

BOINC: A System for Public-Resource Computing and Storage

David P. Anderson
Space Sciences Laboratory
University of California at Berkeley
davea@ssl.berkeley.edu

Abstract

BOINC (Berkeley Open Infrastructure for Network Computing) is a software system that makes it easy for scientists to create and operate public-resource computing projects. It supports diverse applications, including those with large storage or communication requirements. PC owners can participate in multiple BOINC projects, and can specify how their resources are allocated among these projects. We describe the goals of BOINC, the design issues that we confronted, and our solutions to these problems.

1 Public-Resource Computing

The world's computing power and disk space is no longer primarily concentrated in supercomputer centers and machine rooms. Instead it is distributed in hundreds of millions of personal computers and game consoles belonging to the general public. **Public-resource computing** (also known as "Global Computing" or "Peer-to-peer computing") uses these resources to do scientific supercomputing. This paradigm enables previously infeasible research. It also encourages public awareness of current scientific research, it catalyzes global communities centered around scientific interests, and it gives the public a measure of control over the directions of scientific progress.

Public-resource computing emerged in the mid-1990s with two projects, GIMPS and Distributed.net. In 1999, our group launched SETI@home [1], which has attracted millions of participants worldwide. SETI@home now runs on about 1 million computers, providing a sustained processing rate of over 70 TeraFLOPS (in contrast, the largest conventional supercomputer, the NEC Earth Simulator, provides about 35 TeraFLOPs).

The potential resource is much larger. The number of Internet-connected PCs is growing rapidly, and is projected to reach 1 billion by 2015. Together, these PCs could provide many PetaFLOPs of computing power. The public-resource approach applies to storage as well as computing.

If 100 million computer users each provide 10 Gigabytes of storage, the total (one Exabyte, or 10^{18} bytes) would exceed the capacity of any centralized storage system.

In spite of this resource, and an abundance of promising applications, relatively few large-scale public-resource projects have emerged. This is due in part to the lack of appropriate middleware (client and server software, management tools, user-centered web features, and so on). Some open-source systems have been developed, such as Cosm, jxta, and XtremWeb [4], but these systems provide only part of the necessary functionality. Commercial systems such as Entropia [2] and United Devices are more full-featured but not free.

1.1 Contrast with Grid computing

Public-resource computing and Grid computing share the goal of better utilizing existing computing resources. However, there are profound differences between the two paradigms, and it is unlikely that current Grid middleware [5] will be suitable for public-resource computing.

Grid computing involves organizationally-owned resources: supercomputers, clusters, and PCs owned by universities, research labs, and companies. These resources are centrally managed by IT professionals, are powered on most of the time, and are connected by full-time, high-bandwidth network links. There is a symmetric relationship between organizations: each one can either provide or use resources. Malicious behavior such as intentional falsification of results would be handled outside the system, e.g. by firing the perpetrator.

In contrast, public resource computing involves an asymmetric relationship between projects and participants. Projects are typically small academic research groups with limited computer expertise and manpower. Most participants are individuals who own Windows, Macintosh and Linux PCs, connected to the Internet by telephone or cable modems or DSL, and often behind network-address translators (NATs) or firewalls. The computers are frequently turned off or disconnected from the Internet. Participants

are not computer experts, and participate in a project only if they are interested in it and receive "incentives" such as credit and screensaver graphics. Projects have no control over participants, and cannot prevent malicious behavior.

Accordingly there are different requirements on middleware for public resource computing than for Grid computing. For example, BOINC's features such as redundant computing, cheat-resistant accounting, and support for user-configurable application graphics are not necessary in a Grid system; in fact, most of the features described in the remainder of this paper apply primarily to public-resource computing.

Conversely, Grid computing has many requirements that public-resource computing does not. A Grid architecture must accommodate many existing commercial and research-oriented academic systems, and must provide a general mechanism for resource discovery and access. In fact, it must address all the issues of dynamic heterogeneous distributed systems, an active area of Computer Science research for several decades. This has led to architecture such as Open Grid Services Architecture [7], which achieve generality at the price of complexity and, to some extent, performance.

2 BOINC

BOINC (Berkeley Open Infrastructure for Network Computing) is an platform for public-resource distributed computing. BOINC is being developed at U.C. Berkeley Spaces Sciences Laboratory by the group that developed and continues to operate SETI@home. BOINC is open source and is available at <http://boinc.berkeley.edu>.

