



Checkpoint and Replication based Fault Tolerance for Map Reduce Framework in Cloud Environment

Monika Raghuwanshi, Prof. Deepak Kumar Gour
Department of Computer Science and Engineering
Technocrats Institute of Technology, Bhopal

Abstract:

This research work focuses on providing an efficient proactive fault tolerant technique for scientific application in which above mentioned fault is predicted proactively and handled efficiently such that all cloudlets finish their tasks before or on deadline and no virtual machine has been left out for execution. The used techniques for the fault tolerance comes in the category of reactive (fault is handled after it has occurred) and proactive (fault is predicted before it actually happens). Much work has been done in the area of reactive techniques.

Keywords: checkpoint-restart, cloud computing, distributed application, infrastructure-as-a-service, virtualization, scalability, self-healing,

I. Introduction

MapReduce programs do not explicitly relate to the division of work and distribution neither of data, nor to the many low-level complications that distributed execution entails. As long as the application obeys certain restrictions imposed by the programming model, the MapReduce engine ensures fault-tolerant and scalable execution. A MapReduce program may thus be developed and tested on a single machine or small cluster, and be deployed with a high degree of confidence on a much larger scale. This makes MapReduce a particularly good match for data-intensive applications executing in the cloud. Amazon's commercial cloud computing platform, Amazon Web Services, includes MapReduce support both as a separate high-level web service [12], and in the form of customized virtual machine images.

The adoption of MapReduce and similar high-level programming models was encouraged by another concurrent development: a new paradigm for software deployment known as cloud computing. The aim of cloud computing is to turn computing power into a common utility that is always available on demand and in abundant supply, much like electrical power delivered through the electricity grid. This long-held dream has recently become economically viable, with the construction and operation of extremely large-scale, commodity-computer data centers at low-cost locations.

Cloud providers allow customers to purchase on-demand access to computing power in various forms, for example as a cluster of virtual machines, provisioned from a much larger underlying cluster of physical machines. Customers may log in remotely and execute their own software on the virtual machines, and purchase access to additional machines whenever the need arises. As a result, computing clusters are now more accessible than ever: there is no need for small companies to make large up-front investments in hardware to provision for peak loads or anticipated future growth—virtual computing clusters hosted in the cloud may simply be expanded on demand, minimizing financial risks.

According to pay-as-you-go model, subscription-based services that are offered in the cloud are infrastructure, platform, and software (applications).

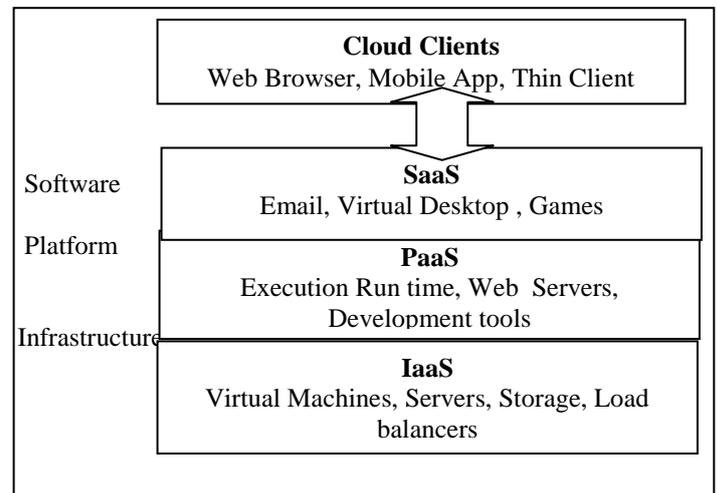


Fig.1 Different service models in cloud computing

Hence these services are provided according to the three different service models Infrastructure as a Service (IaaS), Platform as a Service (PaaS), and Software as a Service (SaaS). Depending upon the type of resource customer has requested, provider responds with the associated service model. These models further comprise of different layers essential for specific service. User can access these services anywhere and at any time with the help of internet. Details of these models are as follows and detailed figure is shown as Figure 1.

