

A Report on:

*Energy efficient algorithm design
in Hadoop clusters*

By

*Kartik Sathyanarayanan
(2013A7PS037G)*



BITS Pilani
Pilani | Dubai | Goa | Hyderabad

BIRLA INSTITUTE OF TECHNOLOGY & SCIENCE, PILANI,
K.K BIRLA GOA CAMPUS

Introduction

Hadoop is an open-source implementation of MapReduce enjoying wide adoption and is often used for short jobs where low response time is critical. Hadoop's performance is closely tied to its task scheduler, which implicitly assumes that cluster nodes are homogeneous and tasks make progress linearly, and uses these assumptions to decide when to speculatively re-execute tasks that appear to be stragglers. In practice, the homogeneity assumptions do not always hold. MapReduce uses speculative execution to improve fault tolerance. Current Hadoop implementation decides whether to run speculative tasks based on the progress rates of running tasks, which does not take into consideration the absolute progress of each task.

The modified Hadoop framework was deployed in 6 t2.medium EC2 instances in a master-slave configuration.

Modified Task Scheduler

The BASE scheduler (which itself was a little improvement from LATE scheduler) was slightly modified to limit unnecessary speculative executions.

The need for speculating a task is determined by the following factors:

- The task should not be done or almost done.
- The task should not have been speculatively executed already.
- Launching a speculative task should not exceed the upper limit on maximum number of currently active tasks.
- The task should have a burn in period – this period is used to obtain usable task-specific CPU metrics for estimating speculation cost.
- The version of HDFS used cannot support multiple simultaneous writes. As a restrictive safety measure, map tasks are not speculated.
- A weighted factor in estimating the speculation benefit is (estimated completion time – speculated completion time) [Nidhi's proposal]. This difference is estimated as the difference mean completion time of the other tasks and the current progress rate (based on the LATE/BASE scheduler ideas) [2].

Ultimately, the decision of speculating a task is taken by checking if the calculated speculation benefit is higher than a threshold. This threshold is decided by a multiple of the standard deviation of the completion time of tasks. Empirically, this is between $2 * \text{standard deviation}$ and $3 * \text{standard deviation}$. This multiple is parameterized as the slow task threshold[1].

Tasks which pass this filtering process are added to a list. This list is then sorted based on the speculative gain (ratio). The estimated time to completion is measured [1] as the percentage of task left to complete – that is, $(1 - \text{progress}) / (\text{progress rate of the task})$. The one with the highest speculative gain is then sent for speculation.

Algorithm

- SpeculativeCap – The max fraction (0-1) of running tasks that can be speculatively re-executed at any time. We have used a value of 1.
- SlowTaskThreshold – The number of standard deviations by which a task's average progress-rates must be lower than the average of all running tasks' for the task to be considered too slow. We have used a value of 0.01.
- SlowNodeThreshold – The number of standard deviations by which a Task Tracker's average map and reduce progress-rates must be lower than the average of all successful map/reduce task's for the TT to be considered too slow to give a speculative task to. We have used a value of 0.01.
- If a task slot becomes available in a node and there are less than SpeculativeCap(fraction) speculative tasks running :
 - Ignore the request if the TaskTracker's progressRate < SlowNodeThreshold
 - Choose candidate tasks : those tasks whose progress rates are below [slowTaskThreshold * mean(progress rates)] - *Nidhi's proposal*
 - Rank the candidate tasks that are not currently being speculated by estimate time left
 - Speculate the task that's expected to complete last

Benchmarks

1. Unmodified Hadoop 0.20

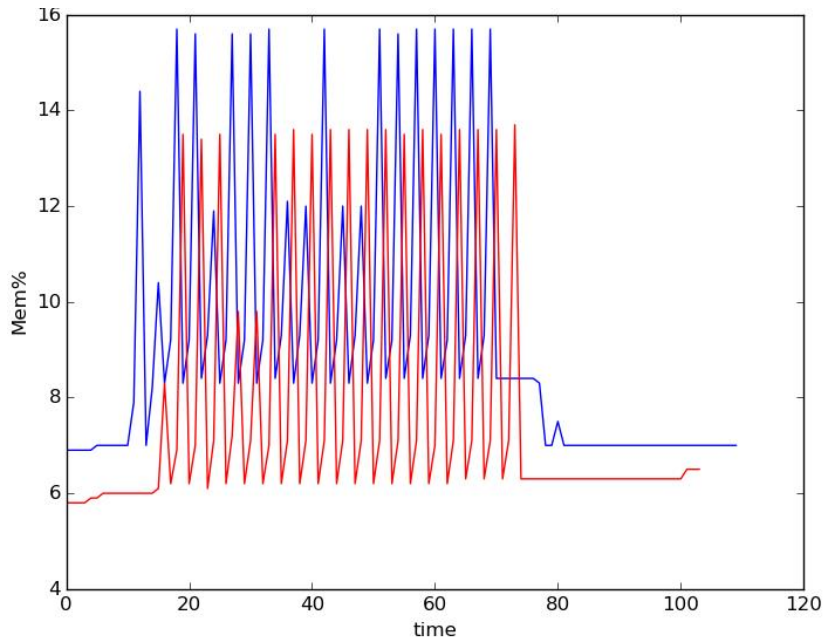
Job Name	Command Line arguments	Job Response time	No. of maps	No. of reduces	Speculative maps launched	Speculative reduces launched
PiEstimator	200, 200	88 secs	200	1	0	0
TeraGen	10 ⁶	19 secs	2	0	0	0
TeraSort	Output of TeraGen	26 secs	2	1	0	0
TeraValidate	Output of TeraSort	22 secs	1	0	0	0
Distributed Pentomino Solver	-	2hrs 39min 38 secs	2001	1	6	0
WordCount	1.2MB file	24 secs	2	1	0	0

2. Modified Hadoop 0.20

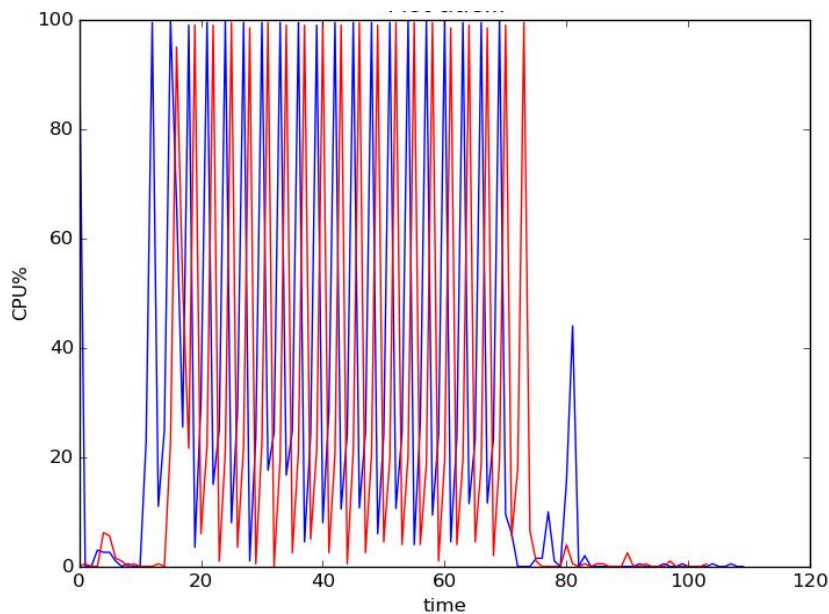
Job Name	Command Line arguments	Job Response time	No. of maps	No. of reduces	Speculative maps launched	Speculative reduces launched
PiEstimator	200, 200	78 secs	200	1	0	0
TeraGen	10 ⁶	13 secs	2	0	0	1
TeraSort	Output of TeraGen	24 secs	2	1	0	0
TeraValidate	Output of TeraSort	21 secs	1	0	0	0
Distributed Pentomino Solver	-	2hrs 38min 45 secs	2001	1	0	1
WordCount	1.2 MB file	21 secs	2	1	0	0

