

# Information Aggregation Techniques in Different Networks

Nini Elsa Shaji

Post-Graduate Student

Department of Computer Science and Engineering  
Karunya University, India

Shamila Ebenezer A

Assistant Professor

Department of Computer Science and Engineering  
Karunya University, India

## ABSTRACT

Information aggregation is a method for reducing the information being exchanged between Grid networks. Resource manager takes the scheduling decisions by using this aggregated information. Aggregated information is kept across each node and the detailed information is kept private, but the resources are available publicly for use. This paper does a comparative study for some of the information aggregation techniques used in different networks.

## General Terms

Grid Networks, Information Aggregation

## Keywords

Grid Networks, Information Aggregation, Scheduling

## 1. INTRODUCTION

Due to the different evolutions in science and engineering, problems in computation fields are becoming complicated. In order to solve these problems, we need to combine and use the computer resources scattered around the world. Hence we use the concept of Grid computing. Grid computing is a way to aggregate the geographically distributed and heterogeneous resources that belongs to different domains. The resource manager in Grid receives different user request and have to assign different task to resources. Static and dynamic resource information, including storage and computation capacity, number of tasks queued; resource availability, etc are some of the scheduling decisions.

Some characteristics of Grid like dynamicity, autonomy and heterogeneity makes it difficult for resource scheduling in a Grid environment. Another reason to deal with is information that being transferred to central monitor. Resources with good performance may get assigned to jobs, when allocation of job comes, thus overloading the resource and reducing the overall balance of the system. This may lead to reduced overall system performance. This paper provides a study on the different information aggregation techniques that proposes as a solution to the issue of scheduling in different networks.

This paper is organized as follows. Section 2 includes a discussion on key concepts presented in this paper. Section 3 is a comparative study of information aggregation in different networks. Section 4 gives the conclusion of this paper.

## 2. KEY CONCEPTS

### 2.1 Grid Computing

Grid computing basically means the use of multiple resources to solve a single problem. At the same time this requires large number of computer cycles and huge amount of data. Application of grid computing is commonly used in scientific

and research projects were volunteer computers joined together to do research projects by running client programs in the participating computers. Grid system has several characteristics [1] like dynamicity, heterogeneity, large-scale, resource sharing etc. It is said to be dynamic because of the constant changes in resources and tasks. Again grid is an interconnected set of resources and tasks of wide variety, which makes it heterogeneous and grid system must cater to the requirements of users in large-scale, i.e. it is capable enough to deal with large number of resources and tasks.

### 2.2 Information Aggregation

Information aggregation [2] relates to summarization of resource information in a Grid network and this information is given to resource manager in order to make scheduling decisions. As size of grid network grows, resource-related information size and dynamicity also grows rapidly, thus making the aggregation and use of this massive amount of information become a challenge for resource management system. The computation and storage tasks are guided non-locally with finer degree of granularity, so the flow of information among different systems across multiple domains will increase. Aggregation techniques are important in order to reduce amount of information being exchanged and frequency of these exchanges, at the same time maximizing its value to Grid resource manager.

Another important motivation of information aggregation is confidentiality and interoperability. This means that, when more resources or domains of resources participate in the Grid, it is desirable to keep sensitive and detailed resource information private, while resources will be still available publically for use. By this way the task scheduler will be able to efficiently and transparently use the resources. In all case, the key to information aggregation is the degree to which the summarized information helps in taking scheduling decision so as to make efficient use of resources.

Quality of information aggregation schemes [3] is measured by its effect on the efficiency of scheduler's decision and also by the reduction it brings in the total resource information that is being transferred to the central monitor.

### 2.3 Clustering

Clustering is the method of grouping computers, workstations, storage devices, interconnections etc in order to provide high availability of resources. It provides the users with a highly available resource and appears to the users as a single large computer. These are distributed for improved performance. Clusters can be categorized as hierarchical clusters [4], compute clusters, etc.

### **3. ANALYSIS OF DIFFERENT AGGREGATION TECHNIQUES**

This section contains a study on some of the aggregation schemes in different networks.

