

# Experiences with Virtual Users' Account System in Polish National Cluster

M.Kupczyk, M.Lawenda, N.Meyer, M.Stroinski, P.Wolniewicz

Poznań Supercomputing and Networking Center (PSNC)

ul. Noskowskiego 10, 61-704 Poznań, Poland

e-mail: { miron, lawenda, meyer, stroins, pawelw }@man.poznan.pl

**ABSTRACT:** In this paper we present the configuration of the Polish national cluster using the Virtual Users' Account System. In order to give users a possibility to access computer resources we configured a set of LSF queues that send jobs to other sites. To avoid the problem with managing users' accounts on remote supercomputers we are using special middleware that allows to run and manage users' jobs without the necessity of creating accounts for all users.

## 1. Introduction

Most often the computing power given to the users in one computing center is not enough for them. Therefore there is a need to connect the supercomputers into grid environments [1,2,3]. Also in Poland most of the Polish supercomputing centers agreed to form a national computing grid. Because of the optical technologies used to build the Polish Optical Network [4], it is possible to form an efficient grid, because the time of access to remote machines is comparable with the time of accessing the local machines.

The simplest solution to create the cluster is the usage of a job processing system that manages jobs. In the first phase we just used the Load Sharing Facility (LSF) queuing system and configured several queues that could send jobs to machines in other centers. But it soon appeared that the good solution for local site conditions - for systems installed in one place, where there is a constant or slow change in the number of users and computing systems does not work well in a distributed environment. While connecting distributed, geographically distant systems, belonging to different institutions, we were faced with the problem of managing users' accounts of the whole nation-wide structure. The problems are due to policies of user management, which are

different in each centre and due to the problems with maintaining users' account coherency on all machines. Therefore we employed the Virtual Users' Account System [5], developed in Poznan Supercomputing and Networking Center, that simplifies users' accounts management in the distributed environment.

In this article we describe the implementation of Polish grid environment based on our Virtual Users' Account System. We present benefits of it and the problems encountered during configuration. We also briefly describe the accounting mechanism we have implemented.

## 2. Virtual Users' Account System

The Virtual Users' Account System (VUS) allows running jobs on remote machines without having an account on this machine. This allows reducing administration overhead connected with maintaining users' accounts for the whole cluster. Instead of using real users' account, a set of generic account is used. This system can co-operate with any job distribution system. In connection with job distribution systems, it allows a better usage of computing resources by sending jobs to currently less utilized machine.

The Virtual Users' Account System has the following benefits in comparison with a typical job distribution system:

1. **Simplified users' accounts administration.** The user does not have to have accounts on remote systems to use them. The pool of virtual users' account is used.
2. **Full accounting information.** Information regarding used CPU time is collected after completing the task and sent to the virtual users' server. Thus, for every user, full information is stored about the used CPU on separated systems and globally.
3. **Automatic file transfer.** There is no need to share files via NFS on all systems connected into a cluster. Files indicated by the user are transferred automatically.
4. **Co-operation with other queuing systems.** The system can be configured in such a way that the jobs can be sent also to other queuing systems.
5. **Easy authorization.** There is no need to define an authorized set of users to a queue in each center. Access to the queue can be granted by adding a user to the queue access list in the central server

