# ANALYTIXLABS

## Hadoop MapReduce (Architecture)
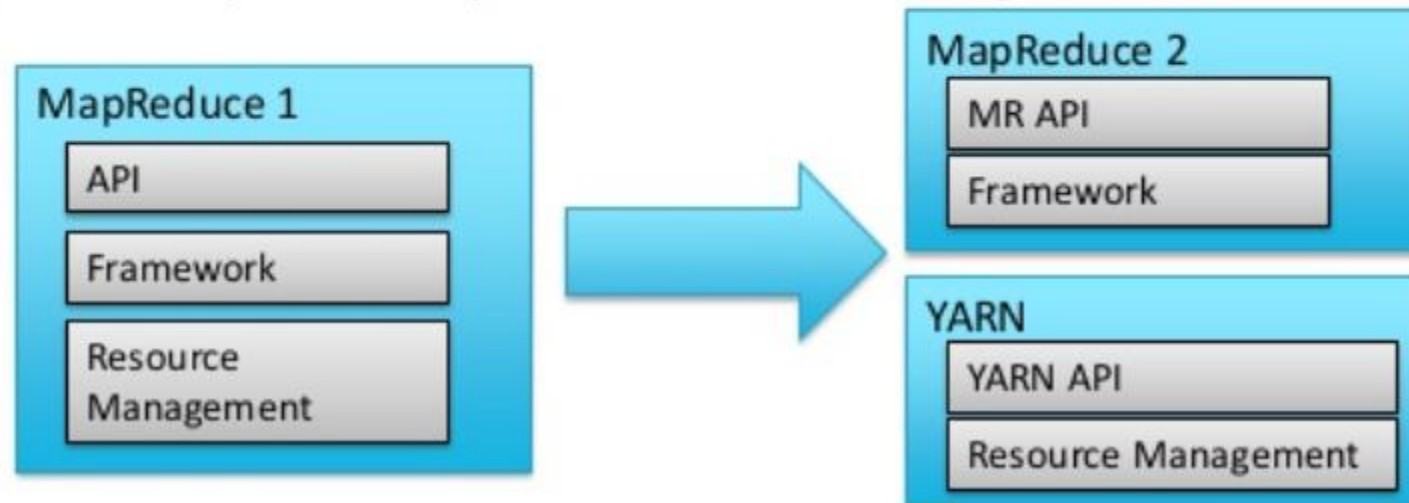
*Learn to Evolve*

## Introduction to YARN and MapReduce 2

- **Overview of MapReduce 1 and 2**

- YARN Architecture

- MapReduce v2

- Managing a YARN Cluster
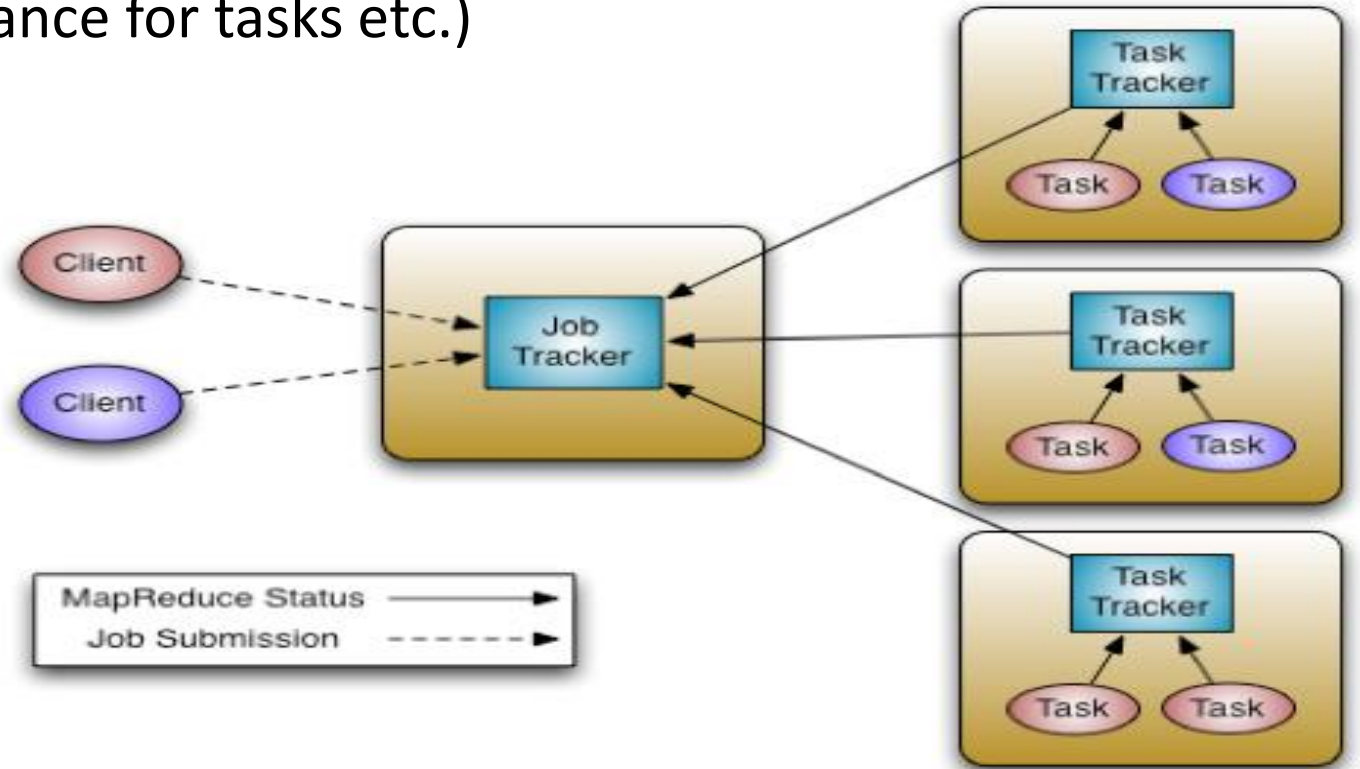
- Cloudera and MRv2

- Conclusion

ANALYTIXLABS

# MRv1 and MRv2

- **MapReduce 1 ("Classic") has three main components**
  - API – for user-level programming of MR applications
  - Framework – runtime services for running Map and Reduce processes, shuffling and sorting, etc.
  - Resource management – infrastructure to monitor nodes, allocate resources, and schedule jobs

- **MapReduce 2 ("NextGen") moves Resource Management into YARN**

# MapReduce V1: Job Tracker & Task Tracker

✓ The JobTracker is responsible for
1. *resource management* (managing the worker nodes i.e. TaskTrackers)
2. *Tracking resource consumption/availability*
3. *Job life-cycle management* (scheduling individual tasks of the job, tracking progress, providing fault-tolerance for tasks etc.)

✓ The TaskTracker has simple responsibilities
1. Launch/teardown tasks on orders from the JobTracker
2. Provide task-status information to the JobTracker periodically.

# Issues with MapReduce (v1)?

- JobTracker has issues related to
  - scalability
  - cluster utilization
  - ability for customers to control upgrades to the stack
  - supporting workloads other than MapReduce itself
  - resiliency to HDFS issues

- MapReduce is essentially batch-oriented and does not support real-time and near real-time processing such as stream processing

- There are utilization issues because map slots might be 'full' while reduce slots are empty (and vice-versa).

