



SATHYABAMA

**INSTITUTE OF SCIENCE AND TECHNOLOGY
(DEEMED TO BE UNIVERSITY)**

Accredited "A" Grade by NAAC | 12B Status by UGC | Approved by AICTE

www.sathyabama.ac.in

SCHOOL OF COMPUTING

DEPARTMENT OF COMPUTER SCIENCE AND ENGINEERING

UNIT V Big Data Analytics – SCSA1603

UNIT 5 MAP REDUCE

9 Hrs

Algorithms using map reduce - Matrix-Vector – Multiplication – Word Count - Understanding inputs and outputs of MapReduce, Data Serialization – Introduction to YARN – MapReduce Vs YARN – YARN Architecture – Scheduling in YARN- – Fair Scheduler – Capacity Scheduler.

Map Reduce

MapReduce is a framework that is used by Hadoop to process the data residing with HDFS. HDFS store data in each block which size is 64 MB or 128 MB. And MapReduce can interaction with HDFS and operates the data in HDFS. The configuration which we set when we set the environment of hadoop. The below diagram will show the structure of hadoop.

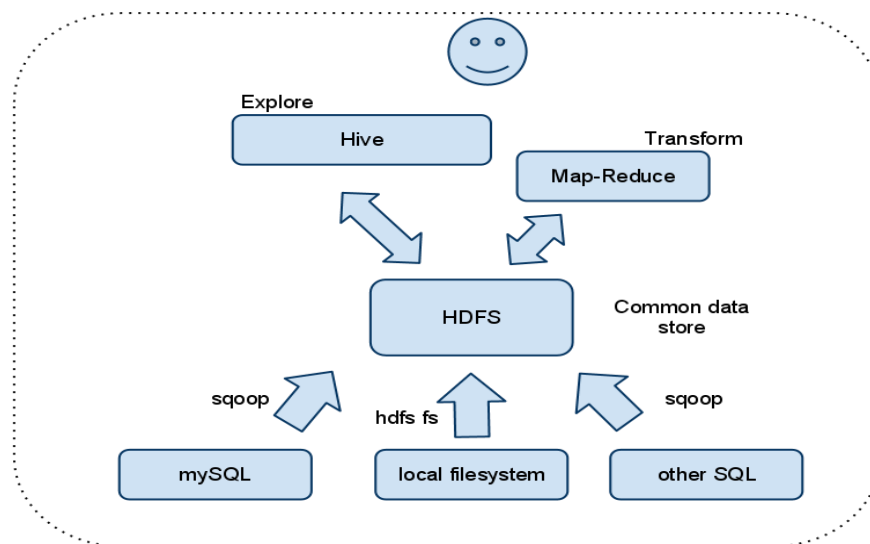


Figure 51.: Structure of Hadoop

MapTask(HDFS data localization):

The unit of input for a map task is an HDFS data block of the input file. The map task functions most efficiently if the data block it has to process is available locally on the node on which the task is scheduled. This approach is called HDFS data localization.

An HDFS data locality miss occurs if the data needed by the map task is not available locally. In such a case, the map task will request the data from another node in the cluster: an operation that is expensive and time consuming, leading to inefficiencies and, hence, delay in job completion.

[

Clients, Data Nodes, and HDFS Storage:

Input data is uploaded to the HDFS filesystem in either of the following two ways:

An HDFS client has a large amount of data to place into HDFS. An HDFS client is constantly streaming data into HDFS.

Both these scenarios have the same interaction with HDFS, except that in the streaming case, the client waits for enough data to fill a data block before writing to HDFS. Data is stored in HDFS in large blocks, generally 64 to 128 MB or more in size. This storage approach allows easy parallel processing of data.

HDFS-SITE.XML

```
<property>
<name>dfs.block.size</name>
<value>134217728</value>ç128MB Block size
</property>
```

OR

```
<property>
<name>dfs.block.size</name>
<value>67108864</value>ç64MB Block size (Default is this value is not set)
</property>
```

Block Replication Factor:

During the process of writing to HDFS, the blocks are generally replicated to multiple data nodes for redundancy. The number of copies, or the replication factor, is set to a default of 3 and can be modified by the cluster administrator as below:

HDFS-SITE.XML

```
<property>
```

```
<name>dfs.replication</name>
<value>3</value>
</property>
<property>
```

When the replication factor is three, HDFS's placement policy is to:

- Put one replica on one node in the local rack,
- Another on a node in a different (remote) rack,
- Last on a different node in the same remote rack.

When a new data block is stored on a data node, the data node initiates a replication process to replicate the data onto a second data node. The second data node, in turn, replicates the block to a third data node, completing the replication of the block.

With this policy, the replicas of a file do not evenly distribute across the racks. One third of replicas are on one node, two thirds of replicas are on one rack, and the other third are evenly distributed across the remaining racks. This policy improves write performance without compromising data reliability or read performance.

5.1 Algorithms using Map Reduce

The MapReduce algorithm contains two important tasks, namely Map and Reduce.

- The map task is done by means of Mapper Class
- The reduce task is done by means of Reducer Class.

Mapper class takes the input, tokenizes it, maps and sorts it. The output of Mapper class is used as input by Reducer class, which in turn searches matching pairs and reduces them.

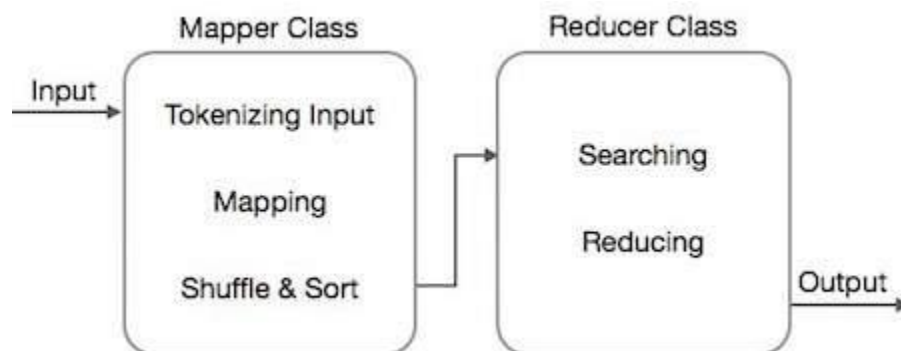


Figure 5.2: The Mapper class and the Reducer class

MapReduce implements various mathematical algorithms to divide a task into small parts and assign them to multiple systems. In technical terms, MapReduce algorithm helps in sending the Map & Reduce tasks to appropriate servers in a cluster.

These mathematical algorithms may include the following –

- Sorting
- Searching
- Indexing
- TF-IDF

Sorting

Sorting is one of the basic MapReduce algorithms to process and analyze data. MapReduce implements sorting algorithm to automatically sort the output key-value pairs from the mapper by their keys.

- Sorting methods are implemented in the mapper class itself.
- In the Shuffle and Sort phase, after tokenizing the values in the mapper class, the **Context** class (user-defined class) collects the matching valued keys as a collection.
- To collect similar key-value pairs (intermediate keys), the Mapper class takes the help of **RawComparator** class to sort the key-value pairs.
- The set of intermediate key-value pairs for a given Reducer is automatically sorted by Hadoop to form key-values (K2, {V2, V2, ...}) before they are presented to the Reducer.

Searching

Searching plays an important role in MapReduce algorithm. It helps in the combiner phase (optional) and in the Reducer phase. Let us try to understand how Searching works with the help of an example.

Example

The following example shows how MapReduce employs Searching algorithm to find out the details of the employee who draws the highest salary in a given employee dataset.

- Let us assume we have employee data in four different files – A, B, C, and D. Let us also assume there are duplicate employee records in all four files because of importing the employee data from all database tables repeatedly. See the following illustration.

name, salary	name, salary	name, salary	name, salary
satish, 26000	gopal, 50000	satish, 26000	satish, 26000
Krishna, 25000	Krishna, 25000	kiran, 45000	Krishna, 25000
Satishk, 15000	Satishk, 15000	Satishk, 15000	manisha, 45000
Raju, 10000	Raju, 10000	Raju, 10000	Raju, 10000

- **The Map phase** processes each input file and provides the employee data in key-value pairs (<k, v> : <emp name, salary>). See the following illustration.

<satish, 26000>	<gopal, 50000>	<satish, 26000>	<satish, 26000>
<Krishna, 25000>	<Krishna, 25000>	<kiran, 45000>	<Krishna, 25000>
<Satishk, 15000>	<Satishk, 15000>	<Satishk, 15000>	<manisha, 45000>
<Raju, 10000>	<Raju, 10000>	<Raju, 10000>	<Raju, 10000>

- **The combiner phase** (searching technique) will accept the input from the Map phase as a key-value pair with employee name and salary. Using searching technique, the combiner will check all the employee salary to find the highest salaried employee in each file. See the following snippet.

<k: employee name, v: salary>

Max= the salary of an first employee. Treated as max salary

```
if(v(second employee).salary > Max){
    Max = v(salary);
}
```

```
else{
    Continue checking;
}
```

The expected result is as follows –

<satish, 26000>	<gopal, 50000>	<kiran, 45000>	<manisha, 45000>
--------------------	-------------------	-------------------	---------------------

- **Reducer phase** – Form each file, you will find the highest salaried employee. To avoid redundancy, check all the <k, v> pairs and eliminate duplicate entries, if any. The same algorithm is used in between the four <k, v> pairs, which are coming from four input files. The final output should be as follows –

<gopal, 50000>

Indexing

Normally indexing is used to point to a particular data and its address. It performs batch indexing on the input files for a particular Mapper.

The indexing technique that is normally used in MapReduce is known as **inverted index**. Search engines like Google and Bing use inverted indexing technique. Let us try to understand how Indexing works with the help of a simple example.

Example

The following text is the input for inverted indexing. Here T[0], T[1], and T[2] are the file names and their content are in double quotes.

T[0] = "it is what it is"

T[1] = "what is it"

T[2] = "it is a banana"

After applying the Indexing algorithm, we get the following output –

"a": {2}

"banana": {2}

"is": {0, 1, 2}

"it": {0, 1, 2}

"what": {0, 1}

Here "a": {2} implies the term "a" appears in the T[2] file. Similarly, "is": {0, 1, 2} implies the term "is" appears in the files T[0], T[1], and T[2].

TF-IDF

TF-IDF is a text processing algorithm which is short for Term Frequency – Inverse Document Frequency. It is one of the common web analysis algorithms. Here, the term 'frequency' refers to the number of times a term appears in a document.

Term Frequency (TF)

It measures how frequently a particular term occurs in a document. It is calculated by the number of times a word appears in a document divided by the total number of words in that document.

$$TF(\text{the}) = (\text{Number of times term the 'the' appears in a document}) / (\text{Total number of terms in the document})$$

Inverse Document Frequency (IDF)

It measures the importance of a term. It is calculated by the number of documents in the text database divided by the number of documents where a specific term appears.

While computing TF, all the terms are considered equally important. That means, TF counts the term frequency for normal words like “is”, “a”, “what”, etc. Thus we need to know the frequent terms while scaling up the rare ones, by computing the following –

$\text{IDF}(\text{the}) = \log_e(\text{Total number of documents} / \text{Number of documents with term 'the' in it})$.

The algorithm is explained below with the help of a small example.

Example

Consider a document containing 1000 words, wherein the word **hive** appears 50 times. The TF for **hive** is then $(50 / 1000) = 0.05$.

Now, assume we have 10 million documents and the word **hive** appears in 1000 of these. Then, the IDF is calculated as $\log(10,000,000 / 1,000) = 4$.

