

How MapReduce works

Date: _____ Page: 42

- submit() - Run MR Job on Job object.
- waitForCompletion() - Submit job then wait for finish.
 - conceal great deal of processing
- Figure 7.1 (How MR run on Hadoop)
- 5 independent entities
 - 1) Client
 - 2) YARN RM
 - 3) NM
 - 4) AM
 - 5) Distributed file system

Job Submission

- Submit() on Job Create JobSubmission instance - Calls submitJobInternal() (Step 1)
- waitForNotification polls job progress once per second
repeat progress to Console (if changed) } Job counters
one displayed when job complete successfully otherwise
error logged to console.
- ASK RM for new App ID (2)
- Copy resources to HDFS directory named jobId. (JAR,
Configs file, input splits) (3)
- Submit job by calling submitApplication() on RM (4)

Job Initialization

- When RM receives a Call submitApplication() - hands off req
to YARN Scheduler
- Scheduler allocate Container (5a)
- RM then launch AM process under NM mgmt (5b)

extramarks

- AM for MR Job is a Java App (main class - MRAppMaster)
- It initializes job by creating book keeping objects to keep track of job's progress, receive progress, complete report from task (6)
- Retrieve input specs from HDFS (7), Create map tasks object, no of reduce task objects, Tasks are given ID.
- AM decide how to run tasks - If job is small - AM choose to run tasks in same JVM (Uber task)
Small job - less than 10 mappers only one reducer
- AM call setupTask() - Create final off dir for tasks and temporary working space for task off.

Task Assignment

- If not Uber task, AM request containers for all map and Reduce Task in job from RM (8)
- Req for map task made first (Priority) than Reducetask (not made until 5% of map task completed)
- Reducetask can run anywhere but req for maptask have data locality constraints
- optimal-task is data local | rack local | diff racks
- Req for memory and CPU for task - M/R task 1024 ms and 1 virtual core (configurable)
- mapreduce.map.memory.mb | CPU.vcores
----- .reduce. ----- .mb | CPU.vcores

Task Execution

- Once task has been assigned resources for a container, AM starts container by contacting NM (9a, 9b)
- Task is executed by Java App (main class YarnChild)
- Before it runs task, it localize resources that that task needs (Job Config, JAR file) (10)
- Finally, it runs the map or reduce task (11)
- YARN child is runs in a dedicated JVM, so any bugs (in map Reduce) don't affect the NM.

Streaming

- It runs special map/reduce tasks (Fig 7.2)
- It comes with process using standard I/O streams.
- Java passes input key-value pairs to external process.
- Run user defined map/reduce func and passes output key-value pair back to Java process.
- NM view - as if child process ran the m/r code itself

Progress and Status updates

- MR jobs are long running jobs (10s to hours), so imp for user to get feedback how job is progressing.
- Job and its tasks have status (includes)
 - State of job/task (running, completed, fail)
 - Progress of M/R
 - Value of job counter
 - Status msg or description

When task is running - it keeps track of progress

- map task - proportion of Input that has been processed

extramarks

- Redundancy - Complexity

- Dividing total progress into 3 parts (phases of shuffle)
- if task run reduces $1/2$ of its input task progress

$$\beta \leq \frac{1}{6} = \frac{1}{3} + \frac{1}{3} + \frac{1}{3} \times \frac{1}{2}$$

What constitutes progress in MR

- Reading an input record
- writing an output record
- Setting the status description
- Incrementing a counter
- Calling Report's progress() method
- child process comm with parent AM through Umbrella interface
- Task reports its progress and status back to its AM
- Client can receive latest status by polling the AM every sec.

Fig 7.3

Job Completion

- When AM receive notification that last task for a job is complete - it changes status - "Successful".
- It then prints a msg. to tell user; return from waitForCompletion()
- Job statistics and counters are printed to console.
- finally on Job completion AM and taskContainers clean up working state | job info is archived by JobHistory Server

Failure

- Real word user code is buggy, process crash, m/c fail.
- Bangs of Hadoop - ability to handle failure and allow job to complete successfully.
- failure of following entities.