2.1 Goals of BOINC

BOINC's general goal is to advance the public resource computing paradigm: to encourage the creation of many projects, and to encourage a large fraction of the world's computer owners to participate in one or more projects. Specific goals include:

Reduce the barriers of entry to public-resource computing. BOINC allows a research scientist with moderate computer skills to create and operate a large public-resource computing project with about a week of initial work and an hour per week of maintenance. The server for a BOINC-based project can consist of a single machine configured with common open-source software (Linux, Apache, PHP, MySQL, Python).

Share resources among autonomous projects. BOINC-based projects are autonomous. Projects are not centrally authorized or registered. Each project operates its own servers and stands completely on its own. Nevertheless, PC owners can seamlessly participate in multiple

projects, and can assign to each project a "resource share" determining how scarce resource (such as CPU and disk space) are divided among projects. If most participants register with multiple projects, then overall resource utilization is improved: while one project is closed for repairs, other projects temporarily inherit its computing power. On a particular computer, the CPU might work for one project while the network is transferring files for another.

Support diverse applications. BOINC accommodates a wide range of applications; it provides flexible and scalable mechanism for distributing data, and its scheduling algorithms intelligently match requirements with resources. Existing applications in common languages (C, C++, FORTRAN) can run as BOINC applications with little or no modification. An application can consist of several files (e.g. multiple programs and a coordinating script). New versions of applications can be deployed with no participant involvement.

Reward participants. Public-resource computing projects must provide "incentives" in order to attract and retain participants. The primary incentive for many participants is **credit**: a numeric measure of how much computation they have contributed. BOINC provides a credit-accounting system that reflects usage of multiple resource types (CPU, network, disk), is common across multiple projects, and is highly resistant to "cheating" (attempts to gain undeserved credit). BOINC also makes it easy for projects to add visualization graphics to their applications, which can provide screensaver graphics.

2.2 Projects using BOINC

A number of public-resource computing projects are using BOINC. The requirements of these projects have shaped the design of BOINC.

SETI@home, a continuation of the original SETI@home project [1], performs digital signal processing of radio telescope data from the Arecibo radio observatory. A BOINC-based version of this project has been developed, and we are currently shifting the existing SETI@home user base (over 500,000 active participants) to the BOINC-based version. The BOINC-based SETI@home will use client disks to archive data, eliminating the need for its central tape archive.

Predictor@home: [11] This project, based at The Scripps Research Institute, studies protein behavior using CHARMM, a FORTRAN program for macromolecular dynamics and mechanics. It is operational within Scripps, and is being readied for a public launch.

Folding@home [10]. This project is based at Stanford University. It studies protein folding, misfolding, aggregation, and related diseases. It uses novel computational methods and distributed computing to simulate time scales thou-

sands to millions of times longer than previously achieved. A BOINC-based project has been implemented and is being tested.

Climateprediction.net [14]. The aim of this project (based at Oxford University) is to quantify and reduce the uncertainties in long-term climate prediction based on computer simulations. This is accomplished by running large numbers of simulations with varying forcing scenarios (initial and boundary conditions, including natural and man-made components) and internal model parameters. The Climateprediction.net application (a million-line FORTRAN program) produces a 2 GB detailed output file for each 50-year simulation run (which takes about 3 PC-months). These output files need to be uploaded and examined in a small fraction of cases - for example, when the smaller summary output file indicates a possible bug in the model.

Climate@home. This project is a collaboration of researchers at NCAR, MIT, UCAR, Rutgers, Lawrence Berkeley Lab, and U.C. Berkeley. Its scientific goals are similar to those of Climateprediction.net, but it will be using the NCAR Community Climate System Model (CCSM). It will collaborate with Climateprediction.net to maximize compatibility and minimize redundant effort, and to enable a systematic comparison of different climate models.

CERN projects. CERN (in Geneva, Switzerland) is deploying a BOINC-based project on 1,000 in-house PCs, and plans to launch the project publicly in coordination with its 50th anniversary in October 2004. The project's current application is a FORTRAN program that simulates the behavior of the LHC (Large Hadron Collider) as a function of the parameters of individual superconducting magnets. CERN researchers are investigating other applications.