II. Fault Tolerance in Cloud Computing

Fault tolerance submits to perfect and permanent operation even when there are defective components. Continuing to work satisfactorily in the presence of obstacles is the science that builds the computing system. The efficient fault tolerant system can allow fault of different types like momentary, discontinuous or stable hardware faults, blunders in software design, operative errors, or corporal damage or externally induced upset. In a real-time application of cloud, it is performed remotely on the computing node, and the probability of error occurrence is high.

These error happening events enhance the desire for fault tolerance technology to achieve the efficiency of real-time computing in cloud environment. If a fault occurs in either the component (hardware) or the design method (software), an error will occur. These events enhance the need for fault tolerance technology.

II.1. Types of Fault Tolerance

Fault tolerance is mainly divided into two types. Fault tolerance of hardware and software.

A. Hardware fault tolerance

Mostly computer systems nowadays get recovered by themselves from failures caused due to hardware components. In systems, each component has a back up with protected duplication or redundancy. If one component fails, the other components can perform the function. Example: Mirroring technology to repair faults in storage media. A common hardware fault tolerance approach is fault masking and dynamic recovery.

i. Fault Masking:

It is a kind of structural duplication method that absolutely identifies the errors within a block of similar components. All identical components perform the same function, their output is ranked, and errors caused by the defective module are removed. Triple module redundancy (TMR) is a very form of fault masking in which components are tripled and their output is voted. Voting process for selecting redundant components can also be tripled such that voting process can rectify the failure of personage voters. If two components in the redundant triplet get failed and the vote is no longer valid, then TMR technique fails. Hybrid redundancy is an extension of TMR where triplexes components are always backed up with extra components and are used to replace failing components that can handle more faults. A voting system requires more than three times the hardware of a non redundant system, but it has the advantage of continuing the calculation without interruption in the event of a failure, and you can use an existing operating system.

ii. Dynamic recovery:

Dynamic recovery techniques are used when only single copy of a task or any calculation is executed at one time. This technology performs self-healing method. Also it performs additional backup operations using additional spare components (proactive redundancy) along with fault masking technique. This requires a special mechanism of recovery to detect faults in the module. There is no need of switching to a redundant component in this technique because faults in modules or components are parsed and recovered by performing tasks such as rollback, initialization, retry, restart, restore operation and continue the process.

B. Software fault tolerance

Software failures can be handled by using static as well as dynamic approaches just like hardware fault tolerance approach. In N version programming, static redundancy is used in which a separately written program is used that executes the same operation. Each module in this technique is developed with up to N different implementations or algorithms. Every programmer does the similar task. Every version submits its output to the voter or decision-maker. The voter or the decision maker then decides the correct output and returns it by calling it as the result of the module. Other dynamic approach is a recovery block method. In this method, a program is separated into blocks, and the acceptance test is executed after the execution of each block. In case the

acceptance test fails, then redundant code block is used for the execution. In the devise diversity approach, merging of both hardware and software fault tolerance is done to make a failure free computer system using different hardware and software on redundant channels. Main goal of this technology is to withstand both hardware and software failures.

III. Proposed Methodology

In proposed work when cloudlets arrive to the processor then the best fit approach will be applied. The best fit approach here refers to the process of allocation of suitable VM to the cloudlet which will execute it in minimum time. In case, if the best fit is not found or if any of the machine is currently not matching with the job's specifications then the proposed algorithm would eliminate this failure by processing that cloudlet in two ways, either it will split the cloudlet in two parts i.e dividing the cloudlet into two parts on the basis of its size or it will make a request for the new virtual machine. The splitting of the cloudlet is a method in which the cloudlet is divided into two parts and both parts are processed simultaneously which consumes less time and also works efficiently. In case of generating a request for the new virtual machine, memory requirement for that VM is checked. If there is no sufficient memory for the new VM, then that cloudlet will be considered as fail, otherwise a new VM is created for that cloudlet.

