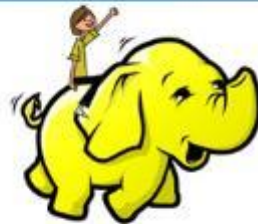
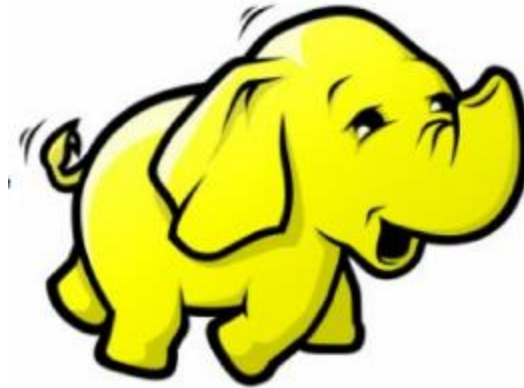


edureka!

Hadoop Administration



Hadoop Administration



Module 6: Advanced Topics: QJM, HDFS Federation and Security

✓ **Module 1**

- ✓ Understanding Big Data
- ✓ Hadoop Components

✓ **Module 2**

- ✓ Different Hadoop Server Roles
- ✓ Hadoop Cluster Configuration

✓ **Module 3**

- ✓ Hadoop Cluster Planning
- ✓ Job Scheduling

✓ **Module 4**

- ✓ Securing your Hadoop Cluster
- ✓ Backup and Recovery

✓ **Module 5**

- ✓ Hadoop 2.0 New Features
- ✓ HDFS High Availability

✓ **Module 6**

- ✓ **Quorum Journal Manager (QJM)**
- ✓ **Hadoop 2.0 - YARN**

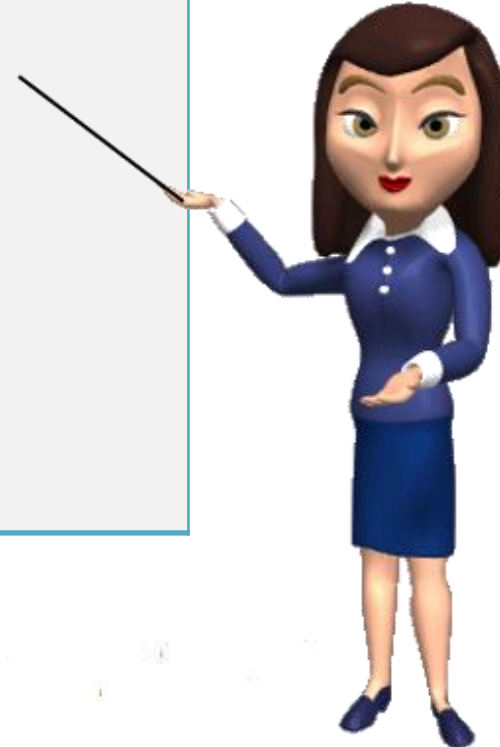
✓ **Module 7**

- ✓ Oozie Workflow Scheduler
- ✓ Hive and Hbase Administration

✓ **Module 8**

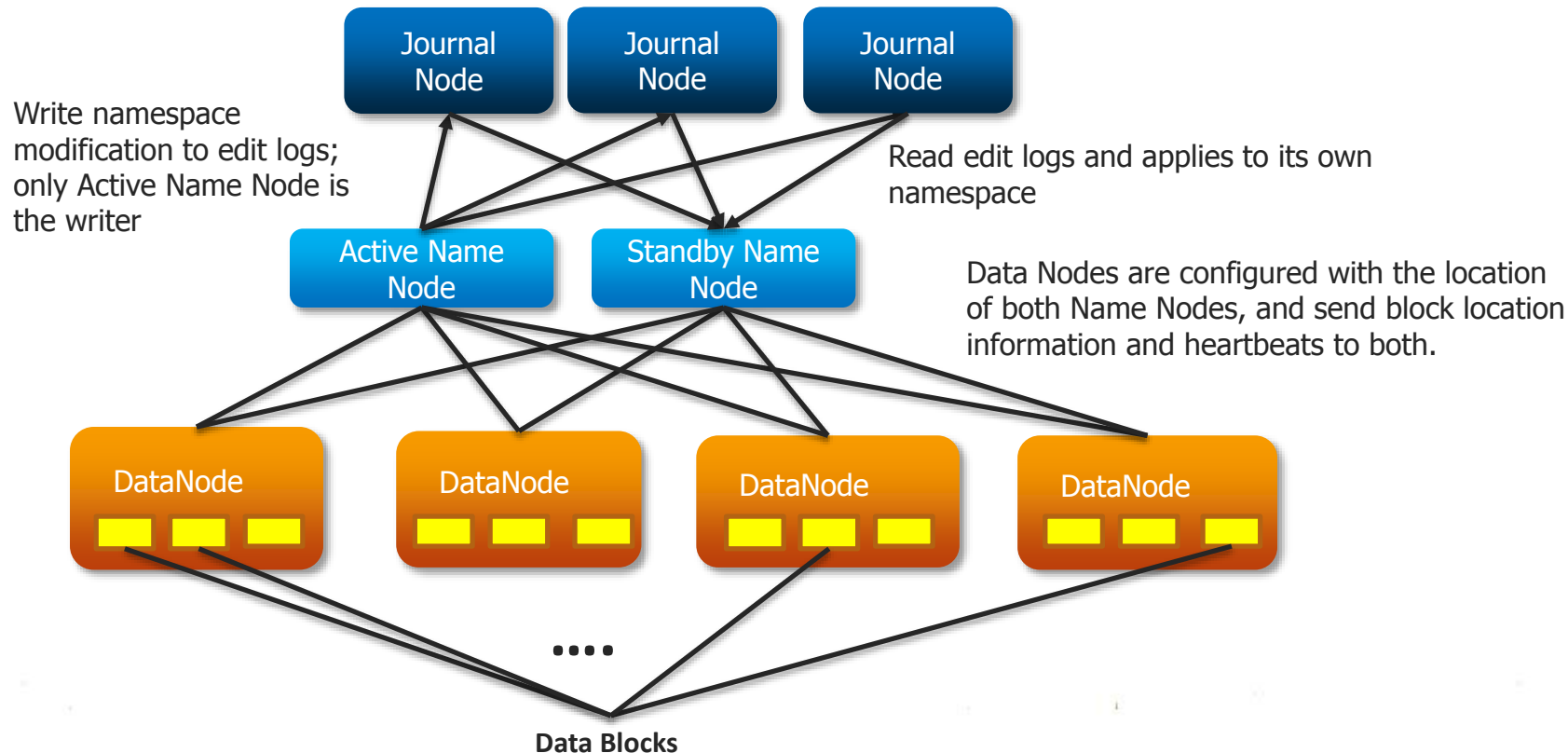
- ✓ Hadoop Cluster Case Study
- ✓ Hadoop Implementation

- 📄 **Hadoop 2.0 New Features**
 - 📄 **HDFS High Availability - QJM**
 - 📄 **YARN and MRv2**
- 📄 **YARN and Hadoop ecosystem**
- 📄 **YARN Components**
- 📄 **Job Tracker and Job Submission**
- 📄 **MR Application Execution in YARN**