CPU and Memory Stats

- Red line : Unmodified hadoop, and Blue line : Modified hadoop
1. PiEstimator 200,200

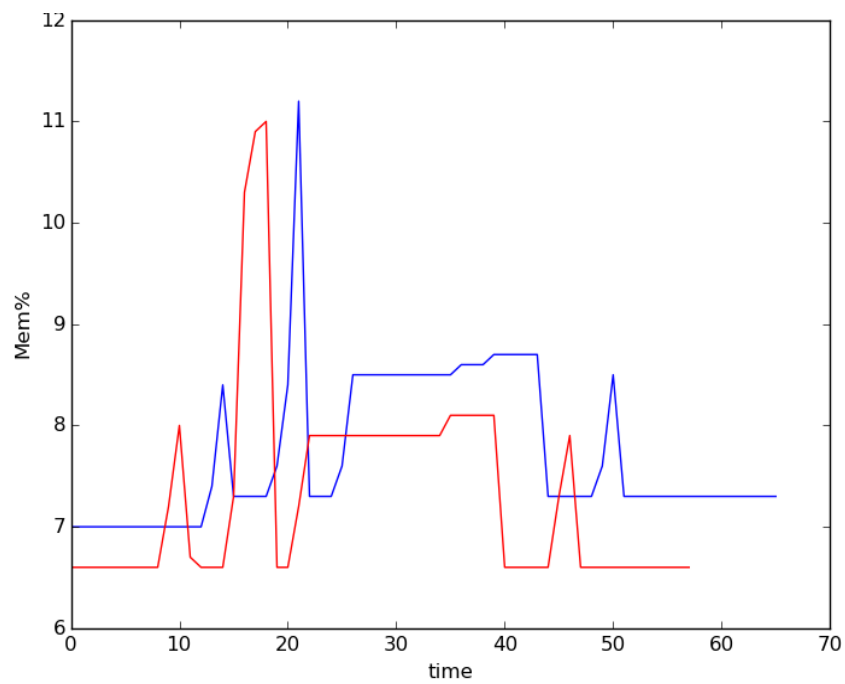
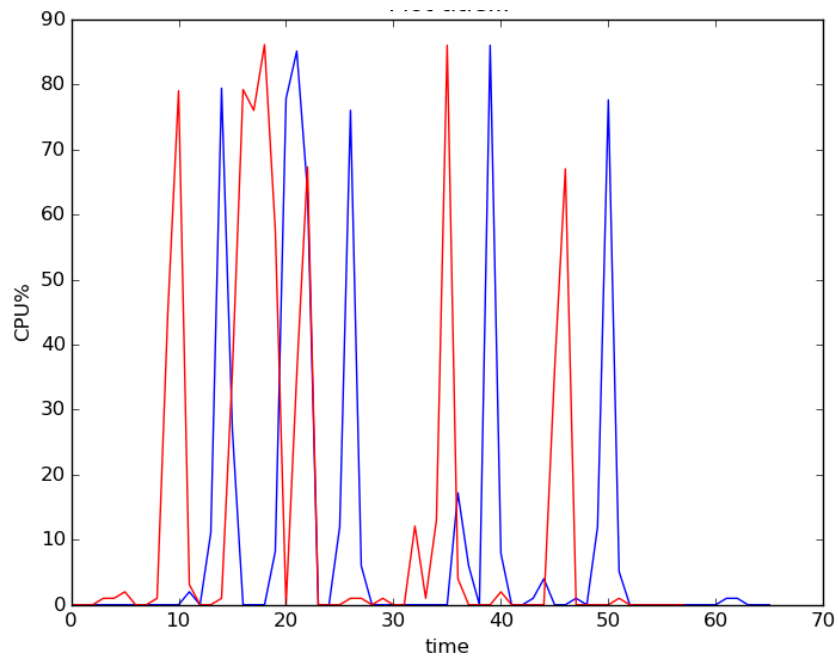


This graph explains the memory usage of a data node during the PiEstimator job. The spikes in the graph represent map or reduce tasks.



This graph explains the CPU usage of a data node during the PiEstimator job. The spikes in the graph represent map or reduce tasks.

2. WordCount 1.2MB text file



The above graphs represent CPU and memory usage of a data node with modified and unmodified hadoop versions for a Word Count job.

Conclusion

The modified hadoop 0.20 gives approximately the same job response times on different benchmarks run as shown in the previous section. The CPU and Memory usage for a few jobs are shown above. The CPU usage spikes are similar in both hadoop versions, whereas Memory usage seems to be slightly on the higher side in the modified hadoop framework.

We were able to observe a difference in number of speculative execution only in the dancingElephant(Distributed Pentomino Solver) job, where 6 maps were speculated in the existing hadoop version. Whereas, only 1 reduce task was speculated in the modified hadoop version. In the Pi Estimator job using Monte-Carlo simulation, with 200 maps and 200 inputs/map we observed that there was one speculative reduce task which started due to the slow progress of the reduce task. A reason for this might be that, Pi Estimator job has overlapping maps and reduce tasks.

A latent assumption here is that the tasks are linear in their progress. This is not a reliable assumption to make in general since the map, shuffle and reduce stages are pipelined and have spill-over mechanisms. Hence, calculating the progress for reduce by evenly weighting shuffle, sort and map is unreliable. To fix this, phase-specific metrics like completion time, mean and average number of spills are needed. Empirical evidence on large-scale cluster management systems, however, show that slightly complicated progress rate heuristics give a usable estimate.

Consider hadoop jobs where maps fetch data from external systems, and emit the data. The reducers in this are identity reducers. The data processed by these jobs is huge. There could be slow nodes in this cluster and some of the reducers run twice as slow as their counterparts. This could result in a long tail. Speculative execution would help greatly in such cases. However given the current hadoop, we have to select speculative execution for both maps and reducers. In this case hurting the map performance as they are fetching data from external systems thereby overloading the external systems.

Future Work

- Incorporate probability of node failure in the benefit calculation of launching a speculative task
- Calculate per-task resource utilization so that CPU, Memory, Disk and Network I/O utilization can be incorporated into the cost calculation
- Come up with a better metric to calculate the estimated completion times of tasks

References

- 1) Zhenhua Guo, Geoffrey Fox, Mo Zhou, Yang Ruan. Improving Resource Utilization in MapReduce. {it Indiana University Report}. May 2012.
- 2) Matei Zaharia, Andy Konwinski, Anthony D. Joseph, Randy Katz, Ion Stoica. Improving MapReduce performance in heterogeneous environments. In {it Proceedings of the 8th USENIX conference on Operating systems design and implementation}, p.29-42, December, 2008.
- 3) Apache Jira issues :-
 - a. <https://issues.apache.org/jira/browse/MAPREDUCE-3404>
 - b. <https://issues.apache.org/jira/browse/HADOOP-1127>
 - c. <https://issues.apache.org/jira/browse/HADOOP-2131>
 - d. <https://issues.apache.org/jira/browse/MAPREDUCE-220>