#### **3.1 Tree Based Aggregation in Hierarchical Networks**

Nodes in a sensor network are formed into a tree [5]. Data aggregation occurs at intermediate nodes in the tree and brief information is passed to the root node of the tree. This type of aggregation is suitable for applications that demands in-network data aggregation. One of the main faces of tree based network is that, it is possible to build an energy efficient data aggregation tree. Data aggregation tree in sensor network starts with a sink which transmits control message. Role of root node of the aggregation tree is assumed by the sink. Control message consists of five fields: ID, parent, power, status and hopcnt which indicates the sensor ID, its parent, its residual power, the status and the number of hops from the sink respectively.

Main advantage of tree based data aggregation is that there is a higher chance for the nodes with higher residual power to become a non-leaf tree node. In order to maintain the tree, a threshold residual power is linked with each of the node. When the residual power of a node is below the threshold power, the node will send a help message for a time period. When the child node receives this help message, it switches to a new parent or it enters to a danger state.

#### **3.2 Distributed Aggregation Trees on a Structured P2P Networks**

Distributed information aggregation has a general application on various distributed system, such as grid resource monitoring [6], P2P reputation aggregation [7], etc. The DAT [8] trees are built implicitly from native Chord routing paths. This is done by leveraging Chord topology and routing mechanisms. DAT trees have to be balanced, for this purpose a balanced routing algorithm on Chord is used. This algorithm will dynamically select the node's parent from its finger nodes by calculating its distance to the root. Almost all aggregate functions are reducible, such as count, min, max, sum, average, distinct count etc. There are two important properties for these reducible aggregation functions. Firstly, the function either return a single value from set of all values (example is min and max), or calculate some property of all the values (example is count and sum). The output value will be much smaller in size than the set of input values in both cases. Secondly, aggregation functions are applied to a large set of inputs recursively.

In order to solve the above aggregation problem, DAT trees are used. In DAT, the aggregate function is applied to each of the node on the values of its child node and the aggregated information is send to the parent node. This is a recursive process and is in a bottom-up fashion. Global aggregated value is calculated by the root node efficiently, since only the values from its direct child are collected. In order to organize the nodes in a tree, DAT algorithm uses already existing neighboring information of Chord. The Chord protocol will automatically update the neighbour nodes using the finger stabilization algorithm, whenever a node joins or leaves the network.

Main advantage of this algorithm is that the tree maintenances overhead is reduced significantly. Maximum number of nodes an aggregated message has to traverse before reaching the root

is determined by height of the tree. Branching factor of a node is determined by the number of children of that particular node. In basic DAT, each node is responsible for aggregating information from its children; the node's branching factor indicates aggregation load of the node.

#### **3.3 Trust Vector Aggregation Algorithm in P2P Networks**

Trust vector based scheme i.e. VectorTrust [9] for aggregation of distributed trust scores is emerged with the emerging internet-scale open content and resource sharing, social networks, complex cyber-physical systems etc. The feature of VectorTrust is localized and distributed simultaneous communication. By nature a VectorTrust enabled system is decentralized and doesn't rely on any centralized server or trust aggregation. VectorTrust system is build on a trust overlay network that is on top of P2P networks.

The trust overlay network is similar to a directed graph. Vertices of the network correspond to peer in a system and the directed edge will be present from one vertex to another vertex, if and only if the first vertex is a client of the second vertex in a direct transaction. A real number ranging between 0 and 1 is present for each of the directed edge and this number shows how much one vertex trust the other one. If it is 0, then the vertex never trust the other one and if it is 1, then the vertex trust the other vertex 100%. This real number is the trust rating and the link with trust rating is called Trust Vector. In the system with each direct transaction, a direct trust link is generated by the participating peers. It also assigns a trust rating in order to represent the quality of particular transaction. Each transaction in the system can either adds a new directed edge in the trust graph, or re-labels the value of an existing edge with its new trust rate.