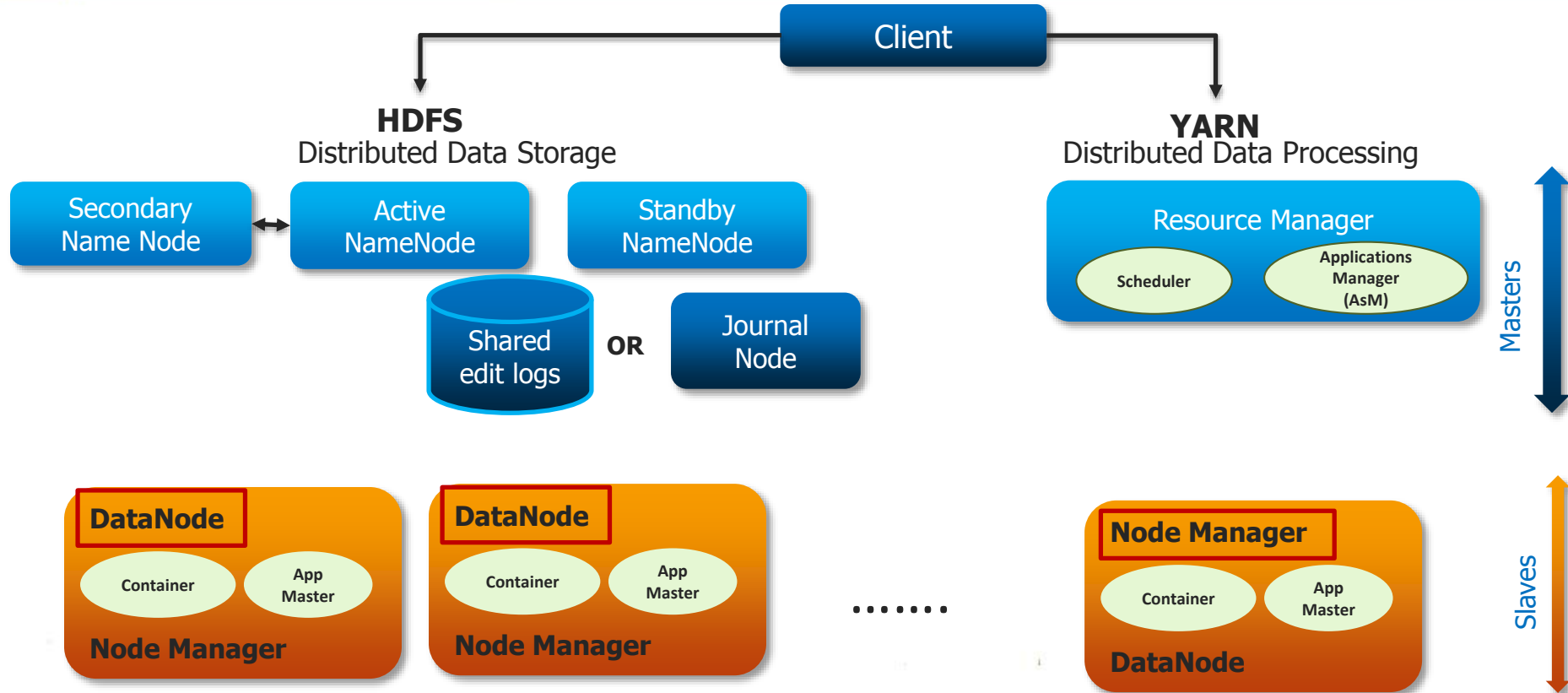


- ✓ Hadoop 1.0 Vs. Hadoop 2.0
- ✓ Hadoop 2.0 Cluster daemons





Hadoop 2.0 Server Roles

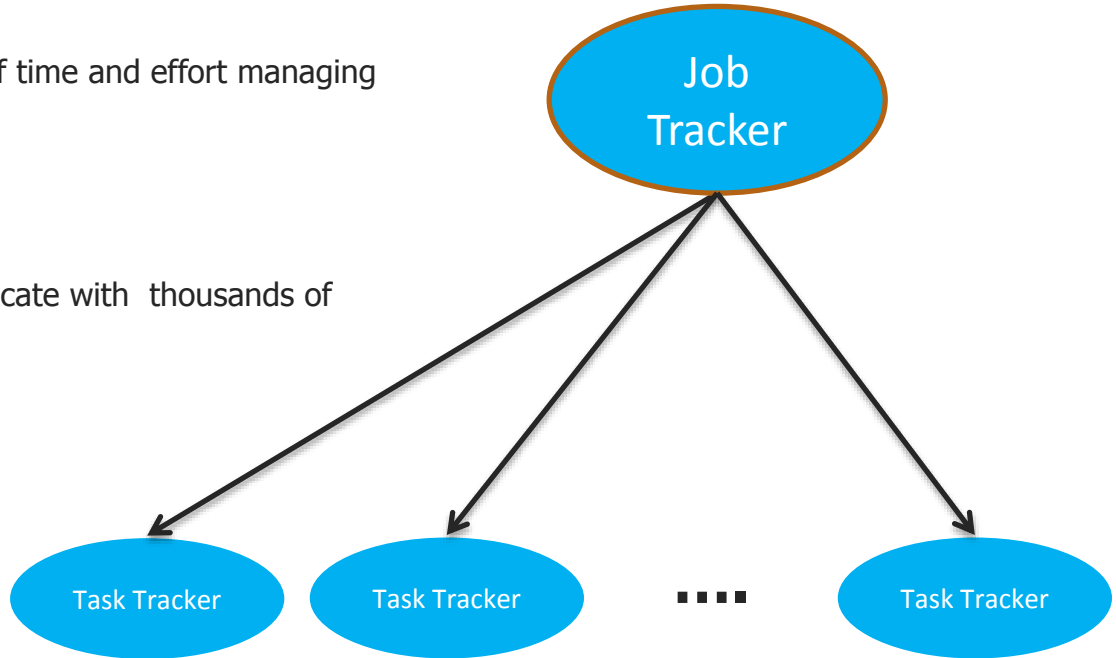


CPU

- ✓ Spends a very significant portion of time and effort managing the life cycle of applications

Network

- ✓ Single Listener Thread to communicate with thousands of Map and Reduce Jobs



As the cluster size grow and reaches to 4000 Nodes

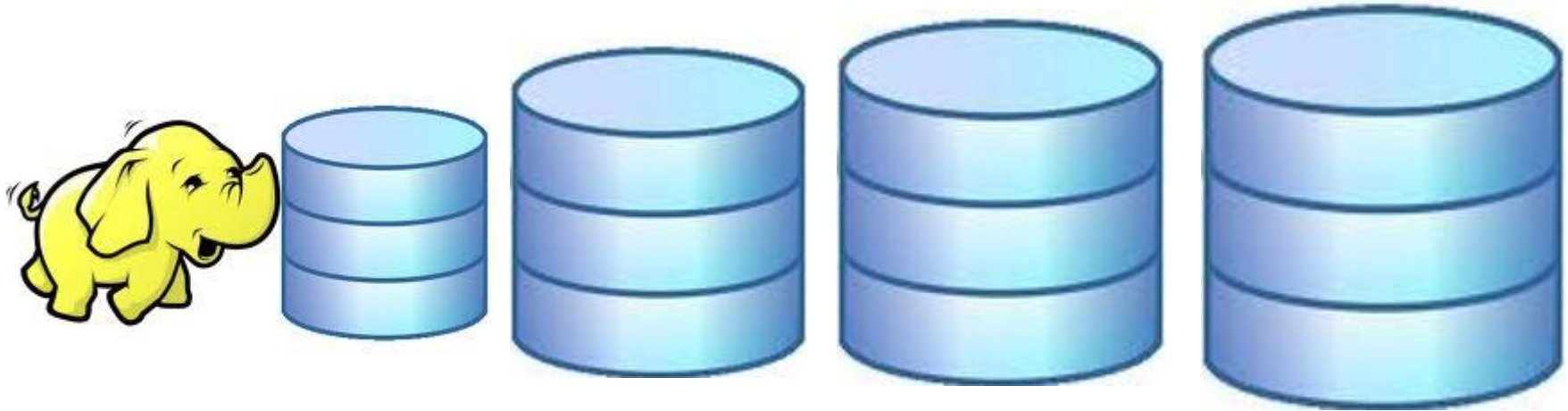
✓ Cascading Failures

- ✓ The DataNode failures results in a serious deterioration of the overall cluster performance because of attempts to replicate data and overload live nodes, through network flooding.

✓ Multi-tenancy

- ✓ As clusters increase in size, you may want to employ these clusters for a variety of models. MRv1 dedicates its nodes to Hadoop and cannot be re-purposed for other applications and workloads in an Organization. With the growing popularity and adoption of cloud computing among enterprises, this becomes more important.