Appendix

Summary of changes in each java file

1. [src/contrib/capacity-scheduler/src/java/org/apache/hadoop/mapred/CapacityTaskScheduler.java](#)
tip.hasSpeculativeTask(currentTime, progress) was replaced with tip.canBeSpeculated(currentTime) because this function was changed in TaskInProgress.java
2. [src/contrib/capacity-scheduler/src/test/org/apache/hadoop/mapred/TestCapacityScheduler.java](#)
tip.hasSpeculativeTask(currentTime, progress) was replaced with tip.canBeSpeculated(currentTime) because this function was changed in TaskInProgress.java
3. [src/contrib/fairscheduler/src/java/org/apache/hadoop/mapred/DefaultTaskSelector.java](#)
tip.hasSpeculativeTask(currentTime, progress) was replaced with tip.canBeSpeculated(currentTime) because this function was changed in TaskInProgress.java
4. [src/mapred/mapred-default.xml](#)
Included variables SpeculativeCap, slowTaskThreshold and slowNodeThreshold for speculative cost calculation
5. [src/mapred/org/apache/hadoop/mapred/JobInProgress.java](#)
 - a) Included the required variables and DataStatistics class for speculative cost calculation.
 - b) Changed findSpeculativeTask() , to incorporate speculative task selection based on the metrics mentioned in the Algorithm
 - c) Introduced getSpeculativeMap() and getSpeculativeReduce() to make calls to findSpeculativeTask() and get the task ID to be speculated.
 - d) EstimatedTimeLeftComparator class introduced to compare the estimated completion time of two tasks, which is used to sort the candidates in findSpeculativeTask()
 - e) isSlowTracker() – to determine if the TaskTracker node is too slow to run speculative tasks on

6. [src/mapred/org/apache/hadoop/mapred/TaskInProgress.java](#)
 - a) Modified hasSpeculativeTask() to canBeSpeculated() so that an extra condition as mentioned in the algorithm is included to filter speculative tasks
 - b) Calculate the bestProgressRate of a task in TaskTracker node with the help of getCurrentProgressRate()
7. [src/mapred/org/apache/hadoop/mapred/JobConf.java](#)

Speculative map and reduce were enabled in this file. Ideally, speculative maps should be disabled and speculative reduces should be enabled. The reasons are stated in the Conclusion. Another good reason we found is that **multiple writes to HDFS is generally discouraged because slows down the write rates of other running map tasks.**
8. [src/core/org/apache/Hadoop/util/LinuxResourceCalculatorPlugin.java](#)

Included this file at the last moment, to calculate per task resource utilization as mentioned by <https://issues.apache.org/jira/browse/MAPREDUCE-220> . But wasn't able to use it correctly. Instantiating the class in JobInProgress.java and using its method gave figures which were similar to "top" output. In my honest opinion, if at all it calculates per task resource utilization it should have been associated with TaskInProgress class. However, figuring out its correct usage would be of paramount importance in calculating Speculative Cost as mentioned in Nidhi's proposal.

The following contains the changes in the Hadoop source code. Code portions colored in red were deleted and those in green were included. Few java files have no color in them in entirety; these files are new.

```
//Replaced tip.hasSpeculativeTask(currentTime, progress) with  
//tip.canBeSpeculated(currentTime) in the following files :  
// 1. src/contrib/capacity-  
//scheduler/src/java/org/apache/hadoop/mapred/CapacityTaskScheduler.java  
// 2. src/contrib/capacity-  
//scheduler/src/test/org/apache/hadoop/mapred/TestCapacityScheduler.java  
// 3. src/contrib/fairscheduler/src/java/org/apache/hadoop/mapred/DefaultTaskSelector.java
```

//Changed src/mapred/mapred-default.xml to include Speculative default variables.

```
504 <property>  
505   <name>mapred.speculative.execution.speculativeCap</name>  
506   <value>0.1</value>  
507   <description>The max percent (0-1) of running tasks that  
508   can be speculatively re-executed at any time.</description>  
509 </property>
```

```
505 511 <property>  
512   <name>mapred.speculative.execution.slowTaskThreshold</name>  
513   <value>1.0</value>The number of standard deviations by which a task's  
514   ave progress-rates must be lower than the average of all running tasks'  
515   for the task to be considered too slow.  
516   <description>  
517   </description>  
518 </property>  
519  
520 <property>  
521   <name>mapred.speculative.execution.slowNodeThreshold</name>  
522   <value>1.0</value>  
523   <description>The number of standard deviations by which a Task  
524   Tracker's ave map and reduce progress-rates (finishTime-dispatchTime)  
525   must be lower than the average of all successful map/reduce task's for  
526   the TT to be considered too slow to give a speculative task to.  
527   </description>  
528 </property>  
529
```

// Changed src/mapred/org/apache/hadoop/mapred/*.java

// 1. JobInProgress.java

```
24 import java.util.Collections;
25 import java.util.Comparator;
26 import java.util.HashSet;
27 import java.util.HashMap;
```

```
199 // Don't lower speculativeCap below one TT's worth (for small clusters)
200 private static final int MIN_SPEC_CAP = 10;
201
202 private static final float MIN_SLOTS_CAP = 0.01f;
```

```
210 //thresholds for speculative execution
211 private float slowTaskThreshold;
212 private float speculativeCap;
213 private float slowNodeThreshold; //standard deviations
```

```
215 //Statistics are maintained for a couple of things
216 //mapTaskStats is used for maintaining statistics about
217 //the completion time of map tasks on the trackers. On a per
218 //tracker basis, the mean time for task completion is maintained
219 private DataStatistics mapTaskStats = new DataStatistics();
220 //reduceTaskStats is used for maintaining statistics about
221 //the completion time of reduce tasks on the trackers. On a per
222 //tracker basis, the mean time for task completion is maintained
223 private DataStatistics reduceTaskStats = new DataStatistics();
224 //trackerMapStats used to maintain a mapping from the tracker to the
225 //the statistics about completion time of map tasks
226 private Map<String,DataStatistics> trackerMapStats =
227     new HashMap<String,DataStatistics>();
228 //trackerReduceStats used to maintain a mapping from the tracker to the
229 //the statistics about completion time of reduce tasks
230 private Map<String,DataStatistics> trackerReduceStats =
231     new HashMap<String,DataStatistics>();
232 //runningMapStats used to maintain the RUNNING map tasks' statistics
233 private DataStatistics runningMapTaskStats = new DataStatistics();
234 //runningReduceStats used to maintain the RUNNING reduce tasks' statistics
235 private DataStatistics runningReduceTaskStats = new DataStatistics();
```

```
213 this.jobtracker = null;
```

```
247 this.jobtracker = tracker;
```

```

249
250     hasSpeculativeMaps = conf.getMapSpeculativeExecution();
251     hasSpeculativeReduces = conf.getReduceSpeculativeExecution();
252     this.nonLocalMaps = new LinkedList<TaskInProgress>();
253     this.nonLocalRunningMaps = new LinkedHashSet<TaskInProgress>();
254     this.runningMapCache = new IdentityHashMap<Node, Set<TaskInProgress>>();
255     this.nonRunningReduces = new LinkedList<TaskInProgress>();
256     this.runningReduces = new LinkedHashSet<TaskInProgress>();
257     this.resourceEstimator = new ResourceEstimator(this);
258     this.status = new JobStatus(jobid, 0.0f, 0.0f, JobStatus.PREP);
259     this.taskCompletionEvents = new ArrayList<TaskCompletionEvent>
260     (numMapTasks + numReduceTasks + 10);
261
262     this.slowTaskThreshold = Math.max(0.0f,
263         conf.getFloat("mapred.speculative.execution.slowTaskThreshold", 1.0f));
264     this.speculativeCap = conf.getFloat(
265         "mapred.speculative.execution.speculativeCap", 0.1f);
266     this.slowNodeThreshold = conf.getFloat(
267         "mapred.speculative.execution.slowNodeThreshold", 1.0f);