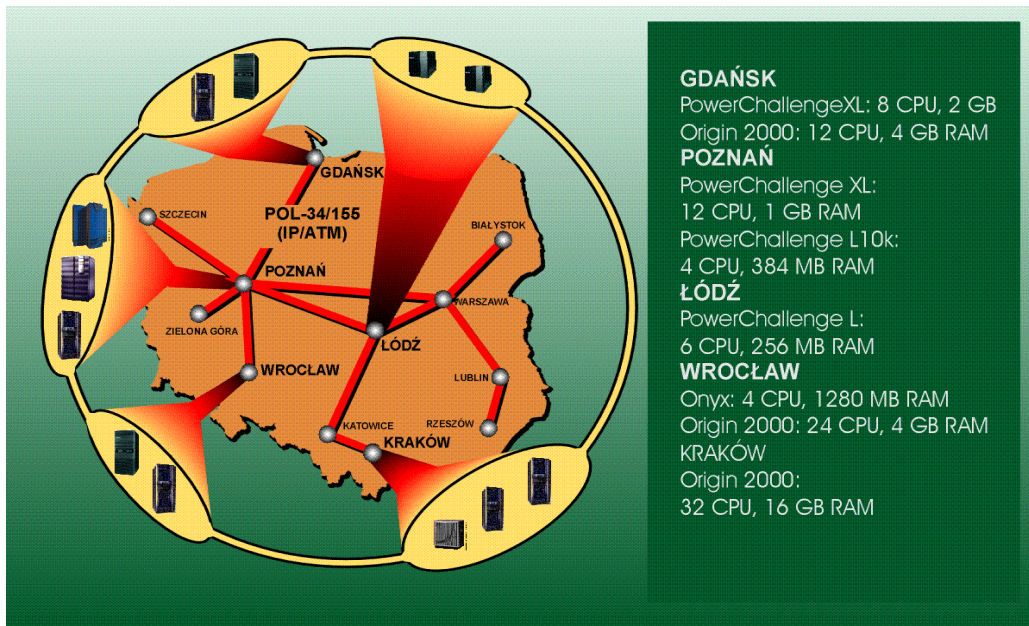


Fig. 1. National Computing Grid

### 3. Polish National Cluster

In Poland there are over a dozen supercomputers located in several supercomputing centers. Most of them are SGI systems (Challenge, Power Challenge, Origin2000, Onyx2, Origin3000) but there are also IBM's, Cray's, HP's. All Polish supercomputers are used mainly to perform chemical, physical, mathematical and engineering computations. The most popular applications are Gaussian98, Gamess, MSI, Abaqus, Fidap, Matlab, Maple, Amber. A large part of the CPU time is also used by users' own developed applications. We decided that it is sensible to dedicate some machines to run only a specific application [6].

Most Polish supercomputing centers agreed to form the Polish national supercomputer cluster. It includes Gdańsk, Poznań, Wrocław, Kraków and Łódź. Additionally, there are

some more Polish centers, which are going to join the national cluster as clients and to have the possibility to access the computing resources from the other centers. The current stage of this cluster construction is built on LSF queuing system. The supercomputing centers and the system dedicated to grid environment are shown in Fig. 1 and listed in Table 1.

Currently in Poland the informational infrastructure program is being started. It is called PIONIER: the Polish Optical Internet – Advanced Application, Services and Technologies for the Informational Society. During the realization of this program there are plans to create a HPC infrastructure. The Polish Grid will consist of the existing and new SGI computers and will run basing on the future Polish Optical Internet. High Performance Computing (HPC) and High Performance Visualization (HPV) services will be used to build the virtual laboratory and tele-immersion applications in 2003 and 2004.

The planned SGI HPC cluster will consist of new SGI systems with Itanium processors. The HPV cluster will be formed of the remaining MIPS systems because MIPS provide better support for the graphics subsystem than Itanium.

Table 1. Supercomputing centers involved in the national cluster configuration

Center	Hardware Platform
Gdańsk	SGI Power Challenge XL, 8CPU, 2GB Origin 2000, 24 CPU, 16 GB RAM
Łódź	SGI Power Challenge L, 6 CPU, 256 MB RAM
Poznań	SGI Power Challenge XL, 12 CPU r8k, 1GB RAM Origin3200, 32 CPU r12k, 8 GB RAM Cray SV1, 8 CPU, 16 GB RAM Cray J90, 16 CPU, 4 GB RAM
Wrocław	Onyx, 4 CPU r10k, 1280MB Origin 2000, 32 CPU r10k, 8 GB RAM
Kraków	Origin 2000, 128 CPU r10k, 16 GB RAM
Szczecin	SGI Power Challenge XL, 4 CPU r8k, 512 MB RAM
Gliwice	Sun Enterprise 6000, 12 CPU, 6GB RAM

#### 4. Implementation and tests

The queues configuration was deliberately designed to let the users take advantage of the interactive applications, especially in the graphical environment. Apart from many queues dedicated to the tasks defined by the users, separate queues for the third party applications were created, e.g. for Gaussian98 or MATLAB. The same solution is planned for MSI and ABAQUS as well.