The TF-IDF weight is the product of these quantities – $0.05 \times 4 = 0.20$.

5.2 Performing Matrix Vector Multiplication using MapReduce

Suppose we have an $n \times n \times n$ matrix M , whose element in row i and column j will be denoted m_{ij} . Suppose we also have a vector v of length n , whose j th element is v_j . Then the matrix-vector product is the vector x of length n , whose i th element x_i is given by $x_i = \sum_{j=1}^n m_{ij} \times v_j$. If $n = 100$, we do not want to use a DFS or MapReduce for this calculation. But this sort of calculation is at the heart of the ranking of Web pages that goes on at search engines, and there, n is in the tens of billions.³ Let us first assume that n is large, but not so large that vector v cannot fit in main memory and thus be available to every Map task. The matrix M and the vector v each will be stored in a file of the DFS. We assume that the row-column coordinates of each matrix element will be discoverable, either from its position in the file, or because it is stored with explicit coordinates, as a triple (i, j, m_{ij}) . We also assume the position of element v_j in the vector v will be discoverable in the analogous way.

The Map Function: The Map function is written to apply to one element of M . However, if v is not already read into main memory at the compute node executing a Map task, then v is first read, in its entirety, and subsequently will be available to all applications of the Map function performed at this Map task. Each Map task will operate on a chunk of the matrix M . From each matrix element m_{ij} it produces the key-value pair $(i, m_{ij} \times v_j)$. Thus, all terms of the sum that make up the component x_i of the matrix-vector product will get the same key, i .

The Reduce Function: The Reduce function simply sums all the values associated with a given key i . The result will be a pair (i, x_i) .

we can divide the matrix into vertical stripes of equal width and divide the vector into an equal number of horizontal stripes, of the same height. Our goal is to use enough stripes so that the portion of the vector in one stripe can fit conveniently into main memory at a compute node. Figure suggests what the partition looks like if the matrix and vector are each divided into five stripes.

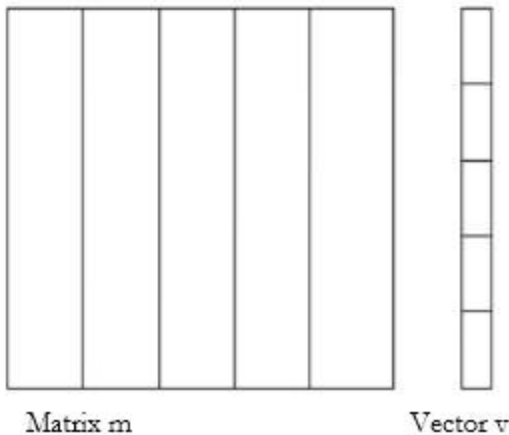


Figure 5.3: Matrix vector representation

Map : $\text{for } i = 0 \text{ to } n-1$ input m_{ij}
 Emit ($i, p_s = \sum m_{ij} * v_j$)

Reduce : Compute $X_i = \sum p_s$

5.3 Matrix Multiplication With MapReduce

MapReduce is a technique in which a huge program is subdivided into small tasks and run parallelly to make computation faster, save time, and mostly used in distributed systems. It has 2 important parts:

- **Mapper:** It takes raw data input and organizes into key, value pairs. For example, In a dictionary, you search for the word “Data” and its associated meaning is “facts and statistics collected together for reference or analysis”. Here the Key is *Data* and the **Value** associated with is *facts and statistics collected together for reference or analysis*.
- **Reducer:** It is responsible for processing data in parallel and produce final output. Let us consider the matrix multiplication example to visualize MapReduce. Consider the following matrix:

$$\begin{bmatrix} 1 & 2 \\ 3 & 4 \end{bmatrix} \times \begin{bmatrix} 5 & 6 \\ 7 & 8 \end{bmatrix}$$

Figure 5.4: 2×2 matrices A and B

Here matrix A is a 2×2 matrix which means the number of rows(i)=2 and the number of columns(j)=2. Matrix B is also a 2×2 matrix where number of rows(j)=2 and number of columns(k)=2. Each cell of the matrix is labelled as A_{ij} and B_{ij} . Ex. element 3 in matrix A is

called A21 i.e. 2nd-row 1st column. Now One step matrix multiplication has 1 mapper and 1 reducer. The Formula is:

Mapper for Matrix A $(k, v) = ((i, k), (A, j, A_{ij}))$ for all k
Mapper for Matrix B $(k, v) = ((i, k), (B, j, B_{jk}))$ for all i

Therefore computing the mapper for Matrix A:

k, i, j computes the number of times it occurs. # Here all are 2, therefore when $k=1$, i can have # 2 values 1 & 2, each case can have 2 further # values of $j=1$ and $j=2$. Substituting all values # in formula

$k=1 \ i=1 \ j=1 \ ((1, 1), (A, 1, 1))$

$j=2 \ ((1, 1), (A, 2, 2))$

$i=2 \ j=1 \ ((2, 1), (A, 1, 3))$

$j=2 \ ((2, 1), (A, 2, 4))$

$k=2 \ i=1 \ j=1 \ ((1, 2), (A, 1, 1))$

$j=2 \ ((1, 2), (A, 2, 2))$

$i=2 \ j=1 \ ((2, 2), (A, 1, 3))$

$j=2 \ ((2, 2), (A, 2, 4))$

Computing the mapper for Matrix B

$i=1 \ j=1 \ k=1 \ ((1, 1), (B, 1, 5))$

$k=2 \ ((1, 2), (B, 1, 6))$

$j=2 \ k=1 \ ((1, 1), (B, 2, 7))$

$j=2 \ ((1, 2), (B, 2, 8))$

$i=2 \ j=1 \ k=1 \ ((2, 1), (B, 1, 5))$

$k=2 \ ((2, 2), (B, 1, 6))$

$j=2 \ k=1 \ ((2, 1), (B, 2, 7))$

$k=2 \ ((2, 2), (B, 2, 8))$

The formula for Reducer is:

Reducer(k, v)= $(i, k) \Rightarrow$ Make sorted Alist and Blist
 (i, k) \Rightarrow Summation ($A_{ij} * B_{jk}$) for j
 Output $\Rightarrow ((i, k), \text{sum})$

Therefore computing the reducer:

We can observe from Mapper computation # that 4 pairs are common (1, 1), (1, 2), (2, 1) and (2, 2) # Make a list separate for Matrix A & # B with adjoining values taken from # Mapper step above:

(1, 1) => Alist = {(A, 1, 1), (A, 2, 2)}

Blist = {(B, 1, 5), (B, 2, 7)}

Now $A_{ij} \times B_{jk}$: $[(1 \times 5) + (2 \times 7)] = 19$ -----(i)

(1, 2) => Alist = {(A, 1, 1), (A, 2, 2)}

Blist = {(B, 1, 6), (B, 2, 8)}

Now $A_{ij} \times B_{jk}$: $[(1 \times 6) + (2 \times 8)] = 22$ -----(ii)

(2, 1) => Alist = {(A, 1, 3), (A, 2, 4)}

Blist = {(B, 1, 5), (B, 2, 7)}

Now $A_{ij} \times B_{jk}$: $[(3 \times 5) + (4 \times 7)] = 43$ -----(iii)

(2, 2) => Alist = {(A, 1, 3), (A, 2, 4)}

Blist = {(B, 1, 6), (B, 2, 8)}

Now $A_{ij} \times B_{jk}$: $[(3 \times 6) + (4 \times 8)] = 50$ -----(iv)

From (i), (ii), (iii) and (iv) we conclude that

((1, 1), 19)

((1, 2), 22)

((2, 1), 43)

((2, 2), 50)

Therefore the Final Matrix is:

$$\begin{bmatrix} 19 & 22 \\ 43 & 50 \end{bmatrix}$$

Figure 5.5.: Final output of Matrix multiplication.

5.4 Word Count Problem using Map reduce and Understanding inputs and outputs of MapReduce

5.4.1 Understanding the inputs and outputs

Input:

What do you mean by Object
What do you know about Java
What is Java Virtual Machine
How Java enabled High Performance

Output:

What	3
Do	2
You	2
Mean	1
By	1
Object	1
Know	1
about	1
Java	3
Is	1
Virtual	1
Machine	1
How	
1	
enabled	1
High	1
Performance	1

5.4.2 Word Count Problem:

The word count problem undergoes two phases, a mapper phase and a reducer phase

Word count Mapper:

```

import java.io.IOException;
import java.util.StringTokenizer;
import org.apache.hadoop.io.IntWritable;
import org.apache.hadoop.io.LongWritable;
import org.apache.hadoop.io.Text;
import org.apache.hadoop.mapred.MapReduceBase;
import org.apache.hadoop.mapred.Mapper;
import org.apache.hadoop.mapred.OutputCollector;
import org.apache.hadoop.mapred.Reporter;
public class WC_Mapper extends MapReduceBase implements Mapper<LongWritable,Text,Text,IntWritable>{
    private final static IntWritable one = new IntWritable(1);
    private Text word = new Text();
    public void map(LongWritable key, Text value,OutputCollector<Text,IntWritable> output,
        Reporter reporter) throws IOException{
        String line = value.toString();
        StringTokenizer tokenizer = new StringTokenizer(line);
        while (tokenizer.hasMoreTokens()){
            word.set(tokenizer.nextToken());
            output.collect(word, one);
        }
    }
}

```

Word count Reducer

```

import java.io.IOException;
import java.util.Iterator;
import org.apache.hadoop.io.IntWritable;
import org.apache.hadoop.io.Text;
import org.apache.hadoop.mapred.MapReduceBase;
import org.apache.hadoop.mapred.OutputCollector;
import org.apache.hadoop.mapred.Reducer;

```

```

import org.apache.hadoop.mapred.Reporter;

public class WC_Reducer extends MapReduceBase implements Reducer<Text,IntWritable,Text,IntWritable> {
    public void reduce(Text key, Iterator<IntWritable> values,OutputCollector<Text,IntWritable> output,
        Reporter reporter) throws IOException {
        int sum=0;
        while (values.hasNext()) {
            sum+=values.next().get();
        }
        output.collect(key,new IntWritable(sum));
    }
}

```

Main method class

```

import java.io.IOException;
import org.apache.hadoop.fs.Path;
import org.apache.hadoop.io.IntWritable;
import org.apache.hadoop.io.Text;
import org.apache.hadoop.mapred.FileInputFormat;
import org.apache.hadoop.mapred.FileOutputFormat;
import org.apache.hadoop.mapred.JobClient;
import org.apache.hadoop.mapred.JobConf;
import org.apache.hadoop.mapred.TextInputFormat;
import org.apache.hadoop.mapred.TextOutputFormat;
public class WC_Runner {
    public static void main(String[] args) throws IOException{
        JobConf conf = new JobConf(WC_Runner.class);
        conf.setJobName("WordCount");
        conf.setOutputKeyClass(Text.class);
        conf.setOutputValueClass(IntWritable.class);
        conf.setMapperClass(WC_Mapper.class);
        conf.setCombinerClass(WC_Reducer.class);
    }
}

```

```

    conf.setReducerClass(WC_Reducer.class);
    conf.setInputFormat(TextInputFormat.class);
    conf.setOutputFormat(TextOutputFormat.class);
    FileInputFormat.setInputPaths(conf,new Path(args[0]));
    FileOutputFormat.setOutputPath(conf,new Path(args[1]));
    JobClient.runJob(conf);
}
}

```

MapReduce is a processing technique and a program model for distributed computing based on java. The MapReduce algorithm contains two important tasks, namely Map and Reduce.

