

Enhancing the Data Oriented Grid Scheduling Using Dynamic Error Detection

B. Radha and V. Sumathy

Abstract—Traditional distributed computing systems closely couple data handling and computation. The key features of the first batch scheduler specialized in data placement and data movement is Stork. Stork is especially designed to understand the semantics and characteristics of data placement tasks, which can include data transfer, storage allocation and de-allocation, data removal, metadata registration and replica location. The Stork also has its own drawbacks in detecting the failures, resulting from back-end system level problems, like connectivity failure which is technically untraceable by users. Error messages are not logged efficiently, and sometimes are not relevant/useful from users' point-of-view. Our study explores the possibility of efficient error detection and reporting system for such environments. Besides, early error detection and error classification have great importance in organizing data placement jobs. It is necessary to have well defined error detection and error reporting methods to increase the usability and serviceability of existing data transfer protocols and data management systems.

Index Terms—Distributed systems, data aware scheduling, error detection, grid computing, performance of systems, scheduling.

I. INTRODUCTION

The latency and the throughput are the two main factors for performance in the closely coupled distributed environment. Failure during the data transfer is very common for example the user may not be aware of the background network connectivity failure [1]. The importance of error propagation and categorization of errors in Grid computing has been mentioned clearly in [2]. The users of the system or the distributed environment may not be aware of what has been went wrong during their data transfer. This paper focuses on how to transfer data efficiently without any disturbance in transfer. The errors are detected prior to data transfer and an alternative service to data transfer is suggested. Here we focus on what sort of an error has occurred and whether the destination node can be reached for transfer of data and how to classify those errors based on their classification, and the performance of the system with and without prior error detection. There have been many efforts to implement file transfer protocols over distributed environments conforming to the security framework of the overall system. These solutions should ideally exploit communication channel to tune-up network and to satisfy high throughput and minimum transfer time

[3], [4]. Parallel data transfers, concurrent connections, and tuning network protocols such as setting TCP buffer are some of the techniques applied [5].

II. DATA SCHEDULING

There are many algorithms like genetic algorithms and heuristic techniques which are available to schedule data. In the previous work [6] of my research a comparison was made on those algorithms for their performance in the data scheduling and the best optimal algorithm was chosen. Even those algorithms work fine they have their own drawbacks which was rectified by using the techniques like Stork [7] which is specialized in data placement and data movement. This scheduler uses the job description language for data placement jobs. It also can interact with high level planners and workflow managers in order to verify when and in which network to send the data. The Stork has also its drawbacks like congestion in the network through which the data is transmitted to the destination node [8]. This sort of problem is rectified using our methodology.

The work flow model of the data and grid resource scheduling is given in the figure1 clearly [9].

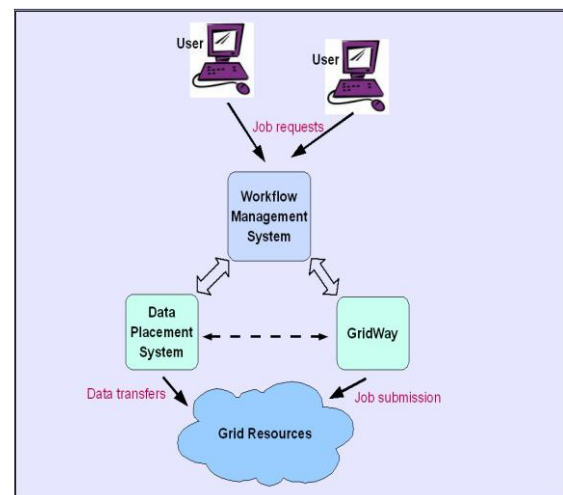


Fig. 1. A workflow management system to schedule job and data.

As per the diagram the users' job request is processed by the workflow management system which manages the job submission to the grid resource through the Grid way and the transfer of data through the data placement system [10].

A. Failure Detection

When data is been transmitted through a network there may occur an error while transmission of data and a message statement may be sent to the sender stating the actual problem that happened while transmission. The

Manuscript received May 11, 2012; revised June 17, 2012.