Algorithm for Trust Vector Aggregation works as follow: firstly each peer's trust table in the network should be initialized. Then for each peer in the network, see if it is a client peer in new direct transaction, if so assign and re-label the trust rating of the peer and insert this information to the local trust table of the peer. Again for each peer in the network, firstly it sends trust table and receive trust table to all previous client peer's and from its neighbour respectively. Once the trust table is received, existing trust rating is replaced with the higher ones in local trust table. The VectorTrust scheme is efficient, accurate, scalable and robust in its nature. Computational complexity in this algorithm is less and it also converges fast. To malicious peers and malicious behavior, VectorTrust remains robust and tolerant.

#### **3.4 Simple Aggregation Algorithm in High Performance Computing**

In a practical grid system, which is composed of different domains, resource model is crucial. One solution is to use a common resource model among the different domains. Since interoperable grid systems are composed of numerous domains, there is a scalability issue with the amount of resource information exchanged between brokers. So in order to save the data transferred, latency time, and communication bandwidth interchange of the resource information is done in an aggregated form. Problem with the aggregated data is the loss of details related to each resource description, but this summarized or aggregated information is ample for the selection of best broker to submit a job. Resource model is defined by a set of resources and their relationships.

Simple aggregation algorithm [10] aggregates the data as much as possible. This algorithm looks for maximum compression for scalability, because of this, it loses more detailed information. Input for this algorithm is a set of resources and its relationship that defines the computer and some fixed attributes. For example, let's take three fixed attributes for aggregating the information. First is the processor type for Computing System resource, secondly the operating system type for Operating System resource, and the file system type for File System resource. Output of this aggregation algorithm is a set of resources in aggregated form and a set of relationships that depict the original resource. Some examples of this aggregated form is the count of resources in same category, maximum or minimum values and the sum of all values i.e., the total. Depending on the kind of resources, the aggregated information differs. Once the resource aggregation is done, algorithm looks for correspondences with the original form in order to make resource relationship.

### **3.5 Categorized Aggregation Algorithm using Grid Broker Selection**

This algorithm [10] tries to find a good balance between the resource data and scalability. Categorized aggregation algorithm considers different attributes and threshold values, in addition to the input set of resources, relationships and fixed categories. As name of the algorithm says, resources are aggregated into categories and subcategories. There will be fixed categories and the attributes and thresholds are defined as subcategory with in this category. Thus these subcategories increase the accuracy of the aggregated information. It is possible to use any number of attributes so that the detailed information will be sufficient. Level of detail can be increased by defining more threshold values, thus it is possible to avoid loss of important resource characteristics and also it is possible to maintain the benefits of aggregation.

Depending on the attributes values and a set of thresholds that define the distinct categories, the algorithm computes the category and subcategory of resources. Later on the algorithm calculates the information that contains the resource in aggregated form within the resources of same category and subcategory. At the end, the algorithm builds the relationship between the aggregated resources and the original resources. Precision of aggregated information is better in categorized aggregation algorithm rather than the simple aggregation algorithm. Complexity of this algorithm increases as the number of categories and subcategories increases.

### **3.6 Topology Aggregation Mechanisms for Delay Bandwidth Sensitive Networks**

Main purpose of topology aggregation [11] is to simplify routing by summarizing and compressing information of lower levels and giving them to the logical higher levels. Balance between accuracy and performance should be taken in to consideration because topology aggregation gives rise to distortion. Usually a network consists of nodes and links, where nodes are the originator and receiver of information while link is the one that serves as transporter by connecting the nodes. Topology aggregation mechanism has two steps. First step is to find an appropriate staircase, which will help in representing the properties of different physical paths between the border nodes. Using this staircase and the border node, each domain is transformed to a full mesh topology. Second step is to transform this full mesh topology to a maximum weighted spanning tree and star structure.

For most of the aggregation schemes, full mesh approach is considered as the basis. The main idea is to convert the topology of a routing domain into a mesh, which contains only the border node of the domain, then these border nodes in the domain is connected through logical links. Among the border nodes of the domain, the tree representation is Maximum Weighted Spanning Tree (MST) approach. There will be only one path between each pair of the border nodes, i.e. MST doesn't have any loops. For restrictive parameters, MST is distortion free. Star approach is a compromise between the full mesh and single node approach. Here the border nodes are connected to nucleus. Nucleus is a virtual center node. The links are called spokes. Spokes include the links of going from the border nodes to the nucleus and vice-versa.