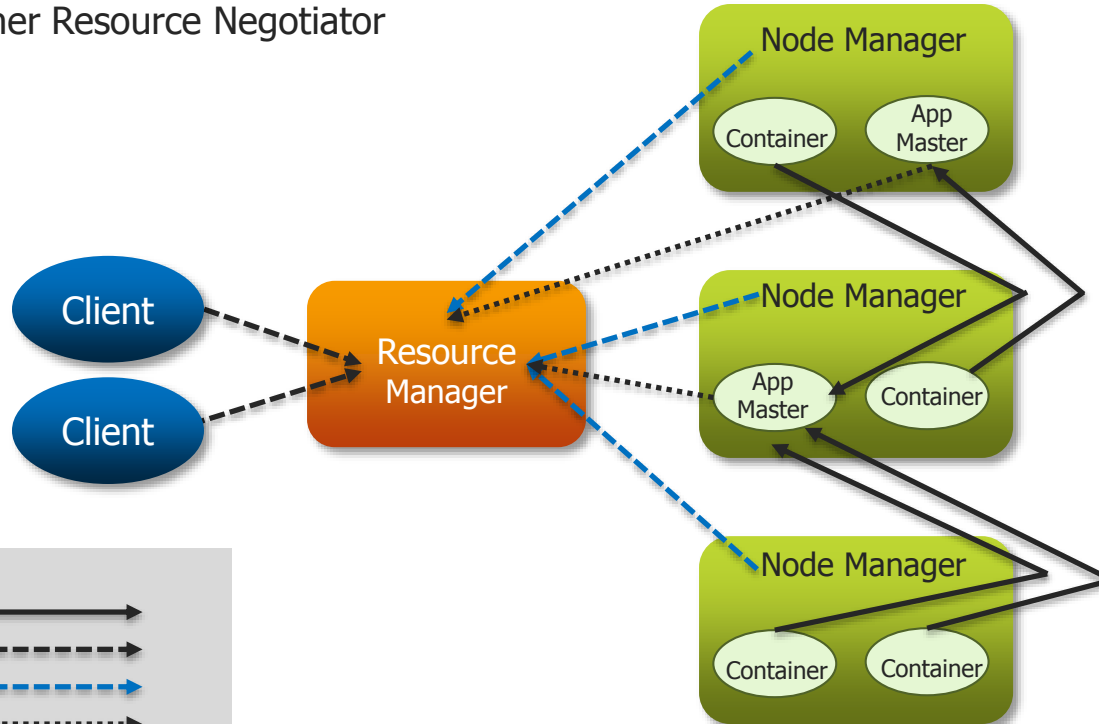




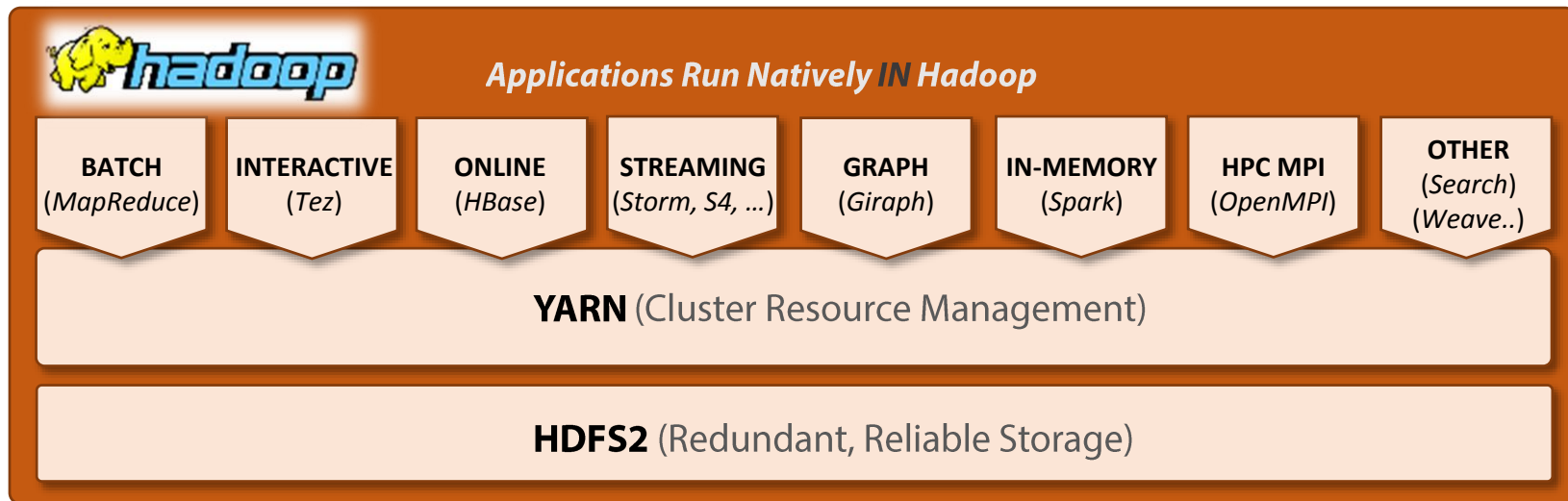
- ✓ Terabytes and Petabytes of data in HDFS can be used only for MapReduce processing

Hadoop 2.0 : YARN

YARN = Yet Another Resource Negotiator



MapReduce Status →
Job Submission →
Node Status →
Resource Request →



<http://hadoop.apache.org/docs/stable2/hadoop-yarn/hadoop-yarn-site/YARN.html>

- ✓ Organizes jobs into queues
- ✓ Queue shares as %'s of cluster
- ✓ FIFO scheduling within each queue
- ✓ Data locality-aware Scheduling

- ✓ **Hierarchical Queues**

To manage the resource within an organization.

- ✓ **Capacity Guarantees**

A fraction to the total available capacity allocated to each Queue.

- ✓ **Security**

To safeguard applications from other users.

- ✓ **Elasticity**

Resources are available in a predictable and elastic manner to queues.

- ✓ **Multi-tenancy**

Set of limit to prevent over-utilization of resources by a single application.

- ✓ **Operability**

Runtime configuration of Queues.

- ✓ **Resource-based scheduling**

If needed, Applications can request more resources than the default.

- ✓ Edit **yarn-site.xml** to enable the Capacity Scheduler

Property	Value
yarn.resourcemanager.scheduler.class	org.apache.hadoop.yarn.server.resourcemanager.scheduler.capacity.CapacityScheduler

- ✓ Configure **'queues'** in capacity-scheduler.xml

The Resource Manager has two main components

- a) NameNode and SNN
- b) Scheduler and Applications Manager
- c) Manager and Application manager4



Answer: Scheduler and Applications Manager



YARN enables a user to interact with all data in multiple ways simultaneously, making Hadoop a true multi-use data platform and allowing it to take its place in a modern data architecture.