The application queues allow running only specific applications. It is realized by means of so called job starters. LSF allows defining a program, which will be run before the user script is processed. For example, as a job starter for the Gaussian queue, g98 is defined. Users cannot submit jobs to a queuing system directly. It is realized by a set of scripts but users still can specify parameters and queuing system limits.

Not all supercomputers in the Polish centers have installed LSF. But our system allows sending jobs to the systems with other queuing systems. We have configured queues that send jobs to NQE queuing system on Cray in PSNC. LSF itself has modules responsible for communication with NQE, but it does not allow transferring files. We also tested our system in communication between LSF and LoadLeveler installed on IBM SP2 and our tests were successful. But because all Polish SP2 are a little bit outdated, we did not decide to include them into the cluster until their upgrade.

We have tested our grid configuration with a set of tests. In each test we submitted several thousands of jobs with frequency one job per minute. All jobs were successfully run, results were properly returned and full accounting information is stored. We assume that the typical utilization of the system will be much lighter. Because most of the jobs run in our systems are rather long-term jobs the average number of jobs on host per month is 100-200. We have also checked the CPU overhead introduced by our system and the results show that the average time used by all daemons per job is about 0.2s. Then the VUS overhead is a few minutes per month. It is not much comparing it to more than 50 hours of system CPU time per month and we do not care about this small loss.

We are aware that the Virtual Users' Account Server can be a bottleneck because it is responsible for all database operation. Currently we have only one such daemon that communicates with our Oracle database. However we have not observed any problem with the response time even for very frequent job submission. But it is possible to duplicate this daemon, e.g. so that there is one copy in each center and all machines communicate with the nearest daemon.

We also implemented the Virtual Users System Monitor that shows the current state of the system. The example view of this application is shown in the Fig. 2. The application shows the current and maximum possible number of virtual jobs on each machine or in each center. The list of current virtual jobs can also be shown.

