

CS535 BIG DATA

PART 1. BATCH COMPUTING MODELS FOR BIG DATA ANALYTICS
1. DISTRIBUTED MODEL FOR SCALABLE BATCH COMPUTING - MapReduce

Sangmi Lee Pallickara
 Computer Science, Colorado State University
<http://www.cs.colostate.edu/~cs535>

9/1/2016 CS535 Big Data - Fall 2016 W2.B.1

FAQs

- Programming Assignment 1
 - Due Sept. 28
 - Submission via Canvas
 - Please check the course Web Page at least twice a week

9/1/2016 CS535 Big Data - Fall 2016 W2.B.2

Today's topics

- YARN
- MapReduce with examples

9/1/2016 CS535 Big Data - Fall 2016 W2.B.3

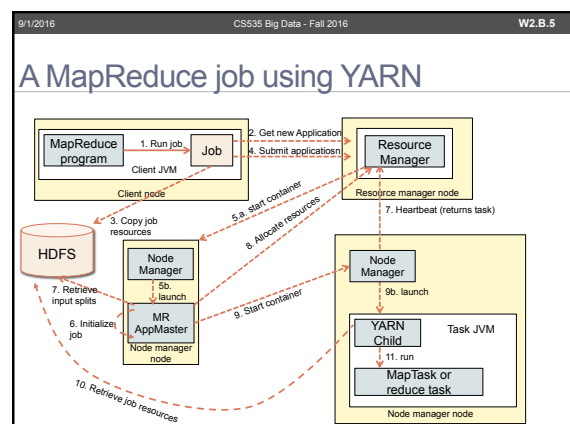
YARN (MapReduce 2)

- To provide the scalability to MapReduce
 - Splitting responsibility of the jobtracker
 - Scheduling
 - Task progress monitoring
- MapReduce is one type of YARN application

9/1/2016 CS535 Big Data - Fall 2016 W2.B.4

YARN (MapReduce 2)

- Resource manager
 - Manages the use of resources across the cluster
- Node manager
 - Launches and monitors the compute containers on machines in the cluster
- Application master
 - Manages the lifecycle of applications running on the cluster
 - Application master negotiates with the resource manager for cluster resources
 - Number of container and certain memory limit
 - Node managers oversee containers not to use more resources than allocated



9/1/2016 CS535 Big Data - Fall 2016 W2.B.6

Job Initialization

- `submitApplication()`
 - Resource manager will hands off the request to the scheduler
 - The scheduler allocates a container
 - The resource manager then launches the application master's process in the container
- Application master
 - The application master for MapReduce jobs
 - MRAppMaster
 - Initializes the job by creating bookkeeping objects
 - To keep track of the job's progress
 - Retrieves the input splits
 - Creates a map task object for each split
 - Creates reduce task object
 - `Mapreduce.job.reduces` property

9/1/2016 CS535 Big Data - Fall 2016 W2.B.7

- Application master
 - Plans job execution
 - If the jobs is small, Application Master will run the tasks in **the same JVM as itself**
 - **Uber task**
 - The overhead of allocating and running tasks in new container outweighs the gain to be had if running them in parallel, compared to running them sequentially on one node

What is a small job?
A small job is one that has less than 10 mappers, only one reducer, and an input size that is less than the size of one HDFS block

9/1/2016 CS535 Big Data - Fall 2016 W2.B.8

Task assignment [1/2]

- Application Master requests container for all the map and reduce tasks in the job
 - From the resource manager (Step 8)
- All the requests includes information about each map tasks' **locality**
 - Host and corresponding racks that the input split resides on
- Scheduler attempts to place tasks on data-local nodes in the ideal case
 - If it is not possible, the scheduler prefers rack-local placement
 - Job is running on a node in the same rack

9/1/2016 CS535 Big Data - Fall 2016 W2.B.9

Task assignment [2/2]

- Requests specify required **memory**
 - 1024MB (by default)
 - This is configurable
 - `mapreduce.map.memory.mb`
 - `mapreduce.reduce.memory.mb`
- In YARN, resources are managed more fine-grained
 - Applications may request a memory capability that is anywhere between **the minimum allocation and a maximum allocation**
 - `yarn.scheduler.capacity.minimum-allocation-mb`
 - Default minimum: 1024MB
 - `yarn.scheduler.capacity.maximum-allocation-mb`
 - Default maximum: 10240 MB
 - Tasks can request any memory allocation between 1 and 10GB(default) in multiple of 1GB
 - `mapreduce.map.memory.mb` and `mapreduce.reduce.memory.mb`

