

REPLICA OPTIMIZATION TECHNIQUE IN DATA GRIDS



Computer Science

KEYWORDS : Optimization technique; Data Grids; Replica Selection Technique; Data movement; Grids.

Rafah M. Almuttairi

University of Babylon, Babylon, Iraq

ABSTRACT

Grid technologies are developed to be used for executing parallel applications over grid's sites in different geographical locations. For executing a grid job that contains a parallel application on a set of grid sites chosen by grid's broker, the files of data needed by the application are distributed in the data grid sites. The data grid files might be used subsequently by different grid's applications leading to multiple replicas of files in various data grid servers. Data grid files needed for a parallel application are transferred from the replica providers onto the computational site chosen by the grid broker for job execution. In this work, a novel optimization method is devised, to determine "set of cheapest" replica sites containing segments of the needed files. The objective of the proposed optimization technique, is to minimize the total transferring time needed for transferring the file's segments from the different replica sites to the computing site. Optimization technique is tested on different kinds of data grid experimental setups. We find that the best algorithm varies according to the configuration of replica providers, computational sites and clients.

Introduction

Data and computational grids are found to be powerful research-beds for executing different kinds of parallel and distributed applications [4]. To execute a parallel distributed application, broker or scheduler chooses a set of grid resources [5, 26]. In data grids, the required inputs data needed by the application are partitioned and distributed on the grid resources via data distribution strategy [28].

Grid's broker selects different sets of resources to execute parallel operations. Similarly, in high-energy physics experiments at CERN, the large amounts of data will be generated, same portions of the generated data might be used for various kinds of processing by different users/ application. Since the amount of data is large, some of the processing may involve parallel computations on the data whereby the data is partitioned and distributed among the resources used for parallel computations as shown in Figure 1.

In data grids, replication techniques are used for caching and replication policies, so that, parallel applications which are dealing with the same data can concurrently use set of data resources. Thus, replicas have to be created with different distributions on different sets of resources when the number of parallel computations on the same data is increased. Increasing replicas can help to reduce data transfer and access times for a computation and accordingly, various replica placement strategies have been proposed [21]. In this article, the following challenge that are related to replica selection are addressed :

- Select set of replica providers
- Selected replica providers share transferring job of the required files
- Cost of file transferring is reasonable
- Total file transferring time is low

The rest of the paper is organized as follows: section 2 is about the internal design of the replica broker, section 3 introduces the proposed Intelligent Optimizer highlights, section 4 covers simulation inputs and section 5 declares simulation results. Section 6 is for discussion conclusion.

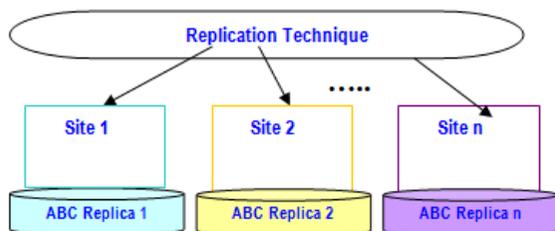


Figure1. Replication three files on different storage sites

Internal design of the Replica Broker

One of the most important operation of the previous sequenced operations is, selecting the best replicas using a Replica Broker. In this section the general description of our proposed Replica Broker is given. We followed the design of the GridBus Grid Service Broker [26] to design our Broker. In general the design of the Broker is composed of three main sub-systems:

User /Application interface sub-system

The interface layer is responsible of forwarding the input files which are:

A files description list:

It contains jobs need to be executed for users.

A replicas description list:

It contains current state of the available replica sites providers.

Core-sub-system

It converts the above inputs, lists of interface sub-system, into "jobs" and "replicas", however, "job" is the abstraction for a unit of work contains the names of the required files which are needed to be moved to a specific computational node, whereas, "replicas" contains the abstraction of replica providers.

The execution sub-system

Once the jobs are prepared and the replicas are discovered, the optimizer is started in the execution sub-system. As it is shown in Figure 2 below, the proposed optimizer selects a best replica provider(s) based on its selection algorithm. Execution sub-system has components which are:

Actuator component : it is a middleware specific component, dispatches the job to the remote storage of grid node.

Job-monitor : it updates the book-keeper by periodically monitoring the jobs using the services of the execution sub-system.

Job-manager : it takes care of cleaning up and gathering the output of the jobs when the job gets completed, Figure 2 illustrates the architecture of the internal design of GridBus Broker and shows the proposed Optimizer .