```

```

342
343     this.nonLocalMaps = new LinkedList<TaskInProgress>();
344     this.nonLocalRunningMaps = new LinkedHashSet<TaskInProgress>();
345     this.runningMapCache = new IdentityHashMap<Node, Set<TaskInProgress>>();
346     this.nonRunningReduces = new LinkedList<TaskInProgress>();
347     this.runningReduces = new LinkedHashSet<TaskInProgress>();
348     this.slowTaskThreshold = Math.max(0.0f,
349         conf.getFloat("mapred.speculative.execution.slowTaskThreshold", 1.0f));
350     this.speculativeCap = conf.getFloat(
351         "mapred.speculative.execution.speculativeCap", 0.1f);
352     this.slowNodeThreshold = conf.getFloat(
353         "mapred.speculative.execution.slowNodeThreshold", 1.0f);
354

```

```

943     int target = findNewMapTask(tts, clusterSize, numUniqueHosts, anyCacheLevel,
944                             status.mapProgress());

```

```

1010     int target = findNewMapTask(tts, clusterSize, numUniqueHosts,
1011                             anyCacheLevel);

```

```

1004     int target = findNewMapTask(tts, clusterSize, numUniqueHosts, maxLevel,

```

```
1005 status.mapProgress());
```

```
1071 int target = findNewMapTask(tts, clusterSize, numUniqueHosts, maxLevel);
```

```
1028 NON_LOCAL_CACHE_LEVEL, status.mapProgress());
```

```
1094 NON_LOCAL_CACHE_LEVEL);
```

```
1206 int target = findNewReduceTask(tts, clusterSize, numUniqueHosts,
```

```
1207 status.reduceProgress());
```

```
1272 int target = findNewReduceTask(tts, clusterSize, numUniqueHosts);
```

```
1270 if (tip.getActiveTasks().size() > 1)
```

```
1335 if (tip.isSpeculating()) {
```

```
1337 LOG.debug("Chosen speculative task, current speculativeMap task count: "
```

```
1338 + speculativeMapTasks);
```

```
1339 }
```

```
1277 if (tip.getActiveTasks().size() > 1)
```

```
1345 if (tip.isSpeculating()) {
```

```
1347 LOG.debug("Chosen speculative task, current speculativeReduce task count: "
```

```
1348 + speculativeReduceTasks);
```

```
1349 }
```

```
1688 public boolean hasSpeculativeMaps() {
```

```
1689 return hasSpeculativeMaps;
```

```
1690 }
```

```
1691
```

```
1692 public boolean hasSpeculativeReduces() {
```

```
1693 return hasSpeculativeReduces;
```

```
1694 }
```

```
1695
```

```
1616 * Find a speculative task
```

```
1617 * @param list a list of tips
```

```
1618 * @param taskTracker the tracker that has requested a tip
```

```
1619 * @param avgProgress the average progress for speculation
```

```
1620 * @param currentTime current time in milliseconds
```

```
1621 * @param shouldRemove whether to remove the tips
```

```
1622 * @return a tip that can be speculated on the tracker
```

```

1697     * Retrieve a task for speculation.
1698     * If a task slot becomes available and there are less than SpeculativeCap
1699     * speculative tasks running:
1700     * 1) Ignore the request if the TT's progressRate is < SlowNodeThreshold
1701     * 2) Choose candidate tasks - those tasks whose progress rate is below
1702     *    slowTaskThreshold * mean(progress-rates)
1703     * 3) Speculate task that's expected to complete last
1704     * @param list pool of tasks to choose from
1705     * @param taskTrackerName the name of the TaskTracker asking for a task
1706     * @param taskTrackerHost the hostname of the TaskTracker asking for a task
1707     * @return the TIP to speculatively re-execute

```

```

1624     private synchronized TaskInProgress findSpeculativeTask(
1625         Collection<TaskInProgress> list, TaskTrackerStatus ttStatus,
1626         double avgProgress, long currentTime, boolean shouldRemove) {

```

```

1709     protected synchronized TaskInProgress findSpeculativeTask(
1710         Collection<TaskInProgress> list, String taskTrackerName,
1711         String taskTrackerHost) {
1712         if (list.isEmpty()) {
1713             return null;
1714         }
1715         long now = jobtracker.getClock().getTime();
1716         if (isSlowTracker(taskTrackerName) || atSpeculativeCap(list)) {
1717             return null;
1718         }
1719         // List of speculatable candidates, start with all, and chop it down
1720         ArrayList<TaskInProgress> candidates = new ArrayList<TaskInProgress>(list);

```

```

1628         Iterator<TaskInProgress> iter = list.iterator();

```

```

1722         Iterator<TaskInProgress> iter = candidates.iterator();

```

```

1632         // should never be true! (since we delete completed/failed tasks)
1633         if (!tip.isRunning()) {

```

```

1725             if (tip.hasRunOnMachine(taskTrackerHost, taskTrackerName) ||
1726                 !tip.canBeSpeculated(now)) {
1727                 //remove it from candidates

```

```

1637
1638         if (!tip.hasRunOnMachine(ttStatus.getHost(),
1639                                ttStatus.getTrackerName())) {
1640             if (tip.hasSpeculativeTask(currentTime, avgProgress)) {
1641                 // In case of shared list we don't remove it. Since the TIP failed
1642                 // on this tracker can be scheduled on some other tracker.
1643                 if (shouldRemove) {
1644                     iter.remove(); //this tracker is never going to run it again

```

```

1646         return tip;
1731         //resort according to expected time till completion
1732         Comparator<TaskInProgress> LateComparator =
1733             new EstimatedTimeLeftComparator(now);
1734         Collections.sort(candidates, LateComparator);
1735         if (candidates.size() > 0) {
1736             TaskInProgress tip = candidates.get(0);
1737             if (LOG.isDebugEnabled()) {
1738                 LOG.debug("Chose task " + tip.getTIPId() + ". Statistics: Task's : " +
1739                    tip.getCurrentProgressRate(now) + " Job's : " +
1740                    (tip.isMapTask() ? runningMapTaskStats : runningReduceTaskStats));
1741             }
1742             return tip;

```

```

2008     private synchronized TaskInProgress getSpeculativeReduce(
2009         String taskTrackerName, String taskTrackerHost) {
2010         TaskInProgress tip = findSpeculativeTask(
2011             runningReduces, taskTrackerName, taskTrackerHost);
2012         if (tip != null) {
2013             LOG.info("Choosing reduce task " + tip.getTIPId() +
2014                 " for speculative execution");
2015         } else {
2016             LOG.debug("No speculative map task found for tracker " + taskTrackerHost);
2017         }
2018         return tip;
2019     }
2020
2021     /**
2022      * Check to see if the maximum number of speculative tasks are
2023      * already being executed currently.
2024      * @param tasks the set of tasks to test

```

```

2025     * @return has the cap been reached?
2026     */
2027     private boolean atSpeculativeCap(Collection<TaskInProgress> tasks) {
2028         float numTasks = tasks.size();
2029         if (numTasks == 0) {
2030             return true; // avoid divide by zero
2031         }
2032
2033         //return true if totalSpecTask < max(10, 0.01 * total-slots,
2034         //                                0.1 * total-running-tasks)
2035
2036         if (speculativeMapTasks + speculativeReduceTasks < MIN_SPEC_CAP) {
2037             return false; // at least one slow tracker's worth of slots(default=10)
2038         }
2039         ClusterStatus c = jobtracker.getClusterStatus(false);
2040         int numSlots = c.getMaxMapTasks() + c.getMaxReduceTasks();
2041         if ((float)(speculativeMapTasks + speculativeReduceTasks) <
2042             numSlots * MIN_SLOTS_CAP) {
2043             return false;
2044         }
2045         boolean atCap = (((float)(speculativeMapTasks+
2046             speculativeReduceTasks)/numTasks) >= speculativeCap);
2047         if (LOG.isDebugEnabled()) {
2048             LOG.debug("SpeculativeCap is "+speculativeCap+", specTasks/numTasks is " +
2049                 ((float)(speculativeMapTasks+speculativeReduceTasks)/numTasks)+
2050                 ", so atSpecCap() is returning "+atCap);
2051         }
2052         return atCap;
2053     }
2054
2055     /**
2056     * A class for comparing the estimated time to completion of two tasks
2057     */
2058     private static class EstimatedTimeLeftComparator
2059     implements Comparator<TaskInProgress> {
2060         private long time;
2061         public EstimatedTimeLeftComparator(long now) {
2062             this.time = now;
2063         }
2064     }
2065
2066     /**
2067     * Estimated time to completion is measured as:
2068     * % of task left to complete (1 - progress) / progress rate of the task.
2069     */