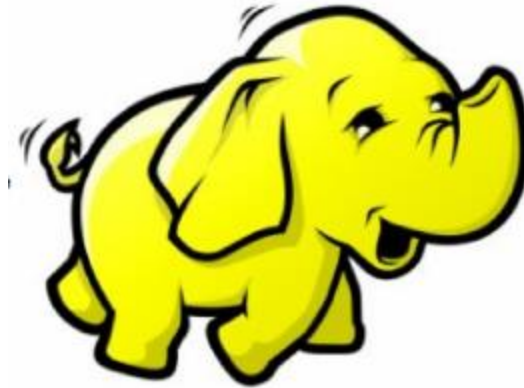
- a) True
- b) False



Answer: True

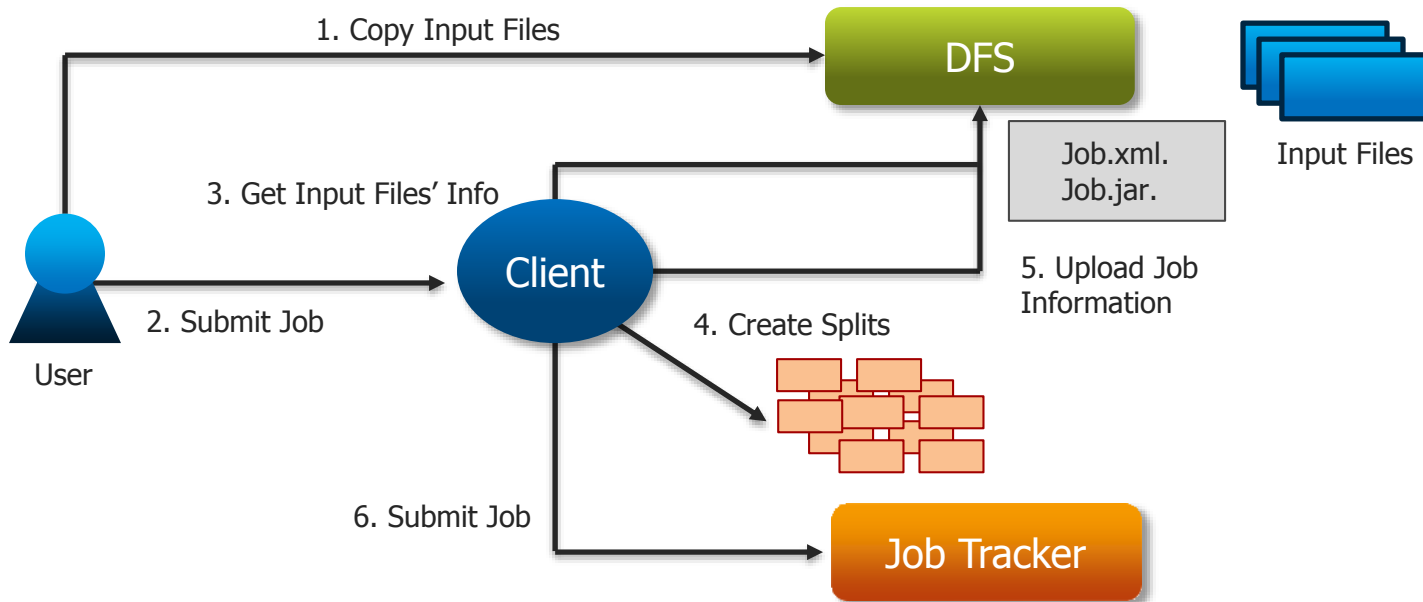


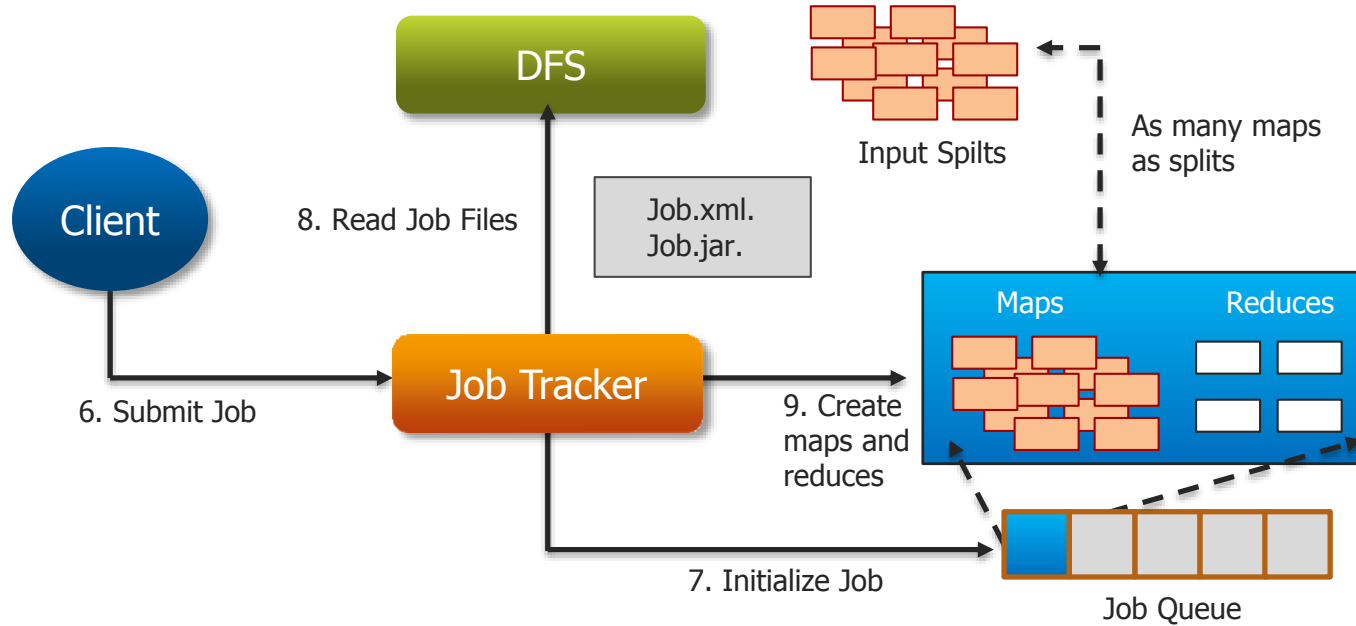
Hadoop Administration

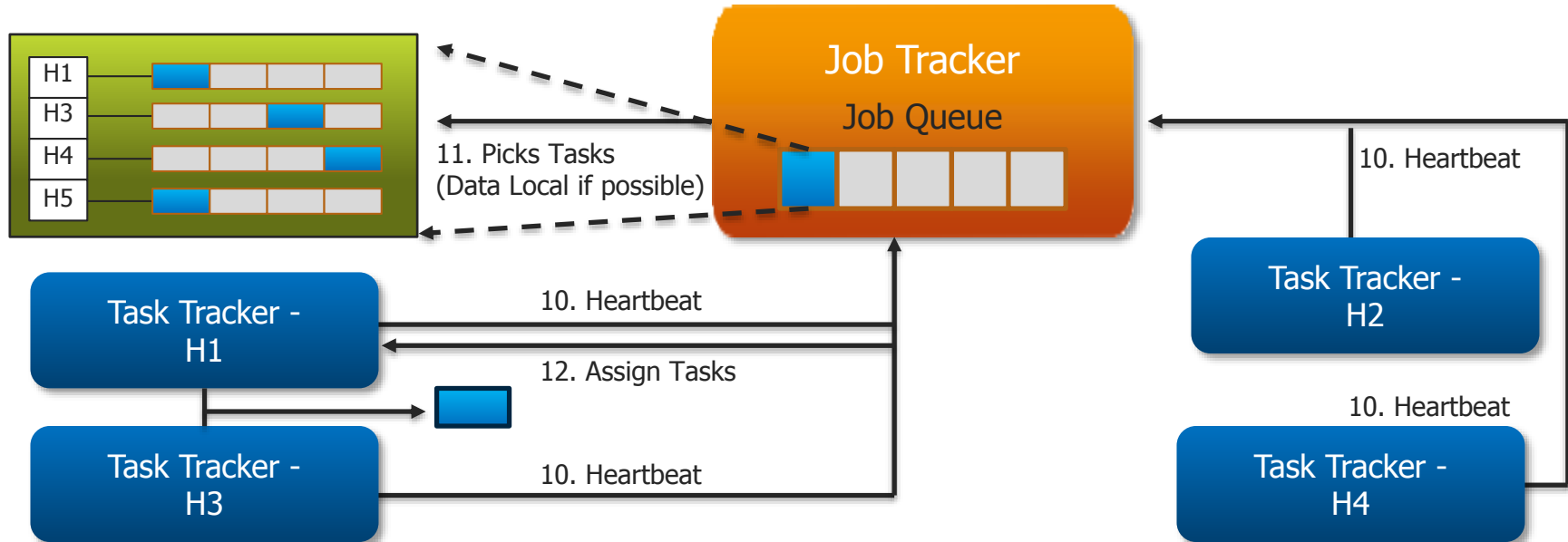


Executing MapReduce Application on YARN

Job Execution







YARN – Beyond MapReduce

YARN = Yet Another Resource Negotiator

✓ **Resource Manager**

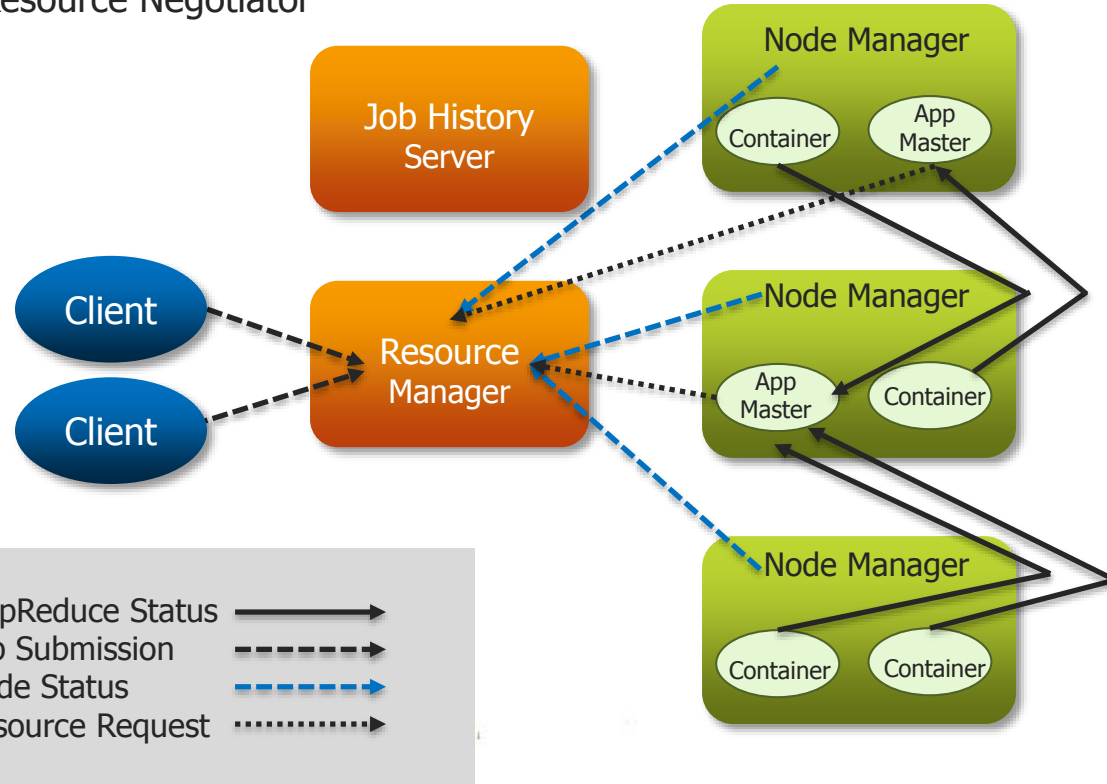
- ✓ Cluster Level resource manager
- ✓ Long Life, High Quality Hardware

✓ **Node Manager**

- ✓ One per Data Node
- ✓ Monitors resources on Data Node

✓ **Application Master**

- ✓ One per application
- ✓ Short life
- ✓ Manages task /scheduling