9/1/2016 CS535 Big Data - Fall 2016 W2.B.10

Task execution

- Application master starts the container by contacting node manager
 - The task is executed by `YarnChild`
 - `YarnChild` runs in a dedicated JVM

9/1/2016 CS535 Big Data - Fall 2016 W2.B.11

Progress and status updates

- Task reports its progress and status back to its application master
 - Every 3 seconds over the umbilical interface
- The client polls the application master every second
 - `mapreduce.client.progressmonitor.pollinterval`

9/1/2016 CS535 Big Data - Fall 2016 W2.B.12

Failures

9/1/2016 CS535 Big Data - Fall 2016 W2.B.13

Failures observed in Google Data Centers

"In each cluster's first year, it's typical that 1,000 individual machine failures will occur; thousands of hard drive failures will occur; one power distribution unit will fail, bringing down 500 to 1,000 machines for about 6 hours; 20 racks will fail, each time causing 40 to 80 machines to vanish from the network; 5 racks will "go wonky," with half their network packets missing in action; and the cluster will have to be rewired once, affecting 5 percent of the machines at any given moment over a 2-day span. And there's about a 50 percent chance that the cluster will overheat, taking down most of the servers in less than 5 minutes and taking 1 to 2 days to recover."

<http://www.datacenterknowledge.com/archives/2008/05/30/failure-rates-in-google-data-centers/>

9/1/2016 CS535 Big Data - Fall 2016 W2.B.14

Failures in Classic MR

- The child task fails
 - Runtime exception from the user code
 - The child JVM reports the error back to its parent before it exits
 - Written in the user logs
 - Tasktracker marks the task attempt as failed
 - Frees a slot to run another task
- Sudden exit of the child JVM
 - Tasktracker notices that the process has exited and marks the attempt as failed
- Hanging tasks
 - If there is no progress update for a while
 - Mark the task as failed
 - Timeout period is normally 10 minutes
 - `mapred.task.timeout`

9/1/2016 CS535 Big Data - Fall 2016 W2.B.15

Tasktracker failure in Classic MR

- Tasktracker stops sending heartbeats
 - Jobtracker will notice if it hasn't received one for 10 minutes (configurable)
 - Remove it from the pool of tasktracker
- Jobtracker arranges tasks including the completed jobs
 - Because the output may not be accessible
- Tasktracker can also be blacklisted if more than four tasks from the same job fail (set by `mapred.max.tracker.failures`)
 - Blacklisted tasktrackers are not assigned tasks.
 - Until faults expire

9/1/2016 CS535 Big Data - Fall 2016 W2.B.16

Jobtracker failure in Classic MR

- The most serious failure mode
- Hadoop has **no mechanism** for dealing with jobtracker failure
 - Single point of failure
- **All running jobs fail**
- After restarting a jobtracker
 - Job should be **resubmitted**
- **This is improved with YARN**

9/1/2016 CS535 Big Data - Fall 2016 W2.B.17

Task failure in YARN

- Failure of the running task is similar to the classic case
 - Runtime exception and sudden exit of the JVM are propagated back to the application master
 - The task attempt is marked as failed
 - Hanging tasks are noticed by the application manager by the absence of a ping over the umbilical channel

9/1/2016 CS535 Big Data - Fall 2016 W2.B.18

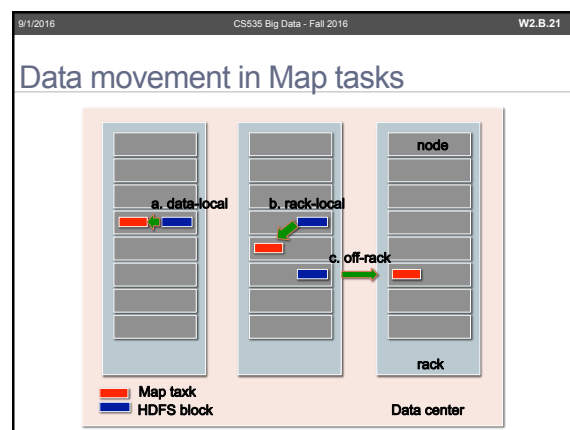
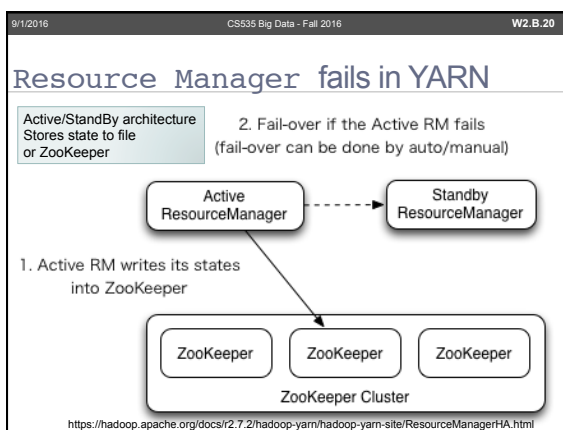
Application master failure in YARN

