# Parallelizing OMNeT++ Simulations using Xgrid

Robin Seggelmann
Münster University of Applied
Sciences
Fachbereich Elektrotechnik
und Informatik
Stegerwaldstrasse 39
D-48565 Steinfurt, Germany
seggelmann@fh-
muenster.de

Irene Rüngeler
Münster University of Applied
Sciences
Fachbereich Elektrotechnik
und Informatik
Stegerwaldstrasse 39
D-48565 Steinfurt, Germany
i.ruengeler@fh-
muenster.de

Michael Tüxen
Münster University of Applied
Sciences
Fachbereich Elektrotechnik
und Informatik
Stegerwaldstrasse 39
D-48565 Steinfurt, Germany
tuexen@fh-muenster.de

Erwin P. Rathgeb
University of Duisburg-Essen
Institute for Experimental
Mathematics
Ellernstrasse 29
D-45326 Essen, Germany
erwin.rathgeb@iem.uni-
due.de

## ABSTRACT

Working with simulations, testing and validating theories often requires a large number of simulation runs. The discrete event simulation environment OMNeT++ already provides functionality for distributed computing, yet the simulated model and modules to be parallelized have to be declared manually. Apple Mac OS X comes with Xgrid support, which allows easily setting up an ad hoc grid for parallel computing. In this paper we will describe the features we had to add to OMNeT++ to be able to use Xgrid to parallelize even several thousand runs. We will point out how to use Xgrid to distribute runs of a simulation not only to multiple CPU cores but also to multiple machines without modifying the simulation itself. Our analysis will reveal that Xgrid allows to reduce the computing time almost proportional to the added parallel computing power.

## Categories and Subject Descriptors

I.6.8 [**Simulation and Modeling**]: Types of Simulation—*Distributed*; I.6.6 [**Simulation and Modeling**]: Simulation Output Analysis; C.2.4 [**Computer-Communication Networks**]: Distributed Systems—*Distributed applications*

## General Terms

Distributed simulation, Grid computing

## Keywords

OMNeT++, Xgrid, Cluster

## 1. INTRODUCTION

Since we are working on the further development of the Stream Control Transmission Protocol (SCTP) [14], we have integrated SCTP in the INET framework [13], [4] of the OMNeT++ simulation environment to evaluate new features. Simulations with several thousand runs, depending on the convergence, have often been necessary, that can require several days even on a high-end workstation. As today's CPUs gain their performance mainly by more and more cores rather than by raising the frequency, the maximum performance can only be obtained by parallelizing the work load. The main issue is that OMNeT++ per default only uses one process per simulation, which results in fully loading only one core of the CPU. A solution would be to parallelize the runs of a simulation without changing the simulation itself.

The easiest way to parallelize the runs is to start them manually one by one, the operating system should then assign the processes to different cores. This is not very convenient because the simulation has to be supervised at all times to start new processes as soon as computing power has become available. A better solution would be to use Xgrid, which can distribute a batch of processes not only to multiple cores but also to multiple machines. Hence, every run can be listed with adequate parameters in a batch and assigned to any idle core of any machine within the grid. This solution is generic and works with every simulation model, even large frameworks like INET.

In this paper, a parallel computation of OMNeT++ simulations using Xgrid will be introduced. In Section 2 an overview of Xgrid, its structure and features is given. Alternative approaches are pointed out in Section 3. In Section 4 we describe the generation of batch jobs with OMNeT++,

and in Section 5 we explain how Xgrid has to be set up to process these jobs efficiently. The results are presented in Section 6, while limitations of this solution are discussed in Section 7. The conclusion is given in Section 8.

## 2. OVERVIEW OF XGRID

In 2004, Apple introduced Xgrid [10] to set up ad hoc grids for parallel computing with as little administration costs as possible. It consists of three roles: controller, client and agent. The controller manages the jobs and the distribution of tasks to the available machines, called agents, as shown in Figure 1. A job consists of an arbitrary number of tasks, each described by a command and associated parameters. Any host running Mac OS X 10.4 Tiger or later can provide an Xgrid controller, agent, or be a client which submits jobs to and retrieves results from the controller. There is also an agent for Mac OS X 10.3 Panther [9] and an open source project of a Java-based agent [8] available. An overview is given in Table 1.

| | Authentication | | |
| | none | Password | Kerberos |
|---|---|---|---|
| **Agent** | | | |
| Mac OS X 10.3 - 10.5 | + | + | + |
| Java | + | − | − |
| | | | |
| **Controller** | | | |
| Mac OS X 10.3 | − | − | − |
| Mac OS X 10.4 | + | + | − |
| Mac OS X 10.4 Server | + | + | + |
| Mac OS X 10.5 | − | − | − |
| Mac OS X 10.5 Server | + | + | + |
| Java | − | − | − |

Table 1: Availability of Xgrid

The minimal setup is a single machine combining all the three roles, but a cluster of up to 128 agents connected over a local network or the internet is also supported [7]. The agent can be dedicated to the grid but this is not mandatory, it can also be configured to only accept tasks when it is idle, that is no user input has occurred for some time. This allows to use user workstations to be part of the cluster when not in use, e.g. over night or on weekends. Agents accept as many tasks as they have CPU cores at a time, so assuming it is single-threaded, each task gets its own core.

