

# The umbrella project of volunteer distributed computing Optima@home \*

Ilya I. Kurochkin  
IITP RAS, Moscow, Russia  
kurochkin@iitp.ru

## 1 Introduction

As far as distributed computing is concerned, which can be defined as an approach to solve large scale computing problems using personal computers (PC) that organized as a computing system, the most crucial part is related to desktop grids and volunteer computing. The definition of *volunteer distributed computing is computing using voluntarily provided computing resources organized in a desktop grid*.

There are several platforms for distributed computing (including desktop grid): Globus [Foster& Kesselman 1997], HTCondor [Litzkow et al., 1988], Legion, but the most widely used is BOINC [Anderson 2004], [boincstats.com] at this moment. There are about one hundred of public international projects in volunteer computing based on BOINC platform [boincstats.com]. These projects involve about 16 million PCs [boincstats.com]. The BOINC platform has a client-server architecture.

The vast majority of volunteer distributed computing are the scientific projects of the world's leading universities and scientific organizations. The total computing power of volunteer computers exceeds the computing power of modern supercomputers. The peak computing power reaches 150 PetaFLOPS, and the current real power is about 15 PetaFLOPS (~ 10 % of the peak power).

One or several experiments can be performed with the help of one volunteer computing project. Experiments in the project are united by a common theme, in addition, they are conducted by one scientific group.

The exception is so-called, umbrella projects. Experiments that are not related to each other by a common theme and experimental purposes can be performed in the umbrella project. Besides, each experiment is carried out by an independent scientific group. The beginning and the end of the experiments occur asynchronously. It is assumed that the umbrella project works at any time and the probability of downtime is much lower. When organizing umbrella projects, efforts can be saved to deploy and support a separate project. As a result, there is an opportunity to conduct experiments with low computational complexity.

Examples of existing international umbrella projects are World Community Grid and PrimeGrid. They are the largest projects in terms of the number of active computing nodes at the moment [boincstats.com].

## 2 How does the BOINC project work?

The majority of voluntary computing projects are organized to solve one scientific problem. A series of numerical experiments is required to solve the problem. As a rule, the problem can be divided into a set of independent tasks [Benoit, et al., 2010]. Each task is calculated on a separate computing node of the distributed system. Most projects have only one application for the conduction of one or a series of computational experiments. Different sets of input data are used for different tasks. This type of problems is called "bag of tasks" in the literature [Choi S. J. et al., 2007] or a problem shared by data. Examples of such problems are problems of combinatorics

---

This work was funded by Russian Science Foundation (project 16-11-10352)

Copyright © by the paper's authors. Copying permitted for private and academic purposes.

In: E. Ivashko, A. Rumyantsev (eds.): Proceedings of the Third International Conference BOINC:FAST 2017, Petrozavodsk, Russia, August 28 - September 01, 2017, published at <http://ceur-ws.org>

[Vatutin et al., 2017] and exhaustive search, SAT approach [Zaikin et al., 2015], simulation modeling problems [Ivashko & Golovin 2015], some of machine learning problems and others.

The BOINC server-side (BOINC server) must be deployed in order to start the project on the BOINC platform. The generated input data for tasks calculation and a computational application are stored on BOINC server. The results of calculations are also stored on the BOINC project server. To connect a computing node to the BOINC project, it is necessary to install the client-side of BOINC (the BOINC client) and join the project.

The computational application is uploaded to a desktop grid node once, and input data sets (work units) are uploaded as needed. Due to the unreliability of desktop grid nodes, each work unit is sent to several users (initial replication copies) to increase the probability of obtaining a work unit on time (before the deadline) or earlier. The deadline parameter equals 10 days and significantly exceeds the work unit execution time (from several minutes to several hours) for most BOINC projects. The results are sent back to the server after the work unit is completed. To check the results, several results should be obtained by default for comparison (quorum). The results of the execution of one work unit obtained at different computational nodes are compared bit-by-bit and the results are recognized as valid in the case of coincidence. Users receive credits for calculating the work unit results that are recognized as correct. Credits are granted in the project according to a single principle and this is called system for granting credits. The speed of the work unit calculation (from the moment the work unit is received to the moment the result is returned), the computational complexity of the work unit, the power of the computer, etc. can be taken into account when granting credits.