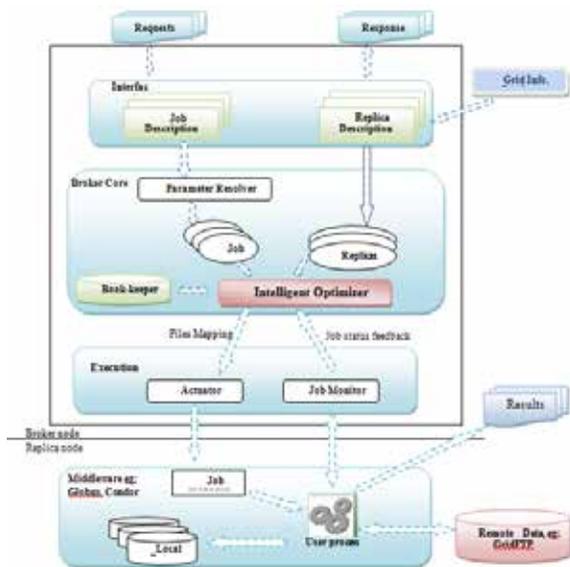


Figure 2: Enhanced Resource Broker architecture

Intelligent Optimizer

This section clarifies our approach, how it works, how it differs from the traditional model and what are the helpful advantages for us to cover the limitations of the traditional model. Our proposed optimizer uses two methods to optimize replica selection process.

discovering sets of associated replica sites

The concept of association rules of data mining approach is used to upgrade the capability of optimizer in the proposed selection technique. It means, the optimizer becomes able to choose multiple replica provider sites. The selected set of sites can concurrently share transferring files to minimize total transferring time which leads to speeding up executing data grid job as shown in Figure 3.

The selected set of replica sites should have similar characteristic of network condition, that is the main reason to use Pincer-Search algorithm. Before going ahead and explaining the steps of our proposed model, there is an important point of difference between the traditional model and our model that must be declared. In our model, we don't depend on the number of hops or the bandwidth criteria to select the best replica site. We use the stability of the network link as a criterion. It means the replica sites having the most stability links will be chosen even though their bandwidth or hops are not optimal. So the retransmission is going to be far less using our model than the traditional method.

To know the stability of links, we used a new testing route term called Single Trip Time (STT). STT is the time taken by the small packet to travel from replica's site to computing site. The STT delays include packet-transmission delays (the transmission rate out of each router and out of the replica site), packet-propagation delays (the propagation on each link), packet-queuing delays in intermediate routers and switches, and packet-processing delays (the processing delay at each router and at the replica site) for a single trip starting from replica site to computing site. It means that STT is the summation of all these delays which can reflect the stability of the links [14]. Before selection process starts, the computing site receives a periodically STTs of all replicas' sites and stores the most recent in a log file called Network History Database (NHD). To extract the best replicas from sites with the stable links, Pincer-Search algorithm is used. Pincer-Search algorithm is popular algorithm used for the as-

sociation rules discovery to extract the hidden knowledge of the large data base.

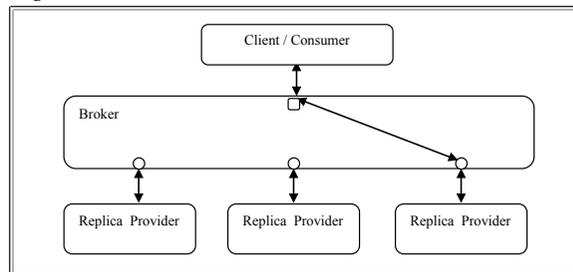


Figure 3: Multiple sites concurrently send different files

The notations and definitions of mining association rules of Pincer-Search algorithm as it has been introduced by Agrawal in 1993 [16] are introduced below:

Definitions:

- I: set of items, where $I = \{i_1, i_2, \dots, i_m\}$
- T: set of transactions, each transaction t is included in I. A transaction represents the values of RTTs between computing site and all replicas sites
- X: set of items from I
- $X \rightarrow Y$: An association rule, where $X \subseteq I, Y \subseteq I, X \cap Y = \emptyset$
- c: confidence of the rule $X \rightarrow Y$, if c% of the transactions in T that contain the set X, contain also the set Y with confidence c% of transactions in T
- s: support of rule $X \cup Y$, if s% of the transactions in T contains the set $X \cup Y$
- k : number of times of reading the whole database
- F: frequent Item with support s and minimum user support is s_1 , if $s \geq s_1$
- Infrequent Item (E): an item which is not frequent is infrequent
- MFCS: Maximum Frequent Candidate Set
- MFS: Maximum Frequent item Set
- DGTT: Data Grid Transactions Table
- ES: Efficient Set
- NHF: Network History File (column represent the replicas sites and rows represent transactions)