### **3.7 Directed Diffusion in Wireless Sensor Networks**

In wireless sensor networks, directed diffusion [12] is a popular data aggregation scheme. Directed diffusion is data centric that means all communication is for named data. All nodes in this network will be application aware. The main advantage of directed diffusion is that it achieves energy savings.

Several elements are involved in directed diffusion such as interests, data messages, gradients and reinforcements. Interest message is a query which specifies the need of a user. It also contains a description of a sensing task for acquiring data that is supported by a sensor network. In sensor networks data is said to be the processed information. These data can be an event which in turn is a short description of the sensed phenomenon and is also represented or named using attribute-value pairs. A sensing task is broadcasted throughout the network as an interest for named data. This broadcast will set up gradients within the network. Gradient is a direction state that is created in each node that receives an interest. The direction is set towards the neighboring node from which the interest is received. Then the data start flowing towards the originating nodes of interests along multiple gradient paths. The sensor network reinforces or strengthens one or a small number of these gradient paths.

In a network, an interest is usually injected to some node in the network and this node is termed as sink. For every task, the sink will broadcast interest to its neighboring nodes. Every node maintains an interest cache, which has several fields like timestamp, duration, expiresAt, data rate field etc. It is possible to have interest aggregation also, if the interest is identical to each other. Whenever a node receives an interest, first the node will see whether that particular interest is already present in interest cache or not. If a match already exists, then the data message is dropped silently. Otherwise, the received message is added to the data cache and the data message is resend to the neighboring nodes.

### **3.8 Node and Link Aggregation in Hierarchical Networks**

Usually in hierarchical network, there are two steps in aggregation [13] one is the node aggregation within peer groups and this is followed by link aggregation between peer groups. Hierarchical structure is highly desirable and efficient. Node aggregation is one in which one or more nodes in the network are replaced with a single complex node and this node should be made transparent to the nearest neighboring node. A node connecting to the complex node is called ingress or egress nodes. An input / output relation in an

information network is characterized by QoS measures. Thus QoS measure can be categorized into additive QoS measure, min (max) functions of QoS measures along the path and the combination of these two measures.

In order to solve the complex node, nucleus node is found out and the ingress node to this nucleus is called as spoke. For each spoke a QoS measure is computed and thus the QoS value for each of link is known. Then we will choose the best QoS measure from all possible paths between the nuclei of peer groups and this value will be the aggregated one. In case, if the QoS measure is additive, the aggregated link will be the one that minimizes the sum of the QoS values. Procedure for link aggregation will get a best link as the aggregated logical link between two complex nodes. This aggregated logical link offers maximum condensation. In addition to all these, this method enhances uniformity in data structure of the hierarchy, which will help in simplifies routing.

**Table 1. Comparative Study of Different Aggregation Techniques**

Aggregation Technique	Network	Advantage	Disadvantage
Tree Based Data Aggregation	Hierarchical Network	Nodes with higher residual power can be a non-leaf node.	(i) Not that optimal. (ii) Throughout the network, overhead is involved in cluster.
Distributed Aggregation Tree	Structured Peer-to-peer Networks	(i) Tree maintenances overhead is less. (ii) Aggregation is done in an efficient and load balanced way.	Poor performance of DAT under extreme node dynamics.
Trust Vector Aggregation	P2P Networks	It is faster and computational complexity is less. Is also robust, efficient, accurate and scalable.	Malicious peer detection rate depends 90% to 99 % on network complexity.
Simple Aggregation	High Performance Computing	Information that is aggregated is accurate.	The percentage of aggregated resources and their relationships is less.
Categorized Aggregation	High Performance Computing	Precision of resource information is good.	As categories and subcategories increases, complexity of the algorithm also increases.