- **Map** takes a set of data and converts it into another set of data, where individual elements are broken down into tuples (key/value pairs).
- **Reduce** task takes the output from a map as an input and combines those data tuples into a smaller set of tuples.

As the sequence of the name MapReduce implies, the reduce task is always performed after the map job.

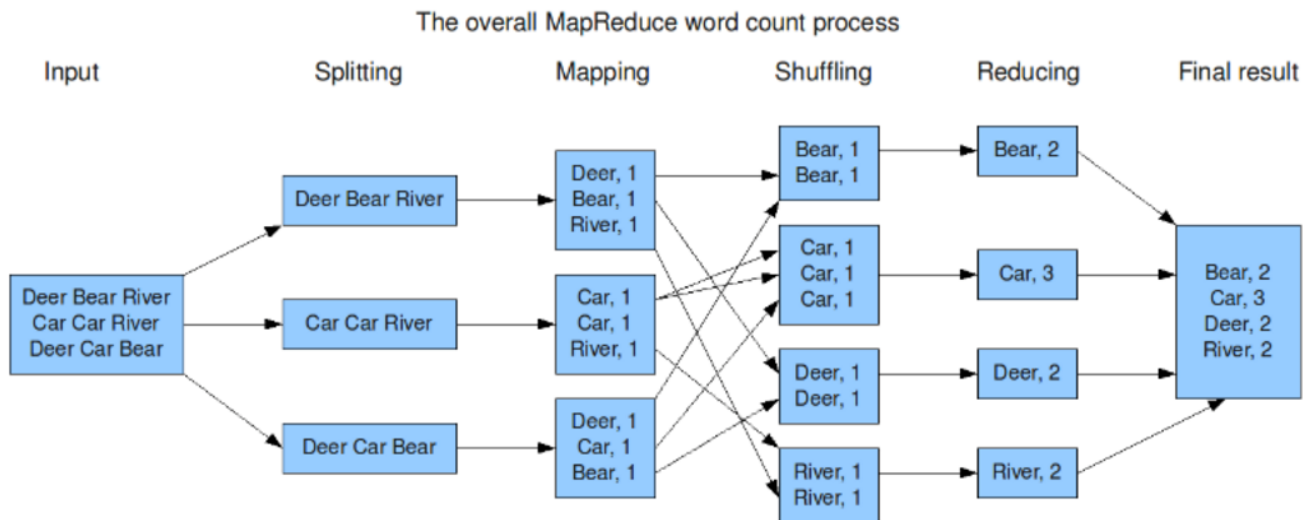


Figure 5.6: The word count problem using Map reduce

5.5 Data Serialization:

- Serialization is the process of converting object data into byte stream data for transmission over a network across different nodes in a cluster or for persistent data storage.
- Data serialization is a process that converts structure data manually back to the original form.
- Serialize to translate data structures into a stream of data. Transmit this stream of data over the network or store it in DB regardless of the system architecture.
- Isn't storing information in binary form or stream of bytes is the right approach.
- Serialization does the same but isn't dependent on architecture.

Consider CSV files contains a comma (,) in between data, so while Deserialization, wrong outputs may occur. Now, if metadata is stored in XML form, a self- architected form of data storage, data can easily deserialize.

Deserialization:

- Deserialization is the reverse process of serialization and converts byte stream data into object data for reading data from HDFS. Hadoop provides Writables for serialization and deserialization purpose.
- Writable and WritableComparable Interfaces to provide mechanisms for serialization and deserialization of data, Hadoop provided two important interfaces Writable and WritableComparable
- **Hadoop File System (HDFS)?**
- Hadoop File System (HDFS) is a distributed file system. Store all types of files in the Hadoop file system. It supports all standard formats such as Gifs, text, CSV, tsv, xls, etc. Beginners of Hadoop can opt for tab_delimiter (data separated by tabs) files because it is -
- Easy to debug and readable
- The default format of Apache Hive
- **What is the Data Serialization Storage format?**
- Storage formats are a way to define how to store information in the file. Most of the time, assume this information is from the extension of the data. Both structured and unstructured data can store on HADOOP-enabled systems. Common Hdfs file formats are -
- Plain text storage
- Sequence files
- RC files
- AVRO
- Parquet

Need for storage formats in Hadoop

- File format must be handy to serve complex data structures

- HDFS enabled applications to find relevant data in a particular location and write back data to another location.
- Dataset is large
- Having schemas
- Having storage constraints
-

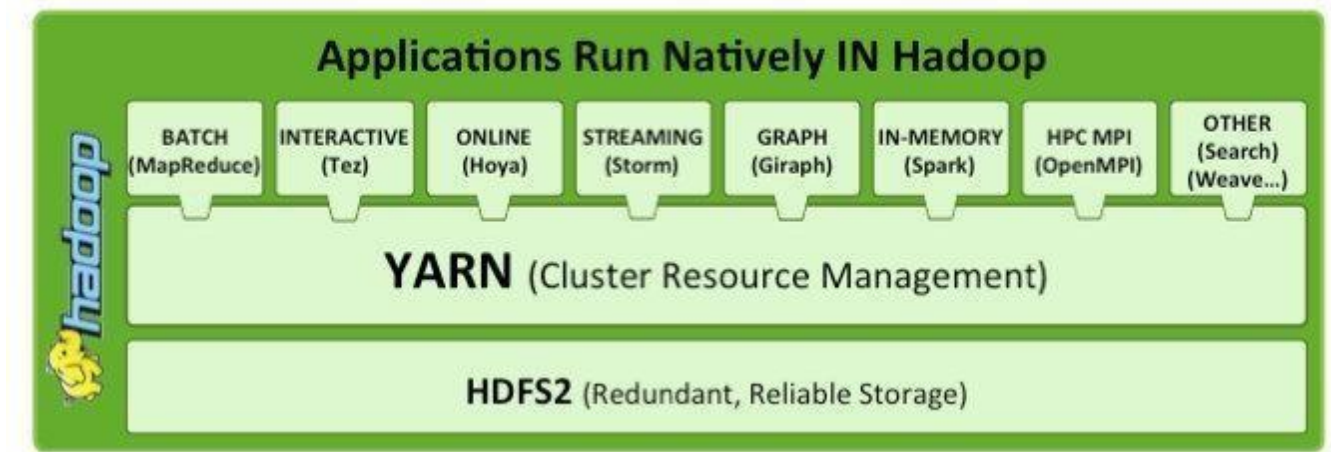
Serialization in Hadoop has two areas -

- **Interproc**
ess communication
When a client calls a function or subroutine from one pc to the pc in-network or server, that calling is a remote procedure call.
-
- **Persisten**
t storage
It is better than java's inbuilt serialization as java serialization isn't compact. Serialization and Deserialization of data help maintain and manage corporate decisions for effective use of resources and data available in Data warehouse or any other database -writable - language specific to java.

5.6 Introduction to YARN and Hadoop2.0

Apache Hadoop 2.0 represents a generational shift in the architecture of Apache Hadoop. With YARN, Apache Hadoop is recast as a significantly more powerful platform – one that takes Hadoop beyond merely batch applications to taking its position as a ‘data operating system’ where HDFS is the file system and YARN is the operating system. YARN is a re-architecture of Hadoop that allows multiple applications to run on the same platform. With YARN, applications run—in Hadoop, instead of—on Hadoop:

Figure 5.7: Hadoop2.0 YARN



The fundamental idea of YARN is to split up the two major responsibilities of the JobTracker and TaskTracker into separate entities. In Hadoop 2.0, the JobTracker and TaskTracker no longer exist and have been replaced by three components:

- **ResourceManager:** a scheduler that allocates available resources in the cluster amongst the competing applications.
- **NodeManager:** runs on each node in the cluster and takes direction from the ResourceManager. It is responsible for managing resources available on a single node.
- **ApplicationMaster:** an instance of a framework-specific library, an ApplicationMaster runs a specific YARN job and is responsible for negotiating resources from the ResourceManager and also working with the NodeManager to execute and monitor Containers.

The actual data processing occurs within the Containers executed by the ApplicationMaster. A Container grants rights to an application to use a specific amount of resources (memory, CPU, etc.) on a specific host.

YARN is not the only new major feature of Hadoop 2.0. HDFS has undergone a major transformation with a collection of new features that include:

- **NameNode HA:** automated failover with a hot standby and resiliency for

theNameNode masters service.

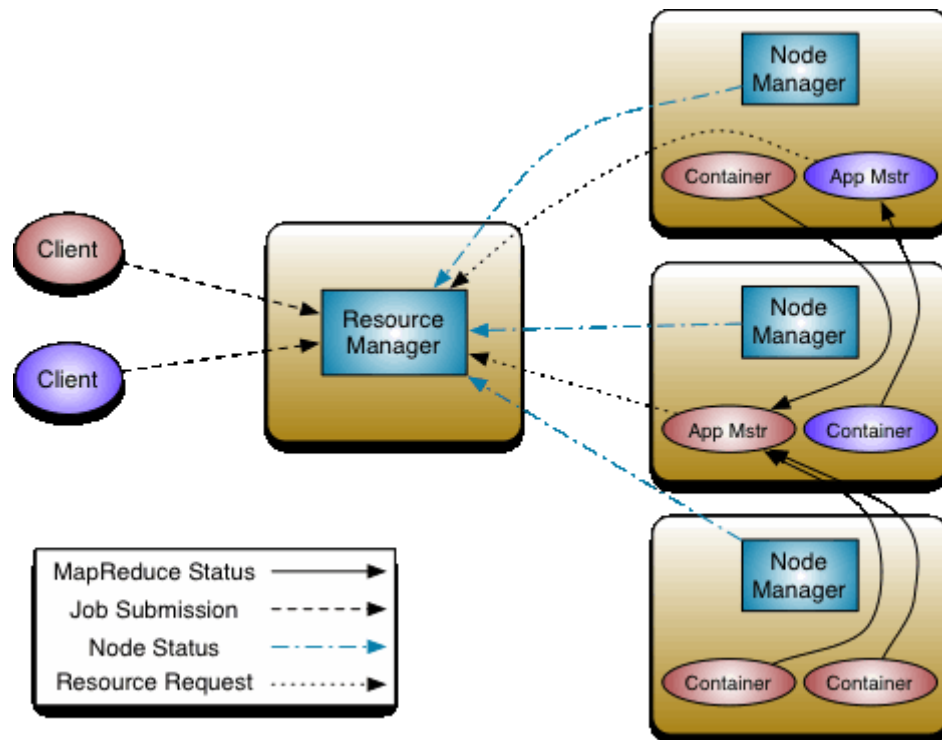
- **Snapshots:** point-in-time recovery for backup, disaster recovery and protection against user errors.
- **Federation:** a clear separation of namespace and storage by enabling generic block storage layer.

MRv2-YARN

The new architecture introduced in Hadoop-

0.23, divides the two major functions of the JobTracker: resource management and job life-cycle management into separate components. The new Resource Manager manages the global assignment of compute resources to applications and the per-application Application Master manages the application's scheduling and coordination.

