

Employing Peer-to-Peer Services for Robust Grid Computing

Jik-Soo Kim

UMIACS and Department of Computer Science
University of Maryland, College Park, MD 20742
jiksoo@cs.umd.edu

1. Introduction

The recent growth of both the Internet and the hardware capabilities of personal computers and workstations enables distributed computing to achieve tremendous computing power by harnessing tens of thousands to millions of machines. These systems are often called *desktop grid* computing systems and leverage unused capacity on high-performance desktop PCs [1, 3]. Desktop grid computing systems mainly target complex scientific applications requiring massive computing power and resources that might exceed those available in a single supercomputing platform. However, existing platforms for desktop grid computing typically employ a client-server architecture, which has inherent shortcomings with respect to robustness, reliability and scalability since the server can be a single point of contention and failure.

Our goal is to design and build a scalable infrastructure for executing Grid applications on a widely distributed set of resources. Such infrastructure must be *decentralized, robust, highly available, and scalable*, while efficiently *mapping* application instances to available resources throughout the system. Fortunately, these are precisely the characteristics promised by new techniques and approaches in Peer-to-Peer (P2P) systems. Using P2P services can provide a robust, reliable, and scalable job submission and execution system that is able to efficiently utilize widely distributed available computational resources. Such a confluence of P2P and distributed computing is a natural step in the progression of Grid computing, and has indeed been described as inevitable [2, 4, 9].

Applications that are suited for our proposed system have both large computational requirements and relatively low I/O requirements. With our astronomy collaborators at the University of Maryland, we have identified multiple problem areas with these characteristics, mainly related to physical simulations and data analysis, including *finding habitable planets* through N-body simulations, *formation of asteroid binaries* through gravity simulations and analysis and modeling of data from the NASA *Deep Impact mission*. Additional astronomy applications may be explored in the later stages of our system development. While we are using the astronomy applications as the initial set for testing the system, applications from many other scientific and engineering disciplines, among others, can make use of the system, as evidenced by the widespread use of computational resources managed by Condor [10] or BOINC-based systems [1].

The rest of the presentation is structured as follows. Section 2 describes our overall system architecture for executing jobs using a P2P overlay network. Section 3 discusses efficient algorithms for matching jobs to resources, while Section 4 presents related work. We conclude in Section 5.

2. System Architecture

We describe a system composed from a relatively loosely coupled set of distributed, cooperating users (peers). Our goal is to use scalable P2P services to allow users to submit jobs to be run in the system and to run jobs submitted by other users on any resources available in the system, essentially allowing a group of users to form an ad-hoc set of shared resources. The overall system, from the point of view of a user, can be thought of as a combination of a centralized, Condor-like grid system for submitting and running arbitrary jobs [10], and a system such as BOINC [1] for farming out jobs from a server to be run on a (potentially very large) collection of machines in a completely distributed environment. However, to execute jobs in this decentralized and distributed environment we have to address several issues as follows:

1. *Job submission* - How can we submit a job into the P2P network?
2. *Matchmaking* - How can we find a resource that *meets* the minimum resource requirements of a job without any centralized control and information about the system for better scalability?
3. *Load balance* - How can we distribute the load (jobs) across the nodes in the system?
4. *Secure job execution* - Compute hosts should be protected from malicious jobs.
5. *Resilience to failures* - The overall system must be resilient to failures of individual resources.

For all that follows, we assume an underlying *Distributed Hash Table* (DHT) infrastructure [13, 14]. DHTs use computationally secure hashes to map arbitrary identifiers to random nodes in a system. This randomized mapping allows DHTs to present a simple insertion and lookup API that is highly robust, scalable, and efficient. A system can build upon these basic services to allow users to place idle computational resources into a general pool and draw upon the resources provided by others when needed. We insert both nodes and jobs into a single DHT, performing matchmaking by mapping a job to a node via the insertion process, and then relying on that node to find candidates that are able and willing to execute the job. By leveraging such an architecture, we are effectively *reformulating* the problem of matchmaking to one of routing in the P2P network.