- No heartbeats to the resource manager from the application master
 - The resource manager will detect the failure and start a new instance of the application master running in a new container
 - All tasks will be rerun (default)
 - Recovery can be enabled
- Client will access resource manager to get the new address of the application master

9/1/2016 CS535 Big Data - Fall 2016 W2.B.19

Node manager failure in YARN

- Resource manager will stop getting heartbeats
 - Remove the failed node manager from the pool of available nodes



9/1/2016 CS535 Big Data - Fall 2016 W2.B.22

Data locality optimization

- Hadoop tries to run the map task on a node where the input data resides in HDFS
 - Minimizes usage of cluster bandwidth
- If all replication nodes are running other map tasks
 - The job scheduler will look for a free map slot on a node in the same rack

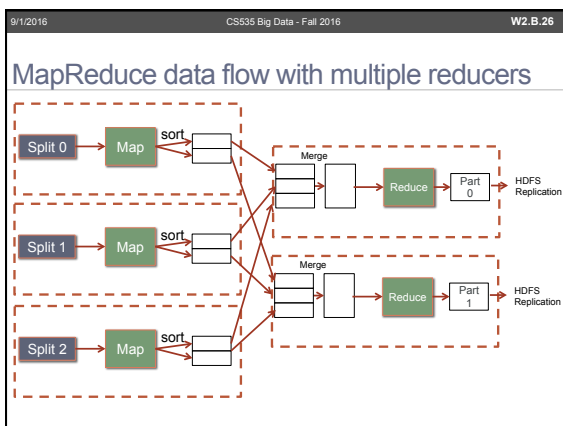
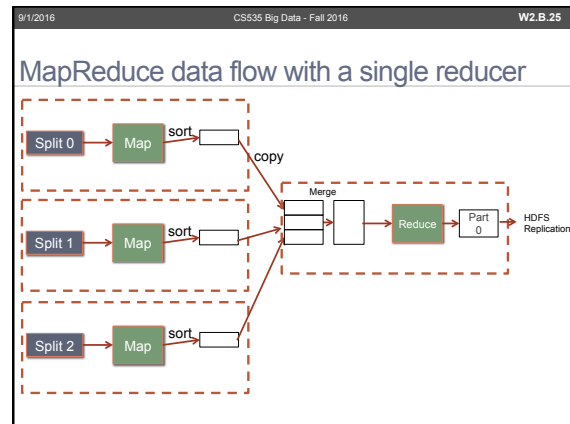
9/1/2016 CS535 Big Data - Fall 2016 W2.B.23

Shuffle and Combiner

9/1/2016 CS535 Big Data - Fall 2016 W2.B.24

Shuffle

- The process by which the system performs the sort and transfers the map outputs to the reducers as inputs
 - MapReduce guarantees that the input to every reducer is sorted by key



9/1/2016 CS535 Big Data - Fall 2016 W2.B.27

Combiner functions

- Minimize data transferred between map and reduce tasks
- Users can specify a **combiner function**
 - To be run on the map output
 - To replace the map output with the combiner output

9/1/2016 CS535 Big Data - Fall 2016 W2.B.28

Combiner example

- Example (from the previous NCDL max temp example)
 - The first map produced,
 - (1950, 0), (1950, 20), (1950, 10)
 - The second map produced,
 - (1950, 25), (1950, 15)
 - The reduce function is called with a list of all the values,
 - (1950, [0, 20, 10, 25, 15])
 - Output will be,
 - (1950, 25)
- We may express the function as,
 - `max(0, 20, 10, 25, 15)`
 - `= max(max(0, 20, 10), max(25, 15))`
 - `= max(20, 25) = 25`

9/1/2016 CS535 Big Data - Fall 2016 W2.B.29

Combiner function

- Run a **local** reducer over Map output
- Reduce the amount of data shuffled between the mappers and the reducers
- Combiner cannot replace the reduce function
 - Why?