An application is either a single job in the sense of classic MapReduce jobs or a DAG of such jobs. The Resource Manager and per-machine Node Manager daemon, which manages the



user processes on that machine, form the computation fabric. The per-application Application Master is, in effect, a framework specific library and is tasked with negotiating resources from the Resource Manager and working with the Node Manager(s) to execute and monitor the tasks. The fundamental idea of YARN is to split up the functionalities of resource management and jobscheduling/monitoring into separate daemons. The idea is to have a global Resource Manager (RM) and per-application Application Master (AM). An application is either a single job or a DAG of jobs. The Resource Manager and the Node Manager form the data-computation framework. The Resource Manager is the ultimate authority that arbitrates resources among all the applications in the system. The Node Manager is the per-machine framework agent who is responsible for containers, monitoring their resource usage (cpu, memory, disk, network) and reporting the same to the Resource Manager/Scheduler. The per-application

ApplicationMaster

is, in effect, a framework specific library and is tasked with negotiating resources from the Resource Manager and working with the Node Manager(s) to execute and monitor the tasks.

Figure 5.8 YARN workflow

The Resource Manager has two main components: Scheduler and Applications Manager. The Scheduler is responsible for allocating resources to the various running applications subject to familiar constraints of capacities, queues etc. The Scheduler is a pure scheduler in the sense that it performs no monitoring or tracking of status for the application. Also, it offers no guarantees about restarting failed tasks either due to application failure or hardware failures. The Scheduler performs its scheduling function based on the resource requirements of the applications; it does so based on the abstract notion of a resource *Container* which incorporates elements such as memory, CPU, disk, network etc.

NameNode High Availability

Prior to Hadoop 2.0.0, the NameNode was a single point of failure (SPOF) in an HDFS

cluster. Each cluster had a single NameNode, and if that machine or process became unavailable, the

cluster as a whole would be unavailable until the NameNode was either restarted or brought upon a separate machine.

This impacted the total availability of the HDFS cluster in two major ways:

- In the case of an unplanned event such as a machine crash, the cluster would be unavailable until an operator restarted the NameNode.
- Planned maintenance events such as software or hardware upgrades on the NameNode machine would result in windows of cluster downtime.

The HDFS High Availability feature addresses the above problems by providing

the option of running two redundant NameNodes in the same cluster in an Active/Passive configuration with a hot standby. This allows a fast failover to a new NameNode in the case that a machine crashes, or a graceful administrator-initiated failover for the purpose of planned maintenance.

5.7 YARN Architecture

In a typical HA cluster, two separate machines are configured as NameNodes. At any point in time, exactly one of the NameNodes is in an *Active* state, and the other is in a *Standby* state. The Active NameNode is responsible for all client operations in the cluster, while the Standby simply acts as a slave, maintaining enough state to provide a fast failover if necessary. In order for the Standby node to keep its state synchronized with the Active node, the current implementation requires that the two nodes both have access to a directory on a shared storage device (eg an NFS mount from a NAS). This restriction will likely be relaxed in future versions. When any namespace modification is performed by the Active node, it durably logs a record of the modification to an edit log file stored in the shared directory. The Standby node is constantly watching this directory for edits, and as it sees the edits, it applies them to its own namespace. In the event of a failover, the Standby will ensure that it has read all of the edits from the shared storage before promoting itself to the Active state. This ensures that the namespace state is fully synchronized before a failover occurs.

In order to provide a fast failover, it is also necessary that the Standby node have up-to-date information regarding the location of blocks in the cluster. In order to achieve this, the DataNodes are configured with the location of both NameNodes, and send block location information and heartbeats to both.

Hardware resources

In order to deploy an HA cluster, you should prepare the following:

- **NameNode machines-**
the machines on which you run the Active and Standby NameNodes should have equivalent hardware to each other, and equivalent hardware to what would be used in a non-HA cluster.

- **Shared storage** - you will need to have a shared directory which both NameNode machines can have read/write access to. Typically this is a remote file system which supports NFS and is mounted on each of the NameNode machines. Currently only a single shared edit directory is supported. Thus, the availability of the system is limited by the availability of this shared edit directory, and therefore in order to remove all single points of failure there needs to be redundancy for the shared edit directory.

HDFS Federation

HDFS Federation improves the existing **HDFS** architecture through a clear separation of namespace and storage, enabling a generic block storage layer. It enables support for multiple namespaces in the cluster to improve scalability and isolation.

This guide provides an overview of the HDFS Federation feature and how to configure and manage the federated cluster.

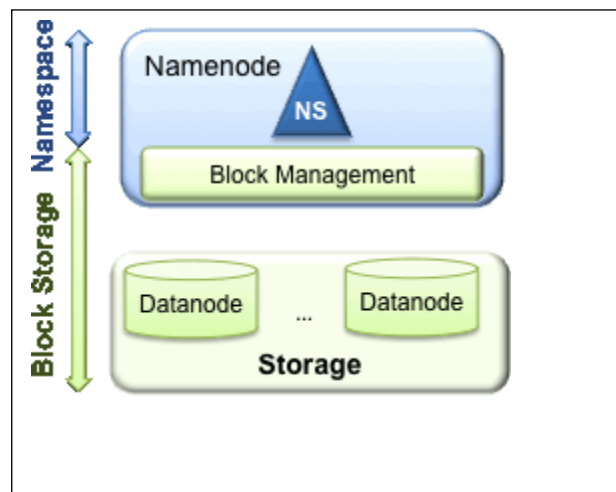


Figure 5.9 Hadoop Storage Representation

HDFS has two main layers:

Namespace

- Consists of directories, files and blocks.

It supports all the namespace related filesystem operations such as create, delete, modify and list files and directories.

Block Storage Service, which has two parts:

- Block Management (performed in the NameNode)

Provides DataNode cluster membership by handling registrations, and periodic heartbeats.

Processes block reports and maintains location of blocks.

Supports block related operations such as create, delete, modify and get block location.

Manages replica placement, block replication for under replicated blocks, and deletes blocks that are over replicated.

Storage - is provided by DataNodes by storing blocks on the local filesystem and allowing read/write access.

The prior HDFS architecture allows only a single namespace for the entire cluster. In that configuration, a single NameNode manages the namespace. HDFS Federation addresses this limitation by adding support for multiple NameNodes/namespaces to HDFS.

Key Benefits

- **Namespace Scalability** - Federation adds namespace horizontal scaling. Large deployments or deployments using lot of small files benefit from namespace scaling by allowing more NameNodes to be added to the cluster.
- **Performance** - File system throughput is not limited by a single NameNode. Adding more NameNodes to the cluster scales the filesystem read/write throughput.
- **Isolation** - A single NameNode offers no isolation in a multi-user environment. For example, an experimental application can overload the NameNode and slow down production critical applications. By using

multiple Namenodes, different categories of applications and users can be isolated to different namespaces.

Federation configuration is backward compatible and allow existing single Namenode configuration to work without any change. The new configuration is designed such that all the nodes in the cluster have the same configuration without the need for deploying different configurations based on the type of the node in the cluster.

Configuration:

Step 1: Add the `dfs.nameservices` parameter to your configuration and configure it with a list of comma-separated Name Service IDs. This will be used by the Datanodes to determine the Namenodes in the cluster.

Step 2: For each Namenode and Secondary Namenode/BackupNode/Checkpointers add the following configuration parameters suffixed with the corresponding Name Service ID into the common configuration file:

Formatting Namenodes

Step 1: Format a Namenode using the following command:

```
[hdfs]${HADOOP_PREFIX}/bin/hdfsnamenode-format[-clusterId<cluster_id>]
```

Choose a unique `cluster_id` which will not conflict to other clusters in your environment. If a `cluster_id` is not provided, then a unique one is autogenerated.

Step 2: Format additional Namenodes using the following command:

```
[hdfs]${HADOOP_PREFIX}/bin/hdfsnamenode-format-clusterId<cluster_id>
```

Upgrading from an older release and configuring federation

Older releases only support a single Namenode. Upgrade the cluster

to new release in order to enable federation. During upgrade, you can provide a ClusterID as follows:

```
[hdfs]$ $HADOOP_PREFIX/bin/hdfs start namenode --config $HADOOP_CONF_DIR -  
upgrade-clusterId<cluster_ID>
```

If cluster_id is not provided, it is auto generated.

Adding a new Namenode to an existing HDFS cluster

Perform the following steps:

Add dfs.nameservice to the configuration.

Update the configuration with the NameServiceID suffix. Configuration key names changed post release 0.20. You must use the new configuration parameter names in order to use federation.

Add the new Namenode related config to the configuration file.

Propagate the configuration file to all the nodes in the cluster.

Start the new Namenode and Secondary/Backup.

Refresh the Datanodes to pick up the newly added Namenode by running the following command again on all the Datanodes in the cluster:

```
[hdfs]$ $HADOOP_PREFIX/bin/hdfsdfsadmin-refreshNameNodes  
<datanode_host_name>:<datanode_rpc_port>
```

Managing a cluster

To start the cluster, run the following command:

```
[hdfs]$ $HADOOP_PREFIX/sbin/start-dfs.sh To stop the cluster run
```

the following command:

```
[hdfs]$ $HADOOP_PREFIX/sbin/stop-dfs.sh
```

Balancer

The Balancer has been changed to work with multiple Namenodes. The Balancer can be run using the command:

```
[hdfs]$ $HADOOP_PREFIX/sbin/hadoop-daemon.sh start balancer [-policy <policy>]
```

Decommissioning

Decommissioning is similar to prior releases. The nodes that need to be decommissioned are added to the exclude file at all of the Namenodes. Each Namenode decommissions its Block Pool. When all the Namenodes finish decommissioning a Datanode, the Datanode is considered decommissioned.

Step 1: To distribute an exclude file to all the Namenodes, use

the following command: [hdfs]\$ \$HADOOP_PREFIX/sbin/distribute-

exclude.sh <exclude_file>

Step 2: Refresh all the Namenodes to pick up the new exclude

file: [hdfs]\$ \$HADOOP_PREFIX/sbin/refresh-namenodes.sh

The above command uses HDFS configuration to determine the configured Namenodes in the cluster and refreshes them to pick up the new exclude file.

Cluster Web Console

Similar to the Namenode status web page, when using federation a Cluster Web Console is available to monitor the federated cluster at http://<any_nn_host:port>/dfsclusterhealth.jsp. Any Namenode in the cluster can be used to access this web page.

Migrating from MapReduce1(MRv1) to MapReduce2

MapReduce from Hadoop 1 (MapReduce MRv1) has been split into two components. The cluster resource management capabilities have become YARN (Yet Another Resource Negotiator), while the MapReduce-specific capabilities remain MapReduce. In the MapReduce MRv1 architecture, the cluster was managed by a service called the JobTracker. TaskTracker services lived on each host and would launch tasks on behalf of jobs. The JobTracker would serve information about completed jobs.

In MapReduce MRv2, the functions of the JobTracker have been split between three services. The ResourceManager is a persistent YARN service that receives and runs applications (a MapReduce job is an application) on the cluster. It contains the scheduler, which, as previously, is pluggable. The MapReduce-specific capabilities of the JobTracker have been moved into the MapReduce ApplicationMaster, one of which is started to manage each MapReduce job and terminated when the job completes. The JobTracker function of serving information about completed jobs has been moved to the JobHistoryServer. The TaskTracker

has been replaced with the NodeManager, a YARN service that manages resources and deployment on a host. It is responsible for launching containers, each of which can house a map or reduce task. Nearly all jobs written for MRv1 will be able to run without any modifications on an MRv2 cluster.

The new architecture has its advantages. First, by breaking up the JobTracker into a few different services, it avoids many of the scaling issues faced by MapReduce in Hadoop 1. More importantly, it makes it possible to run frameworks other than MapReduce on a Hadoop cluster. For example, Impala can also run on YARN and [share resources](#) with MapReduce.