Application Vs. Job Execution

- ✓ MapReduce **Application** Execution
 - ✓ Submission
 - ✓ Initialization
 - ✓ Tasks Assignment
 - ✓ Tasks' Memory
 - ✓ Status Updates
 - ✓ Failure Recovery

✓ **Client**

- ✓ Submit a MapReduce Application

✓ **Resource Manager**

- ✓ Manage the resource utilization across
- ✓ Hadoop Cluster

✓ **Node Manager**

- ✓ Runs on each Data Node
- ✓ Creates execution container
- ✓ Monitors Container's usage

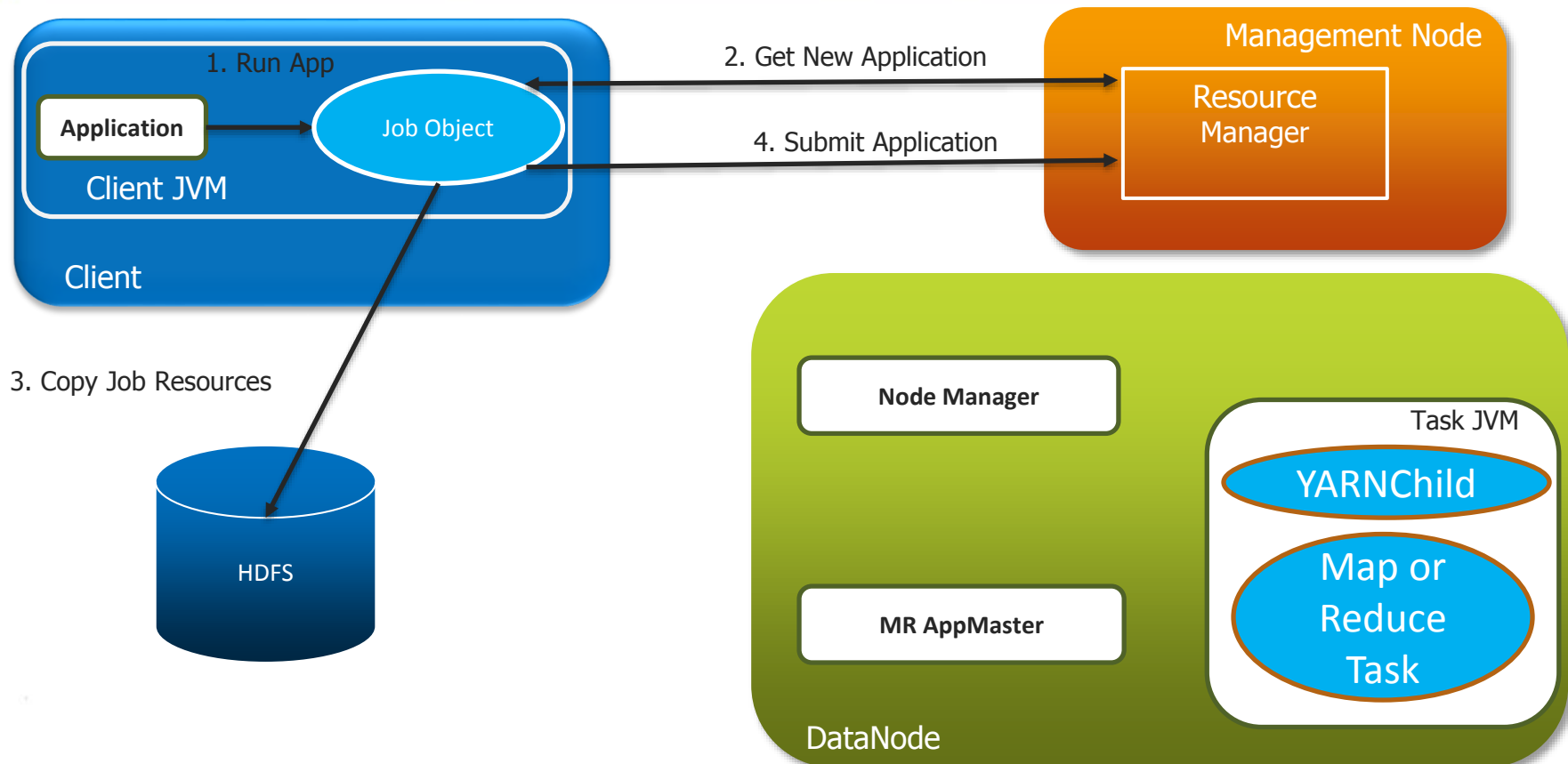
✓ **MapReduce Application Master**

- ✓ Coordinates and Manages MapReduce tasks
- ✓ Negotiates with Resource Manager to schedule asks
- ✓ The tasks are started by Node Manager(s)

