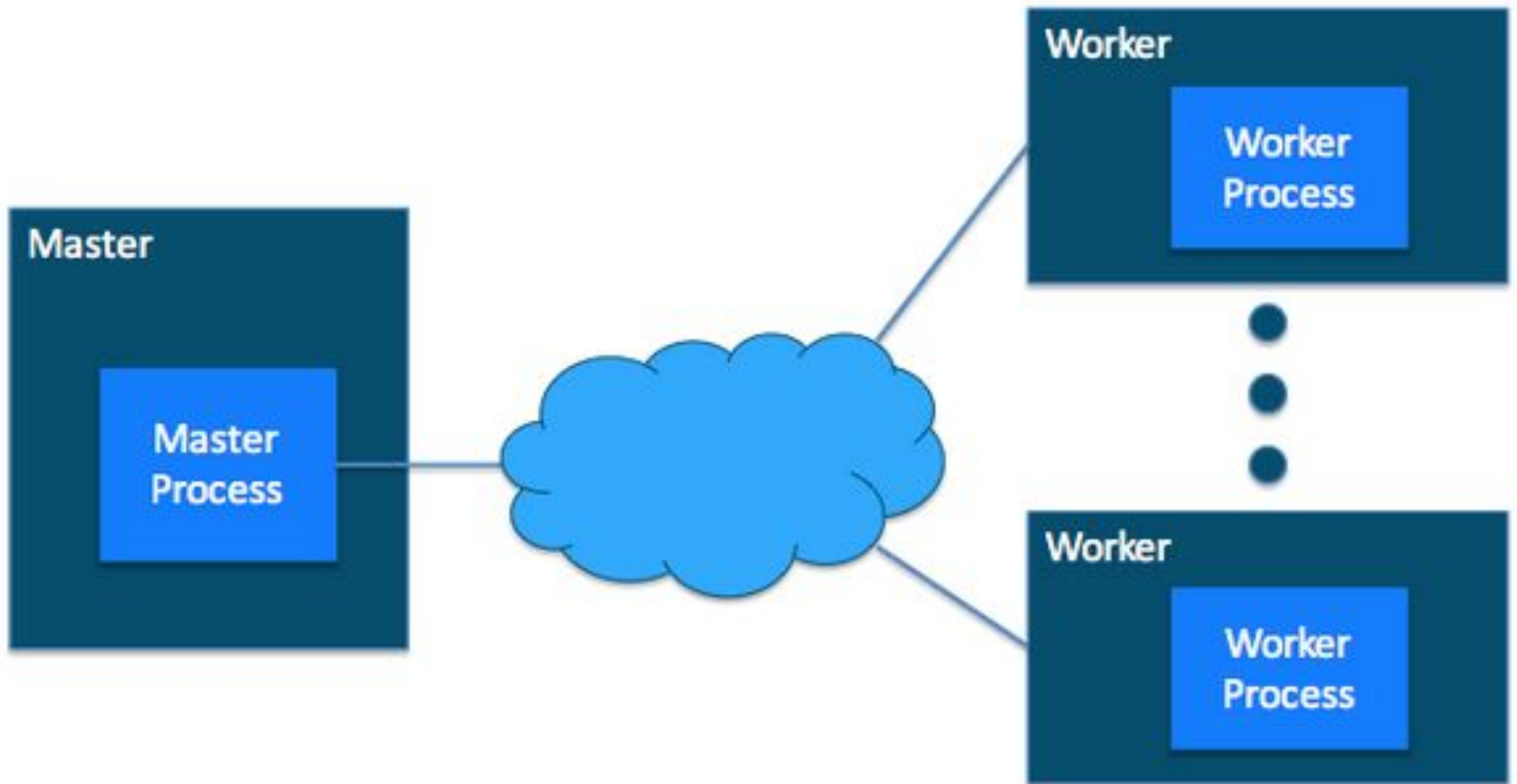| Batch | Script | SQL | NoSQL | Stream | Search | Others |
|-------|--------|-----|-------|--------|--------|--------|
| Map Reduce | Pig | Hive/Tez HCatalog | HBase Accumulo | Storm | Solr | In-Memory Analytics ISV Engines |