With Mac OS X 10.5 Leopard extended security features have been added, which require the controller to be run on the server edition of Mac OS X. The authentication ranges from none over simple passwords to Kerberos with Open Directory (see Table 1). The latter allows tasks to be run on agents with the permission of a directory user, otherwise each task will be sandboxed and will have no permission to access anything beyond its working directory.

Clients can submit tasks individually or within a batch. The files needed for each run can be included in a batch file, which will be transferred to the agents before a task is executed. Including files has the advantage that neither a manual distribution of files to every agent nor a network share is necessary, and that all data can be accessed within the sandbox.

Submitted jobs can be watched and managed in the Xgrid admin tool, shown in Figure 2, while the submission itself
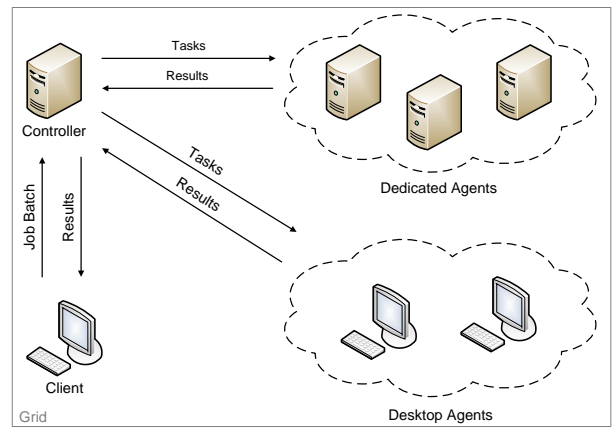


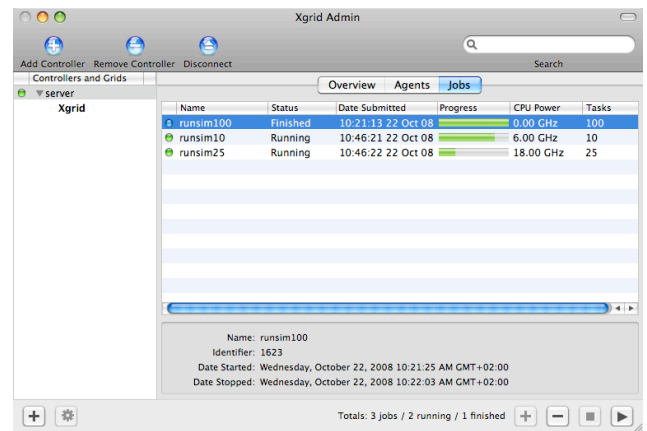Figure 1: The three roles of Xgrid: controller, client and agent



Figure 2: Xgrid administration tool

and the retrieval of results has to be done with a command line tool. The simplest form of a job is a single command, which can be committed directly to the grid and thus no batch file is needed. The command to print a calendar for March 2009 would be:

```
/usr/bin/cal 3 2009
```

This simple command can be submitted to Xgrid easily with the following:

```
xgrid -h controllerhostname \
-auth Kerberos -job run /usr/bin/cal 3 2009
```

With the keyword run the call will block until any available agent executed the command. Alternatively run can be exchanged with submit. The job will be queued to wait for its turn, and the call returns immediately. The latter always applies to batch jobs. To submit a batch, a plist or XML formatted job specification file has to be generated and committed to Xgrid with the keyword batch. That would be:

```
xgrid -h controllerhostname \
-auth Kerberos -job batch job.xml
```

Further on, Xgrid supports the Message Passing Interface (MPI) for inter-process communication. This requires changes to the simulation, discrete modules have to be declared and the MPI interface has to be implemented. As the solution should be as easy and generic as possible we will not apply these features.

# 3. ALTERNATIVE APPROACHES

There are already some alternative solutions to parallelize OMNeT++ simulations. Although they are suitable for the tasks they were designed for, they have some limitations which led us to Xgrid. In the following, five alternative approaches and their characteristics are described.

## 3.1 Parallel discrete event simulation

OMNeT++ comes with parallel discrete event simulation (PDES) [5], which allows the definition of different modules within a simulation in the configuration. These modules are run as logical processes during a simulation and thus can be parallelized. Communication between processes is realized with MPI. This allows sending messages while direct method calls are not possible unless the respective modules are assigned to the same processor.