9/1/2016 CS535 Big Data - Fall 2016 W2.B.30

MapReduce Programming with Examples

9/1/2016 CS535 Big Data - Fall 2016 W2.B.31

Example-1: Selecting rows

- Suppose that you have a fileset with a 1TB dataset

```
Yam:~/data/$ cat samplefile.txt
127.0.0.1 192.168.10.102
127.0.0.1 192.168.10.105
127.0.0.1 192.168.10.101
192.168.10.102 192.168.10.101
...
```

From To

- Retrieve **all of the rows** where “from” and “to” are the **same URL**
- Equivalent to SQL's “select * where from == to” operation

9/1/2016 CS535 Big Data - Fall 2016 W2.B.32

Example-1: Implementation using MapReduce

Map function
For each line, test if it satisfies condition (here, from == to)
If it satisfies the condition, produce the key-value pair (t, t)

Reduce function
Passes each key-value pair to the output

9/1/2016 CS535 Big Data - Fall 2016 W2.B.33

Example-2: Selecting columns

- Suppose that you have a file set with 1TB dataset

```
Yam:~/data/$ cat samplefile.txt
127.0.0.1 192.168.10.102
127.0.0.1 192.168.10.105
127.0.0.1 192.168.10.101
192.168.10.102 192.168.10.101
...
```

From To

- Retrieve all **unique** “from”s

9/1/2016 CS535 Big Data - Fall 2016 W2.B.34

Example-2: Implementation using MapReduce

Map function
For each line, retrieve value of “from”.
Produce the key-value pair (f, f)

Reduce function
The pairs with same fs are shuffled as (f, [f, f, f, ..., f])
Returns (f, f) as the result of eliminating duplications

9/1/2016 CS535 Big Data - Fall 2016 W2.B.35

Example-3: Implementation of Union using MapReduce

Map function
For each line, produce the key-value pair (t, t)

Reduce function
For the shuffled pairs (t, [t, t, ..., t])
Returns (t, t) as the result of eliminating duplications

9/1/2016 CS535 Big Data - Fall 2016 W2.B.36

Example-4: Implementation of Intersection using MapReduce

Map function
For each line in set A, produce the key-value pair (t, "A")
For each line in set B, produce the key-value pair (t, "B")

Reduce function
If there are (t, "A") and (t, "B") in the shuffled results, return (t, t). This eliminates duplications
Otherwise, do not return anything

9/1/2016 CS535 Big Data - Fall 2016 W2.B.37

Example-5: Natural Join

- Suppose that you have two tables with 1TB each

```
Yami:~/data/$ cat samplefile.txt
127.0.0.1 192.168.10.102
127.0.0.1 192.168.10.105
127.0.0.1 192.168.10.101
192.168.10.102 192.168.10.101
-
```

A.txt

```
Yami:~/data/$ cat samplefile-detailed.txt
192.168.10.102 ALIVE
192.168.10.101 DOWN
192.168.10.101 ALIVE
127.0.0.1 DOWN
-
```

B.txt

- If the tuples agree on an attribute that are common to the two schemas, then produce a new tuple that has components for each of the attributes in either schema
- Equivalent to SQL's "join" operation
- second url in the A.txt and first attribute in B.txt are common
- join(A, B) = {(url1, url12, ALIVE), (url12, url13, DOWN)}

9/1/2016 CS535 Big Data - Fall 2016 W2.B.38

Example-5: Implementation of Natural Join using MapReduce

Map function
For each line (a,b) in A, produce the key-value pair (b, ("A", a))
For each tuple (b,c) in B, produce the key-value pair (b, ("B", c))

Reduce function
Construct all pairs comprising one with first component "A" and the other with first component "B". e.g. ("A", a) and ("B", c)
Produce tuple (a,b,c) such that ("A", a) and ("B", c)

9/1/2016 CS535 Big Data - Fall 2016 W2.B.39

Matrix-Vector Multiplication Using MapReduce

9/1/2016 CS535 Big Data - Fall 2016 W2.B.40

Matrix-Vector multiplication using MapReduce (1/3)

- Suppose we have an $n \times n$ matrix M , whose element in row i and column j will be denoted M_{ij}
- v is a $n \times 1$ column vector
- Then the matrix-vector product is the vector x of length n , whose i^{th} element x_i is given by:

$$x_i = \sum_{j=1}^n m_{ij} v_j$$