Topology Aggregation	Delay Bandwidth Sensitive Network	(i) Routing is simplified. (ii) This algorithm is simple, feasible and accurate. (iii) Space complexity is reduced.	-
Directed Diffusion	Wireless Sensor Network	(i) By selecting empirically good paths, it save energy, i.e. energy efficiency is obtained. (ii) Performance is also good.	(i) Careful attention is needed. (ii) Its stable only under some range of network dynamics.
Node and Link Aggregation	Hierarchical Network	(i) To aggregate the nodal topology to a given accuracy, this method is the best one. (ii) There is uniformity in data structure of the hierarchy, which will make simplified routing.	(i) For combined additive min (max) QoS measure, this aggregation is an issue. (ii) It is not possible to generalize on any of the QoS measures.

#### 4. CONCLUSION

In this paper some of the best aggregation algorithms in different networks have been discussed. All these algorithms focus on optimizing different performance measures such as network lifetime, energy consumption etc in different network. Main advantages and disadvantages of these algorithms are also described in the paper. Performance of each of the information aggregation algorithm is strongly related of coupled with infrastructure of the network. This paper will aid to take decisions on which algorithm to be used for information aggregation in different situations or scenarios.

#### 5. REFERENCES

- [1] Miguel L. Bote-Lorenzo, Yannis A. Dimitriadis, Eduardo G'omez- S'anchez, "Grid characteristics and uses: a grid definition", in: Proc. the First European Across Grids Conference, ACG'03, , pp. 291–298, 2004.
- [2] Panagiotis Kokkinos, Emmanouel Varvarigos, "Data Consolidation and Information Aggregation in Grid Networks", Advances in Grid Computing, pp. 95-118, 2008.
- [3] P. Kokkinos, E.A. Varvarigos, "Scheduling efficiency of resource information aggregation in grid networks", Future Generation Computer systems 28, pp. 9–23, 2012.

- [4] Katherine A Heller, Zoubin Ghahramani, “Bayesian Hierarchical Clustering”, Gatsby Computational Neuroscience, U. London, 2000.
- [5] Ramesh Rajagopalan and Pramod K. Varshney, “Data aggregation techniques in sensor networks: A survey”, Department of Electrical Engineering & Computer Science, Syracuse University.
- [6] S. Czajkowski, K. Fitzgerald, I. Foster, and C. Kesselman, “Grid information services for distributed resource sharing”, in Proc. of High Performance Distributed Computing Conference, 2001.
- [7] S. D. Kamvar, M. T. Schlosser, and H. Garcia-Molina, “The EigenTrust algorithm for reputation management in P2P networks”, in Proc. of the 12<sup>th</sup> International conference on World Wide Web, pp. 640.651, 2003.
- [8] M. Cai, K. Hwang, “Distributed aggregation algorithms with load-balancing for scalable grid resource monitoring”, in: IEEE International Parallel and Distributed Processing Symposium, 2007.
- [9] H. Zhao and X. Li, “VectorTrust: Trust Vector Aggregation Scheme for Trust Management in Peer-to-Peer Networks”, The research presented in this paper is supported in part by National Science Foundation (CNS-0709329), pp. 1-6.
- [10] I. Rodero, F. Guim, J. Corbalan, L. Fong, S.M. Sadjadi, “Grid broker selection strategies using aggregated resource information”, Future Generation Computer Systems 26 (1), pp. 72–86, 2010.
- [11] J. Zhang, Y. Han and L. Wang, “New Topology Aggregation Mechanisms for Delay bandwidth Sensitive Networks”, 1-4244-2424-5/08/\$20.00 IEEE, pp. 737–742, 2008.
- [12] C. Intanagonwiwat, R. Govindan, D. Estrin, J. Heidemann and F. Silva, “Directed diffusion for wireless sensor networking”, IEEE Transactions on Networking 11, pp. 2–16, 2003.
- [13] P. Van Mieghem, “Topology information condensation in hierarchical networks”, The International Journal of Computer and Telecommunications 31 (20), pp. 2115–2137, 1999.