The partitioning has some drawbacks, as it limits the simulation design and does not always speed up the computation. If the modules are chosen unfavorably, the simulation may run even slower than without parallelization. Additionally, the setup of an MPI-capable environment has to be done in advance.

## 3.2 SimProcTC

SimProcTC, which is based on the architecture Reliable Server Pooling (RSerPool) [11], can also be used to distribute OMNeT++ simulations to multiple CPUs [12]. An RSerPool architecture consists of an arbitrary number of servers, called pool elements (PE). A group of pool elements is called pool and identified by a pool handle. Every handle is managed by Pool Registrars (PR) to which pool elements register themselves. The registrars keep track of the availability and load of their pool elements, synchronize with other pool registrars and assign available pool elements to pool users (PU) requesting their services.

Similar to the approach with Xgrid, runs of a simulation can be distributed to the pool elements by requesting an available pool element from the registrar for every run. The pool element runs the simulation and sends the results back to the pool user. Unfortunately this solution is not yet integrated on Mac OS X platforms and does not supply any GUI-based tools.

## 3.3 Akaroa

The project Akaroa is an architecture designed for parallel computation of quantitative stochastic simulations [6]. Akaroa consists of two parts, the akslaves which are managed by an akmaster. The akmaster accepts jobs and distributes them to the akslaves. The akslaves report their results to the akmaster which decides if the data already matches a given accuracy. Otherwise the slave has to continue the simulation.

There is an interface to use the Akaroa architecture for OMNeT++ simulations [1]. Every output of the simulation is reported to the akmaster for estimation. When an appropriate accuracy is reached, the simulation will be stopped

gracefully. To use this interface, every simulation has to be modified to include calls to the Akaroa master to judge the current accuracy.

## 3.4 make

The UNIX tool `make` can also be used to parallelize simulations. In the necessary Makefile a target for each run of the simulation has to be created which executes the run. Then a main target is needed which invokes all the previously created run targets. The number of concurrent processes can be specified by calling the `make` command with the parameter `-j`, usually with the number of available CPU cores. It will then execute the given number of targets simultaneously and each process will be assigned to a different core by the operating system. This approach is limited to a single machine, no distribution to multiple hosts can be realized.

## 3.5 Eclipse for OMNeT++ 4.0

In the upcoming 4.0 release of OMNeT++ an IDE based on Eclipse [3] will be included. This environment supports developing and running simulations. It can be configured to start multiple simulation runs at once. This allows to fully load multiple CPU cores of a single machine but like `make` it cannot distribute runs to other machines.

# 4. GENERATING BATCH JOBS WITH OMNET++

Starting an OMNeT++ simulation with just one run in Xgrid can be done without any alterations to OMNeT++. But as simulations normally consist of several runs which correspond to tasks in Xgrid, a batch file is needed to provide the command and the necessary parameters for each task, so that the controller can distribute them. After a job has been started, a job id is returned, that can be used to retrieve the job specification.

```
xgrid -job specification -id n
```

This batch file is in XML format, and therefore its generation can be automated. A simple batch job has the following structure:

```
jobSpecification = {
    applicationIdentifier
        = "com.apple.xgrid.cli";
    inputFiles = {};
    name = "/usr/bin/cal";
    taskSpecifications = {
        0 = {arguments = (3, 2009);
            command = "/usr/bin/cal";
        };
    };
}
```

The section *inputFiles* consists of a list of files that are needed by each agent. For OMNeT++ simulations, this is usually the binary, the `ned`-files, the configuration file, routing files and so on. If the agents do not have a shared medium, all input files have to be written in ASCII coded hexadecimal representation in the XML file. The *taskSpecifications* specify the different runs in OMNeT++. All necessary arguments, e.g. the name of the configuration, the number of the run, and the command are listed. Without

the integration of the generation of the complete file in OMNeT++, the job specification could be written manually by retrieving the specification file for one run and adding the task specifications for all the others. As this is a very tedious work, when it has to be done for hundreds of runs, the aim is to automate the whole process.

To generate the job specification file automatically, two prerequisites have to be fulfilled:

- It has to be figured out, which files have to be included.

- The possible number of runs has to be known.

- Each file has to be converted letter by letter into the ASCII coded hexadecimal representation.

Looking at small examples in OMNeT++, for instance FIFO, the necessary files are all included in the working directory of the example. But with the growing complexity of the frameworks, the input files can be distributed among different directories. To start, for instance, an example of the INET framework [4], about 170 `ned`-files in 50 different directories are loaded in addition to the example related files. To solve this problem, all directories are searched recursively, and the positions of found files are set relative to a base directory.

Figure 3 illustrates the mapping of the original files to the ones in the batch file for the following `-n` switch

```
-n ../..:../../../src
```

starting from the working directory `INET/examples/sctp/fair`.