## YARN : Data Operating System

## HDFS
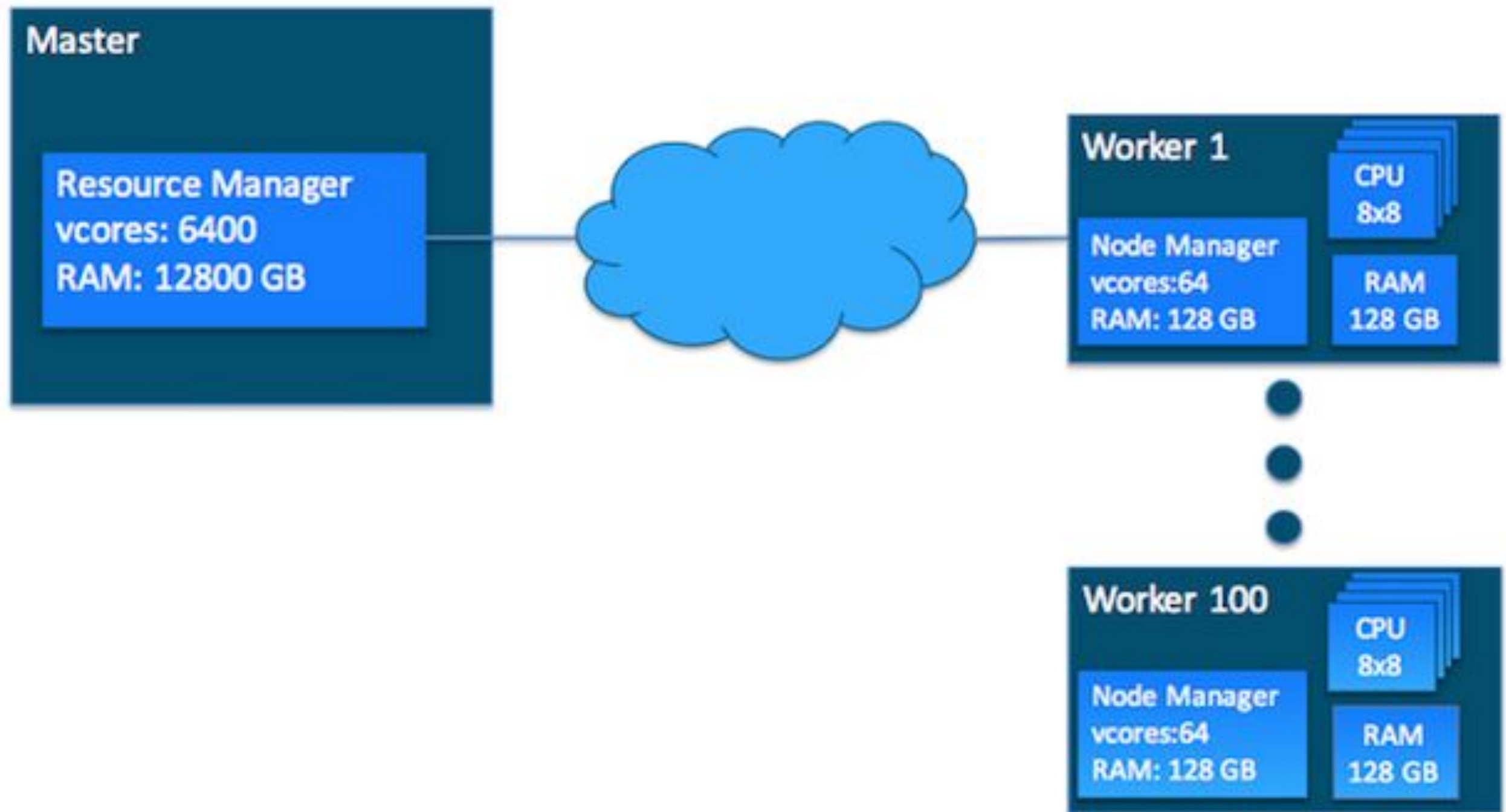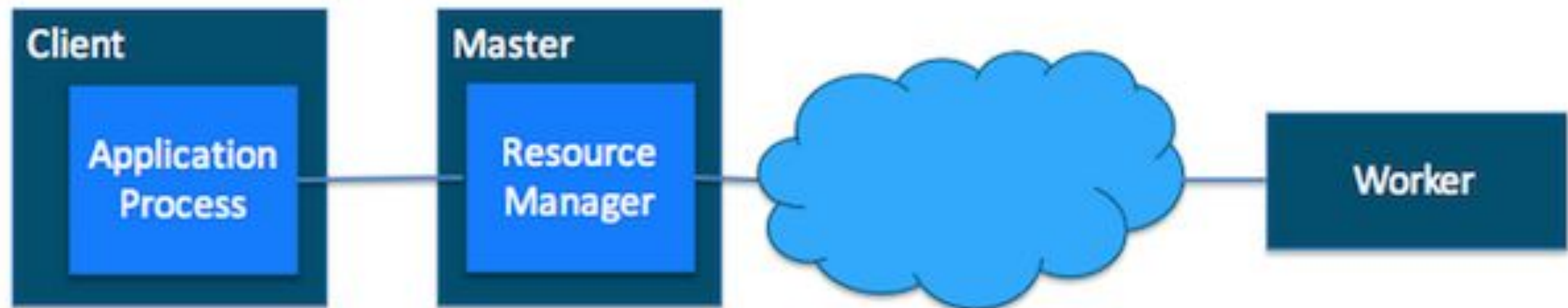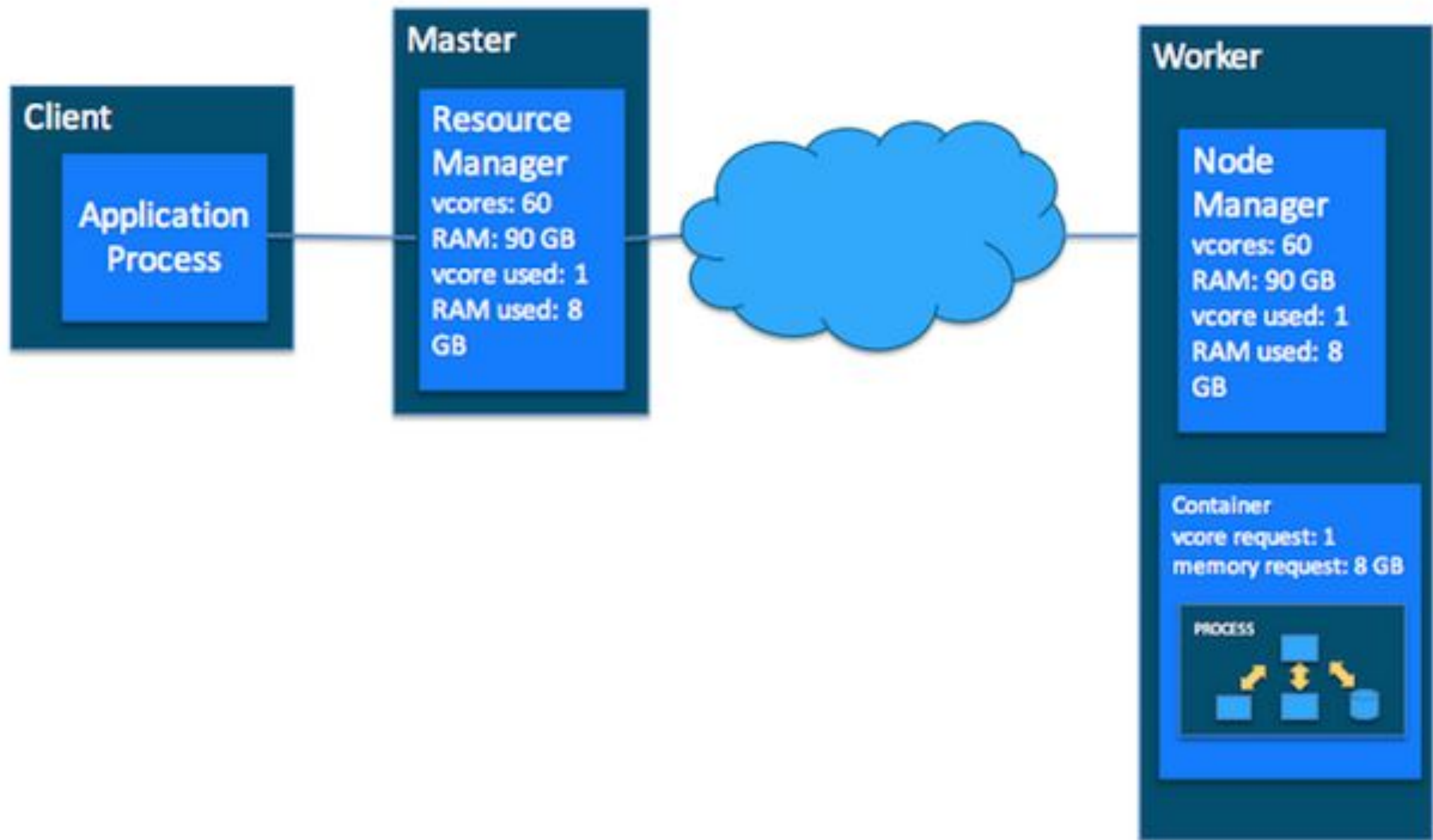### (Hadoop Distributed File System)