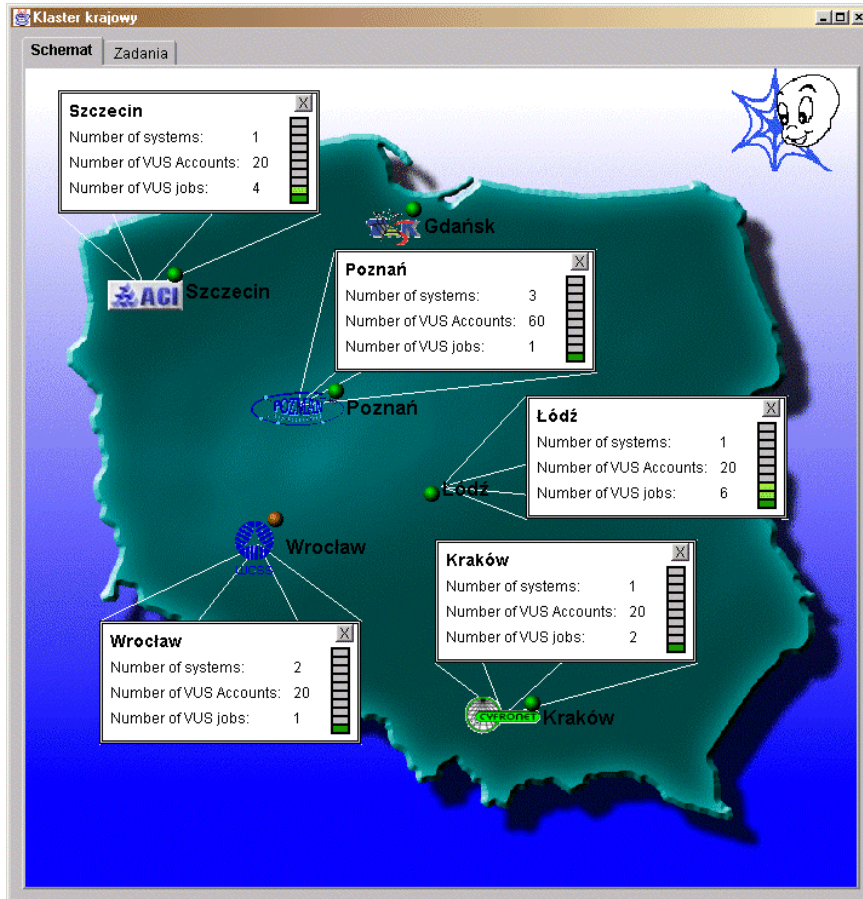


Fig. 2. Virtual User' Account System monitoring application

## 5. Encountered problems

During the creation and operation of our national grid we encountered several following problems:

### 5.1. Setup problems

Initially, most of the supercomputer centers used LSF, but in different version. Additionally some machines used LSF, NQE or PBS. We have chosen the LSF as our main job distribution system, but then we have to upgrade LSF to the newest versions on some machines because of the version incompatibility.

The Installation of the Virtual Users' Account System allowed using different queuing system but still there are some problems with such co-operation. Our system can send jobs to the remote queuing system, but it cannot read the current utilization of the resources. Thus the scheduling can only base on a number of jobs in queues and queues parameters. This can cause the system overloading especially when the remote system run other jobs submitted locally to its queuing system. In our Polish grid machines will be dedicated to run specific kind of jobs and applications and therefore our solution is sufficient to our needs. For example we dedicated our two vector Cray machines to run only Gaussian 98 jobs and users' own vectorized job. Access to these machines will only be through the LSF that submits jobs to the

Cray NQE queuing system. Local load balancing is done by NQE.

During the installation of the Virtual Users' Account System we also encountered configuration problems because of the local security policies of the centers. Most often there were problems with disabled IP ports and the configuration of ATM PVC connections. Some of the centers also allowed the connection only from the specific domains. These are all technical problems but gave us a lot of pains and required a lot of e-mails and phone calls.

## 5.2. Consistency problems

The worst problems are with the LSF configuration consistency. Adding a new LSF cluster into an LSF multicluster requires changes in the configurations files for all clusters. There was a situation that even in a local LSF cluster jobs could not be scheduled because of a simple mistake in a configuration in the remote LSF cluster.

Also the way of invoking application should be the same. Initially applications were installed in different directories on all machines. Also the way of running this application was different. On some systems users were responsible for setting necessary environment variables and on other systems everything was done automatically. We have to unify the way of application invoking. For all main applications there are scripts with the same name that setup all necessary variables and run applications.

## 5.3. The problems with users

The most important problem is that some of our users do not want to run jobs on remote machines. Some of them claim that they run a series of computation experiments and want to be sure that all results will be achieved in exactly the same condition. Especially the users that run the application based on random number generators like Monte Carlo algorithms, insisted on using the same machine with the same random number generator.

The other problem with users is that they do not like to learn how to use a grid environment and how to submit jobs. Especially if they have to specify resource requirements and which files should be transferred to and from the remote machines they prefer to use only local resources. We can help them with automatization of the submission process. For the most popular application we defined a list of files that should be transferred. The submission program in our Virtual Users' Account System recognizes the additional parameter – a kind of application and then automatically transfers these defined files. For example, for Gaussian 98 application the

system has to transfer .inp file to remote system .log file from the remote system and .chk file in both directions. It is possible to simplify the submission even more and write scripts that automatically submit the job to the system in a proper way. For the consistency reason the file list is not stored locally but in the VUS database.