The most important task is to map the `ned`-files. The two paths, separated by a colon, refer to the two "base directories" `INET/examples` and `INET/src/`. Each base directory contains a file called `package.ned`, that includes the root package, i.e. the name of the package, from which the hierarchy of all the other `ned`-files below this directory stem. If the files in the ellipses on the left hand side were just copied, the second package file would overwrite the first one, which would result in errors when matching the expected to the package names provided in the `ned`-files. Therefore, each base directory has to be given an individual name, from which the relative paths can start. In our case we just name the directories `temp1` and `temp2`. Thus the hierarchy of the files is kept and can be copied to the agents.

Hence, an entry for one input file has the following layout:

```
"temp1/transport/sctp/SCTP.ned"={
    fileData = <7061636b 61676520 696e6574 ....>;
}
```

In addition to the `ned`-files, there are example dependent files, that are needed to run the simulation. In INET, these could be routing files or files to be interpreted by the scenario manager. Assuming, that the necessary files are usually kept in the example's directory where the simulation is started, we include all files from this directory. They are not set relative to `temp2`, but stored in the top directory. It is advisable to provide a subdirectory for the specification files to prevent older files from being included in the actual specification.

The configuration files can be specified with the `-f`-switch. If no file is defined, `omnetpp.ini` from the working directory

is taken by default. In the `ini`-files, the network to use is set and the parameters for the modules the network consists of. As it is expected, that the `ini`-file and the network's `ned`-file reside in the same directory, the location of the `ini`-files must also be transferred to the relative hierarchy of the base directories. An alternative would be to store the `ini`-file in the working directory and provide the complete package path for the network.

Our aim was to change as little in OMNeT++ as possible. Therefore, we took advantage of already existing features, that provided us with the number of runs, which we could use to write the task specifications. Each task includes the paths for the `ned`-files and the configuration files, that have to be adjusted to the new hierarchy. In our example, one task specification looks as follows:

```
taskSpecifications = {
        0 = {arguments = (
            "-n",
            "temp1:temp2",
            "-f",
            "temp2/sctp/fair/config.ini",
            "-c",
            "testconfig",
            "-r",
            0
            );
            command = "/home/user/INET/src/inet";
        };
```

In the first stage we added only one new command line switch to set the name of the job specification file. The creation of a specification file is started from the example's directory. For a complex framework like INET the command can be as follows:

```
../../../src/inet -n ../..:../../../src -u Cmdenv \
-f config.ini -c testconfig -s specfile.xml
```

In addition to the parameters for the normal run, that is the location of the `ned` and configuration files, the command environment has to be chosen (`-u Cmdenv`), the configuration with the flag `-c` and the name of the specification file. The job can then be started with the batch job command of Section 2.

As we mentioned before, also the executable belongs to the input files. This file can easily have a size of several megabytes, which will lead to a very large specification file that has to be transferred to the agents. In Section 6 we will prove that Xgrid is very inefficient in transferring large files and will repeat the entire file transfer for every task submitted to an agent. To improve the performance, it is advisable to keep the batch file as small as possible, so we will leave the executable out in further batch jobs. It should rather be stored on a shared medium or copied to the same location on each agent in advance.