# DATA MANAGEMENT

**All Applications**

hadoop

- Cluster
  - About
  - Nodes
  - Applications
    - NEW
    - NEW_SAVING
    - SUBMITTED
    - ACCEPTED
    - RUNNING
    - FINISHED
    - FAILED
    - KILLED
  - Scheduler
- Tools

**Cluster Metrics**

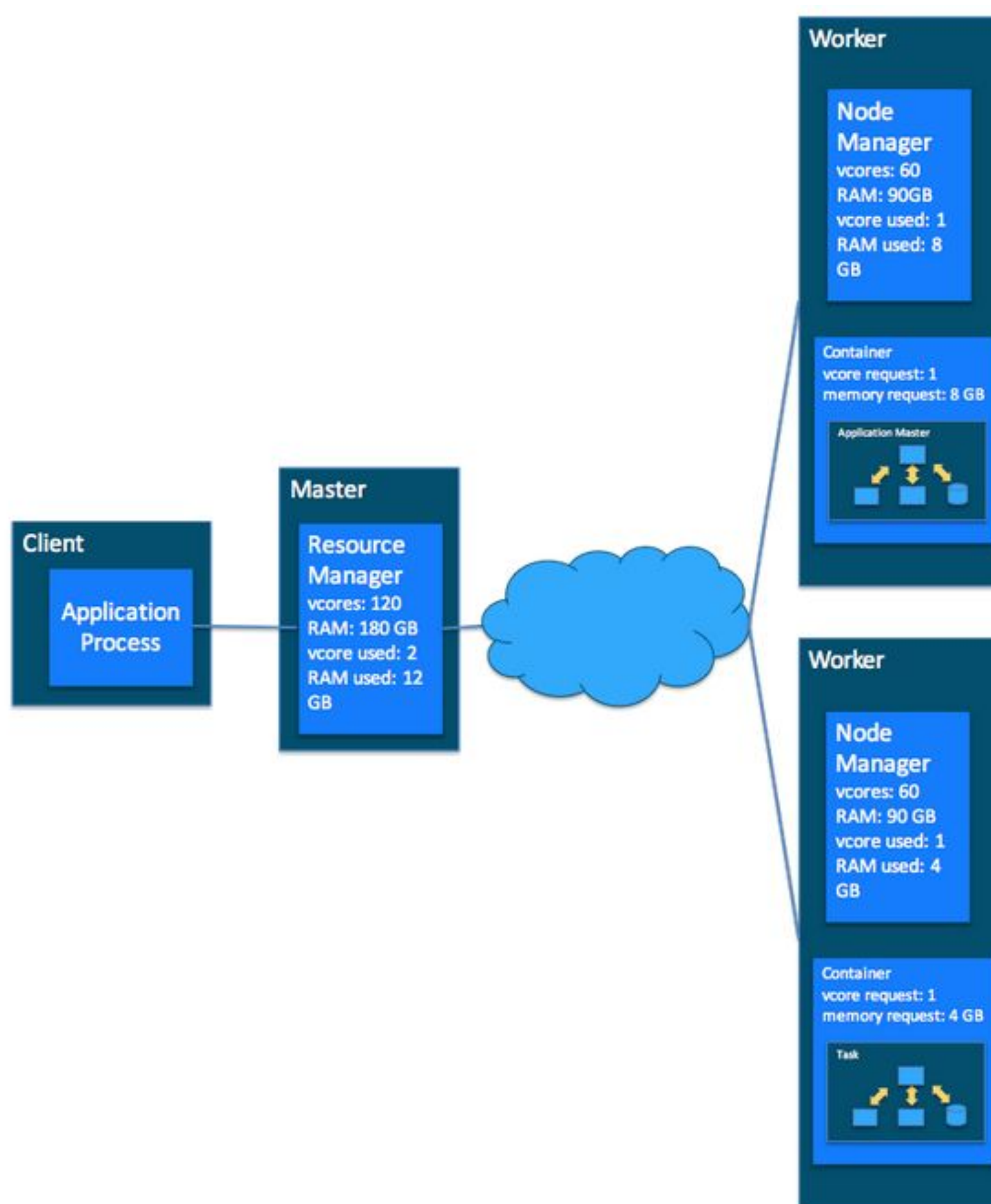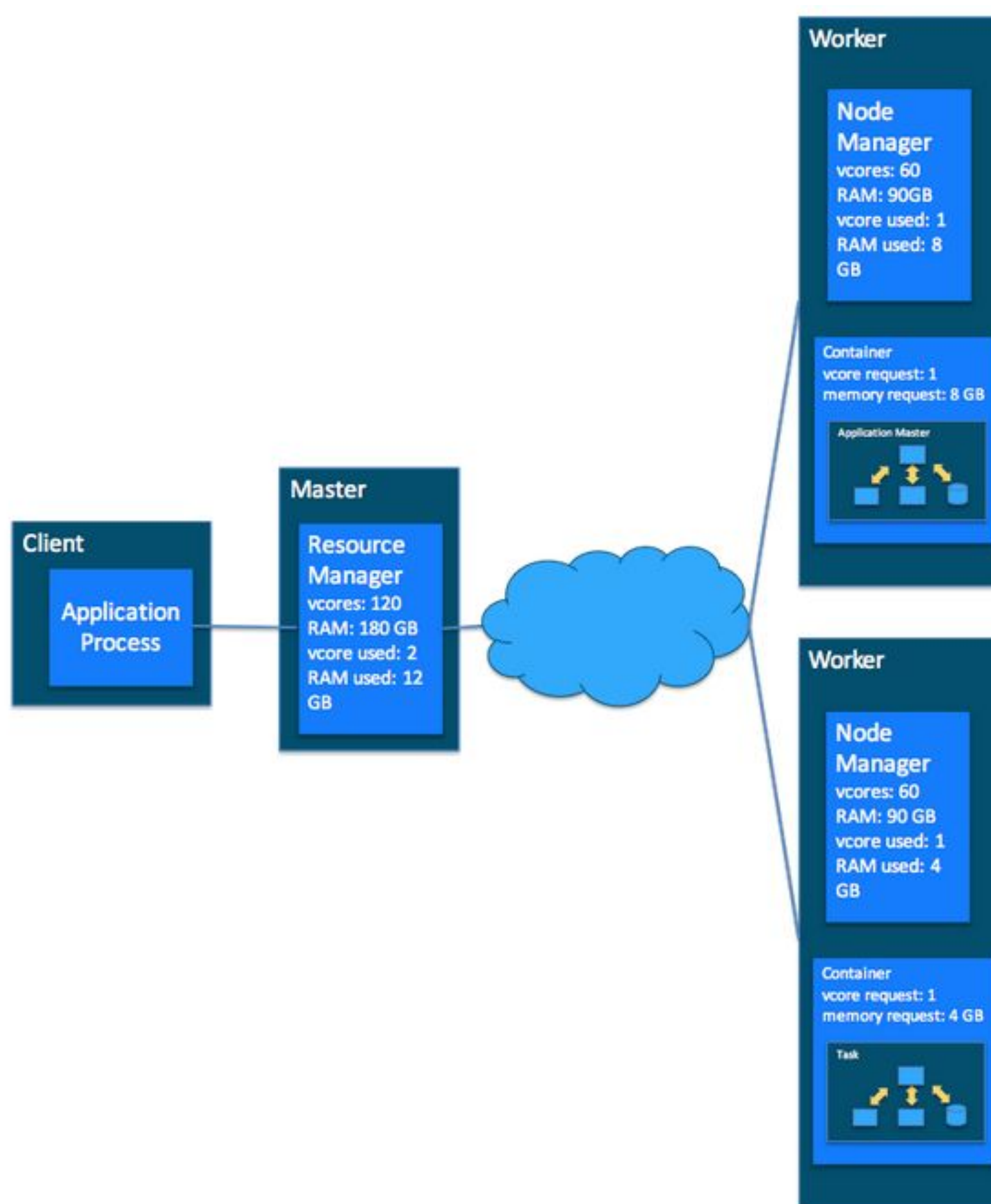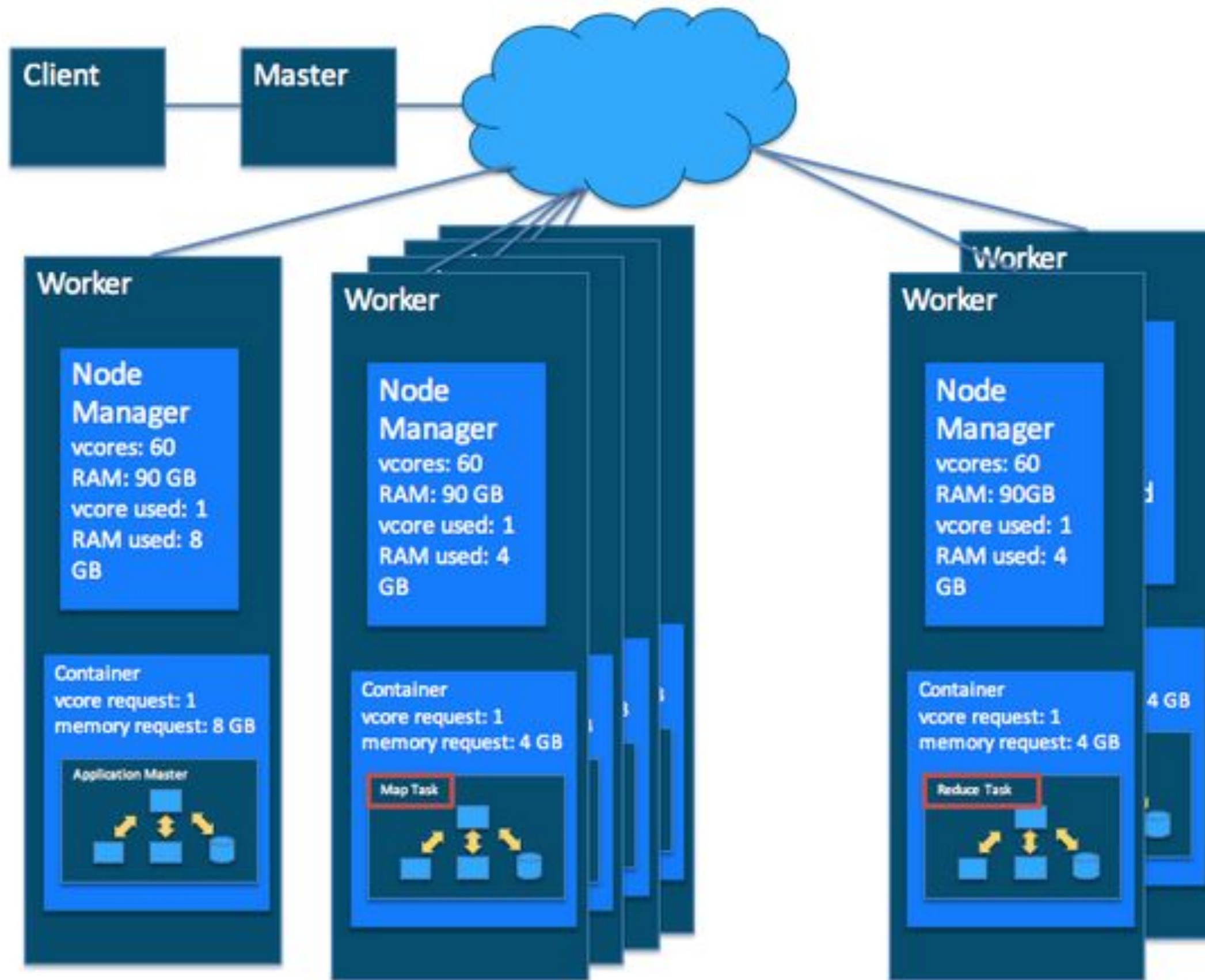| Apps Submitted | Apps Pending | Apps Running | Apps Completed | Containers Running | Memory Used | Memory Total | Memory Reserved | VCores Used | VCores Total | VCores Reserved | Active Nodes | Decommissioned Nodes | Lost Nodes | Unhealthy Nodes | Rebooted Nodes |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| 0 | 0 | 0 | 0 | 0 | 0 B | 4500 GB | 0 B | 0 | 3000 | 0 | 1 | 0 | 0 | 0 | 0 |