A *job* in our system is the data and associated profile that describes a computation to be performed. A job profile contains several characteristics about the job, such as the client that submitted it, its minimum resource requirements, the location of input data, etc. All jobs have modest I/O requirements, with individual input data sets for our initial target applications typically on the order of a few 100 KB or less, with correspondingly small output datasets. However, the jobs for each problem are computationally intensive, since simulation runs consist of advancing physical variables forward in time by solving a set of coupled differential equations, and data analysis runs perform complex operations on the data. Finally, the jobs in the system are *independent*, which implies that no communication is needed between them. This is a typical scenario in a desktop grid computing environment, enabling many independent users to submit their jobs to a collection of node resources in the system.



Figure 1: Overall System Architecture

Figure 1 shows the overall system architecture and flow of job insertion and execution in the P2P network. The steps of job execution are as follows:

1. A client inserts a job into a node in the system (the *injection node*). The DHT provides an external mechanism that can find an existing node in the system [13, 14].
2. The injection node assigns a *Globally Unique Identifier* (GUID) to the job by using its underlying hash function and *routes* the job to the *owner node*.
3. The owner node initiates a matchmaking mechanism to find a *run node* capable of running the job.
4. Once the matchmaking mechanism finds a run node for the job, the owner node sends the job to the run node.
5. The job is inserted into the job queue of the run node, which processes jobs in FIFO order. While processing the jobs, the run node periodically sends *heartbeat* messages to the owner node.
6. When the job is finished, the run node returns the results to the client.

An owner node is responsible for monitoring the execution of the job and ensuring that its results are returned to the client. Heartbeats are communicated directly between run nodes and owner nodes, rather than through DHT routing. This soft-state message plays an important role in failure recovery during the processing of jobs in our system, as job profiles are replicated on both the owner and run nodes. If either the owner node or the run node fails, the other will detect the failure and initiate a recovery protocol so that the job can continue to make progress. If both fail before the recovery protocol completes, the client must resubmit the job. Besides dealing with recovery from failures, the run node must also be able to ensure secure execution of each job in its job queue, to prevent jobs from adversely affecting the state of a node it is running on, and vice versa. More details about our basic framework for the job submission and execution in the P2P network can be found at Kim et al. [8].

As the first step in our concrete system design and implementation, we have concentrated on developing matchmaking algorithms for a decentralized and heterogeneous environment. We next describe the current state of development of our matchmaking algorithms, and provide some preliminary results obtained via simulations.

3. Matchmaking Algorithms

A general-purpose desktop grid system must accommodate heterogeneous clusters of nodes running heterogeneous batches of jobs. The implication is that a matchmaking algorithm must incorporate both node and job information into the process that eventually maps a job onto a specific node. We refined more specific goals for our matchmaking algorithms to achieve efficient matchmaking with low cost and provide good load balancing in decentralized and heterogeneous environments as follows [6, 7]:

1. *Capability* - The matchmaking framework should allow users to specify minimum requirements for any type of resource (CPU speed, memory, etc.).
2. *Load balance* - Load (jobs) must be distributed across the nodes capable of performing them.

3. *Precision* - Resources should not be wasted. All other issues being equivalent, a job should not be assigned to a node that is over-provisioned with respect to that job.
4. *Completeness* - A valid assignment of a job to a node must be found if such an assignment exists.
5. *Low overhead* - The matchmaking must not add significant overhead to the cost of executing a job. This may be challenging, given that the matchmaking is done in a completely decentralized fashion.

In this section, we briefly describe two approaches to address these goals that we have developed: the *Rendezvous Node Tree*, and *CAN-based resource matching*.

3.1 The Rendezvous Node Tree

The Rendezvous Node Tree (RNT) uses a distributed data structure built on top of an underlying Chord DHT [14]. Specifically, the RNT copes with dynamic load balance issues by performing a limited random walk after the initial mapping to an owner node, and addresses *Completeness* by passing information describing the *maximal amount of each resource available* up and down the tree.

An RNT contains all participating nodes in the desktop grid. Each node determines its parent node based only on local information, which enables building the tree in a completely decentralized manner. Due to the uniform distribution of GUIDs of the nodes in the system, the overall height of the RNT is likely to be $O(\log N)$ where N is the total number of live nodes in the system (for details see Kim et al. [5]). Once the parent-child relationship in the RNT is determined, each node periodically sends local subtree resource information (for the subtree rooted by that node) to its parent node, and this information is *aggregated* at each level of the RNT (*hierarchical aggregation*).