### 3 Tweaking the BOINC project

The speed and probability of getting the correct result grows as a result of increasing the number of copies of one work unit. However, the computational capacity of the desktop grid is significantly reduced in total.

The use of replication in the desktop grid follows from the main features of such distributed systems:

- Heterogeneity of nodes in a distributed system, and, consequently, different calculation speed;
- Autonomy of calculations on various nodes, and, consequently, impossibility of constant coordination of calculations among nodes;
- Unreliability of connections and possible disconnection of computing nodes;
- Unstable time of continuous operation of the node and the difficulty of calculating long work units;
- The presence of errors and delays in calculations.

There are several ways to increase the computing power of the entire distributed system:

- Tweaking the replication settings (redundancy setting);
- Presence of checkpoints (saving of intermediate results);
- Ensuring the integrity of input data, checkpoints and final results (possibility of check results for one copy);
- Configuring the load balancing system (increasing the likelihood of obtaining results on time).

For example, to evaluate the correctness of the result for one instance, it is necessary to implement both the results check (on the server side) and the integrity of input data and results on the client-side (computing node).

In addition to the technical ways to increase the computational capacity of a desktop grid, the number of nodes can be increased in the grid system. Since voluntary distributed computing is considered, it is possible to increase the number of nodes only by recruiting new volunteer (called crunchers in slang) and retaining existing ones. Since voluntary computing is considered, it is possible to increase the number of nodes only by recruiting new crunchers and retaining existing ones. In order to do this, it is necessary to understand what crunchers want (the usage of sociological methods), make the project understandable to a wide range of crunchers (compiling a popular scientific description of the project and ongoing experiments), constantly attract crunchers' attention with new information on the project's website and in other ways [Yakimets & Kurochkin (2015)].

In spite of the seeming simplicity of deploying a new project, quite large amount of works must be done. Below is an approximate list of works with a division into 4 sections.

Table 1: List of works for a BOINC project

Section	#	Name of work	Project stage
Technical costs	1.1	Creating the computational application with the preservation of intermediate results	Deployment
	1.2	Creating an input task generator	Deployment
	1.3	Creating a Validator and a Results Aggregator	Deployment
	1.4	Tweaking the parameters of the server part of the BOINC project	Deployment
Organizational costs	2.1	Domain registration	Deployment
	2.2	The information website of the project	Deployment
	2.3	General description of scientific and administrative group of a project	Deployment
	2.4	Organization of competitions in the project	Support
Interaction with the community of crunchers	3.1	Popular scientific description of the scientific component of the project	Deployment
	3.2	Regular popular scientific description of the numerical experiments	Support
	3.3	Regular publication of results on the project website	Support
	3.4	Communicating with crunchers on forums	Support
	3.5	Maintaining the project blog and publishing links to scientific articles	Support
Additional costs	4.1	The development and implementation of the system for granting credits	Deployment
	4.2	Development and setup of the system for issuing virtual prizes	Deployment
	4.3	Design of the website and the information website of the project	Deployment
	4.4	Personalization of the intermediate results obtained	Support
	4.5	Visualization of the results obtained	Deployment

It is common for situations to arise when a research team successfully deploys a BOINC project, but only performs the technical part from the above list of works. At the same time, support is not placed high emphasis or its value is underestimated. As a result, a stage of stagnation comes after an active start when there are no changes in the project, including the project website, the current errors are not corrected, new results are not published. This leads to a natural loss of interest to the project, an outflow of crunchers and, as a result, a significant reduction in the computing power of the project.