Einstein@home. This project involves researchers from University of Wisconsin, U.C. Berkeley, California Institute of Technology, LIGO Hanford Observatory, University of Glasgow, and the Albert Einstein Institute. Its purpose is to detect certain types of gravitational waves, such as those from spinning neutron stars, that can be detected only by using highly selective filtering techniques that require extreme computing power. It will analyze data from the Laser Interferometry Gravitational Observatory (LIGO) and the British/German GEO6000 gravitational wave detector.

UCB/Intel study of Internet resources. This project, a collaboration between researchers at the U.C. Berkeley Computer Sciences Department and the Intel Berkeley Research Laboratory, seeks to study the structure and performance of the consumer Internet, together with the performance, dependability and usage characteristics of home PCs, in an effort to understand what resources are available for peer-to-peer services. This project need to perform actions at specific times of day, or in certain time ranges. While performing these actions other BOINC applications must be suspended. The BOINC API supports these re-

quirements.

2.3 Overview of the BOINC implementation

A BOINC **project** corresponds to an organization or research group that does public-resource computing. It is identified by a single **master URL**, which is the home page of its web site and also serves as a directory of scheduling servers. Participants register with projects. A project can involve one or more **applications**, and the set of applications can change over time.

The server complex of a BOINC project is centered around a relational database that stores descriptions of applications, platforms, versions, workunits, results, accounts, teams, and so on. Server functions are performed by a set of web services and daemon processes: **Scheduling servers** handles RPCs from clients; it issues work and handles reports of completed results. **Data servers** handles file uploads using a certificate-based mechanism to ensure that only legitimate files, with prescribed size limits, can be uploaded. File downloads are handled by plain HTTP.

BOINC provides tools (Python scripts and C++ interfaces) for creating, starting, stopping and querying projects; adding new applications, platforms, and application versions, creating workunits, and monitoring server performance. BOINC is designed to be used by scientists, not system programmers or IT professionals; the tools are simple and well-documented, and a full-featured project can be created in a few hours.

Participants join a BOINC-based project by visiting the project's web site, filling out a registration form, and downloading the **BOINC client**. The client can operate in several modes: as a screensaver that shows the graphics of running applications; as a Windows service, which runs even when no users are logged in and logs errors to a database; as an application that provides a tabular view of projects, work, file transfers, and disk usage, and as a UNIX command-line program that communicates through stdin, stdout and stderr, and can be run from a cron job or startup file.

BOINC provides tools that let participants remotely install the client software on large numbers of machines, and attach the client to accounts on multiple projects.

3 Design issues and solutions

3.1 Describing computation and data

BOINC uses a simple but rich set of abstractions for files, applications, and data. A project defines **application versions** for various platforms (Windows, Linux/x86, Mac OS/X, etc.). An application can consist of an arbitrary set of files.

A **workunit** represents the inputs to a computation: the application (but not a particular version) a set of references input files, and sets of command-line arguments and environment variables. Each workunit has parameters such as compute, memory and storage requirements and a soft deadline for completion. A **result** represents the result of a computation: it consists of a reference to a workunit and a list of references to output files.

Files (associated with application versions, workunits, or results) have project-wide unique names and are immutable. Files can be replicated: the description of a file includes a list of URLs from which it may be downloaded or uploaded. Files can have associated attributes indicating, for example, that they should remain resident on a host after their initial use, that they must be validated with a digital signature, or that they must be compressed before network transfer.

When the BOINC client communicates with a scheduling server it reports completed work, and receives an XML document describing a collection of the above entities. The client then downloads and uploads files and runs applications; it maximizes concurrency, using multiple CPUs when possible and overlapping communication and computation.

BOINC's computational system also provides a distributed storage facility (of computational inputs or results, or of data not related to distributed computation) as a byproduct. This storage facility is much different from peer-to-peer storage systems such as Gnutella, PAST [13] and Oceanstore [9]. In these systems, files can be created by any peer, and there is no central database of file locations. This leads to a set of technical problems (e.g. naming and file location) that are not present in the BOINC facility.

3.2 Redundant computing

Public-resource computing projects must deal with erroneous computational results. These results arise from malfunctioning computers (typically induced by overclocking) and occasionally from malicious participants.