Configuration Migration

Since MapReduce 1 functionality has been split into two components, MapReduce cluster configuration options have been split into YARN configuration options, which go in yarn-site.xml, and MapReduce configuration options, which go in mapred-site.xml. Many have been given new names to reflect the shift. As JobTrackers and TaskTrackers no longer exist in MRv2, all configuration options pertaining to them no longer exist, although many have corresponding options for the Resource Manager, NodeManager, and JobHistoryServer.

A minimal configuration required to run MRv2 jobs on YARN is: yarn-site.xml configuration

```
<?xml version="1.0" encoding="UTF-8"?>
<configuration>
  <property>
    <name>yarn.resourcemanager.hostname</name>
    <value>you.hostname.com</value>
  </property>

  <property>
    <name>yarn.nodemanager.aux-services</name>
    <value>mapreduce_shuffle</value>
  </property>
</configuration>
```

mapred-site.xml configuration

```
<?xmlversion="1.0"encoding="UTF-8"?>
<configuration>
<property>
<name>mapreduce.framework.name</name>
<value>yarn</value>
</property>
</configuration>
```

Below is a table with HA-related configurations used in MRv1 and their equivalents in YARN:

MRv1	YARN/ MRv2	Comment
mapred.jobtrackers.name	yarn.resourcemanager.ha.rm-ids	
mapred.ha.jobtracker.id	yarn.resourcemanager.ha.id	Unlike in MRv1, this must be configured in YARN.
mapred.jobtracker.rpc-address.name.id	yarn.resourcemanager.rpc-address.id	YARN/ MRv2 has different RPC ports for different functionalities. Each port-related configuration must be suffixed with an id. Note that there is no <name> in YARN.
mapred.ha.jobtracker.rpc-address.name.id	yarn.resourcemanager.ha.admin.address	
mapred.ha.fencing.methods	yarn.resourcemanager.ha.fencer	Not required to be specified
mapred.client.failover.*	None	Not required
	yarn.resourcemanager.ha.enabled	Enable HA
mapred.jobtracker.restart.recover	yarn.resourcemanager.recovery.enabled	Enable recovery of jobs after failover
	yarn.resourcemanager.store.class	org.apache.hadoop.yarn
		server.resourcemanager.recovery.ZKRMStateStore
mapred.ha.automatic-failover.enabled	yarn.resourcemanager.ha.automatic-failover.enabled	Enable automatic failover
mapred.ha.zkfc.port	yarn.resourcemanager.ha.automatic-failover.port	

mapred.job.tracker	yarn.resourcemanager.cluster.i	
--------------------	--------------------------------	--

Programming in YARN Framework

Under Hadoop 2.0, MapReduce is but one instance of a YARN application, where YARN has taken center stage as the operating system of Hadoop. Because YARN allows *any* application to run on equal footing with MapReduce, it opened the floodgates for a new generation of software applications with these kinds of features:

More programming models. Because YARN supports any application that can divide itself into parallel tasks, they are no longer shoehorned into the palette of mappers, combiners, and reducers. This in turn supports complex data-flow applications like ETL and ELT, and iterative programs like massively-parallel machine learning and modeling.

Integration of native libraries. Because YARN has robust support for *any* executable – not limited to MapReduce, and not even limited to Java – application vendors with a large mature code base have a clear path to Hadoop integration.

Support for large referenced data. YARN automatically localizes and caches large reference datasets, making them available to all nodes for data local processing. This supports legacy functions like address standardization, which require large reference data sets that cannot be accessed from the Hadoop Distributed File System (HDFS) by the legacy libraries. Of course, MapReduce isn't the only option for processing data at scale using Hadoop. Tools like Pig (a large scale query and analysis system), Hive (a data warehousing application) and others have been available for some time. These tools can express transforms and analysis using more accessible constructs: Hive uses HQL, a language similar to SQL. Pig provides a script language (Pig Latin) to create MapReduce jobs. Business analysts familiar with conventional tools like SQL and SAS should be able to use these tools to write programs to solve large data problems on Hadoop clusters.

The Value of Visual Application Development Tools

A new generation of visual design application development tools could help solve these coding problems. By running as native YARN applications and side-stepping

the need for MapReduce, some of these programs eliminate coding altogether. Other tools reduce coding by generating MapReduce code or by generating scripts like Pig.

Visual designers are powerful for several reasons:

- Increased level of abstraction: Instead of thinking about classes and methods, users see operations, data, and outcomes.
- Fast—what-if!?: The drag-and-connect interface supports quick try/observe/adjust cycles.
- Automatic optimization: Scaling and efficiency are built-in.
- High-level palette: High-level constructs like standardize address, deduplicate consumers, or parse names are often directly on the designer palette.

How does this look in practice? Here's an illustration that shows how three competing approaches differ:

- MapReduce written in Java
- Pig scripts developed from scratch
- A visually-designed process running a native YARN ETL application. The application is from RedPoint Global, but comparable approaches can be seen in Talend and Actian.

Using these three approaches, we conducted a `WordCount` test on 30,000 files (20 gigabytes) of Project Gutenberg books. This test reads lines of text, breaks them into words, and creates a concordance (list of words and the number of times each occurs). Our Hadoop cluster was small—only four nodes—but was large enough to demonstrate the concepts and tradeoffs

MapReduce:

Set-up time: While flexible, MapReduce had the longest learning curve and required significant coding skills—both as a Java programmer and a MapReduce specialist—to prepare the test.

Performance: It took 3 hours 20 minutes to run the test initially due to the `small files problem` that is familiar to seasoned MapReduce programmers. This problem occurs when reading large collections of small files, because MapReduce's default behavior is to assign a mapper task to each file. This results in a huge number of tasks. To address this issue, we created a

customInputFormat class to read multiple files at once. This reduced our run time to 58 minutes. Then we tuned the split sizes and mapper task limit appropriately, which dropped the run time to about six minutes.

Comments: Each performance improvement came at a cost. Overall, nearly a full day of programmer time was spent optimizing the original code.

Pig:

Set-up time: Learning Pig was fairly easy. It was pretty natural to create the coding for this test. However to make a common adjustment in the code—changing the set of whitespace separator to include punctuation—required the addition of a User-defined function or UDF which had to be written in Java.

Pig is generally easy enough to use by people who aren't professional programmers but who know how to write scripting languages like JavaScript or Visual

Basic. Performance: The results weren't stellar: run time was close to 15 minutes.

Comments: While coding took less time, Java programming was ultimately required to meet the test requirements.

Sample Pig script without the User Defined Function:

```
SET pig.maxCombinedSplitSize 67108864
SET pig.splitCombination true
A = LOAD '/testdata/pg/**/*.txt';
B = FOREACH A GENERATE FLATTEN(TOKENIZE((chararray)$0)) AS word;
C = FOREACH B GENERATE UPPER(word) AS word;
D = GROUP C BY word;
E = FOREACH D GENERATE COUNT(C) AS occurrences, group;
F = ORDER E BY occurrences DESC;
STORE F INTO '/user/cleonardi/pg/pig-count';
```

YARN-enabled ETL/ELT designer:

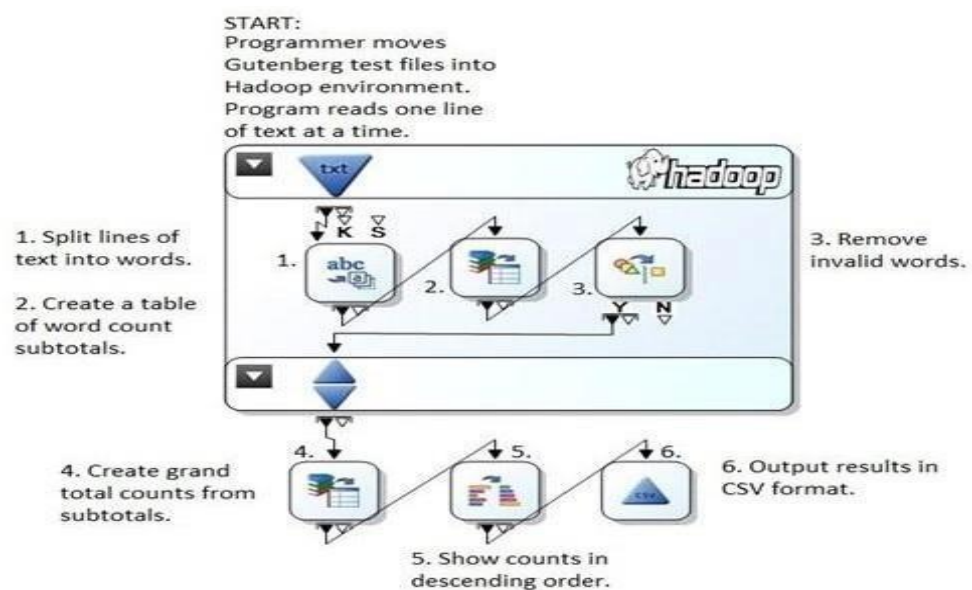
Set-up time: The tool is designed to have a shorter learning curve than even Pig scripting. Dragging tools like -Delimited Input, -Summarize and -Tokenize from the palette and configuring them is designed to be discoverable and intuitive, and the resulting diagram has a one-to-one correspondence between icons and operations. There's no need for coding or learning a language like Java or Pig.

The visual design covers the input file format, tokenizing and counting steps. The resulting dataflow graph contains seven icons along with a grouping construct that shows what executes inside Hadoop. Each icon represents a step in the data transformation.

Performance: The run time for this data flow is just over three minutes with no tuning. *Comments:* Because there is no code to manage, and editing is done visually, running what-if scenarios is quick for non-programmers.

Once the data flow is designed, it can be stored and saved for later use. In addition, the logic can be captured into a macro for sharing and reuse between multiple data flows.

Figure 5.10: Sharing and reusing between multiple data flows



5.8 Scheduling in YARN

5.8.1 Hadoop Fair Scheduler:

FairScheduler, a pluggable scheduler for Hadoop that allows YARN applications to share resources in large clusters fairly.

Fair scheduling is a method of assigning resources to applications such that all apps get, on average, an equal share of resources over time. HadoopNextGen is capable of scheduling multiple resource types. By default, the Fair Scheduler bases scheduling fairness decisions only on memory. It can be configured to schedule with both memory and CPU, using the notion of Dominant Resource Fairness developed by Ghodsi et al. When there is a single app running, that app uses the entire cluster. When other apps are submitted, resources that free up are assigned to the new apps, so that each app eventually gets roughly the same amount of resources. Unlike the default Hadoop scheduler, which forms a queue of apps, this lets short apps finish in reasonable time while not starving long-lived apps. It is also a reasonable way to share a cluster between a number of users. Finally, fair sharing can also work with app priorities - the priorities are used as weights to determine the fraction of total resources that each app should get.

The scheduler organizes apps further into “queues”, and shares resources fairly between these queues. By default, all users share a single queue, named “default”. If an app specifically lists a queue in a container resource request, the request is submitted to that queue. It is also possible to assign queues based on the user name included with the request through configuration. Within each queue, a scheduling policy is used to share resources between the running apps. The default is memory-based fair sharing, but FIFO and multi-resource with Dominant Resource Fairness can also be configured. Queues can be arranged in a hierarchy to divide resources and configured with weights to share the cluster in specific proportions.

In addition to providing fair sharing, the Fair Scheduler allows assigning guaranteed minimum shares to queues, which is useful for ensuring that certain users, groups or production applications always get sufficient resources. When a queue contains apps, it gets at least its minimum share, but when the queue does not need its full guaranteed share, the excess is split between other running apps. This lets the scheduler guarantee capacity for queues while utilizing resources efficiently when these queues don’t contain applications.