A Pincer-Search Method [5,8]

- Let $L_0 = \Phi; K=1; C1=\{\{i\} | i \in I\}; S0 = \Phi$
- Let MFCS = $\{\{1, 2, \dots, n\}\}; MFS = \Phi$
- Do until $Ck = \Phi$ and $Sk-1 = \Phi$
- Read DGTT and count support for Ck and MFCS
- $MFS = MFS \cup \{\text{frequent items in MFCS}\}$
- $Sk = \{\text{infrequent items in Ck}\}$
- call MSCS-gen procedure if $Sk + 1$
- call MFS-pruning procedure
- generate candidates $Ck+1$ from Ck

- if any frequent item in C_k is removed in MFS-pruning procedure then,
- call the Recovery procedure to recover candidates to C_{k+1}
- call MFCS prune procedure to prune candidates in C_{k+1}
- $k = k+1$
- return L_k

MFCS-gen procedure

- for all items $s \in S_k$
- for s is a subset of m , $MFCS = MFCS \cup \{m\}$
- for all items $e \in s$
- if $m \setminus \{e\}$ is not a subset of any item in MFCS then,

$$MFCS = MFCS \cup \{m \setminus \{e\}\}$$

- $ES = MFCS$
- Return ES (it is output from Pincer-Search algorithm is a set of replicas which have a stable links called Efficient Set ES)

Recovery procedure [5]

- for all items $c \in C_k$
- for all items $m \in MFCS$
- if the first $k-1$ items in l are also in m /* suppose $m.item = l.item_{k-1}$ */

- for i from $j+1$ to $|m|$

$$C_{k+1} = C_k \cup \{ \{l.item_1, l.item_2, \dots, l.item_k, m.item_i\} \}$$

MFS-Prune procedure [5]

- for all items $c \in C_k$
- if c is a subset of any item in the current MFS then,

delete c from C_k

MFCS-Prune procedure [5]

- for all items $c \in C_{k+1}$
- if c is not a subset of any item in the current MFCS then,

delete c from C_{k+1}

Hungarian algorithm for selecting cheapest replica set

The Hungarian algorithm is used to solve the linear assignment problem within time bounded by a polynomial expression of the number of agents. The assignment problem is a special case of the transportation problem, which is a special case of the minimum cost flow problem [15].

After getting set of sites with stable links by applying Pincer-Search algorithm, we need to assign a file send task for each replica site from this set of sites so that we get the least costs (prices). To do that, Hungarian algorithm is used.

ES algorithm.

In this section, we declare the steps of our proposed algorithm to get the best set of replica sites working concurrently with minimum cost of getting the requested files.

ES Algorithm.

Step1: Data grid jobs are submitted by User/Resource Broker (RB) to RMS.

Step2: Gathers the replica location information from Replica Location Service (RLS)

Step3: By using Iperf service we probes the routes periodically between the computing site and all replicas sites to build the NHD

Column: The columns will represent the replicas sites; each column represents one site which has a copy of the requested file. Row: The rows will represent the transactions; one transaction represents the values of STTs between computing site and all replica's sites.

Step 4: Calculate the threshold to change the NHD values to binary values

Threshold:

Calculate the means of each column $M = \frac{\sum_{i=1}^N X_i}{N}$

Calculate the Standard deviation for each column as shown below:

$$Stdi = \sqrt{\frac{1}{N} \sum_{i=1}^N (X_i - M)^2}$$

$X = RTT$

$N =$ No of transactions.

Find $Q = (Stdi / M) * 100$

Calculate the Average of all Stdi of all columns. $Av(Stdi)$

Compare $Av(Stdi)$ with $Stdi$

If $(Stdi > Av(Stdi))$ then make that value = 0 Otherwise 1

Step5: Apply Pincer-Search (PS)

PS (NHD,C,S,L) with following inputs:

Inputs: NHD: Network History Data Base which is built in step 4

C: Minimum confidence value.

S: Minimum support value.

Output: L: List of groups of sites.

$L_n = \cup G_n$, n means the Group's order; m means the number of groups.

Step 6: $ES \subseteq L_n$, ES contains sites with stable links.

Step 7: Apply Hungarian algorithm (HA)

HA(ES, Co, HL):

Inputs: Two dimensions matrix with following inputs

Row: (ES) set of sites with stable links which is gained from step7.

Column: Costs of all files in all sites.

Output: Hungarian List of minimum costs, HL.

Step 8: Uses transport services such as GridFTP or UDT to transport the requested files.