**User Metrics for dr.who**

| Apps Submitted | Apps Pending | Apps Running | Apps Completed | Containers Running | Containers Pending | Containers Reserved | Memory Used | Memory Pending | Memory Reserved | VCores Used | VCores Pending | VCores Reserved |
|---|---|---|---|---|---|---|---|---|---|---|---|---|
| 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 B | 0 B | 0 B | 0 | 0 | 0 |

Show 20 entries                                    Search:

Show 20 entries                                    Search:

| ID | User | Name | Application Type | Queue | StartTime | FinishTime | State | FinalStatus | Progress | Tracking UI |
|---|---|---|---|---|---|---|---|---|---|---|

No data available in table

Showing 0 to 0 of 0 entries                    First Previous Next Last

Created by Paint X

http://blog.cloudera.com/blog/2015/10/untangling-apache-hadoop-yarn-part-2/

**DICE** ANALYTICS

```xml
<?xml version="1.0"?>
<allocations>
 <queue name="marketing">
  <weight>30.0</weight>
 </queue>
 <queue name="sales">
  <weight>20.0</weight>
 </queue>
 <queue name="datascience">
  <weight>40.0</weight>
 </queue>

 <queue name="admin">
  <aclSubmitApps>fred,greg</aclSubmitApps>
  <weight>10.0</weight>
 </queue>
</allocations>
```