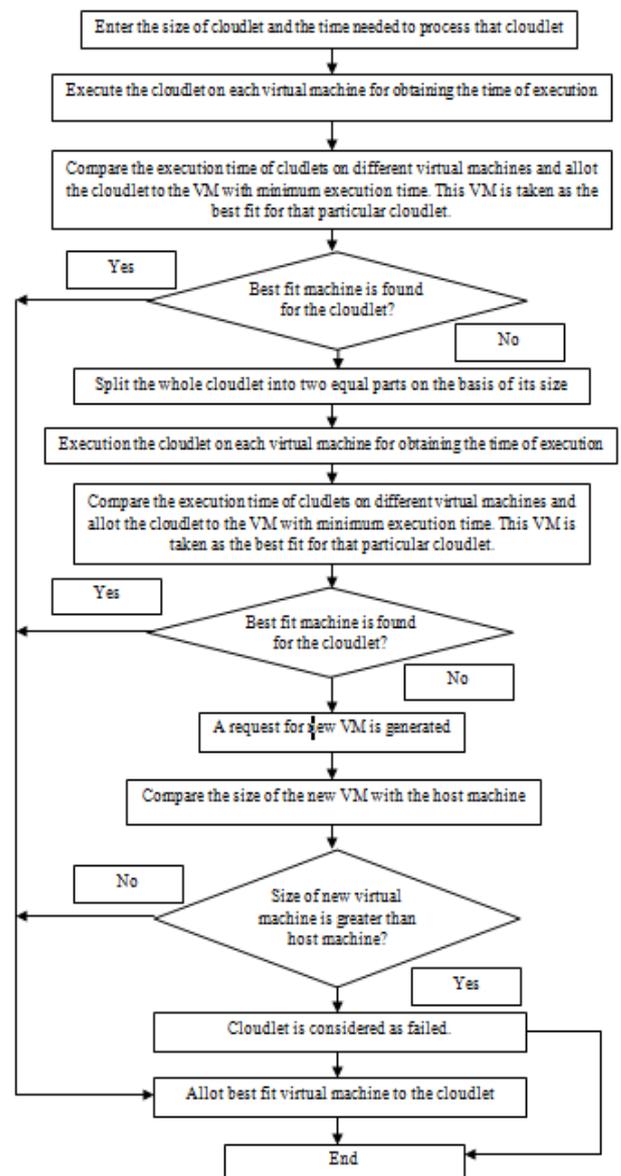


Fig.2 Flow of proposed work

- i. Enter the size of the clouplet (taken from scientific application -Montage) and the time needed to process that clouplet.
- ii. Execute the clouplet on each virtual machine for obtaining the time of execution.
- iii. Compare the execution time of clouplets on different virtual machines and allot the clouplet to the VM with minimum execution time. This VM is taken as the best fit for that particular clouplet.
- iv. If the best fit VM is not found for the execution for a particular clouplet, then split the whole clouplet into two equal parts on the basis of its size and then again repeat the step iii.
- vi. In case after splitting of clouplet into the new clouplets, if any of the clouplet does not get best fit VM, then a request for new VM is generated.
- vii. Compare the size of the new VM with the host machine. If size of new VM is less than the host machine then new VM is created and allocated to the clouplet otherwise clouplet is considered as failed.

III.1. Distributed algorithm

The problem above can be solved by highly efficient approximation algorithms, e.g., branch-and-bound, and fast off-the-shelf solvers, e.g., CPLEX, for moderate-sized input. An additional challenge arises in dealing with the MapReduce job for big data. In such a job, there are hundreds or even thousands of keys, each of which is associated with a set of variables (e.g., xpij and ypk) and constraints in our formulation, leading to a large-scale optimization problem that is hardly handled by existing algorithms and solvers in practice. In this section, we develop a distributed algorithm to solve the problem on multiple machines in a parallel manner. Our basic idea is to decompose the original large-scale problem into several distributively solvable subproblems that are coordinated by a high-level master problem.

III.2. Online algorithm

Until now, we take the data size and data aggregation ratio as input of our algorithms. In order to get their values, we need to wait all mappers to finish before starting reduce tasks, or conduct estimation via profiling on a small set of data. In practice, map and reduce tasks may partially overlap in execution to increase system throughput, and it is difficult to estimate system parameters at a high accuracy for big data applications. These motivate us to design an online algorithm to dynamically adjust data partition and aggregation during the execution of map and reduce tasks. In this section, we divide the execution of a Map Reduce job into several time slots with a length of several minutes or an hour.