# YARN ( Yet Another Resource Negotiator)

- MapReduce has undergone a complete overhaul and CDH 5 now includes MapReduce 2.0 (MRv2).
- The fundamental idea of MRv2's YARN architecture is to split up the two primary responsibilities of the Job Tracker into separate daemons (Global Resource Manager(RM) and Per-Application Masters(AM)
  - 1. Resource management and
  - 2. job scheduling/monitoring
- Yarn provides its core services via two types of long running daemon: a Resource manager and Node manager

http://hadoop.apache.org/docs/stable/hadoop-yarn/hadoop-yarn-site/YARN.html

# YARN ( Yet Another Resource Negotiator)

- **A platform for managing resources in a Hadoop cluster**

- **An Apache Hadoop sub-project**

- **Supports several distributed processing frameworks, including:**
  - MapReduce v2
  - Impala
  - Spark
  - Giraph

- **A full list is at `wiki.apache.org/hadoop/PoweredByYarn`**

# MapReduce2 History

- **Originally architected at Yahoo in 2008**

- **"Alpha" in Hadoop 2 pre-GA**
  - Included in CDH 4

- **YARN promoted to Apache Hadoop sub-project**
  - summer 2013

- **"Production ready" in Hadoop 2 GA**
  - Included in CDH5 (Beta in Oct 2013)

**Hadoop 0.20**

| HDFS | MRv1 |
|------|------|
| Hadoop Common | |

**Hadoop 2.0 (pre-GA)**

| HDFS | MRv2/YARN |
|------|-----------|
| Hadoop Common | |

**Hadoop 2.2 (GA)**

| HDFS | MRv2 |
|------|------|
|      | YARN |
| Hadoop Common | |

ANALYTIXLABS

# Why is YARN needed? (1)

- **MRv1 resource management issues**
  - Inflexible "slots" configured on nodes – Map or Reduce, not both
    - Underutilization of cluster when more map or reduce tasks are running
  - Can't share resources with non-MR applications running on Hadoop cluster (e.g. Impala, Giraph)
  - Scalability – one JobTracker per cluster – limit of about 4000 nodes per cluster

# Why is YARN needed? (2)

- **YARN solutions**
  - No slots
    - Nodes have "resources" – memory and CPU cores – which are allocated to applications when requested
  - Supports MR and non-MR applications running on the same cluster
  - Most Job Tracker functions moved to Application Master – one cluster can have many Application Masters

# YARN Daemons

- **Resource Manager (RM)**
  - Runs on master node
  - Global resource scheduler
  - Arbitrates system resources between competing applications

- **Node Manager (NM)**
  - Runs on slave nodes
  - Communicates with RM

# Running an Application in YARN

- **Containers**
  - Created by the RM upon request
  - Allocate a certain amount of resources (memory, CPU) on a slave node
  - Applications run in one or more containers

NodeManager

| 1 Gb<br>1 core | 3 Gb<br>1 core |

- **Application Master (AM)**
  - One per application
  - Framework/application specific
  - Runs in a container
  - Requests more containers to run application tasks

NodeManager

Application Master

ANALYTIXLABS

# YARN Cluster

```
                                    ┌──────────────────────────────┐
                                    │  NodeManager              🖥  │
                               ┌───►│                              │
                               │    └──────────────────────────────┘
                               │    ┌──────────────────────────────┐
                               │    │  NodeManager              🖥  │
                               ├───►│                              │
                               │    └──────────────────────────────┘
┌──────────────────┐           │    ┌──────────────────────────────┐
│ Resource      🖥 │           │    │  NodeManager              🖥  │
│ Manager          │──────────►├───►│                              │
│                  │           │    └──────────────────────────────┘
└──────────────────┘           │    ┌──────────────────────────────┐
                               │    │  NodeManager              🖥  │
                               └───►│                              │
                                    └──────────────────────────────┘
```

# Resource Manager – High Availability



HDFS HIGH AVAILABILITY

All name space edits logged to shared NFS storage; single writer (fencing)

Read edit logs and applies to its own namespace

NameNode High Availability

*Not necessary to configure Secondary NameNode

Shared Edit Logs

Secondary Name Node

Active NameNode

Standby NameNode

Client

HDFS

YARN

Resource Manager

Next Generation MapReduce

DataNode

DataNode

Node Manager

Node Manager

Node Manager
Container   App Master

Node Manager
Container   App Master

Container   App Master

Container   App Master

DataNode

DataNode

http://hadoop.apache.org/docs/stable2/hadoop-yarn/hadoop-yarn-site/HDFSHighAvailabilityWithNFS.html