The Fair Scheduler lets all apps run by default, but it is also possible to limit the number of running apps per user and per queue through the config file. This can be useful when a user must submit hundreds of apps at once, or in general to improve performance if running too many apps at once would cause too much intermediate data to be created or too much context-switching. Limiting the apps does not cause any subsequently submitted apps to fail, only to wait in the scheduler’s queue until some of the user’s earlier apps finish.

Two things to note about Fair Scheduler in YARN are-

1. By default, the Fair Scheduler bases scheduling fairness decisions only on memory. It can be configured to schedule with both memory and CPU.
2. The scheduler organizes apps further into “queues”, and shares resources fairly between these queues.

As example– If there are two queues **sales** and **finance**. A job is submitted to sales queue, being a sole running job it will get all the resources. Now a job is submitted to finance queue which will result in the new job gradually getting half of the resources. So jobs in both of the queues will have 50% of the resources each. Now another job is submitted to finance queue that will result in half of the resources allocated to finance queue allocated to this new job. So two jobs in the finance queue will now share resources allocated to finance queue (50% of the total resources) in equal proportions where as job in sales queue will use whole 50% of the resources allocated to sales queue.

Hierarchical Queues support:

The fair scheduler in YARN supports hierarchical queues which means an organization can create sub-queues with in its dedicated queue. All queues descend from a queue named “root”. Available resources are distributed among the children of the root queue in the typical fair scheduling fashion. Then, the children distribute the resources assigned to them to their children in the same fashion.

Fair Scheduler Configuration

To use the Fair Scheduler in YARN first assign the appropriate scheduler class in **yarn-site.xml**:

```
<property>
<name>yarn.resourcemanager.scheduler.class</name>
<value>org.apache.hadoop.yarn.server.resourcemanager.scheduler.fair.FairScheduler</value>
</property>
```

Setting up queues

Properties for setting up queues are as follows. These changes are done in the configuration file **etc/hadoop/fair-scheduler.xml**.

<queue> element– Represent queues. Some of the important properties of queue element are as following.

- **minResources**: minimum resources the queue is entitled to, in the form “X mb, Y vcores”. If a queue’s minimum share is not satisfied, it will be offered available resources before any other queue under the same parent.
- **maxResources**: maximum resources a queue is allocated, expressed either in absolute values (X mb, Y vcores) or as a percentage of the cluster resources (X% memory, Y% CPU).
- **weight**: to share the cluster non-proportionally with other queues. Weights default to 1, and a queue with weight 2 should receive approximately twice as many resources as a queue with the default weight.
- **schedulingPolicy**: to set the scheduling policy of any queue. The allowed values are “fifo”, “fair”, “drf” or any class that extends org.apache.hadoop.yarn.server.resourcemanager.scheduler.fair.SchedulingPolicy. Defaults to “fair”.

<defaultQueueSchedulingPolicy> element- Which sets the default scheduling policy for queues; overridden by the schedulingPolicy element in each queue if specified. Defaults to “fair”.

<queueMaxAppsDefault> element- Which sets the default running app limit for queues; overridden by maxRunningApps element in each queue.

<queuePlacementPolicy> element– This element contains a list of rule elements that tell the scheduler how to place incoming apps into queues. Rules are applied in the order that they are listed. All rules accept the “create” argument, which indicates whether the rule can create a new queue. “Create” defaults to true; if set to false and the rule would place the app in a queue that is not configured in the allocations file, we continue on to the next rule.

Valid rules are as follows:

- **specified:** The application is placed into the queue it requested.
- **user:** The application is placed into a queue with the name of the user who submitted it.
- **primaryGroup:** The application is placed into a queue with the name of the primary group of the user who submitted it.
- **secondaryGroupExistingQueue:** The application is placed into a queue with a name that matches a secondary group of the user who submitted it.
- **nestedUserQueue:** The application is placed into a queue with the name of the user under the queue suggested by the nested rule.
- **default:** The application is placed into the queue specified in the ‘queue’ attribute of the default rule. If ‘queue’ attribute is not specified, the app is placed into ‘root.default’ queue.
- **reject:** The application is rejected.

Queue configuration example:

If there are two top level child queues **sales** and **finance** (Descending from root). Within sales queues there are two sub-queues **apac** and **emea** then the queues can be set up to use fair scheduler as given below-

```
<allocations>
<queuename="sales">
<minResources>10000 mb,0vcores</minResources>
<maxResources>50000 mb,0vcores</maxResources>
<weight>2.0</weight>
<schedulingPolicy>fifo</schedulingPolicy>
<queuename="emea"/>
<queuename="apac"/>
</queue>
<queuename="finance">
<minResources>10000 mb,0vcores</minResources>
<maxResources>70000 mb,0vcores</maxResources>
<weight>3.0</weight>
<schedulingPolicy>fair</schedulingPolicy>
</queue>
<queuePlacementPolicy>
<rulename="specified"/>
<rulename="primaryGroup"create="false"/>
<rulename="default"queue="finance"/>
</queuePlacementPolicy>
</allocations>
```

5.8.2 Capacity Scheduler

d YARN defines a minimum allocation and a maximum allocation for the resources it is scheduling for: Memory and/or Cores today. Each server running a worker for YARN has a NodeManager that is providing an allocation of resources which could be memory and/or cores that can be used for scheduling. This aggregate of resources from all Node Managers is provided as the 'root' of all resources the capacity scheduler has available.

The fundamental basics of the Capacity Scheduler are around how queues are laid out and resources are allocated to them. Queues are laid out in a hierarchical design with the topmost parent being the 'root' of the cluster queues, from here leaf (child) queues can be assigned from the root, or branches which can have leafs on themselves. Capacity is assigned to these queues as min and max percentages of the parent in the hierarchy. The minimum capacity is the amount of resources the queue should expect to have available to it if everything is running maxed out on the cluster. The maximum capacity is an elastic like capacity that allows queues to make use of resources which are not being used to fill minimum capacity demand in other queues.

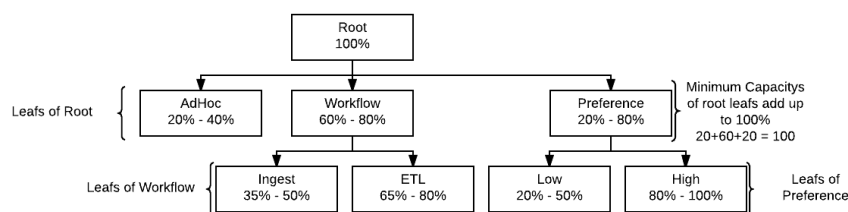


Fig 5.11 Capacity Scheduler

Children Queues like in the figure above inherit the resources of their parent queue. For example, with the Preference branch the Low leaf queue gets 20% of the Preference 20% minimum capacity while the High leaf gets 80% of the 20% minimum capacity. Minimum Capacity always has to add up to 100% for all the leafs under a parent.

Minimum User Percentage and User Limit Factor

Minimum User Percentage and User Limit Factor are ways to control how resources get assigned to users within the queues they are utilizing. The Min User Percentage is a soft limit on the smallest amount of resources a single user should get access to if they are requesting it. For example a min user percentage of 10% means that 10 users would each be getting 10% assuming they are all asking for it; this value is soft in the sense that if one of the users is asking for less we may place more users in queue.

User Limit Factor is a way to control the max amount of resources that a single user can consume. User Limit Factor is set as a multiple of the queues minimum capacity where a user limit factor of 1 means the user can consume the entire minimum capacity of the queue. If the user limit factor is greater than 1 it's possible for the user to grow into

maximum capacity and if the value is set to less than 1, such as 0.5, a user will only be able to obtain half of the minimum capacity of the queue. Should you want a user to be able to also grow into the maximum capacity of a queue setting the value greater than 1 will allow the minimum capacity to be overtaken by a user that many times.

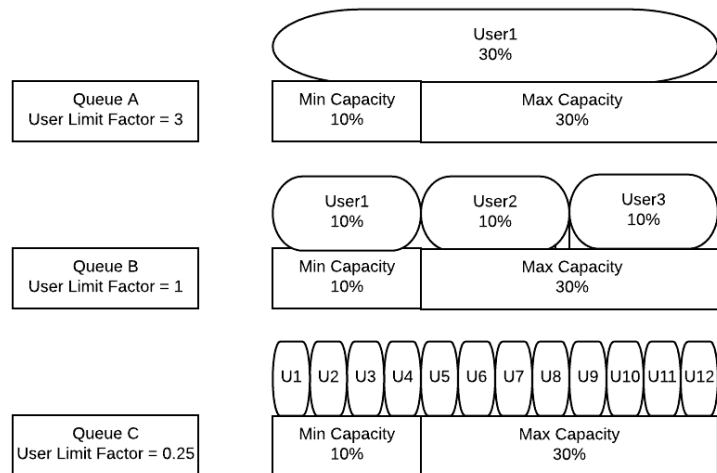


Fig 5.12 Queuing in capacityscheduler

A common design point that may initially be non-intuitive is creation of queues by workloads and not by applications and then using the user-limit-factor to prevent individual takeover of queues by a single user by using a value of less than 1.0. This model supports simpler operations by not allowing queue creation to spiral out of control by creating one per LoB but by creating queues by workloads to create predictable queue behaviors. Once a cluster is at full use and has applications waiting in line to run the lower user-limit-factors will be key to controlling resource sharing between tenants.

Initially this may not provide the cluster utilization at the start of their Hadoop platform journey due to using a smaller user limit to limit user resources, there are many approaches but one to consider is that initially it may be justified to allow a single tenant to take over a tiny cluster (say 10 nodes) but when expanding the platform footprint lower the user limit factor so that each tenant keeps the same amount of allocatable resources as before you added new nodes. This lets the user keep his initial experience, but be capped off from taking over the entire cluster as it grows making room for new users easily to get allocated capacity while not degrading the original user's experience.

Archetypes

Designing queue archetypes to describe effective behavior of the tenant in the queue provides a means to measure changes against to see if they align or deviate from the expectation. While by no means a complete list of workload behaviors the below is a good list to start from. Create your own definitions of the types of queues your organization needs as the types of patterns applications and end users are performing.

- **AD-HOC**

- This is where random user queries, unknown, and new workloads may be ran, there is no expectation as to resource allocation behaviors but can work as a good place to initially run applications to get a feel for per application tuning needs.

- **PREFERENCE**

- These are applications that should get resources before and retain them longer. This could be for many reasons such as catch up applications, emergency runs, or other operational needs.

- **MACHINE LEARNING**

- Machine Learning applications can be typically characterized by their long run times and large or intensive resource requirements. Termination of a task for some machine learning workloads can have long running impacts by greatly extending the duration

- **DASHBOARDING**

- Low concurrency (refresh rate) but high number of queries per day. Dashboards need to be refreshed in an expected time but are very predictable workloads

- **EXPLORATION**

- Exploration users have a need for low latency queries and need throughput to churn through very large datasets. Resources will likely be held for use the entire time the user is exploring to provide an interactive experience

- **BATCH WORKFLOW**

- Designed to provide general computing needs for transformation and batch workloads. Setup is concerned most frequently with throughput of applications not individual application latency

- **ALWAYS ON**

- Applications that are always running with no concept of completion. Applications that keep resources provisioned while waiting for new work to arrive. Slider deployed instances such as LLAP

Container Churn

Churn within a queue is best described as a constant existing and starting of new containers. This behavior is very important to have a well behaving cluster were queues are quickly rebalanced to their minimum capacity and to balance capacity of the queue between its users fairly. The anti-pattern to churn is the long lived container that allocates itself and never releases as it can prevent proper rebalancing of resources in some cases it can completely block other applications from launching or queues from getting their capacity back. When not using preemption if a queue was to grow into elastic space but never release

its containers the elastic space will never be given back to queues that have been guaranteed it. If a application running in your queues has long lived containers make special note and consider ways to mitigate its impact to other users with features like user-limit-factors, preemption, or even dedicated queues without elastic capacity.