IV. Results

In proposed work when clouplets arrive to the processor then the best fit approach will be applied. The best fit approach here refers to the process of allocation of suitable VM to the clouplet which will execute it in minimum time. The simulation results are shown below:

```

user@node:~$ start-all.sh
This script is deprecated. Instead use start-dfs.sh and start-yarn.sh
16/12/22 02:33:15 WARN util.NativeCodeLoader: Unable to load native-hadoop library for your platform... using builtin-java classes where applicable
Starting namenodes on [localhost]
localhost: starting namenode, logging to /usr/local/hadoop-2.6.0/logs/hadoop-user-namenode-node.out
localhost: starting datanode, logging to /usr/local/hadoop-2.6.0/logs/hadoop-user-datanode-node.out
Starting secondary namenodes [0.0.0.0]
0.0.0.0: starting secondarynamenode, logging to /usr/local/hadoop-2.6.0/logs/hadoop-user-secondarynamenode-node.out
16/12/22 02:33:44 WARN util.NativeCodeLoader: Unable to load native-hadoop library for your platform... using builtin-java classes where applicable
Starting yarn daemons
starting resourcemanager, logging to /usr/local/hadoop-2.6.0/logs/yarn-user-resourcemanager-node.out
localhost: starting nodemanager, logging to /usr/local/hadoop-2.6.0/logs/yarn-user-nodemanager-node.out
user@node:~$ jps
14883 SecondaryNameNode
14733 Jps
14689 NodeManager
14874 NameNode
14562 ResourceManager
14197 DataNode
user@node:~$

```

Fig.3 Hadoop Desktop window

Fig.3 shows hadoop desktop window. In this window write the code for the proposed research work. In this specified the data node, node manager, Jps, resource manager etc.

```

user@node:~$ hadoop fs -put /home/user/Desktop/ProjectTrafficFlowPrediction/dataset.csv
16/12/22 02:35:08 WARN util.NativeCodeLoader: Unable to load native-hadoop library for your platform... using builtin-java classes where applicable
user@node:~$ hadoop fs -put /home/user/Desktop/ProjectTrafficFlowPrediction/dataset.csv
16/12/22 02:36:18 WARN util.NativeCodeLoader: Unable to load native-hadoop library for your platform... using builtin-java classes where applicable
user@node:~$ hadoop fs -ls /TrafficFlow
16/12/22 02:37:28 WARN util.NativeCodeLoader: Unable to load native-hadoop library for your platform... using builtin-java classes where applicable
Found 1 items
-rw-r--r-- 1 user supergroup 99659 2016-12-22 02:37 /TrafficFlow/dataset.csv
user@node:~$ cd /home/user/Desktop/ProjectTrafficFlowPrediction/TrafficFlowPrediction/
user@node:~$ hadoop jar /home/user/Desktop/ProjectTrafficFlowPrediction/TrafficFlowPrediction.jar TrafficFlowPrediction /TrafficFlow/dataset.csv /TrafficFlow/output
16/12/22 02:39:06 WARN util.NativeCodeLoader: Unable to load native-hadoop library for your platform... using builtin-java classes where applicable
16/12/22 02:39:08 INFO client.RMProxy: Connecting to ResourceManager at /0.0.0.0:8032
16/12/22 02:39:09 INFO client.RMProxy: Connecting to ResourceManager at /0.0.0.0:8032
16/12/22 02:39:10 WARN mapreduce.JobSubmitter: Hadoop command-line option parsing not performed. Implement the Tool interface and execute your application with ToolRunner to remedy this.