ANALYTIXLABS

# Resource Manager – High Availability

# YARN Cluster: Running an Application



1. A client program submits the application, including the necessary specification to launch the application-specific Application Master
2. The Resource Manager assumes the responsibility to negotiate a specific container in which to start the Application Master and then launches the Application Master

## YARN Cluster: Running an Application

Client

Application: MyApp

Resource Manager

NodeManager

NodeManager

NodeManager

Resource Request

Application Master

Container IDs

NodeManager

3. The Application Master on boot-up, register with the Resource Manager for details which allow it to directly communicate with its own Application Master

4. During normal operation Application Master negotiates appropriate resource container via the resource-request protocol

ANALYTI**X**LABS

# YARN Cluster: Running an Application



5. On successful container allocation, the Application Master launches the container by providing the container launch specification to the Node Manager. The launch specification, typically includes the necessary information to allow the container to communicate with the Application Master itself

# YARN Cluster: Running an Application

**Client** **Client**

Application: YourApp

Application: MyApp

Resource Manager

Launch

NodeManager

NodeManager
- Application Master
- MyApp

NodeManager
- Application Master

NodeManager
- MyApp

6. In case of new application submitted, Resource Manager launches new application specific Application Master

ANALYTIXLABS

# YARN Cluster: Running an Application



7. The application code executing within the container then provides the necessary information(progress, status) to its Application Master via application specific protocol

8. During the application execution, the client that submitted the program communicates directly with the Application Master to get the status, progress update etc. via an application specific protocol

9. Once the application is complete, and all necessary work has been finished, the application master deregisters with the Resource Manager and shuts down, allowing its own container to be repurposed.

# YARN Schedulers (1)

- **Pluggable in Resource Manager**

- **YARN includes** three **schedulers**
  - CapacityScheduler          - FIFO Scheduler          - Preemptive Scheduling*
  - FairScheduler                    - Dominant Resource Fairness*          - Delay Scheduling
                                                                                    - Speculative Execution*

- **How are these different than MRv1 schedulers?**
  - Support any YARN application, not just MR
  - No more "slots" – tasks are allocated based on resources (memory and CPU for now)
  - FairScheduler: *pools* are now called *queues*

ANALYTIXLABS

# YARN Schedulers (2)

- **Hierarchical queues**
  - Queues can contain sub-queues
  - Sub-queues share the resources assigned to queues



Engineering (weight=2)
- Product A
- Product B

Marketing (weight=1)
- Website Logs ETL
- Fast Lane
- Regular
- Data Anal

# Resource Manager Things to Know

**What it does**

- Manages nodes
  - Tracks heartbeats from NodeManagers
- Manages containers
  - Handles AM requests for resources
  - De-allocates containers when they expire or the application completes
- Manages ApplicationMasters
  - Creates a container for AMs and tracks heartbeats
- Manages security
  - Supports Kerberos



Resource Manager

## Node Manager Things to Know

- **What it does**
  - Communicates with the RM
    - Registers and provides info on node resources
    - Sends heartbeats and container status
  - Manages processes in containers
    - Launches AMs on request from the RM
    - Launches application processes on request from AM
    - Monitors resource usage by containers; kills run-away processes
  - Provides logging services to applications
    - Aggregates logs for an application and saves them to HDFS
  - Runs auxiliary services
  - Maintains node level security via ACLs

NodeManager

# Resource Requests

**Resource Request**
- Resource name (hostname, rackname or *)
- Priority (within this application, not between applications)
- Resource requirements
  - memory (MB)
  - CPU (# of cores)
  - more to come, e.g. disk and network I/O, GPUs, etc.
- Number of containers

Resource Manager

NodeManager

Application Master

**Container(s)**
- ID
- Node

# Launch Container

**Container Launch Context**
- Container ID
- Commands (to start application)
- Environment (configuration)
- Local Resources (e.g. application binary, HDFS files)