B. Radha is with Sri Ramkrishna Engineering College, Coimbatore, 641022, India (email: radhakbr10@gmail.com).

V. Sumathy is with Government College of Technology, Coimbatore, India.

failures will be like the remote host server may be down, or file transfer service is not functioning in the host, or file transfer service is not supporting some of the features requested, there may be a mal-functionality in the service protocol, or user credentials are not satisfied, or any other problem occurred in the source server. In addition it is also necessary to verify whether destination host and the service is available or not so that data transfer to that particular destination will not be processed until the errors are rectified. As well the information about the active services in the node will help to choose the alternative protocols for transfer.

III. EVALUATION AND DISCUSSION

In this work we proceed in such a way to prove that during the data transfer mechanism the data is transferred to the destination with prior error detection works effectively compared to the normal data transfer with the scheduler tools available.

A. Data Aware Scheduler

In this work first we have worked with the data transfer mechanism and with the code we developed the scheduler activity which is shown in Fig. 2.

Fig. 2. This picture represents the Normal grid scheduling with the CCR value as 0.1.

The output for the above input is taken from the simulation result as shown in Fig 3

Likewise for the various CCR values we have taken the output and comparison is done.

Fig. 3. This picture gives the simulation result for CCR value 0.1

B. Data Scheduling with Prior Error Notification and Rectification

Here we concentrate on how data is transferred with the Prior error detection mechanism and how the errors are rectified. From the simulation result we give the same value for the CCR and the result is taken and shown in Fig 4.

Fig. 4. The picture shows the simulation where an input for CCR is given as 0.1 for the second case.

The input for the CCR is given as 0.1 for the data scheduler with prior error notification. The output for the given value of CCR is given in the Fig. 5.

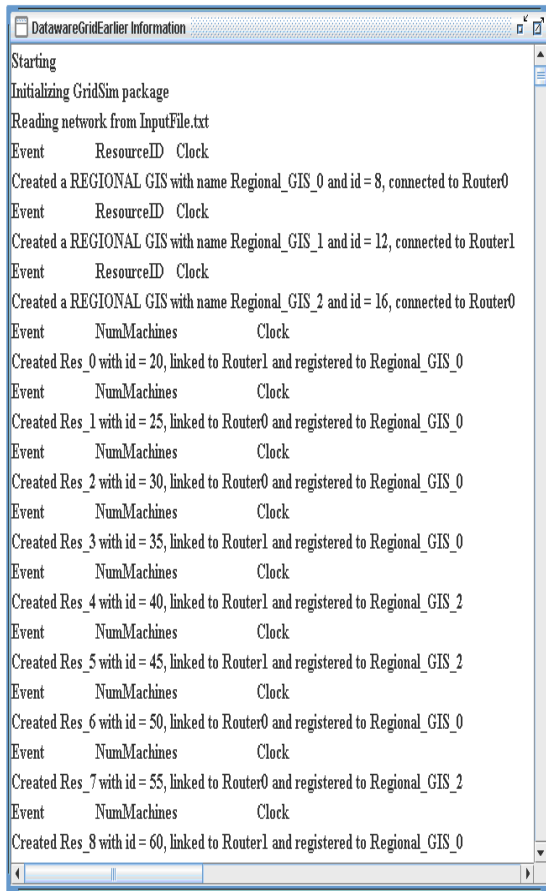


Fig. 5. The picture showing the result for the value 0.1

Different values in correspondence with CCR is given is for the three constraints and is represented here in the form of table input in Table I.

TABLE I: THE VALUES FOR X AXIS AND Y AXIS FOR ALL THE THREE CONSTRAINTS

| Algorithm | X Axis | Y Axis |
|--------------------------------------|--------|--------|
| Data aware grid | 0.1 | 103.99 |
| | 0.2 | 104.99 |
| | 0.3 | 105.98 |
| | 0.4 | 91.99 |
| Data Grid with prior error detection | 0.1 | 90.99 |
| | 0.2 | 95.96 |
| | 0.3 | 96.99 |
| | 0.4 | 91.99 |

IV. RESULTS

Based on the simulation results the values were plotted in graph and the result is given in different graphs based on CCR, PPI, and NRR, NM.

The first graph shows the execution time in milliseconds for CCR vs. Time. The CCR here represents the communication to computation ratio in which both computation time and communication of the node time is

together considered for graph and is given in Fig. 6.