```

Fig.4 Create Directory

Fig.4 shows create directory for process of remove replication files. In this research work create directory in hadoop distributed filling system at given path for remove replication files. The directory is specified the path for giving the data and obtain the data.

```

dataset.csv /TrafficFlow/output
16/12/22 02:39:06 WARN util.NativeCodeLoader: Unable to load native-hadoop library for your platform... using builtin-java classes where applicable
16/12/22 02:39:08 INFO client.RMProxy: Connecting to ResourceManager at /0.0.0.0:8032
16/12/22 02:39:09 INFO client.RMProxy: Connecting to ResourceManager at /0.0.0.0:8032
16/12/22 02:39:10 WARN mapreduce.JobSubmitter: Hadoop command-line option parsing not performed. Implement the Tool interface and execute your application with ToolRunner to remedy this.
16/12/22 02:39:11 INFO mapred.FileOutputFormat: Total input paths to process : 1
16/12/22 02:39:11 INFO mapreduce.JobSubmitter: number of splits:2
16/12/22 02:39:12 INFO mapreduce.JobSubmitter: Submitting tokens for job: job_1482354231061_0001
16/12/22 02:39:13 INFO input.YarnClientImpl: Submitted application application_1482354231061_0001
16/12/22 02:39:14 INFO mapreduce.Job: The url to track the job: http://node:8088/proxy/application_1482354231061_0001/
16/12/22 02:39:14 INFO mapreduce.Job: Running job: job_1482354231061_0001
16/12/22 02:39:31 INFO mapreduce.Job: Job job_1482354231061_0001 running in uber mode : false
16/12/22 02:39:31 INFO mapreduce.Job: map 0% reduce 0%
16/12/22 02:40:24 INFO mapreduce.Job: map 100% reduce 0%
16/12/22 02:40:51 INFO mapreduce.Job: map 100% reduce 100%
16/12/22 02:40:53 INFO mapreduce.Job: Job job_1482354231061_0001 completed successfully
16/12/22 02:40:54 INFO mapreduce.Job: Counters: 49
File System Counters
FILE: Number of bytes read=17795
FILE: Number of bytes written=353940
FILE: Number of read operations=0
FILE: Number of large read operations=0
FILE: Number of write operations=0
HDFS: Number of bytes read=102972
HDFS: Number of bytes written=353940

```

Fig.5 Mapreduce shows in Command window

Fig.5 shows mapreduce. In this command window shows different commands for map reduce. 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.

```

Total megabyte-seconds taken by all map tasks=142726192
Total megabyte-seconds taken by all reduce tasks=34439168
Map-Reduce Framework
Map input records=1038
Map output records=1038
Map output bytes=14026
Map output materialized bytes=16114
Input split bytes=186
Combine input records=0
Combine output records=0
Reduce input groups=25
Reduce shuffle bytes=16114
Reduce input records=1038
Reduce output records=25
Spilled Records=2476
Shuffled Map=2
Failed Shuffles=0
Merged Map outputs=2
GC time elapsed (ms)=283
CPU time spent (ms)=4870
Physical memory (bytes) snapshot=520740956
Virtual memory (bytes) snapshot=2581345280
Total committed heap usage (bytes)=328204288
Shuffle Errors
BAD_ID=0
CONNECTION=0
IO_ERROR=0
MRIOO=0
MRIOO_LENGTH=0
MRIOO_MAP=0
MRIOO_REDUCE=0
File Input Format Counters
Bytes Read=86768
File Output Format Counters
Bytes Written=328

```

Fig.6 Mapreduce framework

Fig.6 shows Mapreduce framework. In this figure shows total megabyte seconds taken by all map tasks and total megabyte seconds taken by all reduce. MapReduce is a framework using which we can write applications to process huge amounts of data, in parallel, on large clusters of commodity in a reliable manner.

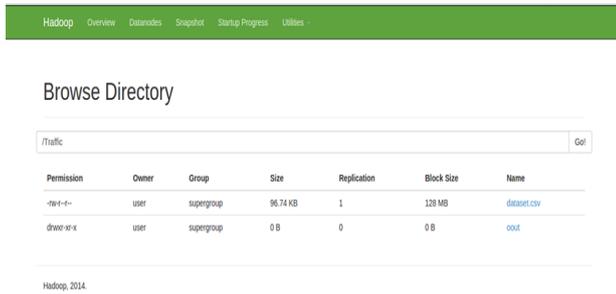


Fig.7 shows the traffic data

Fig.7 shows the traffic data. In this shows block size, file name, group name and no. of replication file. In this show two files browse.