There are frequent mistakes in deploying and supporting voluntary computing projects:

- The absence of checkpoint in computational application;
- Lack of feedback from the administration of a project;
- The lack of a popular scientific description of a numerical experiment on a project website;
- Prolonged absence of new assignments and project timeout;
- The deadlock of cruncher computers while running a computational application;
- The occurrence of errors in the calculations for the most part;
- Long time of tasks calculation (more than one day);
- Absence of estimation of working time for each work unit.

Drawing up a list of works and estimating costs when deploying and maintaining a voluntary computing project is the key to its successful functioning and increasing computing power.

## 4 Umbrella BOINC-project

Umbrella BOINC project is a project that have several independent computational applications. There is a functionality in the client-side that allow user choosing computational applications for launch Figure 1 shows. An example of the umbrella project is the World Community Grid, which is supported by IBM and supports medical experiments.

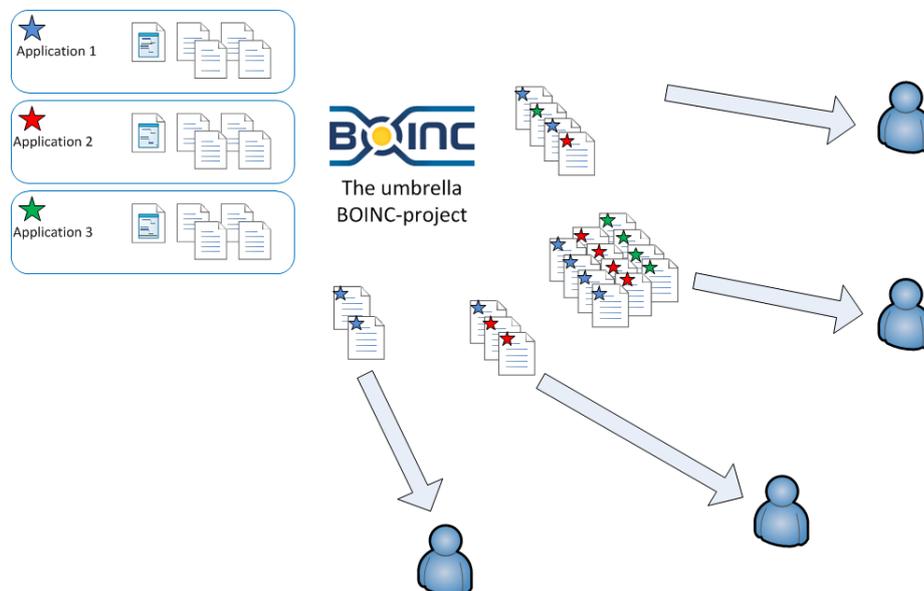


Figure 1: Scheme of tasks distribution in an umbrella project

Using the umbrella project for voluntary distributed computing can significantly reduce the cost of organizing and maintaining the project. In fact, it is only necessary to refine the computational application and make a brief description of the conducted experiment. According to the list of works, it is necessary to implement only the technical part (1) and the popular scientific description of the experiment (3.2). The rest of the work is carried out by the organizers of the umbrella project, as they have more experience in supporting projects of voluntary computing.

There is an opportunity of conducting numerical experiments of different duration: long (more than 6 months), medium (1-6 months) and short (less than 1 month) and also the experiments of various scientific groups. The needs of several scientific groups in the computing resources will exceed the needs of one scientific group. Consequently, the umbrella project will constantly contain tasks for calculation in the interests of one or more experiments. At the same time, computational applications can use different resources for calculation (CPU, GPU, Intel Xeon Phi).

The audience of the already functioning umbrella project will be several times larger than the audience of a separate distributed computing project even after the initial stage, when the number of crunchers in the project is small.

The usage of the umbrella project allows you to conduct small experiments with minimal costs on an existing project with a large computing power.

The development of the concept of voluntary computing and the increasing popularity of BOINC platform led to the fact that many new voluntary computing projects appeared in Russia in 2016-2017 (Uspex@home, AndersonAttack@home, XASONS forCOD, ODLK@home others).