## 5.4. Security

Communication between the machines from different centers is most often conducted on a public network where data is not safe. Users do not quite like to send their jobs and result unencrypted, therefore a system that transmits the data remotely must encrypt all data. We are using SSL for all communications. SSL certificates are also used to ensure that no one wants to intrude upon the cluster.

Another problem is with the security of the tool used for the creation of the grid. We can be responsible for security holes in our Virtual Users' Account System, but we cannot help the security holes in a job distribution system and in operating systems. Currently we know at least three security holes in LSF that can result in acquiring root privileges.

## 5.5. Binary compatibility

Our cluster is heterogeneous and this causes problems with binaries compatibility. Application queues require only submitting data files so it is not important on which architecture jobs run. But when the users submit jobs to general queues they have to supply binaries, which are transferred to the destination system. Users can define an architecture or the operating system as resource requirements. By default the jobs start on systems with the same architecture as the submitting host.

## 6. Accounting

When the job can be run on remote machines and use generic accounts one, of the most important thing is to have full accounting information. It must include full information about the utilization of resources by the real users. It is also important to have full information about jobs and its real owner in every moment of the job execution. This requires that the system be reliable and able to keep this information even in case of the system or network crashes. Let us imagine the problem caused by the loss of mapping information about scheduled or run job. The job then is some kind of a zombie job because there is no information where to return the results and who should pay for the utilization of resources.

To keep the information consistency we have a specially designed daemon that is run on every host and takes care

about sending information. All the important information is temporally stored and in case of any problems it can be retransmitted.

Accounting is made in two steps. First, the Virtual Users Manager records the start and end time of job execution and on which virtual account this job was running. The second step is done during the standard system accounting activity started by cron. This standard accounting runs our virtual accounting script that sends accounting information for each virtual users job to Virtual Users' Account Server. Because information about the real user and the start time of the job was stored earlier it is possible to unambiguously recognize job which account is concerned. Then the Virtual Users' Database is updated with all necessary accounting information. The global accounting information is stored for user@domain instead of users from each single domain host separately.

## 7. Conclusions

In this paper we outlined the current configuration of the Polish National Cluster that consists of several supercomputers from different centers connected by a fast network. We use our Virtual Users' Account System that saves our administrators much work connected with the users' accounts management. Because of that system we can dynamically change the configuration of the cluster, add new machines, change destination of the queues transparently for users. The users do not have to care about applying for an account on all Polish supercomputers. They just run their jobs and receive data and the system takes care of transmitting and running the job. And now jobs from the overloaded sites are sent to currently less utilized sites that, when improved, mean the utilization of all supercomputers. Although users can now run a job on different systems we still have full account information. But still there are problems with the grid operation. Some of them, like consistency and security problems, are caused by our job distribution system – LSF. There is also a need for teaching users how to use grid environment.

## References

- [1] Foster, I., Kesselman, C (eds.) *The GRID – Blueprint for a New Computing Infrastructure*. Morgan Kaufman, San Francisco, 1999
- [2] Grid Forum, <http://www.gridforum.org>
- [3] The Data Grid Project, <http://www.cern.ch/grid>
- [4] J. Rychlewski, J. Węglarz, S. Starzak, M. Stroiński, *PIONIER – Polish Optical Internet*,. conference material ISTHMUS 2000, Poznan, 2000

- [5] W. Dymaczewski, N. Meyer, M. Stroiński, P. Wolniewicz, *Virtual User Account System for distributed batch processing*, HPC 99, Amsterdam, April 1999
- [6] M. Kupczyk, N. Meyer, M. Stroiński, P. Wolniewicz, *Application Servers in Polish Computing Grid*, SGI Users 2000, October 2000