Using a shared medium for all the input files can reduce the size of the specification file further. Therefore, we introduced the command line switch `-t` to give the user the opportunity to decide against a self-contained job and to exclude the input files from the batch file and call them from the shared medium. As a consequence, the job specification is altered to contain only absolute paths to the files of the working directory in the *inputFiles* section and to the `ned` and configuration files in the task specifications.
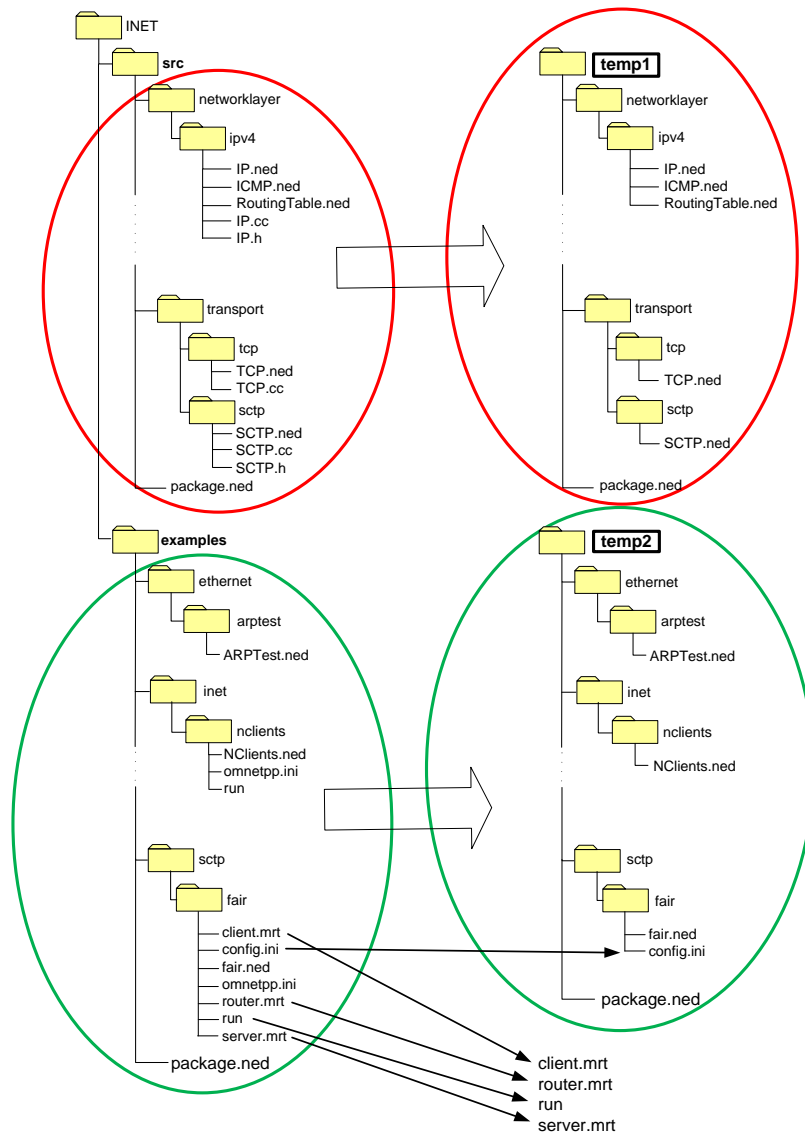
**Figure 3: Mapping of the original files to the Xgrid hierarchy**

## 5. SETTING UP XGRID

A basic setup of Xgrid is just activating the controller and agents with or without password authentication. Every agent in the same network only needs the password and will then use Bonjour [2] to discover the server and connect to it. Bonjour is a Zeroconf service to automatically detect and configure services offered by other hosts in a network.

That is basically everything needed to be done before jobs can be submitted to the controller and then get distributed to the agents. However, since 10.5 Leopard the agents are running received tasks within a sandbox which limits read and write access to a temporary working directory. Every file needed by the tasks has to be included in the submission of the job and has to be transferred to the agents by the controller for each task in advance. Files that have been created by the task are transferred back to the controller after finishing the computation and can be retrieved by the client. Xgrid is not very efficient at transferring files, the maximum supported cumulative size is 128 MB, so it is better to configure the simulations to use a network share for reading and writing large files.

If Xgrid agents authenticate to the server only with a password or do not authenticate themselves at all, they will run every task they receive in a sandbox. This makes it impossible for the tasks to access any network share. The only possibility to grant access to anything outside their working directory is to run the job with normal user permissions. This is achieved by setting up an Open Directory server on the controller's machine to which every host running an Xgrid agent has to be bound to. Open Directory requires a configured DNS service on the controller as well. Both services can be configured with a graphical user interface.