NodeManager

MyApp

NodeManager

Application Master

# Non-MR2 YARN Applications

- **Distributed Shell**

- **Impala**

- **Apache Giraph**

- **Spark**

- **Others**
    - http://wiki.apache.org/hadoop/PoweredByYarn

# YARN and MapReduce

- **YARN does not know or care what kind of application it is running**
  - Could be MR or something else (e.g. Impala)

- **MR2 uses YARN**
  - Hadoop includes a MapReduce ApplicationMaster (MRAppMaster) to manage MR jobs
  - Each MapReduce job is an a new instance of an application

# Running a MapReduce Application in MRv2

# Running a MapReduce Application in MRv2

Running a MapReduce Application in MRv2

Running a MapReduce Application in MRv2

# Running a MapReduce Application in MRv2

# Running a MapReduce Application in MRv2

# Running a MapReduce Application in MRv2

# Running a MapReduce Application in MRv2

# Running a MapReduce Application in MRv2

# Running a MapReduce Application in MRv2

# The MapReduce Framework on YARN

- **In YARN, Shuffle is run as an auxiliary service**
  - Runs in the NodeManager JVM as a persistent service

# Fault Tolerance

- **Any of the following can fail**
  - Task (Container) – Handled just like in MRv1
    - MRAppMaster will re-attempt tasks that complete with exceptions or stop responding (4 times by default)
    - Applications with too many failed tasks are considered failed

# Fault Tolerance

- **Any of the following can fail**
  - Task (Container) – Handled just like in MRv1
    - MRAppMaster will re-attempt tasks that complete with exceptions or stop responding (4 times by default)
    - Applications with too many failed tasks are considered failed
  - Application Master
    - If application fails or if AM stops sending heartbeats, RM will re-attempt the whole application (2 times by default )
    - MRAppMaster optional setting: Job recovery
      - if false, all tasks will re-run
      - if true, MRAppMaster retrieves state of tasks when it restarts; only incomplete tasks will be re-run

ANALYTIXLABS

# Resource Manager UI: Nodes

`http://rmhost:8088/cluster/nodes`



link to Node Manager UI

# Resource Manager UI: Applications

`http://rmhost:8088/cluster/apps`



Link to Application Details…

List of running and recent applications

# Resource Manager UI: Application Detail

# MRAppMaster UI: Jobs

`http://rmhost:8088/proxy/appid`



Link to Job Details...

# MRAppMaster UI: Tasks

`http://rmhost:8088/proxy/appid/mapreduce/job/jobid`

# MR Job History Server

- **YARN does not keep track of job history**

- **MapReduce Job History Server**
  - Archives job's metrics and metadata
  - Can be accessed through Job History UI or Hue

`http://rmhost:19888/jobhistory`

# Cloudera Manager

- **Full support for MR2 added in CM 5**

# Hue

- **Hue supports browsing MRv2 jobs**
  - Connects to Job History Server for "retired" jobs
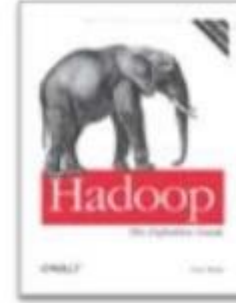
# Where To Learn More

- *Hadoop: The Definitive Guide,* **3<sup>rd</sup> Edition**
  - Chapter 6 – focuses on how MR is implemented on YARN

- **Cloudera Blog posts –** blog.cloudera.com/blog/category/yarn
  - Migrating to MapReduce 2 on YARN
  - Writing Hadoop Programs That Work Across Releases
  - and more...

# Contact Us

Visit us on: http://www.analytixlabs.in/


For more information, please contact us: http://www.analytixlabs.co.in/contact-us/

Or email: info@analytixlabs.co.in


Call us we would love to speak with you: (+91) 9910509849


Join us on:

Twitter - http://twitter.com/#!/AnalytixLabs

Facebook - http://www.facebook.com/analytixlabs

LinkedIn - http://www.linkedin.com/in/analytixlabs

Blog - http://www.analytixlabs.co.in/category/blog/