DICE
ANALYTICS

```xml
<?xml version="1.0"?>
<allocations>
 <queue name="marketing">
  <weight>3.0</weight>
  <queue name="reports">
   <weight>40.0</weight>
  </queue>
  <queue name="website">
   <weight>20.0</weight>
  </queue>
 </queue>
 <queue name="sales">
  <weight>4.0</weight>
  <queue name="northamerica">
   <weight>30.0</weight>
  </queue>
  <queue name="europe">
   <weight>30.0</weight>
  </queue>
 </queue>
 <queue name="datascience">
  <weight>13.0</weight>
  <queue name="short_jobs">
   <weight>100.0</weight>
  </queue>
  <queue name="best_effort_jobs">
   <weight>0.0</weight>
  </queue>
 </queue>
</allocations>
```
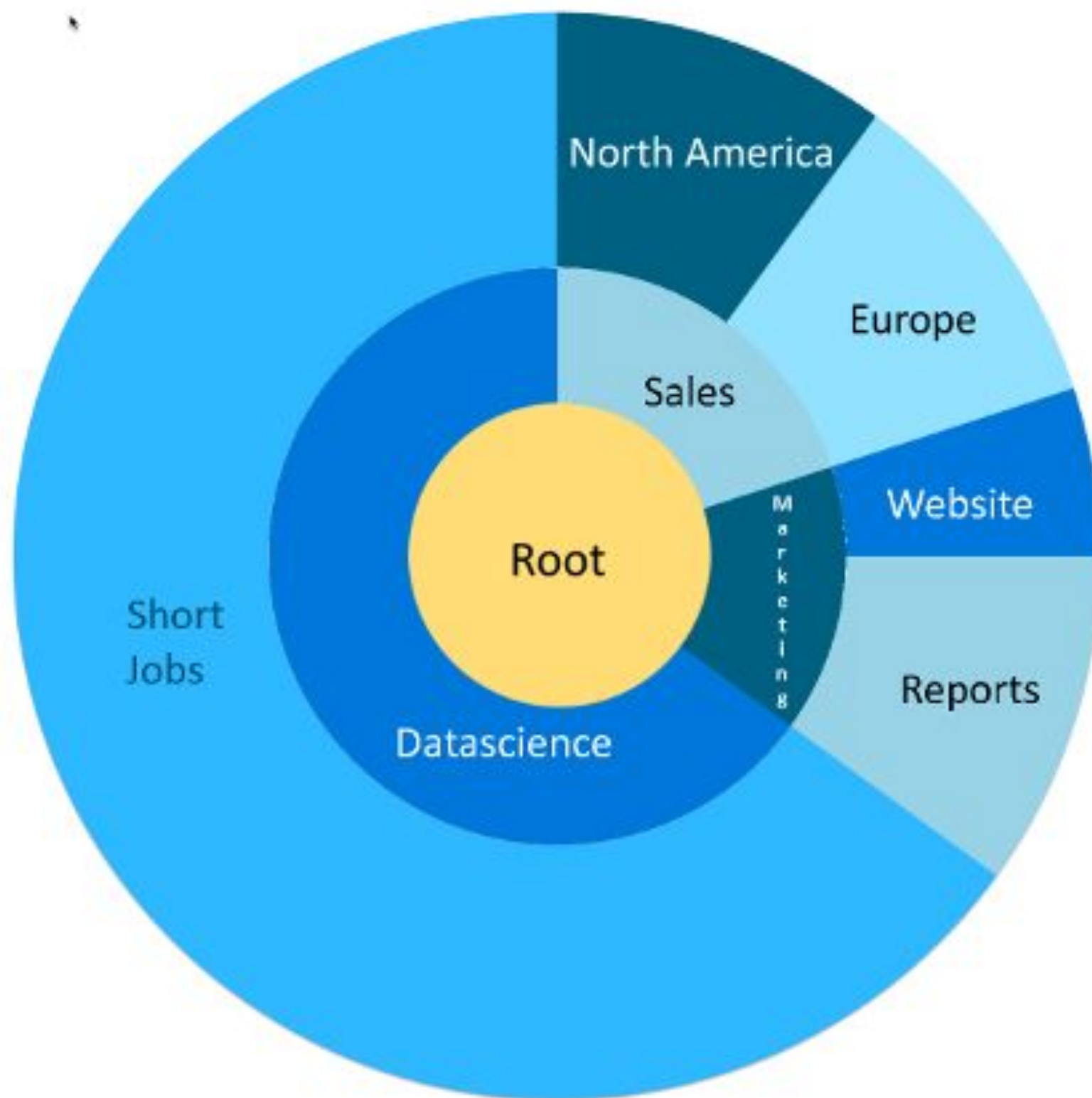
DICE
ANALYTICS

# Assignment

## Fair & Capacity Scheduler

## LIFO & FIFO Concept in YARN

**DICE**
ANALYTICS

# BRAD PITT

*Rolling Stone*
PETER TRAVERS

"ONE OF THE BEST
FILMS OF THE YEAR."