Jobs are injected into the system by mapping a job to a randomly chosen node that becomes the job's owner node, which achieves a good initial load balancing by spreading the jobs across the system. The owner node initiates a search for a node on which to run the job. The search first proceeds through the subtree rooted at the owner node, only searching up the tree into subtrees rooted at the ancestors of the owner node if the subtree does not contain any satisfactory candidates. The search is *pruned* using the maximal resource information carried by the RNT. Rather than stopping at the first candidate capable of executing a given job, the search proceeds until at least k capable nodes are found for better load balancing (*extended search*). More details about the RNT can be found in [5, 7].

3.2 Content-Addressable Network

A Content-Addressable Network (CAN) [13] is a DHT that maps GUIDs of nodes and data to points in a d -dimensional space, so that each node divides up the CAN space into rectangular *zones* and maintains *neighbor* information. Based on this basic CAN, we can formulate the matchmaking problem as a routing problem in a CAN space. By treating each *resource type* as a distinct dimension, nodes and jobs can be mapped into the CAN space by using their capabilities or requirements for each resource type, respectively, to determine their coordinates. Then the matchmaking process becomes somewhat straightforward since we can search for *the closest node whose coordinates in all dimensions meet or exceed the job's requirements*.

A job is inserted into the system by using its requirements as coordinates and defining the owner of the resulting zone as the owner node of the job. The owner node creates a list of candidate run nodes, and chooses the (approximately) least loaded among them based on load information periodically exchanged between neighboring nodes. The candidate nodes are drawn from the owners of neighboring zones, such that each candidate is at least as capable as the original owner node in all dimensions (capabilities), but more capable in at least one dimension.

The basic CAN procedure works in all cases, but may cause some problems for the basic CAN mechanisms when many nodes have similar, or even identical, resource capabilities. Since the coordinates of a node are defined by its resource capabilities, identical nodes are mapped to the same place in the CAN volume. The best way to distribute ownership of a zone across multiple such nodes is not immediately obvious. Conversely, many jobs might have very similar requirements. For example, many jobs will likely be inserted into the system with no resource requirements at all specified. In this case, all of those jobs will be mapped to the single node that owns the zone containing the origin in the CAN space resulting in load imbalance.

We address this problem by supplementing the "real" dimensions (those corresponding to node capabilities) with a *virtual dimension*. Coordinates in the virtual dimension are generated *uniformly at random*. Whenever a new node joins the system, a representative point for the new node is generated by combining the resource capabilities of the node and a randomly generated virtual dimension value. Therefore, even when multiple identical nodes join the system, they are mapped to distinct locations, and CAN zone splitting

is straightforward. Similarly, when a new job is inserted into the system, the new job's coordinates become a combination of the job's requirements and a randomly assigned virtual dimension coordinate. In combination, the randomly assigned node and job coordinates act to break up clusters and spread load more evenly over nodes. More details about CAN matchmaking can be found in [7].

3.3 Experimental Results

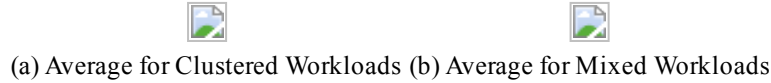


Figure 2: Job Wait Time for Clustered and Mixed Workloads

In this section, we briefly discuss and analyze experimental results obtained via simulations. We have employed an event-driven simulator to investigate the basic behavior of a P2P network, namely creating and maintaining the network and performing lookups into the DHT based on peer IDs.

The results for our matchmaking algorithms for different workload scenarios and under relatively heavy loads, with multiple clients submitting jobs over time at different average rates are shown in Figure 2. Our test workloads differ on two axes. Workloads are categorized as either clustered or mixed. The former divides all nodes and jobs into a small number of equivalence classes (in terms of resource capabilities and constraints, respectively), where all nodes or jobs in a given equivalence class are identical. The latter assigns node capabilities and job constraints randomly. Based on these concepts, the overall problem space for Grid computing environments can be divided along two axes, measuring the degree to which the nodes and jobs are either clustered or mixed. Systems such as Condor [10] mainly target mixed jobs running on clustered nodes, while systems like BOINC [1] often deal with clustered jobs on mixed nodes. Our intent is to effectively support all four scenarios.