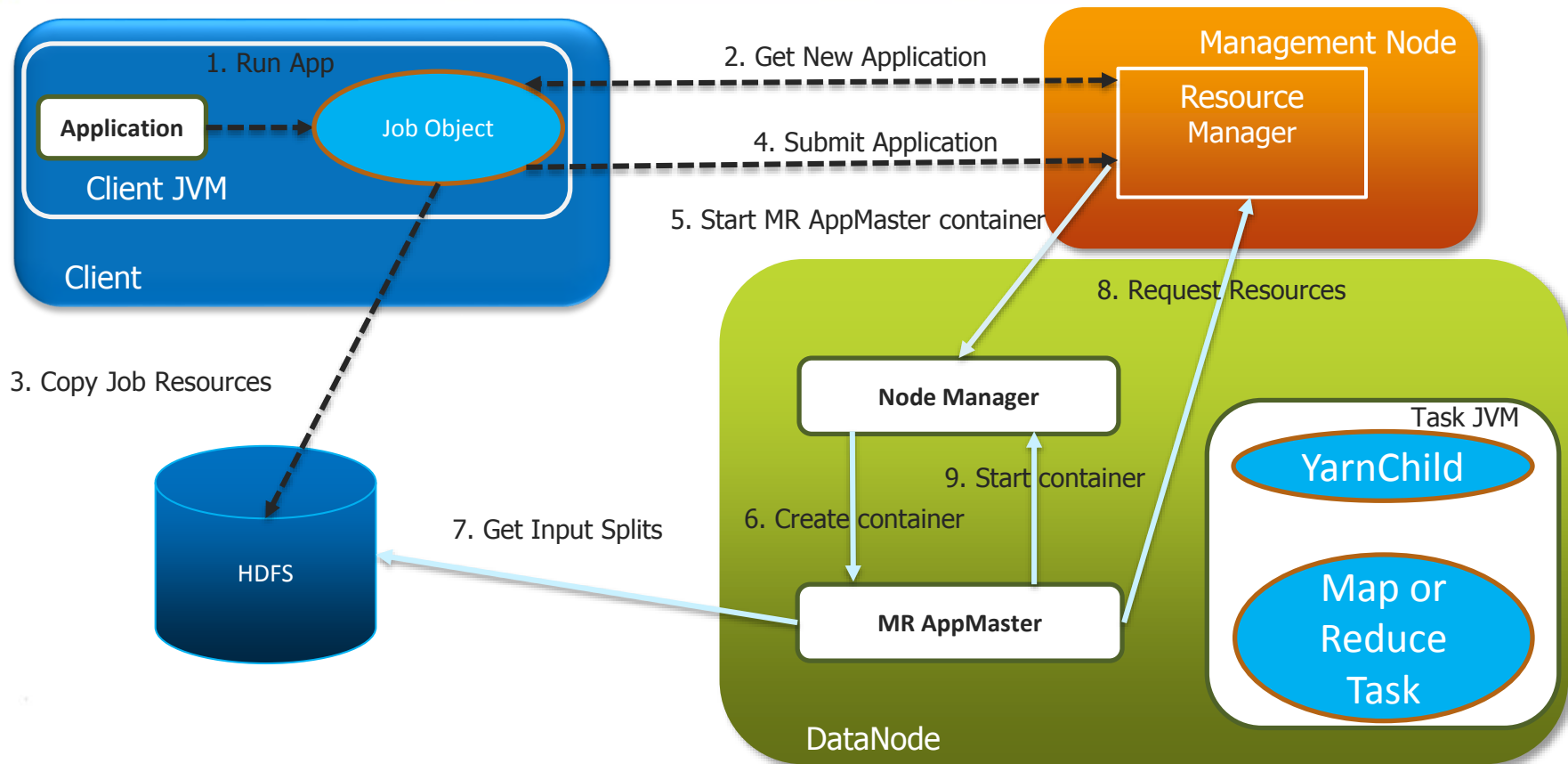
✓ **HDFS**

- ✓ shares resources and task's artefacts among YARN components

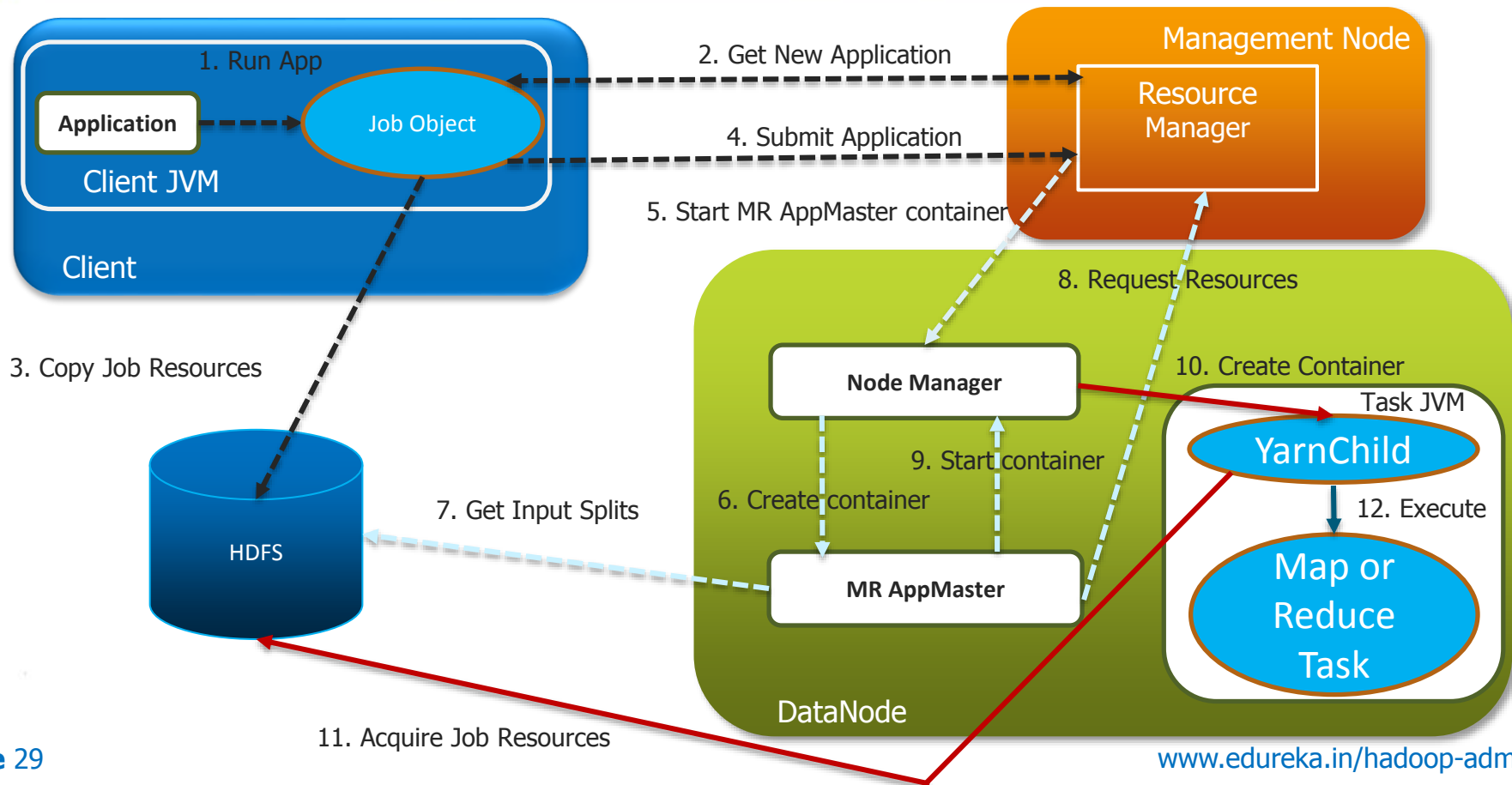
YARN MR Application Execution Flow



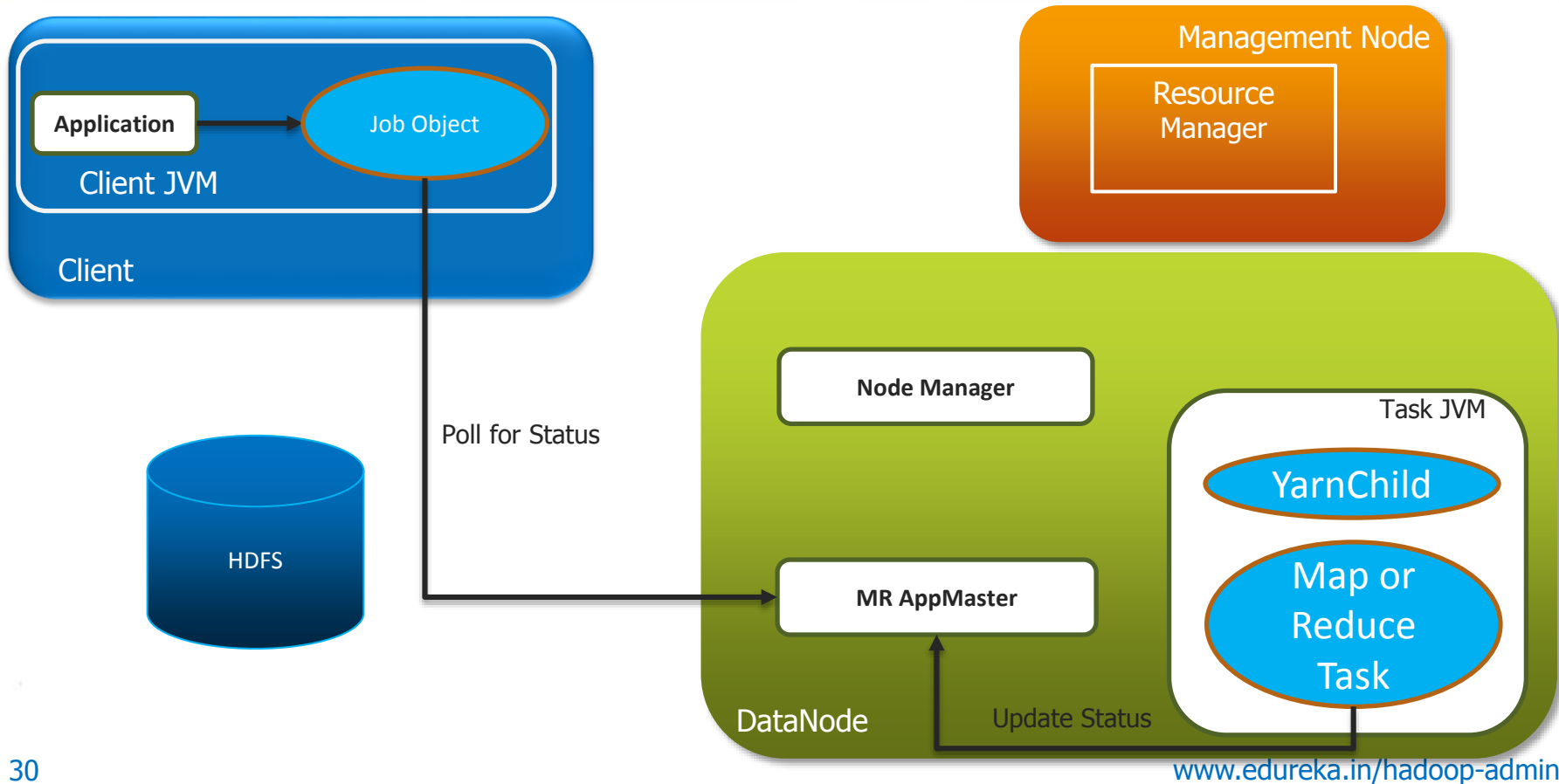
YARN MR Application Execution Flow

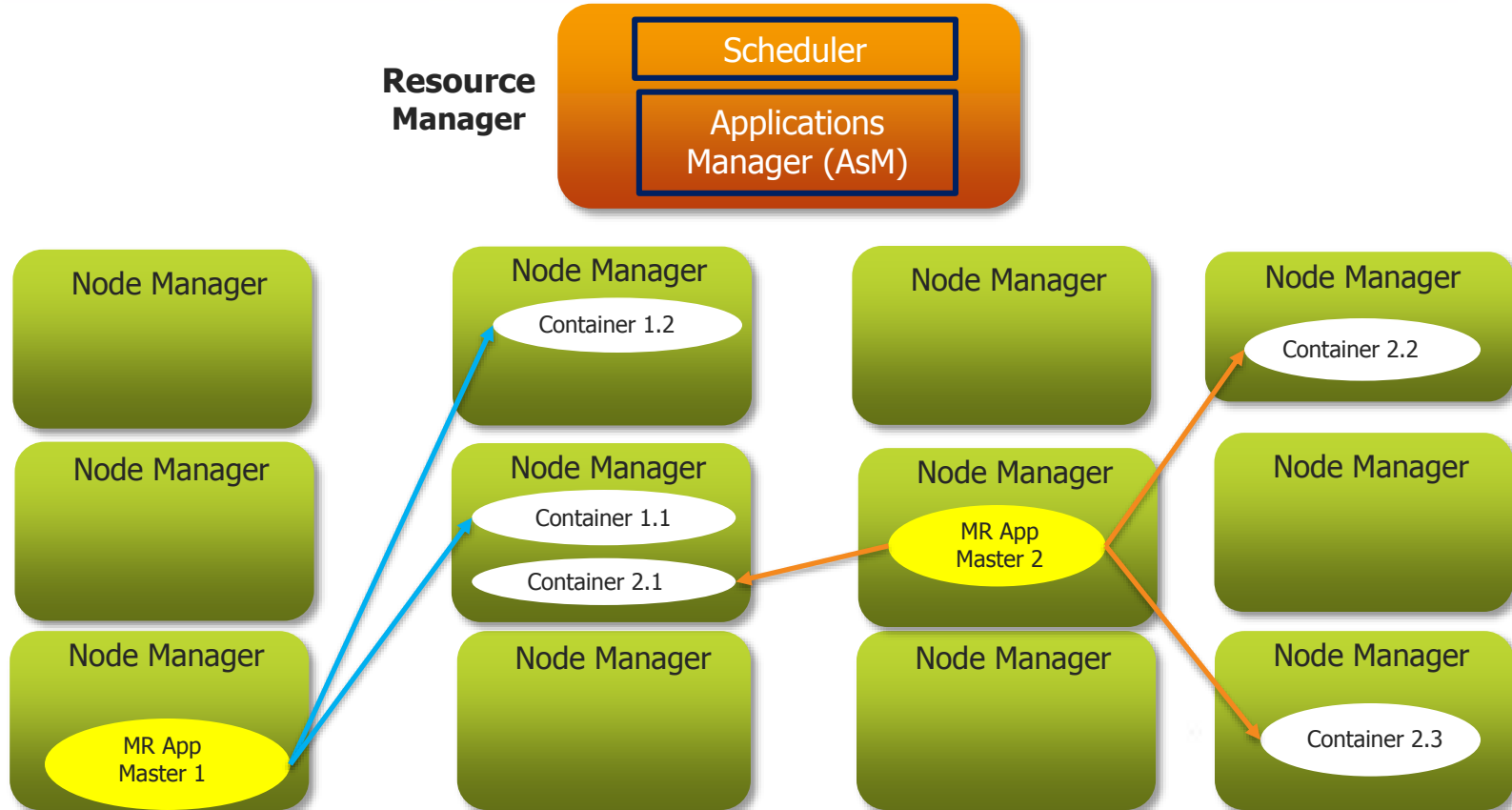


YARN MR Application Execution Flow

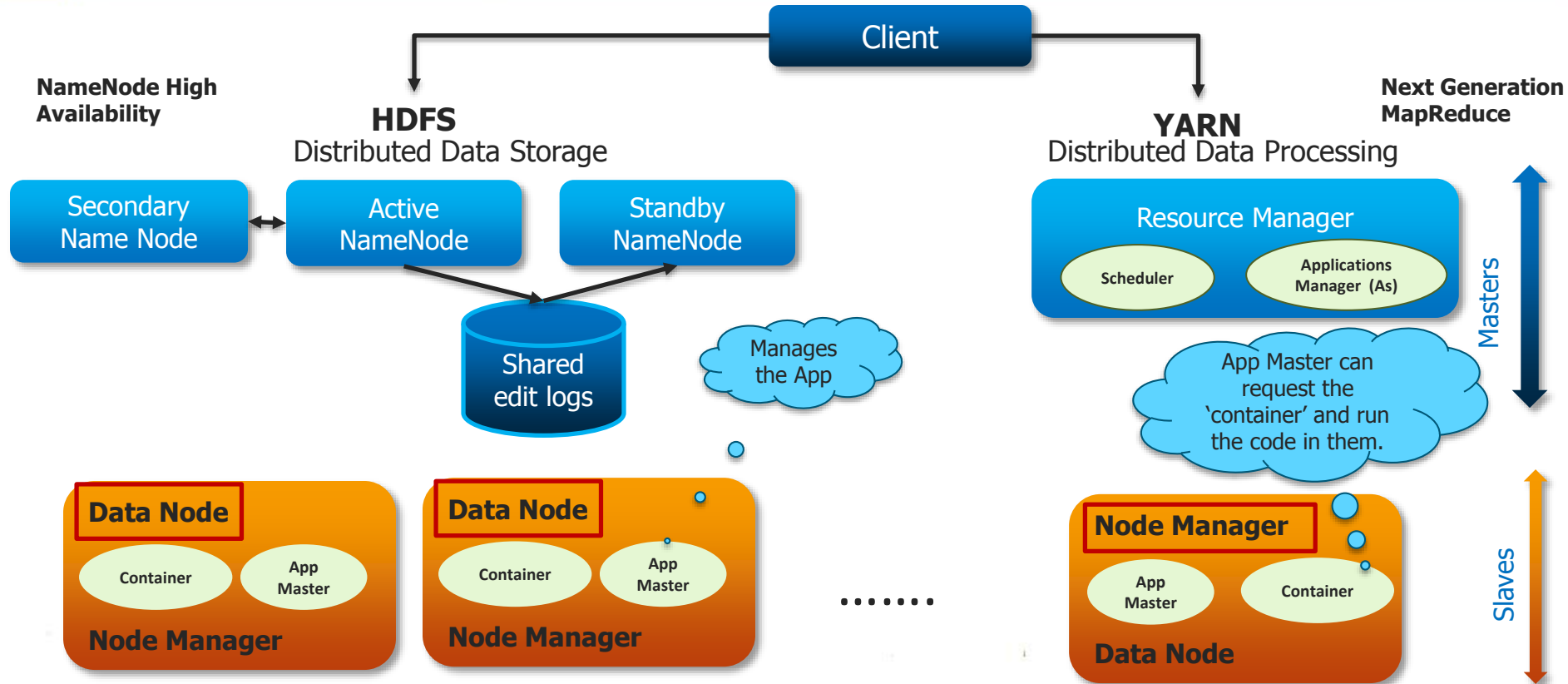


YARN MR Application Execution Flow





Hadoop 2.0 – In Summary



YARN was developed to overcome the following disadvantage in MRv1?

- a) Single Point of Failure of NameNode
- b) Only one version can be run in classic MapReduce
- c) Too much burden on Job Tracker

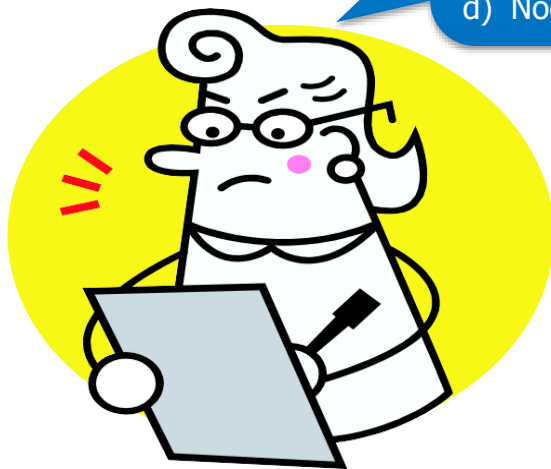


Answer: Too much burden on Job Tracker



In YARN, the functionality of Job Tracker has been replaced by which of the following YARN features:

- a) Job Scheduling
- b) Task Monitoring
- c) Resource Management
- d) Node management



Task Monitoring and Resource Management. The fundamental idea of MRv2 is to split up the two major functionalities of the Job Tracker, i.e. resource management and job scheduling/monitoring, into separate daemons. A global **Resource Manager (RM)** for resources and per-application **Application Master (AM)** for task monitoring.



In YARN, which of the following daemons takes care of the container and the resource utilization by the applications?