Binding a host to an Open Directory server allows users within the directory, called network users, to log on to any of the connected machines. Network users have the same access rights as local standard users but exist on every bound host. This is an essential prerequisite to execute a task as a
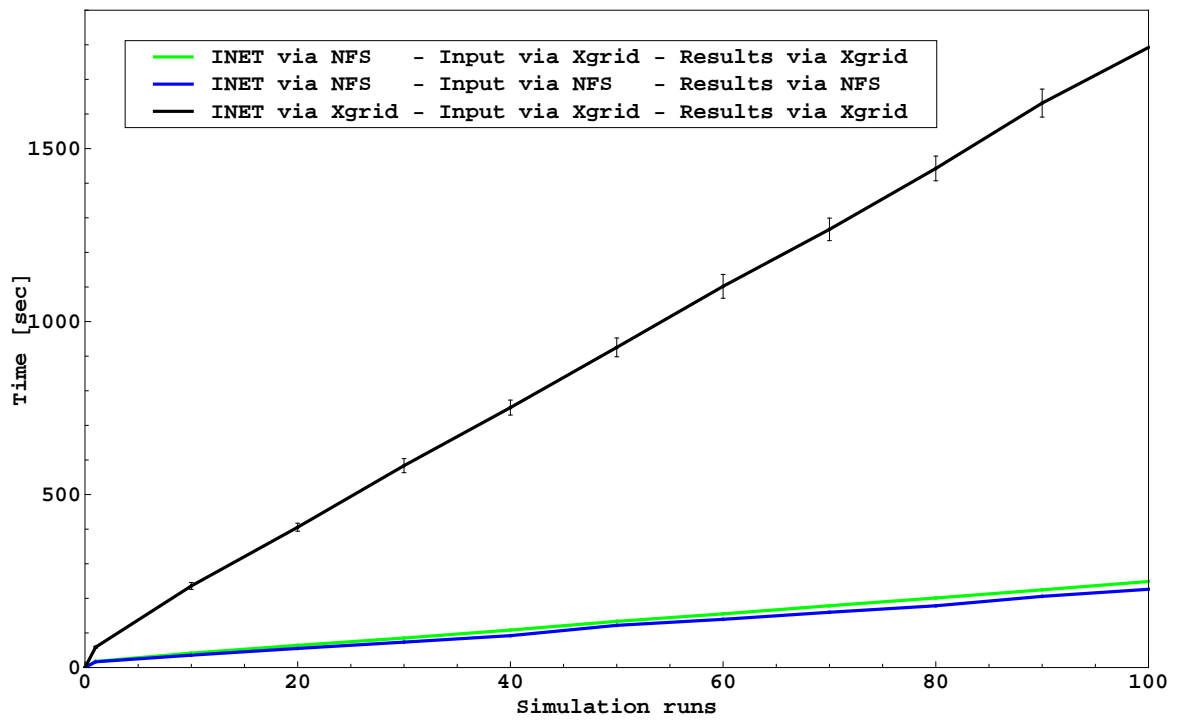
**Figure 4:** Performance comparison of a large INET binary input file either read from an NFS share or transferred by Xgrid
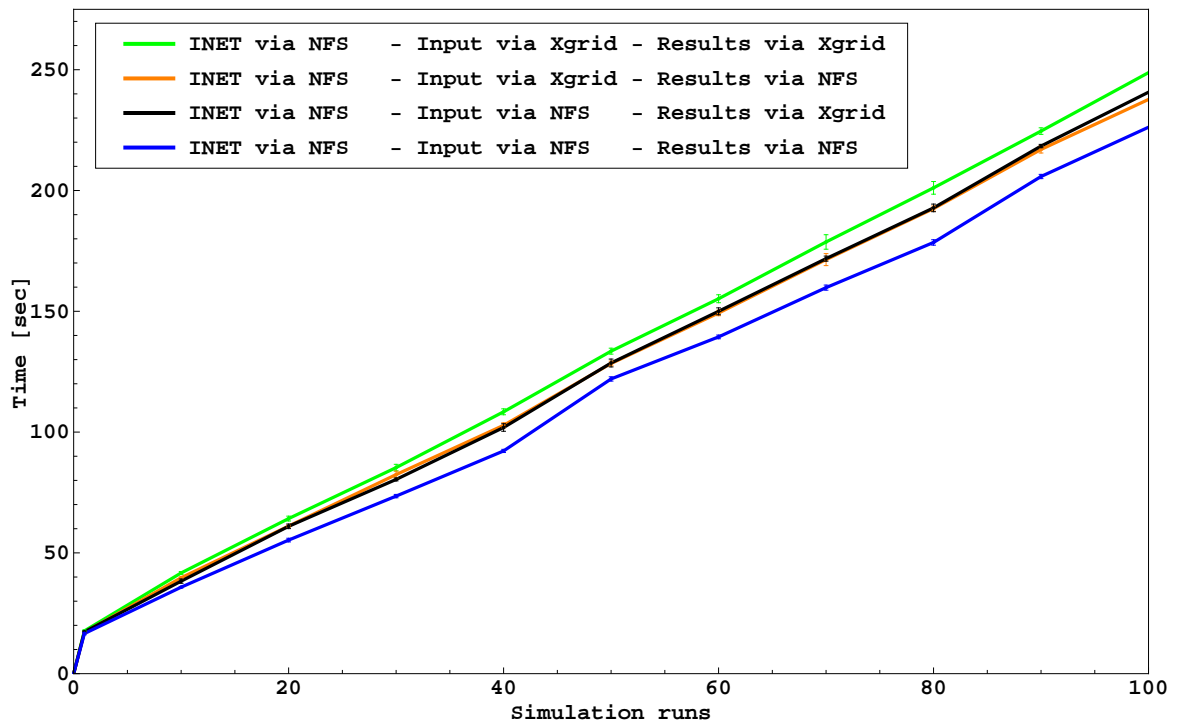