In the experiments presented here, workloads are also distinguished by whether the jobs are lightly or heavily constrained. For a given job, each type of resource has a fixed independent probability of being constrained: "lightly-constrained" jobs have an average of 1.2 constraints (out of the 3) and "heavily-constrained" jobs have an average of 2.4. As a job has more resource requirements (i.e., heavily-constrained workloads), it is likely to be harder to match the job to the available resources, since fewer nodes in the system can meet those multiple constraints. All of the test workloads consist of 1000 nodes and 10000 jobs, each of which has an average running time of about 200 seconds. The job arrival times are based on a Poisson distribution with an average inter-arrival rate of 0.1 seconds. To see how well the workload could be balanced, we also show results for a centralized scheme that uses knowledge of the status of all nodes and jobs. Such a scheme would be very expensive to implement in a decentralized P2P system, but serves as a target for achieving the best possible load balance from an online matchmaking algorithm.

Overall, we found that for most scenarios, the CAN-based matchmaking framework shows very competitive performance in terms of balancing loads, even compared to the centralized scheme, with low matchmaking cost (in results not shown, we have verified that both the RNT and CAN can find an appropriate run node for a job with a small number of hops through the P2P overlay network). However, we found that under some conditions the CAN-based algorithm works very poorly due to serious load imbalance, namely when jobs with few resource requirements are run on nodes with heterogeneous (mixed) resource capabilities (i.e., the lightly-constrained workloads in Figure 2(b)).

In ongoing work, we have improved the basic CAN-based matchmaking mechanism to address this problem by pushing jobs into underloaded regions of the CAN space based on dynamic aggregated load information [6]. The basic concept is that when a new job is inserted into the system and routed to the owner node, the job is pushed into an underloaded region in the CAN space. To determine whether to initiate pushing of a job, a fixed amount of current system load information is propagated along each dimension in the CAN space. If the overall system is lightly loaded, the job can be pushed into the upper regions of the CAN space (farther from the origin) and utilize the more capable nodes in the system. In preliminary experiments not shown here, we have verified that the modified CAN-based matchmaking mechanism dramatically improves the quality of load balancing compared to the basic scheme presented here, still with low matchmaking cost.

4. Related Work

Recently there have been several research efforts to combine P2P and Grid computing techniques to improve the robustness, reliability and scalability of commonly available client-server based desktop grid infrastructure.

Research such as [4, 11] proposes a P2P architecture to locate and allocate resources in a Grid environment by employing a *Time-To-Live* (TTL) mechanism. TTL-based mechanisms are relatively simple but effective ways to find a resource (that meets the job requirements) in a widely distributed environment without incurring too much overhead in the search. However, such mechanisms may fail to find a resource capable of running a given job, even though such a resource exists somewhere in the network (lack of *Completeness*).

Studies on encoding static or dynamic information about computational resources using a DHT hash function for resource discovery have also been conducted [2, 12]. However, there can be a load balancing problem for these encoding techniques, since a small fraction of the nodes can contain a majority of the resource information whenever there are many nodes that have very similar (or identical) resource capabilities in the system (lack of *Load balance*).

The CCOF (Cluster Computing on the Fly) project [15] has conducted a comprehensive study of generic searching methods in a highly dynamic P2P environment to locate idle computer cycles throughout the Internet. More recent work from the CCOF researchers, on a peer-based desktop grid system called WaveGrid, constructed a *timezone-aware* overlay network based on a Content-Addressable Network [13] to use idle night-time cycles geographically distributed across the globe [16]. However, the host availability model in that work is not based on the resource requirements of the jobs nor the varying resource capabilities of nodes in the system (lack of *Capability*).

5. Conclusions and Future Work

We have proposed an architecture that employs P2P services to allow users to submit jobs to be run in the system and to run jobs submitted by other users on any resources available in the system. Our experimental results obtained via simulations show that the system can reliably execute Grid applications on a widely distributed set of resources with good load balancing and low matchmaking cost. We are in the process of building a prototype system using CAN-based matchmaking, and will characterize its behavior on real workloads, via consultation with our application-area collaborators in astronomy and physics. In the future, we will measure and report on the behavior of our system for heterogeneous environments running real applications. We believe that our research contributes to current research efforts to merge P2P and Grid computing techniques, which may lead to a truly global computing system.