```



```

2068      * This assumes that tasks are linear in their progress, which is
2069      * often wrong, especially since progress for reducers is currently
2070      * calculated by evenly weighting their three stages (shuffle, sort, map)
2071      * which rarely account for 1/3 each. This should be fixed in the future
2072      * by calculating progressRate more intelligently or splitting these
2073      * multi-phase tasks into individual tasks.
2074      *
2075      * The ordering this comparator defines is: task1 < task2 if task1 is
2076      * estimated to finish farther in the future => compare(t1,t2) returns -1
2077      */
2078      public int compare(TaskInProgress tip1, TaskInProgress tip2) {
2079          //we have to use the Math.max in the denominator to avoid divide by zero
2080          //error because prog and progRate can both be zero (if one is zero,
2081          //the other one will be 0 too).
2082          //We use inverse of time_remainaing=[(1- prog) / progRate]
2083          //so that (1-prog) is in denom. because tasks can have arbitrarily
2084          //low progRates in practice (e.g. a task that is half done after 1000
2085          //seconds will have progRate of 0.0000005) so we would rather
2086          //use Math.maxnon (1-prog) by putting it in the denominator
2087          //which will cause tasks with prog=1 look 99.99% done instead of 100%
2088          //which is okay
2089          double t1 = tip1.getCurrentProgressRate(time) / Math.max(0.0001,
2090              1.0 - tip1.getProgress());
2091          double t2 = tip2.getCurrentProgressRate(time) / Math.max(0.0001,
2092              1.0 - tip2.getProgress());
2093          if (t1 < t2) return -1;
2094          else if (t2 < t1) return 1;
2095          else return 0;
2096      }
2097  }
2098
2099  /**
2100   * Compares the ave progressRate of tasks that have finished on this
2101   * taskTracker to the ave of all succesfull tasks thus far to see if this
2102   * TT one is too slow for speculating.
2103   * slowNodeThreshold is used to determine the number of standard deviations
2104   * @param taskTracker the name of the TaskTracker we are checking
2105   * @return is this TaskTracker slow
2106   */
2107  protected boolean isSlowTracker(String taskTracker) {
2108      if (trackerMapStats.get(taskTracker) != null &&
2109          trackerMapStats.get(taskTracker).mean() -
2110          mapTaskStats.mean() > mapTaskStats.std()*slowNodeThreshold) {

```

```

2111         if (LOG.isDebugEnabled()) {
2112             LOG.debug("Tracker " + taskTracker +
2113                 " declared slow. trackerMapStats.get(taskTracker).mean() :" + trackerMapStats
2114                 " mapTaskStats :" + mapTaskStats);
2115         }
2116         return true;
2117     }
2118     if (trackerReduceStats.get(taskTracker) != null &&
2119         trackerReduceStats.get(taskTracker).mean() -
2120         reduceTaskStats.mean() > reduceTaskStats.std()*slowNodeThreshold) {
2121         if (LOG.isDebugEnabled()) {
2122             LOG.debug("Tracker " + taskTracker +
2123                 " declared slow. trackerReduceStats.get(taskTracker).mean() :" +
2124                 trackerReduceStats.get(taskTracker).mean() +
2125                 " reduceTaskStats :" + reduceTaskStats);
2126         }
2127         return true;
2128     }
2129     return false;
2130 }
2131
2132 static class DataStatistics{
2133     private int count = 0;
2134     private double sum = 0;
2135     private double sumSquares = 0;
2136
2137     public DataStatistics() {
2138     }
2139
2140     public DataStatistics(double initNum) {
2141         this.count = 1;
2142         this.sum = initNum;
2143         this.sumSquares = initNum * initNum;
2144     }
2145
2146     public void add(double newNum) {
2147         this.count++;
2148         this.sum += newNum;
2149         this.sumSquares += newNum * newNum;
2150     }
2151
2152     public void updateStatistics(double old, double update) {
2153         sub(old);

```

```

2153         add(update);
2154     }
2155     private void sub(double oldNum) {
2156         this.count--;
2157         this.sum -= oldNum;
2158         this.sumSquares -= oldNum * oldNum;
2159     }
2160
2161     public double mean() {
2162         return sum/count;
2163     }
2164
2165     public double var() {
2166         //  $E(X^2) - E(X)^2$ 
2167         return (sumSquares/count) - mean() * mean();
2168     }
2169
2170     public double std() {
2171         return Math.sqrt(this.var());
2172     }
2173
2174     public String toString() {
2175         return "DataStatistics: count is " + count + ", sum is " + sum +
2176             ", sumSquares is " + sumSquares + " mean is " + mean() + " std() is " + std();
2177     }
2178
2179     private void updateTaskTrackerStats(TaskInProgress tip, TaskTrackerStatus ttStatus,
2180         Map<String,DataStatistics> trackerStats, DataStatistics overallStats) {
2181         float tipDuration = tip.getExecFinishTime()-tip.getDispatchTime();
2182         DataStatistics ttStats =
2183             trackerStats.get(ttStatus.getTrackerName());
2184         double oldMean = 0.0d;
2185         //We maintain the mean of TaskTrackers' means. That way, we get a single
2186         //data-point for every tracker (used in the evaluation in isSlowTracker)
2187         if (ttStats != null) {
2188             oldMean = ttStats.mean();
2189             ttStats.add(tipDuration);
2190             overallStats.updateStatistics(oldMean, ttStats.mean());
2191         } else {
2192             trackerStats.put(ttStatus.getTrackerName(),
2193                 (ttStats = new DataStatistics(tipDuration)));
2194             overallStats.add(tipDuration);
2195         }
2196     }
2197     if (LOG.isDebugEnabled()) {

```

```

2378         LOG.debug("Added mean of " + ttStats.mean() + " to trackerStats of type " +
2379             (tip.isMapTask() ? "Map" : "Reduce") +
2380             " on " + ttStatus.getTrackerName() + ". DataStatistics is now: " +
2381             trackerStats.get(ttStatus.getTrackerName()));
2382     }
2383 }
2384
2385 public void updateStatistics(double oldProg, double newProg, boolean isMap) {
2386     if (isMap) {
2387         runningMapTaskStats.updateStatistics(oldProg, newProg);
2388     } else {
2389         runningReduceTaskStats.updateStatistics(oldProg, newProg);
2390     }
2391 }
2392
2393 public DataStatistics getRunningTaskStatistics(boolean isMap) {
2394     if (isMap) {
2395         return runningMapTaskStats;
2396     } else {
2397         return runningReduceTaskStats;
2398     }
2399 }
2400
2401 public float getSlowTaskThreshold() {
2402     return slowTaskThreshold;
2403 }
2404
2405
2406
2407
2408
2409
2410
2411
2412
2413
2414
2415
2416
2417
2418
2419
2420
2421
2422
2423
2424
2425
2426
2427
2428
2429
2430
2431
2432
2433
2434
2435
2436
2437
2438
2439
2440
2441
2442
2443
2444
2445
2446
2447
2448
2449
2450
2451
2452
2453
2454 private void decrementSpeculativeCount(boolean wasSpeculating,
2455     TaskInProgress tip) {
2456     if (wasSpeculating) {
2457         if (tip.isMapTask()) {
2458             speculativeMapTasks--;
2459             LOG.debug("Decrement count. Current speculativeMap task count: " +
2460                 speculativeMapTasks);
2461         } else {
2462             speculativeReduceTasks--;
2463             LOG.debug("Decrement count. Current speculativeReduce task count: " +
2464                 speculativeReduceTasks);
2465         }
2466     }
2467 }
2468