Simulation inputs

ES approach is tested using:

The Network History File NHF Real of real Grid environment called PRAGMA Grid [7]. Uohyd nodes represent a Computing Site where the files should be moved to. The rest of sites represent the replicas where the required files are stored see Figure 4. Iperf service is used to get the history of Round trip time between Uohyd and other replica sites [4]

Cost of the replicas are taken using Grid Information Service that responsible to get information from replica providers

TION MODELS

In this section, a comparison between ESM and Random Model is explained. Both models, ESM and RM are used to select the cheapest replica sites. Let us use a study case, the list of required files is $\{f_1, f_2, f_3, f_4, f_5\}$, and the list of selected replicas using ES is $ES = \{S_1, S_2, S_3, S_4, S_5\}$, so five sites are needed to concurrently transfer five files in case one site sends a single file.

Random Model

In this model, the five sites are selected randomly to get the required files of J_i [28].

So, Random Selection List is, $SL = \{f_1 \rightarrow S_1, f_2 \rightarrow S_2, f_3 \rightarrow S_5, f_4 \rightarrow S_4, f_5 \rightarrow S_3\}$. Now to find list of files using ESM when Hungarian method is applied. The result of selection is: $HL = \{f_1 \rightarrow S_5, f_2 \rightarrow S_1, f_3 \rightarrow S_3, f_4 \rightarrow S_2, f_5 \rightarrow S_4\}$. Then, we compared the total price of files of the Random Model with the price of the files selected using ESM As shown in Figure 7 below the price of the files using ESM is \$89 whereas the price of the same files using Random Model is \$112 [28].

Sequential Model SM

SM is another selection model is used to compare the result with our proposed model ESM. The selection of the Sequential Model is done as follow [28]:

Sequential selection model is $SL = \{f_1 \rightarrow S_1, f_2 \rightarrow S_2, f_3 \rightarrow S_3, f_4 \rightarrow S_4, f_5 \rightarrow S_5\}$, the total price of the Sequential Model is \$101 whereas the total price of ESM is \$89. Our proposed model ESM uses the Hungarian algorithm always gives better way to get files from replica providers with cheap price as shown in Figure 7.



Figure 7. Comparison of three selection methods

6. Discussion and conclusion

In this paper, we proposed a new replica selection model in data grid to optimize the following points:

- Minimizing the total time of executing the job by minimizing file transfer time
- Minimizing the total cost of files
- Our model utilizes two algorithms
- Pincer-Search algorithm for first optimization point [8]
- Hungarian algorithm for second optimization point [15]

The difference between our model and the traditional model is:

Our technique gets a set of sites with stable links work concurrently to transfer requested files. The traditional model selects one site as a best replica's site and getting a set of sites would not reflect the real network condition. i.e., most probably this model will not pay any attention whether these sites uncongested links or not at the transferring moment because the traditional model depends upon the Bandwidth alone or Hop counts alone which do not describe the real network condition, whereas we depend on the STT which reflects the real network conditions.

7. Future work

Being a node of PRAGMA we are looking forward to deploy our technique as an independent service in PRAGMA data grid infrastructure to speed up the execution of data grid job and minimize total cost of requested files.