References

- [1] D. P. Anderson, C. Christensen, and B. Allen. Designing a Runtime System for Volunteer Computing. In *Proceedings of the 2006 IEEE/ACM SC06 Conference*, Nov. 2006.
- [2] A. S. Cheema, M. Muhammad, and I. Gupta. Peer-to-peer Discovery of Computational Resources for Grid Applications. In *Proceedings of the 6th IEEE/ACM International Workshop on Grid Computing (GRID 2005)*, Nov. 2005.
- [3] A. Chien, B. Calder, S. Elbert, and K. Bhatia. Entropia: Architecture and Performance of an Enterprise Desktop Grid System. *Journal of Parallel and Distributed Computing*, 63(5):597-610, May 2003.
- [4] A. Iamnitchi and I. Foster. A Peer-to-Peer Approach to Resource Location in Grid Environments. *Grid Resource Management: State of the Art and Future Trends*, In J. Nabrzyski, J. M. Schopf and J. Weglarz editors, pages 413-429, Kluwer Academic Publishers, 2004.
- [5] J.-S. Kim, B. Bhattacharjee, P. J. Keleher, and A. Sussman. Matching Jobs to Resources in Distributed Desktop Grid Environments. Technical Report CS-TR-4791 and UMIACS-TR-2006-15, University of Maryland, Department of Computer Science and UMIACS, Apr. 2006.
- [6] J.-S. Kim, P. Keleher, M. Marsh, B. Bhattacharjee, and A. Sussman. Using Content-Addressable Networks for Load Balancing in Desktop Grids. In *Proceedings of the 16th IEEE International Symposium on High Performance Distributed Computing (HPDC 2007)*, June 2007. To appear.
- [7] J.-S. Kim, B. Nam, P. Keleher, M. Marsh, B. Bhattacharjee, and A. Sussman. Resource Discovery Techniques in Distributed Desktop Grid Environments. In *Proceedings of the 7th IEEE/ACM International Conference on Grid Computing (GRID 2006)*, Sept. 2006.
- [8] J.-S. Kim, B. Nam, M. Marsh, P. Keleher, B. Bhattacharjee, D. Richardson, D. Wellnitz, and A. Sussman. Creating a Robust Desktop Grid using Peer-to-Peer Services. In *Proceedings of the 2007 NSF Next Generation Software Workshop (NSFNWS 2007)*, Mar. 2007.

- [9] J. Ledlie, J. Schneidman, M. Seltzer, and J. Huth. Scooped, Again. In *Proceedings of the 2nd International Workshop on Peer-to-Peer Systems (IPTPS'03)*, Feb. 2003.
- [10] M. J. Litzkow, M. Livny, and M. W. Mutka. Condor -A Hunter of Idle Workstations. In *Proceedings of the 8th International Conference on Distributed Computing Systems*, June 1988.
- [11] C. Mastroianni, D. Talia, and O. Verta. A Super-Peer Model for Building Resource Discovery Services in Grids: Design and Simulation Analysis. In *Proceedings of the European Grid Conference (EGC2005)*, Feb. 2005.
- [12] D. Oppenheimer, J. Albrecht, D. Patterson, and A. Vahdat. Design and Implementation Tradeoffs for Wide-Area Resource Discovery. In *Proceedings of the 14th IEEE International Symposium on High Performance Distributed Computing (HPDC-14)*, July 2005.
- [13] S. Ratnasamy, P. Francis, M. Handley, R. Karp, and S. Shenker. A Scalable Content Addressable Network. In *Proceedings of the ACM SIGCOMM*, Aug. 2001.
- [14] I. Stoica, R. Morris, D. Karger, M. F. Kaashoek, and H. Balakrishnan. Chord: A Scalable Peer-to-peer Lookup Service for Internet Applications. In *Proceedings of the ACM SIGCOMM*, Aug. 2001.
- [15] D. Zhou and V. Lo. Cluster Computing on the Fly: Resource Discovery in a Cycle Sharing Peer-to-Peer System. In *Proceedings of the 4th International Workshop on Global and Peer-to-Peer Computing*, Apr. 2004.
- [16] D. Zhou and V. Lo. WaveGrid: a Scalable Fast-turnaround Heterogeneous Peer-based Desktop Grid System. In *Proceedings of the 20th International Parallel & Distributed Processing Symposium*, Apr. 2006.