```

```
2588     boolean wasSpeculating = tip.isSpeculating();
```

```
2592     decrementSpeculativeCount(wasSpeculating, tip);
```

//2. TaskInProgress.java

```
35 import org.apache.hadoop.mapred.JobInProgress.DataStatistics;
```

```
77 private double oldProgressRate;
```

```
78 private long startTime = 0;
```

```
79 private long execStartTime = 0;
```

```
79 private long dispatchTime = 0; // most recent time task attempt given to TT
```

```
80 private long execStartTime = 0; // when we started first task-attempt
```

```
236 public long getStartTime() {
```

```
237     return startTime;
```

```
236 public long getDispatchTime() {
```

```
237     return this.dispatchTime;
```

```
559 572 // if task is a cleanup attempt, do not replace the complete status,
```

```
560 573 // update only specific fields.
```

```
561 574 // For example, startTime should not be updated,
```

```
562 575 // but finishTime has to be updated.
```

```
563 576 if (!isCleanupAttempt(taskid)) {
```

```
564 577     taskStatuses.put(taskid, status);
```

```
578     if ((isMapTask() && job.hasSpeculativeMaps()) ||
```

```
579         (!isMapTask() && job.hasSpeculativeReduces())) {
```

```
580         long now = jobtracker.getClock().getTime();
```

```
581         double oldProgRate = getOldProgressRate();
```

```
582         double currProgRate = getCurrentProgressRate(now);
```

```
583         job.updateStatistics(oldProgRate, currProgRate, isMapTask());
```

```
584         //we need to store the current progress rate, so that we can
```

```
585         //update statistics accurately the next time we invoke
```

```
586         //updateStatistics
```

```
587         setProgressRate(currProgRate);
```

```
588     }
```

```
863 * Return whether the TIP has a speculative task to run. We
```

```
864 * only launch a speculative task if the current TIP is really
```

```
865 * far behind, and has been behind for a non-trivial amount of
```

```
866      * time.
```

```
887      * Can this task be speculated? This requires that it isn't done or almost
888      * done and that it isn't already being speculatively executed.
889      *
```

```
868      boolean hasSpeculativeTask(long currentTime, double averageProgress) {
869          //
870          // REMIND - mjc - these constants should be examined
871          // in more depth eventually..
872          //
873
874          if (!skipping && activeTasks.size() <= MAX_TASK_EXECS &&
875              (averageProgress - progress >= SPECULATIVE_GAP) &&
876              (currentTime - startTime >= SPECULATIVE_LAG)
877              && completes == 0 && !isOnlyCommitPending()) {
878              return true;

```

```
893      boolean canBeSpeculated(long currentTime) {
894          DataStatistics taskStats = job.getRunningTaskStatistics(isMapTask());
895          if (LOG.isDebugEnabled()) {
896              LOG.debug("activeTasks.size(): " + activeTasks.size() + " "
897                  + activeTasks.firstKey() + " task's progressrate: " +
898                  getCurrentProgressRate(currentTime) +
899                  " taskStats : " + taskStats);

```

```
880      return false;
901      return (!skipping && isRunnable() && isRunning() &&
902          activeTasks.size() <= MAX_TASK_EXECS &&
903          currentTime - dispatchTime >= SPECULATIVE_LAG &&
904          completes == 0 && !isOnlyCommitPending() &&
905          (taskStats.mean() - getCurrentProgressRate(currentTime) >
906          taskStats.std() * job.getSlowTaskThreshold()));
881 907  }
912      boolean isSpeculating() {
913          return (activeTasks.size() > MAX_TASK_EXECS);
1125      * Compare most recent task attempts dispatch time to current system time so
1126      * that task progress rate will slow down as time proceeds even if no progress
1127      * is reported for the task. This allows speculative tasks to be launched for

```

```

1128     * tasks on slow/dead TT's before we realize the TT is dead/slow. Skew isn't
1129     * an issue since both times are from the JobTrackers perspective.
1130     * @return the progress rate from the active task that is doing best
1131     */
1132     public double getCurrentProgressRate(long currentTime) {
1133         double bestProgressRate = 0;
1134         for (TaskStatus ts : taskStatuses.values()) {
1135             double progressRate = ts.getProgress() / Math.max(1,
1136                 currentTime - dispatchTime);
1137             if ((ts.getRunState() == TaskStatus.State.RUNNING ||
1138                 ts.getRunState() == TaskStatus.State.SUCCEEDED) &&
1139                 progressRate > bestProgressRate) {
1140                 bestProgressRate = progressRate;
1141             }
1142         }
1143         return bestProgressRate;
1144     }
1145
1146     private void setProgressRate(double rate) {
1147         oldProgressRate = rate;
1148     }
1149     private double getOldProgressRate() {
1150         return oldProgressRate;
1151     }

```

//3. JobConf.java

```

public void setMapSpeculativeExecution(boolean speculativeExecution) {
    //setBoolean("mapred.map.tasks.speculative.execution", speculativeExecution);
    setBoolean("mapred.map.tasks.speculative.execution", true); //Changed by Kartik
}

public void setReduceSpeculativeExecution(boolean speculativeExecution) {
    //setBoolean("mapred.reduce.tasks.speculative.execution",
    //    speculativeExecution);
    setBoolean("mapred.reduce.tasks.speculative.execution",
        true); //Changed by Kartik
}