Some of them carry out small computational experiments, and activity in the project alternates with a period of absence of assignments. The usage of the umbrella project for small experiments has significantly allowed to reduce costs of organizing calculations.

It should be noted that the requirements for the umbrella project itself are increasing with a significant reduction in the cost of launch of a new computing experiment.

The umbrella project is characterized by the following features:

- Big number of users;
- Diversity of computing tasks;
- Large amount of input data and results;
- Increased requirements for reliability of the server-side of the BOINC project;
- Scaling simplicity in the case of significant increase of number and/or activity of users;
- Increased requirements for interactions with users, both on the project website and on other cruncher community communication platforms.

## 5 Fault-tolerance and scalability of the BOINC project

The problem of fault tolerance and scalability of the server part of projects is one of the main ones in the development of the BOINC platform [Kurochkin & Saevskiy (2016)]. Failure of the project is necessary due to a large number of technical failures. Due to the large average load on the umbrella project, the requirements for the hardware and software components of the server part of the BOINC project are growing.

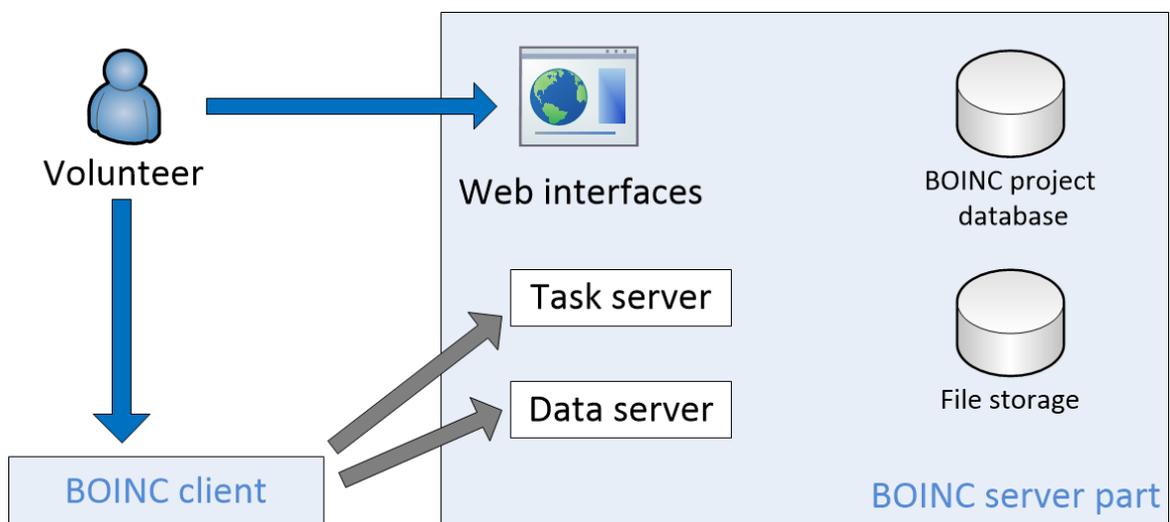


Figure 2: The scheme of the server part of the BOINC project

For the normal operation of the umbrella project and enhancement of fault tolerance, the following steps can be taken:

- Using distributed file systems;
- Database replication;
- Implementing a load balancing system;
- Implementation of monitoring of the components of the server part of the BOINC project.

In addition to the described steps, one should not forget about such measures as

- Regular backup of the project database;
- Uploading results to cloud or external file storages;
- Regular restart of all services (daemons) of BOINC server part.

Distributed file systems can be used as a solution to ensure uninterrupted access to files. Distributed file systems can provide the necessary performance and reliability for the BOINC project. But it should be noted that it will be necessary to optimize the use of the network by these systems. Since in the case of a large stream of data being downloaded to the server and using the replication function at the file system level, rather than at the hardware level, the telecommunications network becomes a bottleneck.