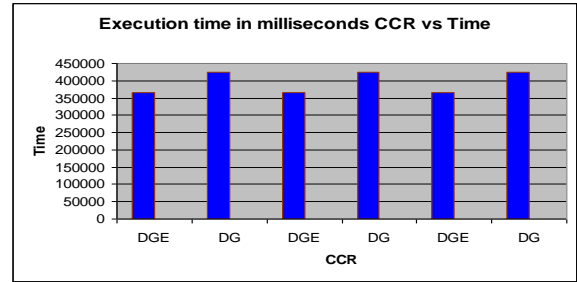


Fig. 6. The graph showing the CCR vs. Time for data grid (DG) and data grid prior error rectification.

Likewise based on the comparison of the performance of both the data grid and the data grid with prior error rectification the execution time of the later is less compared to the former as well the later performs well compared to the data grid.

The normalized value of the resource usage and the performance prediction information is compared in both the cases and is plotted in graph and is given below in figure7. where the data grid with prior error rectification shows better performance than the data scheduler. the performance of data scheduler is much higher as well the normalized resource usage of the data scheduler also seems to be high in compared the second method

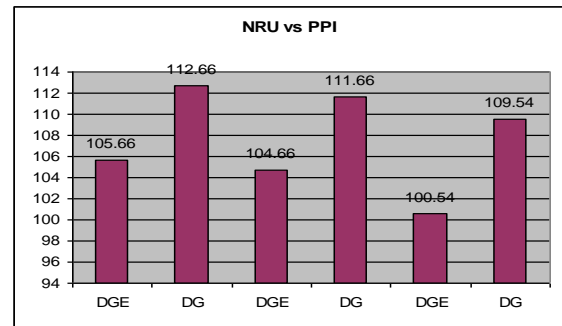


Fig. 7. The graph showing NRU vs. PPI

V. CONCLUSION

In this paper we made a comparative analysis of the data scheduler grid with prior error rectification method and the experiment is implemented in the GridSim and the results are shown in the graph and from the results it is been proved that data scheduled with prior error detection and rectification works more efficiently and fast than the data grid scheduler. All possible comparisons are made based on CCR, PPI and NRU and based on those only the conclusion is given. The future work will be focused on how the faults are rectified dynamically when the data or the resource is scheduled dynamically to the node in the heterogeneous distributed environment. The faults may be like node fault or server problem, traffic in network and the network itself becoming a problem. All the above problems will be focused and the forth coming work will be enhanced in such manner.

REFERENCES

- [1] D. Cieslak, N. Chawla, and D. Thain, "Troubleshooting Thousands of Jobs on Production Grids Using Data Mining Techniques," in *Proc. of IEEE Grid Computing*, September 2008.

- [2] D. Thain and M. Livny, "Error scope on a computational grid: Theory and practice," in *Proc. of the 11th IEEE Symposium on High Performance Distributed Computing (HPDC'02)*, pp. 199–208, 2002.
- [3] Developer's Guide. [Online]. Available: <http://www.globus.org>.
- [4] W. Allcock *et al.*, "The globus striped gridftp framework and server," in *Proc. of ACM/IEEE conference on Supercomputing*, pp. 45, 2005.
- [5] M. Balman and K. T. "Data scheduling for large scale distributed applications," in *Proc. of the 5th DCEIS, International Conference on Enterprise Information Systems (ICEIS'07)*, June, 2007.
- [6] B. Radha and V. Sumathy, "Comparison of ACO and PSO in Grid Job Scheduling," *CIIT International journal of networking and communication Engineering*, 2009.
- [7] T. Kosar and M. Balman, "A new paradigm: Data-aware scheduling in grid computing," *Future Generation Computer Systems*, vol. 6, September 2008.
- [8] T. Kosar and M. Livny "Stork: Making Data Placement a First Class Citizen in the Grid," presented at International Conference on Distributed Computing Systems, vol. 1, March 2004.
- [9] SciDAC, "Center for Enabling Distributed Petascale Science," *Technical report, A Department of Energy*, [Online]. Available: <http://www.scidac.gov/compsci/distpeta.html>
- [10] A. Peris, J. Hernandez, E. Huedo, and I. Llorente, "Data location-aware job scheduling in the grid," in *Proc. of Application to the GridWay metascheduler 17th International Conference on Computing in High Energy and Nuclear Physics (CHEP09)* vol. 209, 2010.