```



```
//Entirely new file – src/core/org/apache/Hadoop/util/*  
// 1. LinuxResourceCalculatorPlugin.java
```

```
package org.apache.hadoop.util;  
  
import java.io.BufferedReader;  
import java.io.FileNotFoundException;  
import java.io.FileReader;  
import java.io.IOException;  
import java.util.regex.Matcher;  
import java.util.regex.Pattern;  
  
import org.apache.commons.logging.Log;  
import org.apache.commons.logging.LogFactory;  
  
/**  
 * Plugin to calculate resource information on Linux systems.  
 */  
public class LinuxResourceCalculatorPlugin extends ResourceCalculatorPlugin {  
    private static final Log LOG =  
        LogFactory.getLog(LinuxResourceCalculatorPlugin.class);  
  
    /**  
     * proc's meminfo virtual file has keys-values in the format  
     * "key:[ \t]*value[ \t]kB".  
     */  
    private static final String PROCFS_MEMFILE = "/proc/meminfo";  
    private static final Pattern PROCFS_MEMFILE_FORMAT =  
        Pattern.compile("^([a-zA-Z]*):[ \t]*([0-9]*)[ \t]kB");  
  
    // We need the values for the following keys in meminfo  
    private static final String MEMTOTAL_STRING = "MemTotal";  
    private static final String SWAPTOTAL_STRING = "SwapTotal";  
    private static final String MEMFREE_STRING = "MemFree";  
    private static final String SWAPFREE_STRING = "SwapFree";  
    private static final String INACTIVE_STRING = "Inactive";  
    private static final int UNAVAILABLE = -1;  
  
    /**  
     * Patterns for parsing /proc/cpuinfo  
     */  
    private static final String PROCFS_CPUINFO = "/proc/cpuinfo";  
    private static final Pattern PROCESSOR_FORMAT =
```

```

        Pattern.compile("^processor[ \\t]:[ \\t]*([0-9]*)");
private static final Pattern FREQUENCY_FORMAT =
    Pattern.compile("^cpu MHz[ \\t]*:[ \\t]*([0-9.]*)");

/**
 * Pattern for parsing /proc/stat
 */
private static final String PROCFS_STAT = "/proc/stat";
private static final Pattern CPU_TIME_FORMAT =
    Pattern.compile("^cpu[ \\t]*([0-9]*)" +
        "[ \\t]*([0-9]*)[ \\t]*([0-9]*)[ \\t].*");

private String procfsMemFile;
private String procfsCpuFile;
private String procfsStatFile;
long jiffyLengthInMillis;

private long ramSize = 0;
private long swapSize = 0;
private long ramSizeFree = 0; // free ram space on the machine (kB)
private long swapSizeFree = 0; // free swap space on the machine (kB)
private long inactiveSize = 0; // inactive cache memory (kB)
private int numProcessors = 0; // number of processors on the system
private long cpuFrequency = 0L; // CPU frequency on the system (kHz)
private long cumulativeCpuTime = 0L; // CPU used time since system is on (ms)
private long lastCumulativeCpuTime = 0L; // CPU used time read last time (ms)
// Unix timestamp while reading the CPU time (ms)
private float cpuUsage = UNAVAILABLE;
private long sampleTime = UNAVAILABLE;
private long lastSampleTime = UNAVAILABLE;
private ProcfsBasedProcessTree pTree = null;

boolean readMemInfoFile = false;
boolean readCpuInfoFile = false;

/**
 * Get current time
 * @return Unix time stamp in millisecond
 */
long getCurrentTime() {
    return System.currentTimeMillis();
}

public LinuxResourceCalculatorPlugin() {
    procfsMemFile = PROCFS_MEMFILE;
    procfsCpuFile = PROCFS_CPUINFO;

```

```

procfsStatFile = PROCFS_STAT;
jiffyLengthInMillis = ProcfsBasedProcessTree.JIFFY_LENGTH_IN_MILLIS;
String pid = System.getenv().get("JVM_PID");
pTree = new ProcfsBasedProcessTree(pid);
}

```

```

/**

```

```

 * Constructor which allows assigning the /proc/ directories. This will be
 * used only in unit tests

```

```

 * @param procfsMemFile fake file for /proc/meminfo

```

```

 * @param procfsCpuFile fake file for /proc/cpuinfo

```

```

 * @param procfsStatFile fake file for /proc/stat

```

```

 * @param jiffyLengthInMillis fake jiffy length value

```

```

 */

```

```

public LinuxResourceCalculatorPlugin(String procfsMemFile,
                                     String procfsCpuFile,
                                     String procfsStatFile,
                                     long jiffyLengthInMillis) {
    this.procfsMemFile = procfsMemFile;
    this.procfsCpuFile = procfsCpuFile;
    this.procfsStatFile = procfsStatFile;
    this.jiffyLengthInMillis = jiffyLengthInMillis;
    String pid = System.getenv().get("JVM_PID");
    pTree = new ProcfsBasedProcessTree(pid);
}

```

```

/**

```

```

 * Read /proc/meminfo, parse and compute memory information only once

```

```

 */

```

```

private void readProcMemInfoFile() {
    readProcMemInfoFile(false);
}

```

```

/**

```

```

 * Read /proc/meminfo, parse and compute memory information

```

```

 * @param readAgain if false, read only on the first time

```

```

 */

```

```

private void readProcMemInfoFile(boolean readAgain) {

    if (readMemInfoFile && !readAgain) {
        return;
    }
}

```

```

// Read "/proc/memInfo" file

```

```

BufferedReader in = null;

```

```

FileReader fReader = null;

```

```

try {
    fReader = new FileReader(procfsMemFile);
    in = new BufferedReader(fReader);
} catch (FileNotFoundException f) {
    // shouldn't happen....
    return;
}

```

Matcher mat = null;

```

try {
    String str = in.readLine();
    while (str != null) {
        mat = PROCFS_MEMFILE_FORMAT.matcher(str);
        if (mat.find()) {
            if (mat.group(1).equals(MEMTOTAL_STRING)) {
                ramSize = Long.parseLong(mat.group(2));
            } else if (mat.group(1).equals(SWAPTOTAL_STRING)) {
                swapSize = Long.parseLong(mat.group(2));
            } else if (mat.group(1).equals(MEMFREE_STRING)) {
                ramSizeFree = Long.parseLong(mat.group(2));
            } else if (mat.group(1).equals(SWAPFREE_STRING)) {
                swapSizeFree = Long.parseLong(mat.group(2));
            } else if (mat.group(1).equals(INACTIVE_STRING)) {
                inactiveSize = Long.parseLong(mat.group(2));
            }
        }
        str = in.readLine();
    }
} catch (IOException io) {
    LOG.warn("Error reading the stream " + io);
} finally {
    // Close the streams
    try {
        fReader.close();
        try {
            in.close();
        } catch (IOException i) {
            LOG.warn("Error closing the stream " + in);
        }
    } catch (IOException i) {
        LOG.warn("Error closing the stream " + fReader);
    }
}

```

readMemInfoFile = true;

```

}

/**
 * Read /proc/cpuinfo, parse and calculate CPU information
 */
private void readProcCpuInfoFile() {
    // This directory needs to be read only once
    if (readCpuInfoFile) {
        return;
    }
    // Read "/proc/cpuinfo" file
    BufferedReader in = null;
    FileReader fReader = null;
    try {
        fReader = new FileReader(procfsCpuFile);
        in = new BufferedReader(fReader);
    } catch (FileNotFoundException f) {
        // shouldn't happen....
        return;
    }
    Matcher mat = null;
    try {
        numProcessors = 0;
        String str = in.readLine();
        while (str != null) {
            mat = PROCESSOR_FORMAT.matcher(str);
            if (mat.find()) {
                numProcessors++;
            }
            mat = FREQUENCY_FORMAT.matcher(str);
            if (mat.find()) {
                cpuFrequency = (long)(Double.parseDouble(mat.group(1)) * 1000); // kHz
            }
            str = in.readLine();
        }
    } catch (IOException io) {
        LOG.warn("Error reading the stream " + io);
    } finally {
        // Close the streams
        try {
            fReader.close();
            try {
                in.close();
            } catch (IOException i) {
                LOG.warn("Error closing the stream " + in);
            }
        }
    }
}

```

```

    } catch (IOException i) {
        LOG.warn("Error closing the stream " + fReader);
    }
}
readCpuInfoFile = true;
}

/**
 * Read /proc/stat file, parse and calculate cumulative CPU
 */
private void readProcStatFile() {
    // Read "/proc/stat" file
    BufferedReader in = null;
    FileReader fReader = null;
    try {
        fReader = new FileReader(procfsStatFile);
        in = new BufferedReader(fReader);
    } catch (FileNotFoundException f) {
        // shouldn't happen....
        return;
    }

    Matcher mat = null;
    try {
        String str = in.readLine();
        while (str != null) {
            mat = CPU_TIME_FORMAT.matcher(str);
            if (mat.find()) {
                long uTime = Long.parseLong(mat.group(1));
                long nTime = Long.parseLong(mat.group(2));
                long sTime = Long.parseLong(mat.group(3));
                cumulativeCpuTime = uTime + nTime + sTime; // milliseconds
                break;
            }
            str = in.readLine();
        }
        cumulativeCpuTime *= jiffyLengthInMillis;
    } catch (IOException io) {
        LOG.warn("Error reading the stream " + io);
    } finally {
        // Close the streams
        try {
            fReader.close();
            try {
                in.close();
            } catch (IOException i) {