Load balancing is the distribution of the load on different components of the BOINC project. DNS servers by default can distribute the load evenly. But not all requests are uniform, and one can take 1 second and another 200 milliseconds. This can cause a different load on the server, despite the same number of requests received. The balancer can distribute requests to the server using various algorithms (eg round-robin, random, etc.).

Database replication this solution will increase the bandwidth of the database for reading, a duplicating element will appear, which eliminates a single point of failure of the system. Separating the database into 2 databases - analytical (large) and operational (small) will entail certain difficulties. As the modifications in the code significantly increase the costs and complicate the further support of the BOINC project when upgrading the platform from the official repository. Another approach is to use several replicas of the database to use a replica in the event of a failure of the main database.

As a result, the following scheme of the server part of the umbrella project is suggested, taking into account the requirements for fault tolerance and scalability.

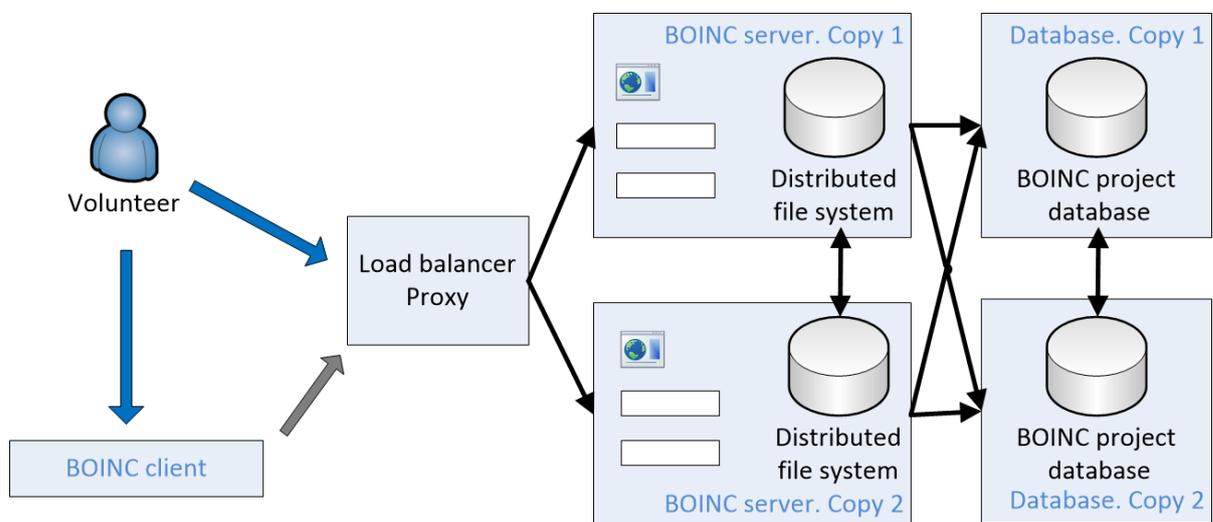


Figure 3: The proposed scheme of the umbrella project

## 6 Examples of possible usage of a umbrella project in distributed computing

Private umbrella project in distributed computing. Instrument for conducting laboratory works. Within the special course for students of technical specialties.

Private umbrella project in distributed computing with possible participation of big number of users. Popularization of distributed computing technologies among high school students (9-11). BOINC, as an example of actual distributed computing technologies. Within an optional course for calculations, studying subject areas of projects, for creation of technical applications, as an information exchange platform, as a tool for conducting competitions.

Private umbrella project in distributed computing SandBox for testing applications.

Open umbrella project of voluntary calculations for conducting series of experiments on various subjects (large and small experiments). As an area of interest for new small experiments and arranging the creation of new projects in distributed computing.

## 7 Conclusions

Drawing up a list of works and estimating costs when deploying and maintaining the volunteer computing project is a key of its successful operation and increase in computing power.

The use of an umbrella project allows the scientific team conducting of small experiments with minimal costs on an existing project with a large computing power.

Attraction and training of specialists in distributed computing through the "school-university-scientific organization" is the basis for the development of high-performance computing in the future.