**Figure 5:** Performance comparison of input and result files either read from an NFS share or transferred by Xgrid
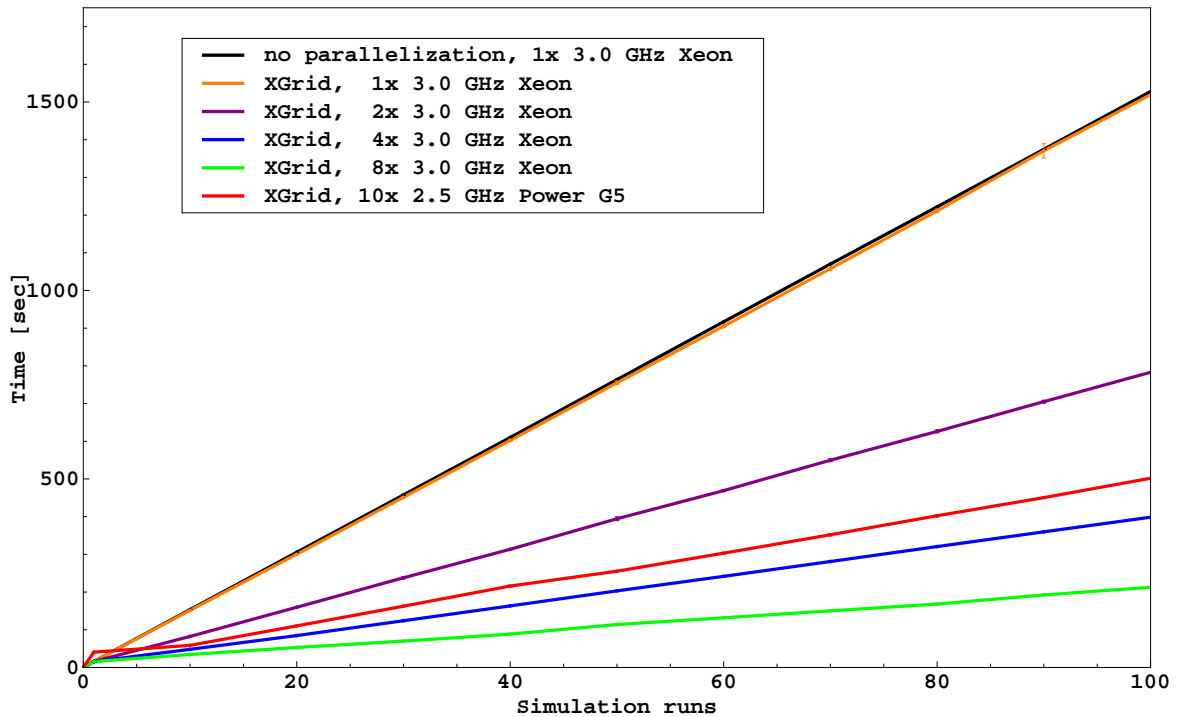
**Figure 6: Conventional computing compared to Xgrid with an increasing number of CPU cores**

standard user because every task within a job has to be run with the same permissions on all agents. After the agents have been bound to the Open Directory, Kerberos authentication can be used. The client has to create a Kerberos ticket for a user which will automatically be valid on every agent. When the job is submitted to the controller the tasks are run with the permissions of the ticket's owner. This allows the tasks to access anything a normally logged on user would be able to, including network shares.

Since the tasks are now able to access the network, shares have to be set up for users who should be allowed to submit jobs. A convenient solution is to enable home directories stored on the Open Directory server which the connected hosts can mount via NFS. This allows network users to access their files from every host in the directory and thus from every agent. Any permission issues should be solved without any additional administrative work.

It is then possible to create tasks for which the agents either get the data from the controller or just read them from the home directory of the submitting user. The results can be written directly to the home directory, whereby the available computing power of the agents can be used most efficiently.

## 6.    MEASUREMENT RESULTS

To compare the performance of Xgrid supported computations, our agents are an 8 core Mac Pro and several dual core G5 Powermacs accepting tasks from a dedicated controller which also provides NFS shares. The simulation used is always the SCTP example 'fair' of the INET framework, which is a quite simple simulation with a short runtime and small configuration and result files. The time needed to

compute an increasing number of runs is measured.

With the first series of measurements we examined the most efficient job specification, more precisely, which files should be included in the batch file and which better be distributed via a network share. The easiest way to handle jobs is to make them self-contained, that is include everything. Unfortunately we discovered that this may be slower even with all 8 cores of the Mac Pro than running the simulation without Xgrid. All included files are transmitted to the agents every time before a task is submitted and the results are sent back to the controller after finishing a task. With growing file sizes this can be a severe problem because the agents do not only have to wait until they have received all data to start computation, they also have to wait until the results are transferred back. Depending on the network this can take several minutes and lead to some agents always waiting for data transfer without doing any computations. Figure 4 shows the topmost graph that, letting Xgrid handle all the file transfers, is up to 10 times slower than using an NFS share, which corresponds to the lower graphs.