REFERENCE

- [1] Rashedur M. Rahman, Ken Barker, Reda Alhaji, Replica selection in grid environment: a data-mining approach, Distributed systems and grid computing (DSGC), pp: 695 - 700 , 2005 | [2] S. Vazhkudai, J. Schopf, Using regression techniques to predict large data transfers, in: Computing Infrastructure and Applications, The International Journal of High Performance Computing Applications, IJHPCA (August) (2003). | [3] R. Kavitha, I. Foster, "Design and evaluation of replication strategies for a high performance data grid", in, Proceedings of Computing and High Energy and, Nuclear Physics, 2001. | [4] <http://goc.pragma-grid.net/cgi-bin/scm-swcb/probe.cgi?source=GOC&grid=PRAGMA> | [5] A. K. Pujar, "Data Mining Techniques", Niversities Press, India, 2001. | [6] W. Bell, et al., "OptorSim - A grid simulator for studying dynamic data replication strategies", Journal of High Performance Computing Applications 17 (4) (2003) | [7] <http://goc.pragma-grid.net/wiki/index.php/UoHyd> | [8] DI Lin, ZM Kedem, "Pincer-Search: A New Algorithm for Discovering the Maximum Frequent Set", - IEEE Transactions on Knowledge and Data, 2002. | [9] S. Vazhkudai, S. Tuecke, I. Foster, "Replica Selection in the Globus Data Grid", Cluster Computing and the Grid, IEEE International Symposium on, p. 106, 1st IEEE International Symposium on Cluster Computing and the Grid (CCGrid'01), 2001 | [10] A. Tirumala, J. Ferguson, "Iperf 1.2 - The TCP/UDP Bandwidth Measurement Tool", 2002. | [11] R. Wolski, "Dynamically forecasting network performance using the Network Weather Service", Cluster Computing (1998). | [12] Yunhong Gu, Robert L. Grossman, "UDT: UDP-based data transfer for high-speed wide area networks, Computer Networks", Volume 51, Issue 7, 16 May 2007, Pages 1777-1799, Elsevier | [13] I. Foster, S. Vazhkudai, and J. Schopf, "Predicting the performance of wide-area data transfers". In Proceedings of the 2008 ACM/IEEE conference on Supercomputing (SC 2008), pages 286-297, Nov. 2008. | [14] J. F. Kurose and K. W. Ross, "Computer networking - A top-down approach featuring the internet". Addison-Wesley-Longman, 2001. | [15] http://en.wikipedia.org/wiki/Assignment_problem | [16] R. Agrawal, I. Tomasz, and A. Swami, "Mining association rules between sets of items in large databases". In Proceedings of the 1993 ACM SIGMOD international conference on Management of data, SIGMOD '93, pages 207-216, New York, NY, USA, 1993. ACM. | [17] R. M. Rahman, R. Alhaji, and K. Barker, "Replica selection strategies in data grid". In Journal of Parallel Distributing and Computing, volume 68, pages 1561-1574, 2008. | [18] A. Jaradat, R. Salleh, and A. Abid, "Imitating k-means to enhance data selection". In Journal of Applied Sciences 2009, volume 19, pages 3569-3574, 2009. | [19] R. M. Almuttairi, R. Wankar, A. Negi, C. R. Rao, and M. S. Almahana, "New replica selection technique for binding replica sites in data grids". In Proceedings of 1st International Conference on Energy, Power and Control (EPC-IQ), pages 187-194, Dec. 2010. | [20] R. M. Almuttairi, R. Wankar, A. Negi, and C. R. Rao, "Intelligent replica selection strategy for data grid". In Proceedings of the 2010 International Conference on Grid Computing and Applications (GCA'10), pages 95-101, 2010. | [21] R. M. Almuttairi, R. Wankar, A. Negi, and C. R. Rao, "Replica selection in data grids using preconditioning of decision attributes by k-means clustering". In Proceedings of the Vaagdevi International Conference on Information Technology for Real World Problems, pages 18-23, Los Alamitos, CA, USA, 2010. IEEE Computer Society. | [22] R. M. Almuttairi, R. Wankar, A. Negi, and C. R. Rao, "Rough set clustering approach to replica selection in data grids". In Proceedings of the 10th International Conference on Intelligent Systems Design and Applications (ISDA), pages 1195-1200, Egypt, 2010. | [23] R. M. Almuttairi, R. Wankar, A. Negi, and C. R. Rao, "Smart replica selection for data grids using rough set approximations". In Proceedings of the International Conference on Computational Intelligence and Communication Networks (CICN2010), pages 466-471, Los Alamitos, CA, USA, 2010. IEEE Computer Society. | [24] R. M. Almuttairi, R. Wankar, A. Negi, and C. R. Rao, "Enhanced data replication broker". In Proceedings of the 5th Multi-disciplinary International Workshop on Artificial Intelligence (MIWAI2011), pages 286/297, 2011. | [25] Praveen Ganghishetti, Rajeev Wankar, Rafah M. Almuttairi, C. Raghavendra Rao, "Rough Set Based Quality of Service Design for Service Provisioning in Clouds", 6th International Conference on Rough Sets and Knowledge Technology (RSKT 2011) pp. 268-273, Lecture Notes in Computer Science 6954, Springer, 2011, ISBN 978-3-642-24424-7, October 2011, Canada. | [26] K. Nadiminti, S. Venugopal, H. Gibbins, T. Ma, and R. Buyya, The grid-bus grid service broker and scheduler (v2.4) user guide. In online at: <http://www.Gridbus.org/broker/2.4/manualv2.4.pdf>, 2005. | [27] S. Vadhiyar, J. Dongarra, GrADSolve, A grid-based rpc system for parallel computing with application-level scheduling, Journal of Parallel and Distributed Computing 64 (2004) 774-783. | [28] M. S. Almahana, R. M. Almuttairi, "Enhanced Replica Selection Technique for binding Replica Sites in Data Grids". In Pro. Of International Conference on Intelligent Infrastructure, 47th Annual National Convention of the Computer Society of India organized by The Kolkata Chapter December 1-2, 2012, Science City, Kolkata. |