The use of an umbrella distributed computing project on BOINC platform allows a high school student and a student to get knowledge in practice with the current technologies in distributed computing, to study various aspects of Desktop grid functioning, to learn about the features of applications development for grid systems and to gain experience in conduction of a computational experiment on a real distributed system.

## References

- [Foster& Kesselman 1997] I Foster, C Kesselman *Globus: A metacomputing infrastructure toolkit*, International Journal of High Performance Computing Applications 11 (2), 1997, pp.115-128.
- [Litzkow at al., 1988] M.J. Litzkow, M. Livny, M.W. Mutka (1988) *Condor-a hunter of idle workstations*, Distributed Computing Systems, IEEE.
- [Anderson 2004] Anderson D. (2004). BOINC: a system for public-resource computing and storage. Grid Computing, IEEE.
- [boincstats.com ] The server of statistics of voluntary distributed computing projects on the BOINC platform. <http://boincstats.com>.
- [Vatutin at al., 2017] Vatutin E.I., Zaikin O.S., Zhuravlev A.D., Manzyuk M.O., Kochemazov S.E., Titov V.S. *Using grid systems for enumerating combinatorial objects on example of diagonal Latin squares* // CEUR Workshop proceedings. Selected Papers of the 7th International Conference Distributed Computing and Grid-technologies in Science and Education. 2017. Vol. 1787. pp. 486490. urn:nbn:de:0074-1787-5.
- [Benoit, et al., 2010] Benoit, et al., *Scheduling Concurrent Bag-of-Tasks Applications on Heterogeneous Platforms*, IEEE Trans. Computers, vol. 59, no. 2, pp. 202-217, Feb. 2010.
- [Choi S. J. et al., 2007] Choi S. J. et al., *Characterizing and classifying desktop grid* //Cluster Computing and the Grid, 2007. CCGRID 2007. Seventh IEEE International Symposium on. – IEEE, 2007. – . 743-748.
- [Zaikin at al., 2015] Oleg Zaikin, Alexander Semenov and Ilya Otpuschennikov. *Solving Weakened Cryptanalysis Problems for the Bivium Keystream Generator in the Volunteer Computing Project SAT@home* // Proceedings of the Second International Conference BOINC-based High Performance Computing: Fundamental Research and Development (BOINC:FAST 2015), Petrozavodsk, Russia, September 14-18, 2015. pp. 22-30.
- [Ivashko & Golovin 2015] Ivashko E., Golovin A. *Partition Algorithm for Association Rules Mining in BOINC-based Enterprise Desktop Grid. Lecture Notes in Computer Science*. Parallel Computing Technologies 13th International Conference, 2015, 268272, Springer.
- [Yakimets & Kurochkin (2015)] Yakimets V.N., Kurochkin I.I. (2015). *The voluntary distributed calculations in Russia: the sociological analysis* In the collection: INFORMATION SOCIETY: EDUCATION, SCIENCE, CULTURE AND TECHNOLOGIES of the FUTURE Works XVIII of the joint conference "Internet and Modern Society" (IMS-2015)., St. Petersburg: ITMO university, pp. 345-352.
- [Afanasiev at al., 2015] Afanasiev, A. P., Bychkov, I. V., Manzyuk, M. O., Posypkin, M. A., Semenov, A. A., Zaikin, O. S. (2015). *Technology for integrating idle computing cluster resources into volunteer computing projects*. In Proc. of The 5th International Workshop on Computer Science and Engineering, Moscow, Russia (pp. 109-114).

[Kurochkin & Saevskiy (2016)] Kurochkin Ilya, Saevskiy Anatoliy *BOINC forks, issues and directions of development* // 5th International Young Scientists Conference in High Performance Computing and Simulation (YSC 2016), 26-28 October 2016, Krakow, Poland, Procedia Computer Science Volume 101, 2016, ISSN 1877-0509, pages 369378 (<http://dx.doi.org/10.1016/j.procs.2016.11.043>).