Features & Behaviors of the Capacity Scheduler

CPU Scheduling (Dominant Resource Fairness)

CPU Scheduling is not enabled by default and allows for the oversubscription of cores with no enforcement of use or preferred allocation taken into account. Many batch driven workloads may experience a throughput reduction if CPU scheduling is used, but may be required for strict SLA guarantees. A method called Dominant Resource Fairness (DRF) is used to decide what resource type to weight the utilization as: DRF takes the most used resource and treats you as using the highest percent for scheduling.

There are 2 primary parts to CPU Scheduling

1. Allocation and Placement
2. Enforcement

Allocation and placement is solved simply by enabling CPU Scheduling so that the Dominant Resource Fairness algorithm and VCores Node Manager's report begin getting used by the scheduler. This solves for some novel problems like an end user that used to schedule his Spark applications with 1 or 2 executors but 8 cores per, and then all other tasks that ran on these nodes due to the excess of available memory were impacted. By enabling simple CPU Scheduling other tasks would no longer be provisioned onto the servers were all the cores are being utilized and find other preferred locations for task placement in the cluster.

Enforcement is solved by utilizing CGroups. This allows for YARN to ensure that if a task asked for a single vcore that they cannot utilize more than what was requested. Sharing of vcores when not used can be enabled or a strict enforcement of only what was scheduled can be done. The Node Manager's can also be configured as to what the max amount of CPU use on the server they will allow all tasks to sum up to, this allows for cores to be guaranteed to OS functions.

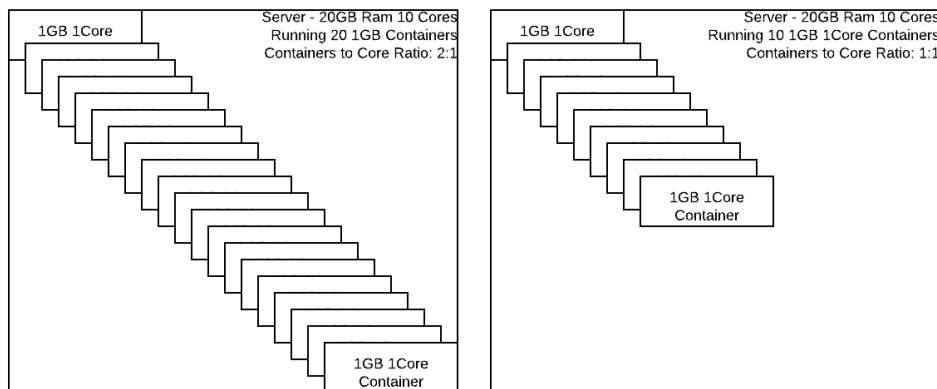


Fig 5.13 Resource allocation to CPU cores

The above figure gives an idea of how many concurrent containers can change if limited to the smallest resource (typically CPU cores.) It's unlikely a true 1:1 Core to Container ratio is what will be required but this aspect of tuning is best left to monitoring historical system metrics and then increasing or decrease the number of NodeManagerVCores available for scheduler to allow more or less containers on a given server group.

Preemption

When applications are making use of elastic capacity in their queues and another application comes along to demand back their minimum capacity (which is being used in another queue as elastic) traditionally the application would have to wait for the task to finish to get its resource allocations. With Preemption enabled resources in other queues can be reclaimed to provide the minimum capacity back to the queue that needs it. Preemption will try not to outright kill a application, and will take reducers last as they have to repeat more work then mappers if they have to re-run. From an ordering perspective Preemption looks at the youngest application and most over-subscribed ones first for task reclamation.

Preemption has some very specific behaviors and some of them don't function as expected for users. The most commonly expected behavior is for the queues to preempt within themselves to balance the resources over all users. That assumption is wrong as preemption only works across queues today, resources that are mis-balanced within a queue between users needs to look into other ways to control this such as User Limit Factors, Improved Queue Churned, and Queue FIFO/FAIR Policies.

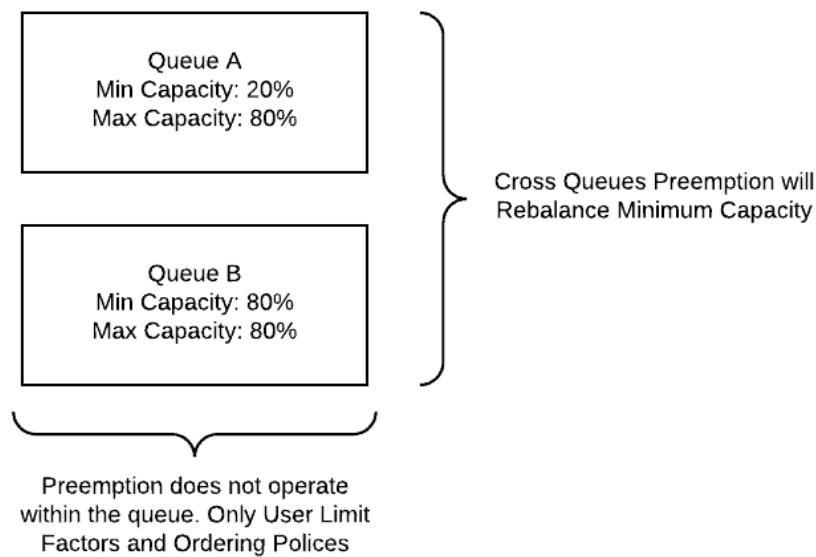


Figure 5.14: FIFO fair policy

Another behavior is that Preemption will not preempt resources if it cannot provide enough to fulfill another resource allocation request. While typically this is not an issue with large cluster for small clusters with large maximum container sizes could run into a scenario where Preemption is not configured to reclaim a container of the largest possible size and therefore will do nothing at all. The main properties used to control this behavior is the Total PreemptionPer Round and the Natural Termination Factor. Total PreemptionPer Round is the percentage of resources on the cluster that can be preempted at once, the Natural Termination Factor is the percent of resources out of the total cluster (100%) requested that will be preempted up to the Total Preemption Per Round.

The last misconception is that Preemption will rebalance max (elastic) capacity use for between queues that wish to grow into it. Maximum Capacity is given on a first come first serve basis only. If a single queue has taken over all the of the clusters capacity, and another application start in a queue that needs its minimum capacity back, only the minimum capacity will be preempted and all of the maximum capacity being used by the other queue will remain until the containers churn naturally.

Queue Ordering Policies

There are currently two types of ordering policies the Capacity Scheduler supports today: FIFO and FAIR. The default that queues start with is FIFO which in my experience is not the behavior that clients expect from their queues. By configuring between FIFO and FAIR at the Queue Leaf level you can create behaviors that lead to throughput driven processing or shared fair processing between applications running. One important thing to know about

ordering policies is that they operated at the application level in the queue and do not care about which user owns the applications.

With the FIFO policy applications are evaluated in order of oldest to youngest for resource allocation. If an application has outstanding resource requests they are immediately fulfilled first come first serve. The result of this is that should an application have enough outstanding requests that they would consume the entire queue multiple times before completing they will block out other applications from starting while they are first allocated resources as the oldest application. It is this very behavior that comes most unexpected to users and creates the most discontent as a users can even block their own applications with their own applications!

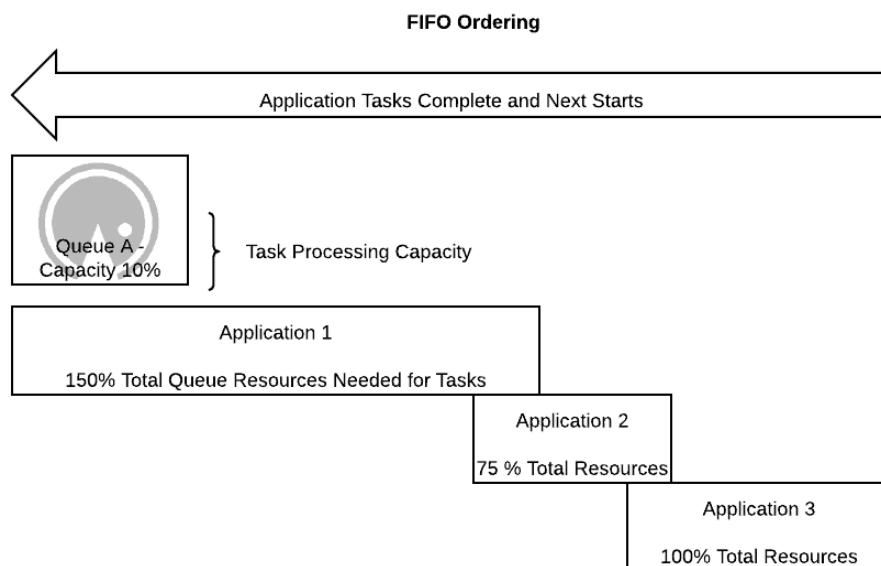


Figure 5.15 FIFO Ordering

That behavior is easily solved for today by using the FAIR ordering policy. When using the FAIR ordering policy on a leaf queue applications are evaluated for resource allocation requests by first the application using the least resources first and most last. This way new applications entering the queue who have no resources for processing will be first asked for their required allocations to get started. Once all applications in the queue have resources they get balanced fairly between all users asking for them.

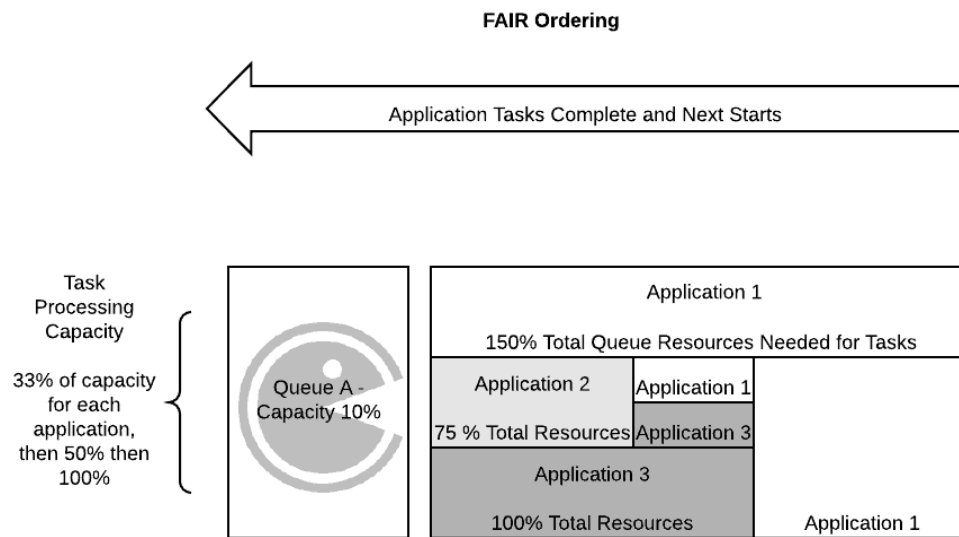


Figure 5.16 Fair ordering

It's important to note that this behavior only occurs if you have good container churn in your queue. Because Preemption does not exist inside of a queue resources cannot be forcefully redistributed within it and the FAIR ordering policy only is concerned with new allocations of resources not current ones; what does this mean? That if the queue is currently utilized with tasks which never complete or run for long periods of time without allowing for container churn to occur in the queue will hold the resources and still prevent applications from executing.