"BRAD PITT IS
SENSATIONAL." TIME
RICHARD CORLISS

# MONEYBALL

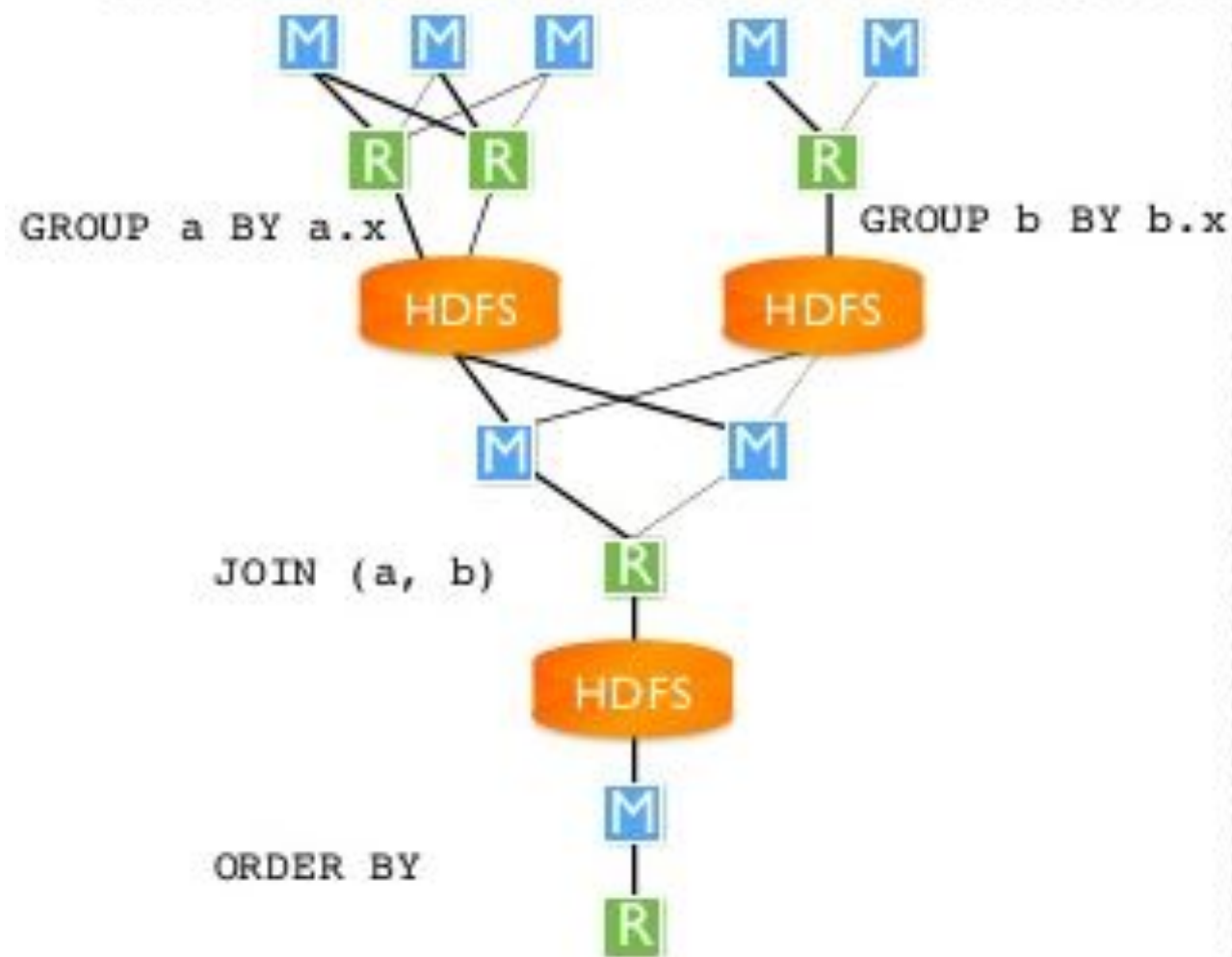## JONAH HILL   PHILIP SEYMOUR HOFFMAN
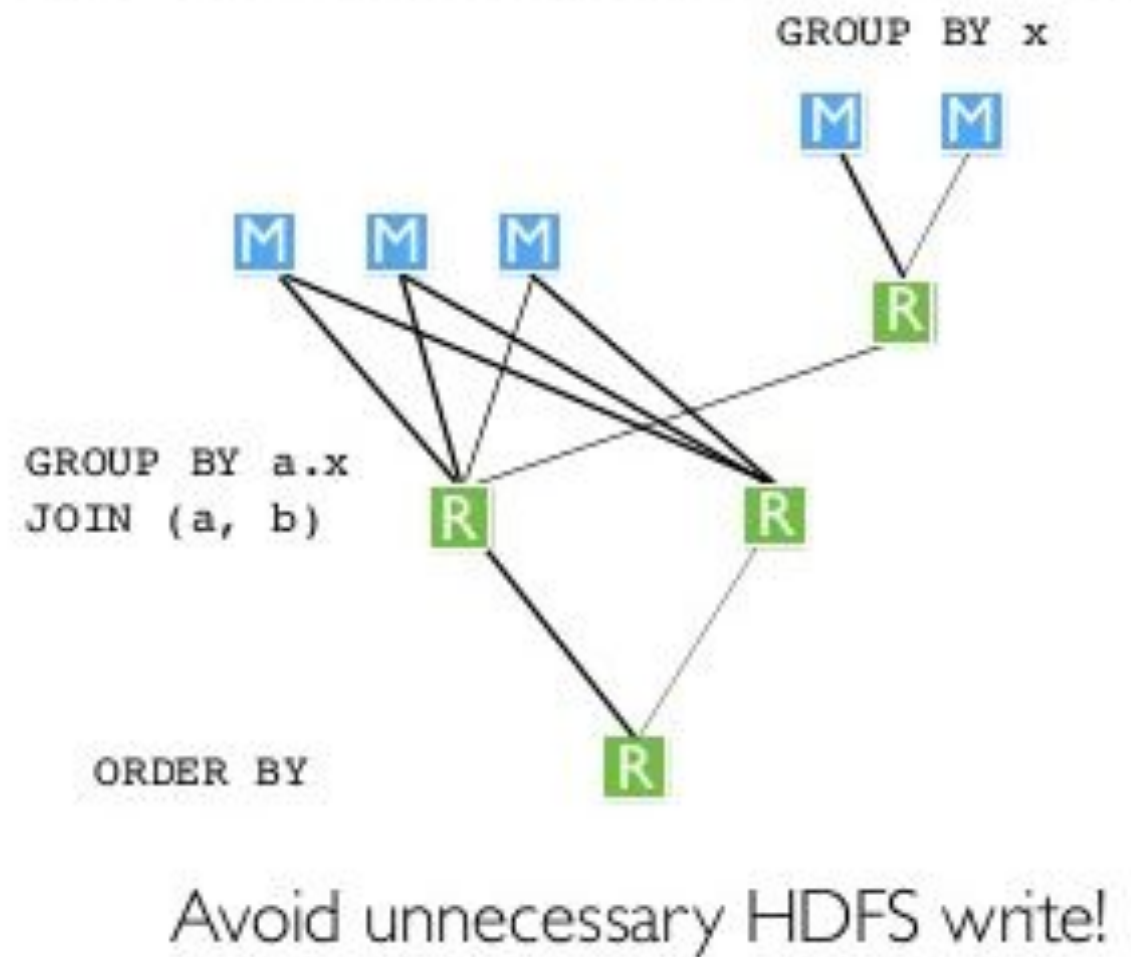
### BASED ON A TRUE STORY

IN THEATERS SEPTEMBER 23

# Tez



SELECT g1.x, g2.avg, g2.cnt
FROM (SELECT a.x AVERAGE(a.y) AS avg FROM a GROUP BY a.x) g1
JOIN (SELECT b.x, COUNT(b.y) AS avg FROM b GROUP BY b.x) g2
ON (g1.x = g2.x) ORDER BY avg;

**MapReduce**

**Tez**

GROUP a BY a.x

GROUP b BY b.x

JOIN (a, b)

ORDER BY

GROUP BY x

GROUP BY a.x
JOIN (a, b)