Fig.8 Line graph of traffic flow



Fig.9 3D Bar graph of traffic flow

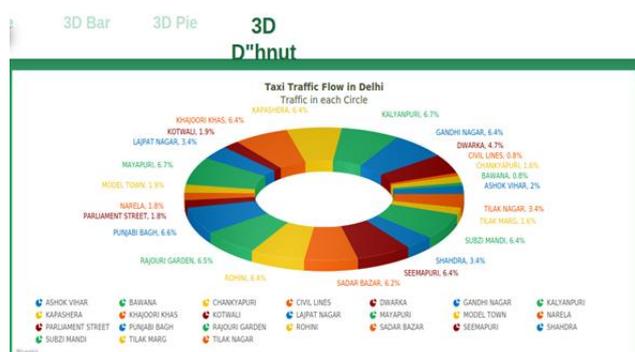


Fig.10 3D D'hunt graph of traffic flow

Fig.8 shows Line graph of traffic flow. In this figure x axis show circles and y axis show number of taxis. In this figure circles shows the traffic. There is different traffic in different places shows in this figure.

Fig.9 shows 3D Bar graph of traffic flow. In this figure x axis show circles and y axis show number of taxis. In this figure traffic shows in 3D bar chart in different places.

Fig.10 shows 3D D'hunt graph of traffic flow. In this figure traffic shows in percentage in different places. This is pie chart.

V. Conclusion

The main goal of this analysis work was to propose an economical proactive fault tolerance technique for scientific applications in Cloud. The planned technique satisfies all the analysis gaps found within the literature survey of fault tolerance techniques. The planned technique handles the fault of allocation of virtual machines proactively specified all cloudlets complete their tasks on or before point and there's successful allocation of virtual machines to all or any the cloudlets. Likelihood of failure within the planned system is less because the fault occurred within the system has been tolerated proactively. The projected work assures the reliability and the efficiency of the system where all the virtual machines are dedicated to their task and no VM has been left out.

References

- [1] Cao, Jiajun, et al. "Checkpointing as a service in heterogeneous cloud environments." Cluster, Cloud and Grid Computing (CCGrid), 2015 15th IEEE/ACM International Symposium on. IEEE, 2015.
- [2] Liu, Yang, and Wei Wei. "A replication-based mechanism for fault tolerance in mapreduce framework." Mathematical Problems in Engineering 2015 (2015).
- [3] Di, Sheng, et al. "Optimization of cloud task processing with checkpoint-restart mechanism." High Performance Computing, Networking, Storage and Analysis (SC), 2013 International Conference for. IEEE, 2013.
- [4] Egwuotuoha, Ifeanyi P., et al. "A fault tolerance framework for high performance computing in cloud." Proceedings of the 2012 12th IEEE/ACM International Symposium on Cluster, Cloud and Grid Computing (ccgrid 2012). IEEE Computer Society, 2012.
- [5] Egwuotuoha, Ifeanyi P., et al. "A survey of fault tolerance mechanisms and checkpoint/restart implementations for high performance computing systems." The Journal of Supercomputing 65.3 (2013): 1302-1326.
- [6] Liu, Jing, et al. "An intelligent job scheduling system for web service in cloud computing." Indonesian Journal of Electrical Engineering and Computer Science 11.6 (2013): 2956-2961.
- [7] J. Ansel, K. Arya, and G. Cooperman, "DMTCP: Transparent checkpointing for cluster computations and the desktop," in 23rd IEEE International Symposium on Parallel and Distributed Processing (IPDPS-09). IEEE, 2009, pp. 1-12.
- [8] J. Cao, G. Kerr, K. Arya, and G. Cooperman, "Transparent checkpointrestart over InfiniBand," in ACM Symposium on High Performance Parallel and Distributed Computing (HPDC'14). ACM Press, 2009.

- [9] D. Borthakur, The Hadoop Distributed File System: Architecture and Design, <http://hadoop.apache.org/docs/r0.18.0/hdfs design.pdf>.
- [10] S. Ghemawat, H. Gobioff, and S. T. Leung, “The Google file system,” SIGOPS Operating Systems Review, vol. 37, no. 5, pp. 29–43, 2003.
- [11] B. Selic, “Fault tolerance techniques for distributed systems,” <http://zh.scribd.com/doc/37243421/Fault-Tolerance-chniquesfor-Distributed-Systems>.
- [12] F. Wang, J. Qiu, J. Yang, B. Dong, X. Li, and Y. Li, “Hadoop high availability through metadata replication,” in Proceedings of the 1st International Workshop on Cloud Data Management (CloudDB '09), pp. 37–44, ACM, November 2009.