Username and Application Driven Calculations

Calculations in the Capacity Scheduler look at two primary attributes when attempting to provide allocations: User Name and Application ID. When it comes to sharing resources between users in a queue aspects like Minimum User Percentage and User Limit Factor both look at the User Name itself; this can obviously cause some conflicting problems if you're using a service account for multiple users to run jobs as only 1 user will appear to the Capacity Scheduler. Within a Queue which application gets resource allocations are driven by the leaf queues ordering policies: FIFO or FAIR which only care about the application and not which user is running it. In FIFO resources are allocated first to the oldest application in the queue and only when it no longer requires any will the next application gets an allocation. With FAIR applications that are using the least about of resources are first asking if there are pending allocations for them and fulfilled if so, if not the next application with least resources is checked; this helps to evenly share queues with applications that would normally consume them.

Default Queue Mapping

Normally to target a specific queue the user provides some configuration information telling the client tools what queue to request. But commonly users utilize tools that have a hard time passing configurations downstream to target specific queues. With Default Queue Mappings we can route an entity by its username or groups it belongs to into specific queues. Take note that the default queue routing configuration matches whichever routing attribute comes first. So if a group mapping is provided before the user mapping that matches the user he is routed to the queue for the group.

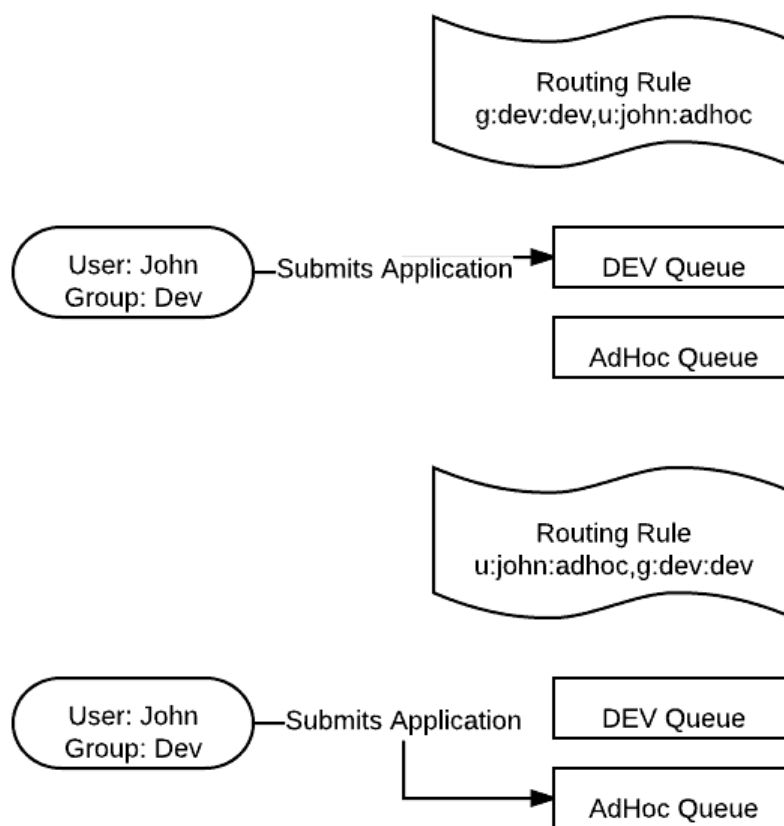


Figure 5.17 Default queue mapping

Priority

When resources are allocated over multiple queues the one with the lowest relative capacity gets resources first. If you're in a scenario where you want to have a high priority queue that receives resources before others then changing to a higher priority is a simple way to do

this. Today using queue priorities with LLAP and Tez allows for more interactive workloads as these queues can be assigned resources at a higher priority to reduce the end latency that an end user may experience.

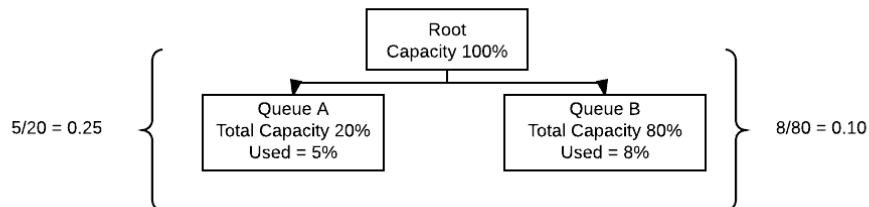


Figure 5.18 Priority queueing

Above the two queues provide an example of how relative capacity is used unmodified by queue priorities. In this case even while Queue A is smaller than Queue B and while Queue B is using more absolute resources it is elected to continue receiving them first because its relative capacity is lower than that of Queue A. If the desired behavior required Queue A to always receive resource allocations first the Queue priority should be increased above that of Queue B. When assigning priority a higher value represents higher priority.

Labels

Labels are better describe as partitions of the cluster. All clusters start with a default label or partition that is exclusive in the sense that as new labeled partitions are added to the cluster they are not shared with the original default cluster partition. Labels be defined as exclusive or shared when created and a node can only have a single label assigned to it. More common uses of labels is for the targeting of GPU hardware in the cluster or to deploy licensed software against only a specific subset of the cluster. Today LLAP also uses Labels to leverage dedicated hosts for long running processes.

Shared labels allow applications in other labeled partitions such as the default cluster to grow into them and utilize the hardware if no specific applications are asking for the label. If an application comes along that targets the label specifically the other applications that were utilizing it will be preempted off the labeled node so the application needing it can make use of it. Exclusive labels are just that, exclusive and will not be shared by any others; Only applications specifically targeting labels will run solely on them. Access to labeled partitions is provided to leaf queues so users able to submit to them are able to target the label. If you wish to autoroute users to labels say creating a GPU queue that automatically uses the GPU label this is possible.

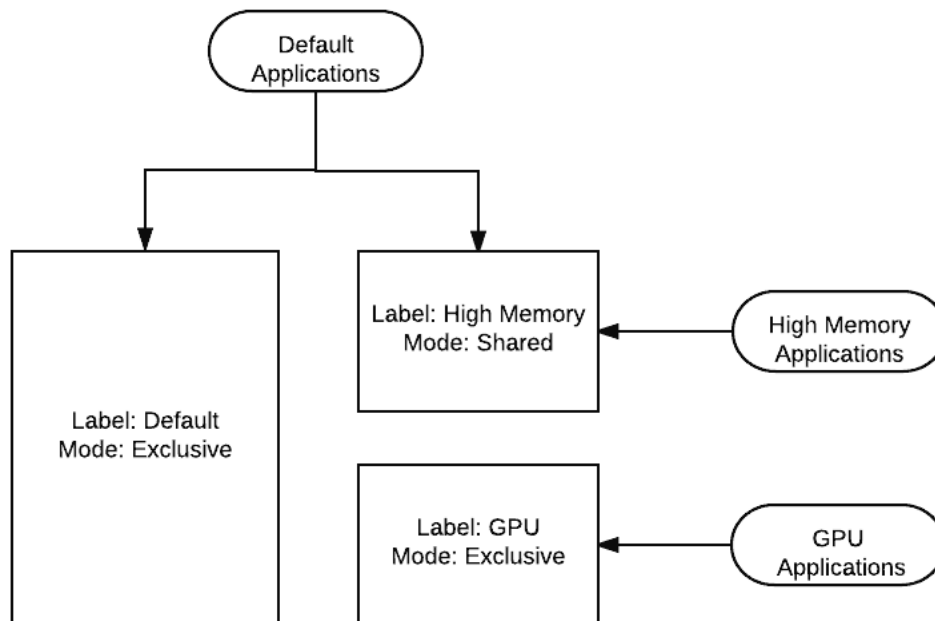


Figure 5.19 Labelling

Queue Names

Queue leaf names have to be unique with the Capacity Scheduler. For example if you created a queue in the capacity scheduler as `root.adhoc.dev` `dev` will have to be unique as a leaf over all queue names and you cannot have a `root.workflow.dev` queue as it would no longer be unique. This is in line with how queues are specified for submission by only using the leaf name and not the entire composite queue name. Parents of leafs are never submitted to directly and have no need to be unique so you can have `root.adhoc.dev` and `root.adhoc.qa` without issue as both `dev` and `qa` are unique leaf names.

Limiting Applications per Queue

Spawning of a queue by launching many applications into it so that none can effectively complete can create bottlenecks and impact SLAs. At the very worst the entire queue becomes deadlocked and nothing is able to process without an admin physically killing jobs to free up resources for compute tasks. This is easily prevented by placing a limit on the total number of applications allowed to run in the leaf queue, alternatively it's possible to control the % of resources of the leaf that can be used by Application Masters. By default

this value is typically rather large over 10,000 applications (or 20% of leaf resources) and can be configured per leaf if needed otherwise the value is inherited from prior parent queues of the leaf.

Container Sizing

An unknown to many people utilizing the Capacity Scheduler is that containers are multiples in size of the minimum allocation. For example if your minimum scheduler mb of memory per container was 1gb and you requested a 4.5gb sized container the scheduler will round this request up to 5gb. With very high minimums this can create large resource wastage problems for example with a 4gb minimum if we requested 5gb we would be served 8gb providing us with 3 extra GB that we never even planned on using! When configuring and minimum and maximum container size the maximum should be evenly divisible by the minimum.

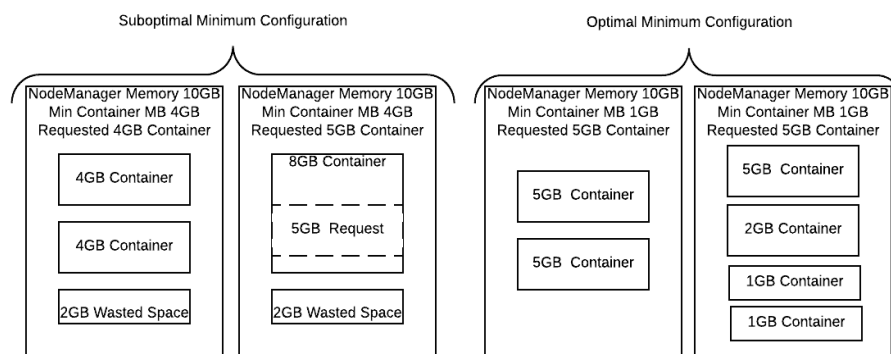


Figure 5.20 Container Seizing

References:

1. Matrix Multiplication
https://www.tutorialspoint.com/map_reduce/map_reduce_algorithm.htm
2. Fair Scheduler: <https://hadoop.apache.org/docs/stable/hadoop-yarn/hadoop-yarn-site/FairScheduler.html>
3. <https://knpcode.com/hadoop/yarn/fair-scheduler-yarn/>
4. Hadoop Capacity scheduler: <https://blog.cloudera.com/yarn-capacity-scheduler/>
5. <https://hortonworks.com/blog/better-slas-via-resource-preemption-in-yarns-capacityscheduler/>
6. Data serialization: <https://www.xenonstack.com/blog/data-serialization-hadoop#:~:text=Architecture%2C%20Tools%2C%20Challenges-,Guide%20to%20Data%20Serialization%20in%20Hadoop,regardless%20of%20the%20system%20architecture.>