ORDER BY

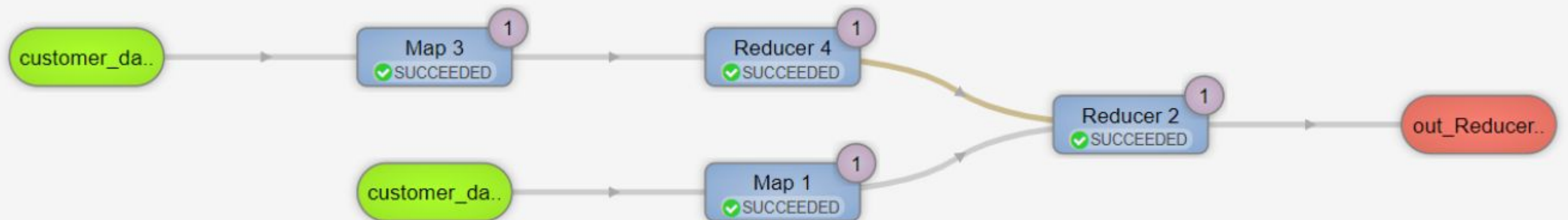Avoid unnecessary HDFS write!

```sql
SELECT
    *
FROM
    (SELECT PROVINCE,COUNT(*) AS CNT FROM customer_data GROUP BY PROVINCE ) AS A
LEFT JOIN
    (SELECT PROVINCE,COUNT(*) AS CNT FROM customer_data GROUP BY PROVINCE ) AS B
ON
A.PROVINCE = B.PROVINCE;
```



*Example by: Sheharyar Alam BD-03*

Select

a.name,count(*) from testdb.mytable

a,testdb.mytable b,testdb.mytable c
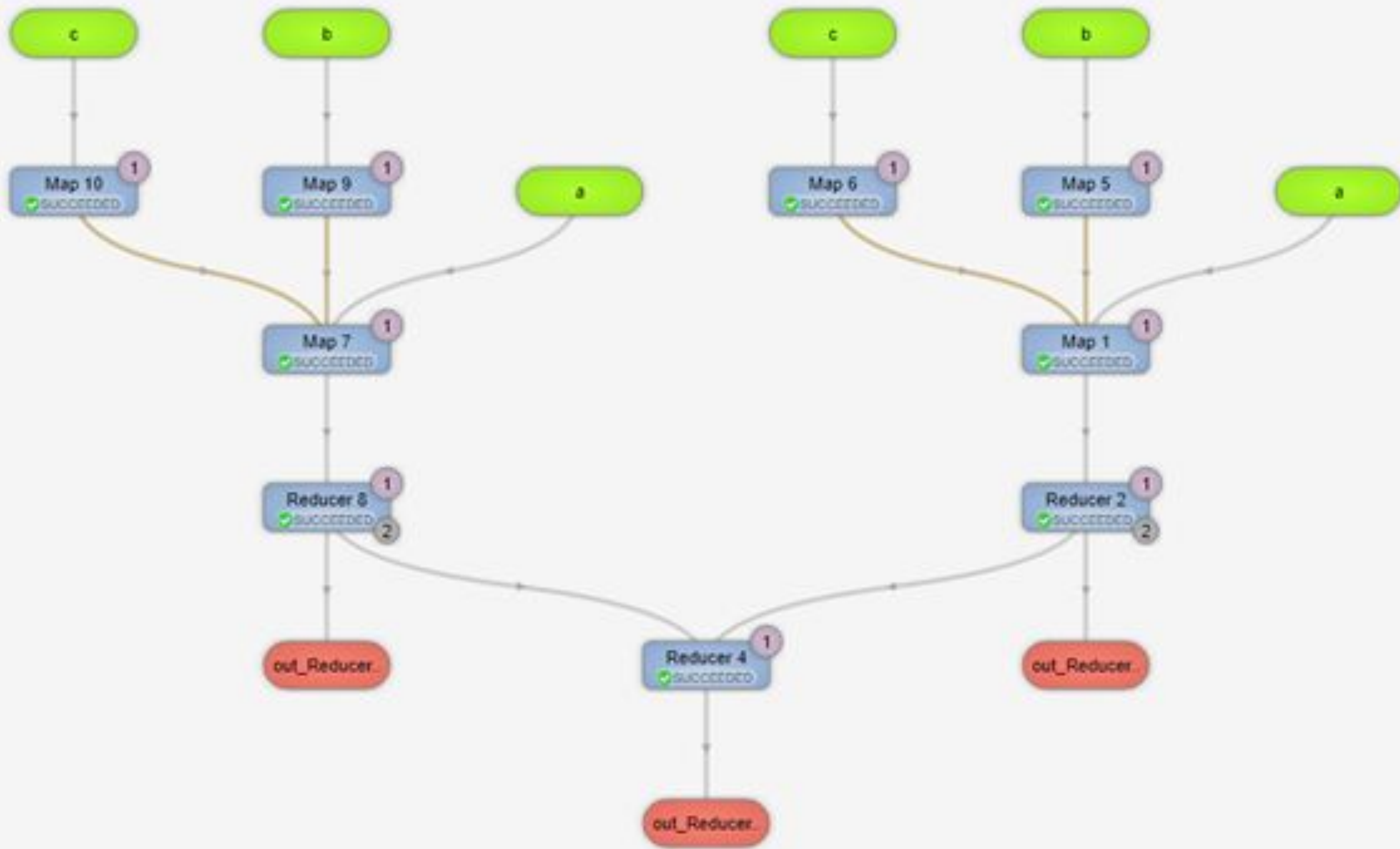
where a.name=b.name

and a.name=c.name

group by a.name

union

select

a.name,count(*) from testdb.mytable

a,testdb.mytable b,testdb.mytable c

where a.name=b.name

# *Hive Optimizations*

*http://chennaihug.org/knowledgebase/673/*

https://www.slideshare.net/Hadoop_Summit/w-235phall1pandey

DICE ANALYTICS