```

```

        LOG.warn("Error closing the stream " + in);
    }
    } catch (IOException i) {
        LOG.warn("Error closing the stream " + fReader);
    }
    }
}

/** { @inheritDoc} */
@Override
public long getPhysicalMemorySize() {
    readProcMemInfoFile();
    return ramSize * 1024;
}

/** { @inheritDoc} */
@Override
public long getVirtualMemorySize() {
    readProcMemInfoFile();
    return (ramSize + swapSize) * 1024;
}

/** { @inheritDoc} */
@Override
public long getAvailablePhysicalMemorySize() {
    readProcMemInfoFile(true);
    return (ramSizeFree + inactiveSize) * 1024;
}

/** { @inheritDoc} */
@Override
public long getAvailableVirtualMemorySize() {
    readProcMemInfoFile(true);
    return (ramSizeFree + swapSizeFree + inactiveSize) * 1024;
}

/** { @inheritDoc} */
@Override
public int getNumProcessors() {
    readProcCpuInfoFile();
    return numProcessors;
}

/** { @inheritDoc} */
@Override
public long getCpuFrequency() {

```

```

    readProcCpuInfoFile();
    return cpuFrequency;
}

```

```

/** { @inheritDoc} */
@Override
public long getCumulativeCpuTime() {
    readProcStatFile();
    return cumulativeCpuTime;
}

```

```

/** { @inheritDoc} */
@Override
public float getCpuUsage() {
    readProcStatFile();
    sampleTime = getCurrentTime();
    if (lastSampleTime == UNAVAILABLE ||
        lastSampleTime > sampleTime) {
        // lastSampleTime > sampleTime may happen when the system time is changed
        lastSampleTime = sampleTime;
        lastCumulativeCpuTime = cumulativeCpuTime;
        return cpuUsage;
    }
    // When lastSampleTime is sufficiently old, update cpuUsage.
    // Also take a sample of the current time and cumulative CPU time for the
    // use of the next calculation.
    final long MINIMUM_UPDATE_INTERVAL = 10 * jiffyLengthInMillis;
    if (sampleTime > lastSampleTime + MINIMUM_UPDATE_INTERVAL) {
        cpuUsage = (float)(cumulativeCpuTime - lastCumulativeCpuTime) * 100F /
            ((float)(sampleTime - lastSampleTime) * getNumProcessors());
        lastSampleTime = sampleTime;
        lastCumulativeCpuTime = cumulativeCpuTime;
    }
    return cpuUsage;
}

```

```

/**
 * Test the { @link LinuxResourceCalculatorPlugin}
 *
 * @param args
 */
public static void main(String[] args) {
    LinuxResourceCalculatorPlugin plugin = new LinuxResourceCalculatorPlugin();
    System.out.println("Physical memory Size (bytes) : "
        + plugin.getPhysicalMemorySize());
    System.out.println("Total Virtual memory Size (bytes) : "

```



```

        + plugin.getVirtualMemorySize());
System.out.println("Available Physical memory Size (bytes) : "
        + plugin.getAvailablePhysicalMemorySize());
System.out.println("Total Available Virtual memory Size (bytes) : "
        + plugin.getAvailableVirtualMemorySize());
System.out.println("Number of Processors : " + plugin.getNumProcessors());
System.out.println("CPU frequency (kHz) : " + plugin.getCpuFrequency());
System.out.println("Cumulative CPU time (ms) : " +
        plugin.getCumulativeCpuTime());
try {
    // Sleep so we can compute the CPU usage
    Thread.sleep(500L);
} catch (InterruptedException e) {
    // do nothing
}
System.out.println("CPU usage % : " + plugin.getCpuUsage());
}

```

```

@Override
public ProcResourceValues getProcResourceValues() {
    pTree = pTree.getProcessTree();
    long cpuTime = pTree.getCumulativeCpuTime();
    long pMem = pTree.getCumulativeRssmem();
    long vMem = pTree.getCumulativeVmem();
    return new ProcResourceValues(cpuTime, pMem, vMem);
}
}

```

```
// 2. ResourceCalculatorPlugin.java
package org.apache.hadoop.util;

import org.apache.hadoop.conf.Configuration;
import org.apache.hadoop.conf.Configured;
import org.apache.hadoop.util.ReflectionUtils;

/**
 * Plugin to calculate resource information on the system.
 */
public abstract class ResourceCalculatorPlugin extends Configured {

    /**
     * Obtain the total size of the virtual memory present in the system.
     *
     * @return virtual memory size in bytes.
     */
    public abstract long getVirtualMemorySize();

    /**
     * Obtain the total size of the physical memory present in the system.
     *
     * @return physical memory size bytes.
     */
    public abstract long getPhysicalMemorySize();

    /**
     * Obtain the total size of the available virtual memory present
     * in the system.
     *
     * @return available virtual memory size in bytes.
     */
    public abstract long getAvailableVirtualMemorySize();

    /**
     * Obtain the total size of the available physical memory present
     * in the system.
     *
     * @return available physical memory size bytes.
     */
    public abstract long getAvailablePhysicalMemorySize();
}
```

```

/**
 * Obtain the total number of processors present on the system.
 *
 * @return number of processors
 */
public abstract int getNumProcessors();

/**
 * Obtain the CPU frequency of on the system.
 *
 * @return CPU frequency in kHz
 */
public abstract long getCpuFrequency();

/**
 * Obtain the cumulative CPU time since the system is on.
 *
 * @return cumulative CPU time in milliseconds
 */
public abstract long getCumulativeCpuTime();

/**
 * Obtain the CPU usage % of the machine. Return -1 if it is unavailable
 *
 * @return CPU usage in %
 */
public abstract float getCpuUsage();

/**
 * Obtain resource status used by current process tree.
 */
public abstract ProcResourceValues getProcResourceValues();

public static class ProcResourceValues {
    private final long cumulativeCpuTime;
    private final long physicalMemorySize;
    private final long virtualMemorySize;
    public ProcResourceValues(long cumulativeCpuTime, long physicalMemorySize,
                              long virtualMemorySize) {
        this.cumulativeCpuTime = cumulativeCpuTime;
        this.physicalMemorySize = physicalMemorySize;
        this.virtualMemorySize = virtualMemorySize;
    }
}

/**
 * Obtain the physical memory size used by current process tree.
 * @return physical memory size in bytes.

```

```

    */
    public long getPhysicalMemorySize() {
        return physicalMemorySize;
    }

    /**
     * Obtain the virtual memory size used by a current process tree.
     * @return virtual memory size in bytes.
     */
    public long getVirtualMemorySize() {
        return virtualMemorySize;
    }

    /**
     * Obtain the cumulative CPU time used by a current process tree.
     * @return cumulative CPU time in milliseconds
     */
    public long getCumulativeCpuTime() {
        return cumulativeCpuTime;
    }
}

/**
 * Get the ResourceCalculatorPlugin from the class name and configure it. If
 * class name is null, this method will try and return a memory calculator
 * plugin available for this system.
 *
 * @param clazz class-name
 * @param conf configure the plugin with this.
 * @return ResourceCalculatorPlugin
 */
public static ResourceCalculatorPlugin getResourceCalculatorPlugin(
    Class<? extends ResourceCalculatorPlugin> clazz, Configuration conf) {

    if (clazz != null) {
        return ReflectionUtils.newInstance(clazz, conf);
    }

    // No class given, try a os specific class
    try {
        String osName = System.getProperty("os.name");
        if (osName.startsWith("Linux")) {
            return new LinuxResourceCalculatorPlugin();
        }
    } catch (SecurityException se) {
        // Failed to get Operating System name.
    }
}

```

```
    return null;
}

// Not supported on this system.
return null;
}
}
```