BOINC provides support for **redundant computing**, a mechanism for identifying and rejecting erroneous results. A project can specify that N results should be created for each workunit. Once $M \leq N$ of these have been distributed and completed, an application-specific function is called to compare the results and possibly select a **canonical result**. If no consensus is found, or if results fail, BOINC creates new results for the workunit, and continues this process until either a maximum result count or a timeout limit is reached.

Malicious participants can potentially game the system by obtaining large numbers of results and detecting groups of results that comprise a quorum. BOINC makes this difficult by a work-distribution policy that sends only at most one result of a given workunit to a given user. Projects can

also limit the total number of results sent to a given host per day.

BOINC implements redundant computing using several server daemon processes:

- The **transitioner** implements the redundant computing logic: it generates new results as needed and identifies error conditions.
- The **validator** examines sets of results and selects canonical results. It includes an application-specific result-comparison function.
- The **assimilator** handles newly-found canonical results. Includes an application-specific function which typically parses the result and inserts it into a science database.
- The **file deleter** deletes input and output files from data servers when they are no longer needed.

In this architecture servers and daemons can run on different hosts and can be replicated, so BOINC servers are scalable. Availability is enhanced because some daemons can run even while parts of the project are down (for example, the scheduler server and transitioner can operate even if the science database is down).

Some numerical applications produce different outcomes for a given workunit depending on the machine architecture, operating system, compiler, and compiler flags. In such cases it may be difficult to distinguish between results that are correct but differ because of numerical variation, and results that are erroneous. BOINC provides a feature called **homogeneous redundancy** for such applications. When this feature is enabled, the BOINC scheduler send results for a given workunit only to hosts with the same operation system name and CPU vendor. In this case, strict equality can be used to compare results. BOINC is compatible with other schemes for ensuring result correctness [8].

3.3 Failure and backoff

Public-resource computing projects may have millions of participants and a relatively modest server complex. If all the participants simultaneously try to connect to the server, a disastrous overload condition will generally develop. BOINC has a number of mechanisms to prevent this. All client/server communication uses exponential backoff in the case of failure. Thus, if a BOINC server comes up after an extended outage, its connection rate will be the long-term average.

The exponential backoff scheme is extended to computational errors as well. If, for some reason, an application fails immediately on a given host, the BOINC client will not

repeatedly contact the server; instead, it will delay based on the number of failures.

3.4 Participant preferences

Computer owners generally participate in distributed computing projects only if they incur no significant inconvenience, cost, or risk by doing so. BOINC lets participants control how and when their resources are used. Using these controls, called **general preferences**, participants specify the hysteresis limits of work buffering on machines (which determines the frequency of network activity); whether BOINC can do work while mouse/keyboard input is active; during what hours can BOINC do work; how much disk space can BOINC use; how much network bandwidth can BOINC use; and so on. These preferences are edited via a web interface, and are propagated to all hosts attached to the account. Participants can create separate sets of preferences for computers at home, work, and school.

Some non-obvious controls are important to certain classes of participants. For example, DSL service in some countries has monthly transfer limits (typically a few hundred MB). BOINC provides a preference for upload, download and combined transfer limits over arbitrary periods. Some BOINC-based applications perform computations that are so floating-point intensive that they cause CPU chips to overheat. BOINC allows users to specify a duty cycle for such applications on a given CPU.

3.5 Credit and accounting

BOINC provides an accounting system in which there is a single unit of "credit", a weighted combination of computation, storage, and network transfer. This can be metered in various ways. By default, the BOINC client runs benchmarks on each CPU, and a result's "claimed credit" is based on this benchmark and elapsed CPU time. Credit "cheating" is made difficult using the redundancy mechanism described above: Each result claims a certain amount of credit, but is granted only the average or minimum (the policy is project-specified) of the claimed credit of correct results.

Our experience with SETI@home has shown that participants are highly motivated by credit, and are particularly interested in their ranking relative to other users. This information is typically displayed on web-based "leaderboards" showing the ranking of participants or teams of participants. There are many ways in which leaderboards can be subdivided, filtered, ordered, and displayed. For example, a leaderboard might show only participants from a particular country, or only those using a single PC; and it might rank entities by total or recent average credit. Rather than supply all these views, BOINC provides a mechanism that

exports credit-related data (at the level of participant, team, and host) in XML files that can be downloaded and processed by **credit statistics sites** operated by third parties. Several of these currently exist.