Task failure

- ① - user code (map/reduce) task throw a runtime exception.
- TaskJVM reports errors back to its parent app Master by exits.
- make it to user logs.
- AM marks task as failed, free up the Container so res available to another task.
- ② Another failure node - sudden exit of taskJVM - JVM bugs cause JVM to exit
- NM notices and inform AM so it can mark as failed.
- ③ Hanging task - AM hasn't received progress update - mark task as failed.
- timeout period (10 min) configurable (mapreduce.task.timeout)
- setting ∞ - long running tasks never marked as fail
never free up its container - slow down cluster
Task - 4 time retries (mapreduce.map.maxattempts)
- Task attempt may be killed - speculative Execution

AM failure

- App in YARN retried in event failure -
- Max no of attempt from MR AM is controlled by mapreduce.am.max-attempts property (def 2)
- YARN impose limit for max no of attempt of your app PM after (AM) extramarks

(yarn.resourcemanager.am.map-attempts)

- 1 Recovery - AM send Heartbeat to RM, will detect failure
Start new instance of master in new container
RM uses job history to recover state of the tasks
- 2 MR Client poll AM for progress report
 - if Am fail client need new AM.
 - During job initialization client ask AM for AM address and Caches⁵⁰ not to overload RM
 - if Am fail client experience timeout - go to RM again for new add.

NM failure

- If NM fail / slow stop sending HB to RM (10min)
- So RM remove from pool of nodes to schedule containers
- AM / Task on failed NM - recompute map tasks
- AM arrange map tasks that were completed successfully on failed NM to rerun if they belong to incomplete job.
(intermediate o/p of failed NM)
- NM - blacklisted - if not failure for app is triggered even if NM has not failed, AM reschedule task on diff node (if more than 3 task fail on NM)

RM failure

- failure - no jobs, task containers can be launched.
- SPOF
- HA - active - standby conf. | info about running app is stored in highly available State Store in (2K/MDFS)
Standby can recover state of failed RM

Transition of RM from Standby to Active handled by failover Controller. (def automatic) - Use 2K leader election (only 1 Active RM)

Shuffle and Sort

- MR guarantee input to every reducer is sorted by key.
- Shuffle - System perform Sort - transfer map o/p to reducer as I/P

Fig 7.4

Map Side

- Buffering writes in memory
- done presenting for efficiency reasons
- Each map task has circular memory buffer and write o/p.
- 100MB (mapreduce-task.io.sort.mb.)
- When threshold (80%) - mapreduce.map.sort.spill.factor background thread starts spill contents to disk.
- If buffer fills - map will block until spill is complete
- By write to disk - thread divide data into partitions
- within each partition - thread performs in memory sort
- Good idea - Compress map o/p - faster, save disk space
- reduce amt of data to transfer
- Det - not compressed (mapreduce.map.outbuf.compress | mapreduce.map.output.compress.codec.)

Reducer Side

- map task o/p is on diff m/c, Reducer can be on another m/c
- map task may finish at diff time, so start reduce task Start copying o/p as each completes. (copy phase)
- Reduce task - Small no of copy threads - fetch o/p in parallel

extramarks

mapreduce.reduce (Def thread - 3)
mapreduce.shuffle parallel (parallel)

- map o/p copied to reducer task Jvm memory (small) else disk.
- when in memory buffer reached threshold size threshold of no of map o/p, it is merged and spilled to disk.
and background thread merge them in large sorted file. (decompressed)
- After copy sort phase - which merge the map o/p
eg 50 map o/p, merge factor = 10, & looks Fig - 7.5

Reduce phase - reduce func is invoke for each key in sorted o/p - o/p written to HDFS. - first block replica written to local disk.

Configuration Tuning

- Principle - Give shuffle as much memory as possible
- Mapside - best performance can be obtained by avoiding m/t spills to disk
- Reduceside - best performance is obtained when intermediate data can reside entirely in memory

Fig - Table 7.1, 7.2