- a) Node Manager
- b) Job Tracker
- c) Task tracker
- d) Application Master
- e) Resource manager



Answer: Application Master



Can we run MRv1 Jobs in a YARN enabled Hadoop Cluster?

- a) Yes
- b) No



Answer: Yes. You need recompile the Jobs in MRv2 after enabling YARN to run the Job successfully on a YARN enabled Hadoop Cluster.



Which of the following YARN daemon is responsible for launching the tasks?

- a) Task Tracker
- b) Resource Manager
- c) Application Master
- d) Application Master Server



Answer: Application Master



Tasks for you



- **Attempt the following Assignments using the documents present in the LMS:**
 - Create an Apache Hadoop 2.0 Cluster using a Virtual Machine in VMPlayer or VirtualBox on AWS EC2 free tier.



The screenshot shows the LMS interface for 'Module 6: Advanced Topics: QJM, HDFS Federation and Security'. The module description states: 'In this module, you will understand basics of Hadoop security, Managing security with Kerberos, HDFS Federation setup and Log Management. You will also understand HDFS High Availability using Quorum Journal Manager (QJM)'. The interface lists three items: 'Module 6 Recording' (with a play icon), 'Module 6 Presentation' (with a presentation icon and a 'Download' link), and 'Hadoop Admin Assignment for Module 6' (with a document icon and a 'Download' link). Annotations with dashed arrows point to these items: 'Recording of the Class' points to the recording, 'Assignment' points to the assignment, and 'Presentation' points to the presentation download link.

Module 6: Advanced Topics: QJM, HDFS Federation and Security

In this module, you will understand basics of Hadoop security, Managing security with Kerberos, HDFS Federation setup and Log Management. You will also understand HDFS High Availability using Quorum Journal Manager (QJM).

- Module 6 Recording
- Module 6 Presentation [Download](#)
- Hadoop Admin Assignment for Module 6 [Download](#)

After completion of this assignment, you should be able to: Create an Apache Hadoop 2.0 Cluster using a Virtual Machine in VMPlayer or VirtualBox on AWS EC2 free tier.

Set-up
Guide



Hadoop HDFS Federation Setup

Download

This document is about Hadoop HDFS Federation setup which gives step-by-step guide to the following.

1. Fresh HA install
2. Properties in hdfs-site.xml for Automatic failover
3. Non-HA to HA



Hadoop Admin Quiz for Module 6 (6 Questions)



30 MINUTES

This quiz is based on topics covered in Module-6; Configuring HDFS Federation, Service Monitoring, Service and Log Management, Auditing and Alerts, Service Monitoring, Basics of Hadoop Platform Security, Securing the Platform, Configuring Kerberos.



Take Quiz

Quiz

edureka!

Thank You

See You in Class Next Week