As part of the accounting system, BOINC provides a **cross-project identification** mechanism that allows accounts on different projects with the same email address to be identified, in a way that doesn't allow email addresses to be inferred. This mechanism allows leaderboard sites to display credit statistics summed over multiple BOINC-based projects.

Participants demand immediate gratification; they want to see their credit totals increase at least daily. Thus projects with long workunits (such as the climate prediction projects) need to grant credit incrementally as the workunit is being processed. BOINC offers a **trickle messages** mechanism, providing bidirectional, asynchronous, reliable, ordered messages, piggybacked onto the regular client/server RPC traffic. This can be used to convey credit or to report a summary of computational state; in the latter case, reply messages can abort wayward computations.

Some projects use non-CPU computational resources. For example, SETI@home is working with NVIDIA to make a version of its application that can do a major portion of the signal processing on an NVIDIA graphics coprocessor chip. This imposes several requirements on BOINC. To maximize resource usage, the local scheduler must run such applications concurrently with CPU-intensive applications. The BOINC API allows applications to report computation (in units of FLOPs, for example) to the core client, since computation can no longer be measured in terms of benchmarks and CPU times. The BOINC API also allows applications to describe the specific architectural features (e.g. the presence of a GPU, and its model number) so that this information can be stored in the database and made available for debugging and resource analysis.

3.6 User community features

BOINC provides participant-oriented web sites features such as

- The ability to form teams.
- The ability to create and browse 'user profiles' including text and images.
- Message boards, including a **dynamic FAQ** system that encourages participants to answer each others' questions. These facilities are integrated with the accounting system: credit and seniority provide a form of reputation system [12].

These features are important in attracting and retaining participants, and in providing a "customer support" mechanism that consumes little project resources.

3.7 Handling large numbers of platforms

Although the bulk of public computing resources use the Windows/Intel platform, there are many other platforms, far more than can easily be accessed by a typical project. BOINC provides a flexible and novel framework for distributing application executables. Normally, these are compiled and distributed by the project itself, for a given set of platforms (those accessible to the project). This mechanism is fine for most participants, but it's inadequate for some, such as:

- Participants who, for security reasons, want to only run executables they have compiled themselves.
- Participants whose computers have platforms not supported by the project.
- Participants who want to optimize applications for particular architectures.

To meet these needs BOINC provides an **anonymous platform mechanism**, usable with projects that make their application source code available. Participants can download and compile the application source code (or obtain executables from a third-party source) and, via an XML configuration file, inform the BOINC client of these application versions. Then, when the client communicates with that project's server, it indicates that its platform is "anonymous" and supplies a list of available application versions; the server supplies workunits (but not application versions) accordingly.

3.8 Graphics and screensaver behavior

The BOINC client software appears monolithic to participants but actually consists of several components:

- The **core client** performs network communication with scheduling and data servers, executes and monitors applications, and enforces preferences. It is implemented as a hierarchy of interacting finite-state machines, can manage unlimited concurrency of network transfers and computations, and is highly portable.
- A **client GUI** provides a spreadsheet-type view of the projects to which the host is attached, the work and file transfers in progress, and the disk usage. It also provides an interface for operations such as joining and quitting projects. It communicates with the core client via XML/RPC over a local TCP connection.
- An **API**, linked with application executables, that interacts with the core client to report CPU usage and fraction done, to handle requests to provide graphics, and to provide heartbeat functionality.

- A **screensaver** program (implemented within the platform-specific screensaver framework) which, when activated, instructs the core client to provide screensaver graphics. The core client generally delegates graphics this to a running graphics-capable application; if there is none, it generates graphics itself.

BOINC also provides a framework for **project preferences**, specific to individual projects; this can be used, for example, to control the appearance of visualization graphics.

3.9 Local scheduling

The BOINC core client, in its decisions of when to get work and from what project, and what tasks to execute at a given point, implements a "local scheduling policy". This policy has several goals:

- To maximize resource usage (i.e. to keep all processors busy);
- To satisfy result deadlines;
- To respect the participant's resource share allocation among projects;
- To maintain a minimal "variety" among projects. This goal stems from user perception in the presence of long workunits. Participants will become bored or confused if they have registered for several projects and see only one project running for several months.

The core client implements a scheduling policy, based on a dynamic "resource debt" to each project, that is guided by these goals.

4 Conclusion

We have described the public-resource computing paradigm, contrasted it with Grid computing, and presented the design of a software system, BOINC, that facilitates it. BOINC is being used by several existing projects (SETI@home, Predictor@home, climateprediction.net) and by several other projects in development.

Many areas of the BOINC design are incomplete. For example, some projects require efficient data replication: Einstein@home uses large (40 MB) input files, and a given input file may be sent to a large number of hosts (in contrast with projects like SETI@home, where each input file is different). In its initial form, Einstein@home will simply send the files separately to each host, using a system of replicated data servers. Eventually we plan to use a mechanism such as BitTorrent [3] for efficiently replicating files using peer-to-peer communication.

We are studying policies and mechanisms for dealing with scarce-resource situations (i.e., when disks become full due to expanding BOINC or non-BOINC disk usage, or when user preferences change). These mechanisms must enforce the participant's resource shares. They must delete files in a rational order based on resource share, replication level, project-specified priority, and so on. We are also studying the issue of multiple disks; BOINC now uses storage only on the disk or filesystem on which it is installed. This mechanism must allow user to specify preferences for whether and how alternate disks are used, and it must deal with situations where hosts share network-attached storage.

References

- [1] D. P. Anderson, J. Cobb, E. Korpela, M. Lebofsky, and D. Werthimer. "SETI@home: An experiment in public-resource computing". *Communications of the ACM*, Nov. 2002, Vol. 45 No. 11, pp. 56-61.
- [2] A. Chien, B. Calder, S. Elbert, K. Bhatia. "Entropy: architecture and performance of an enterprise desktop grid system". *Journal of Parallel and Distributed Computing* 63, 5. P. 597-610, May 2003
- [3] B. Cohen. "Incentives Build Robustness in BitTorrent", *Workshop on Economics of P2P systems*. June 2003.
- [4] Gilles Fedak, Cecile Germain, Vincent Neri, Franck Cappello. "XtremWeb: A Generic Global Computing Platform". *IEEE/ACM - CCGRID'2001 Special Session Global Computing on Personal Devices*, May 2001, IEEE press
- [5] I. Foster and C. Kesselman. *Globus: A metacomputing infrastructure toolkit*. *Int'l Supercomputer Applications*, 11(2), p. 115-128, 1997.
- [6] I. Foster and A. Iamnitchi. "Death, Taxes, and the Convergence of Peer-to-Peer and Grid Computing". In *2nd International Workshop on Peer-to-Peer Systems (IPTPS '03)*, Feb. 2003.
- [7] I. Foster, C. Kesselman, J.M. Nick and S. Tuecke. "Grid Services for Distributed Systems Integration", *IEEE Computer*, 35 (6). 2002.
- [8] C. Germain. "Result Checking in Global Computing Systems". *ACM Int. Conf. on Supercomputing (ICS 03)*. 2003. P. 226-233.
- [9] J. Kubiawicz, D. Bindel, Y. Chen, P. Eaton, D. Geels, R. Gummadi, S. Rhea, H. Weatherspoon, W. Weimer, C. Wells, B. Zhao. *Oceanstore: An Architecture for Global-scale Persistent Storage*. *Proceedings of CAM ASPLOS*, November 2000.
- [10] V. Pande et al., "Atomistic Protein Folding Simulations on the Submillisecond Time Scale Using Worldwide Distributed Computing". *Biopolymers*, Vol. 68. 2003. P. 91-109.
- [11] Predictor@home: <http://predictor.scripps.edu>
- [12] Paul Resnick, Ko Kuwabara, Richard Zeckhauser, Eric Friedman. "Reputation systems", *Communications of the ACM*, v.43 n.12, p.45-48, Dec. 2000
- [13] A. Rowstron and P. Druschel. Storage management and caching in PAST, a large-scale persistent peer-to-peer storage facility. *Symposium on Operating System Principles*. P. 188-201. 2001.
- [14] D. Stainforth et al., "Climateprediction.net: Design Principles for Public-Resource Modeling Research", *Proceedings of the 14th IASTED International Conference on Parallel and Distributed Computing Systems*. (2002).