To avoid this bottleneck, only files of a cumulative size of up to a few megabytes should be included or transferred back as results. We therefore moved the large 34 MB INET binary to an NFS share and only included the small configuration files in the batch file and set the simulation to write its results to an NFS share, too. This is a convenient solution because the INET binary usually does not change and so the jobs are almost self-contained. The results can be obtained from the share with similar effort like retrieving them from the Xgrid controller. Storing all input files on a network share is actually slightly faster but the files cannot be moved or modified during simulation, so only one simu-

lation can be submitted to Xgrid at a time or there must be a separate directory structure for each job. Figure 5 shows the differences between files on NFS shares and transferred by Xgrid, while the more NFS is used instead of Xgrid, the more it tends to result in a higher performance.

With the combination of including files and using a network share as a viable setup, the performance of Xgrid supported can be compared to conventional computation. The number of CPU cores used with Xgrid is increased with every series of measurements. As shown in Figure 6, Xgrid with only one core enabled is as fast as only running OMNeT++. An additional core almost halves the computation time and even more cores reduce the time to complete the simulation nearly linearly. This can be seen in Figure 6, where from top to bottom with every graph the measurement series is done with additional CPU power, except for the third from the bottom, which is just for comparison how many slower Power G5 cores perform against fewer faster Xeon cores.

## 7. LIMITATIONS AND CONSTRAINTS

The transmission of necessary files before running a job is usually no issue, even though they have to be transferred for each job separately, because they tend to be small. However, the collection of the results can cause significant delays for large files and only a cumulative size of 128 MB is supported anyway. To avoid this delay and limit, it is more efficient to have the agents write their results continuously to a network share. This causes a constant network and disk usage which should be lower than trying to write everything at once after finishing a task and no pausing is necessary. Unfortunately this requires a more complex setup with an Open Directory service and network shares.

The parallelization using Xgrid can be used to parallelize multiple runs of the same simulation. However, the depicted solution only allows the variation of parameters. Xgrid supports the use of MPI and therefore simulations could be partitioned for distributed computing, but this is specific to every simulation and requires an elaborate preparation.

In the course of this paper we assumed that OMNeT++ is built without shared libraries, so that no additional libraries have to be transmitted to run the simulations. In a future version the necessary libraries shall be automatically included in the task specifications to allow an even easier use of our extension.

Another limitation is that the `ini`-files are not searched for the inclusion of other files. Therefore, the user has to specify all needed configuration files explicitly.

## 8. CONCLUSION

This paper described the parallelization of OMNeT++ simulations with Xgrid. Changes in OMNeT++ are required to generate job description files for Xgrid automatically, which enable the user to easily set up his simulation for Xgrid. Processing the job specification file, the Xgrid controller then distributes the separate runs of the simulation to all available machines for parallel computation. Smaller files, such as configuration files, can be included in the job description file, whereas larger files should be read from or written to network shares for better performance. Such a setup accelerates the computation of a simulation almost linear to the additional computing power compared to conventional computation with a single process.

It is planned to include this extension as a tool in an upcoming OMNeT++ release.

## 9. REFERENCES

[1] Akaroa2 for OMNeT++. *Retrieved from: http://www-tkn.ee.tu-berlin.de/research/akaroa-omnetpp/index.html.*

[2] Bonjour Protocol Specifications. *Retrieved from: http://developer.apple.com/networking/bonjour/specs.html.*

[3] Eclipse Open Source Integrated Development Environment. *Retrieved from: http://www.eclipse.org.*

[4] INET Framework Documentation. *Retrieved from: http://www.omnetpp.org/staticpages/index.php?page=20041019113420757.*

[5] OMNET++ User Manual Version 3.2. *Retrieved from: http://www.omnetpp.org/doc/manual/usman.html.*

[6] Project Akaroa. *Retrieved from: http://www.cosc.canterbury.ac.nz/research/RG/net_sim/simulation_group/akaroa/about.chtml.*

[7] Xgrid Administration and High Performance Computing. *Retrieved from: http://images.apple.com/server/macosx/docs/Xgrid_Admin_and_HPC_v10.5.pdf.*

[8] Xgrid Agent for Java. *http://sourceforge.net/projects/xgridagent-java.*

[9] Xgrid Agent for Mac OS X 10.3. *http://www.apple.com/support/downloads/xgridagentformacosx103.html.*

[10] Xgrid: High Performance Computing for the Rest of Us. *Retrieved from: http://developer.apple.com/hardwaredrivers/hpc/xgrid_intro.html.*

[11] T. Dreibholz and E. Rathgeb. Reliable Server Pooling A Novel IETF Architecture for Availability-Sensitive Services. *Proceedings of the 2nd IEEE International Conference on Digital Society (ICDS)*, February 2008.

[12] T. Dreibholz and E. Rathgeb. A Powerful Tool-Chain for Setup, Distributed Processing, Analysis and Debugging of OMNeT++ Simulations. *International developer's Workshop on OMNeT++ (OMNeT++ 2008)*, March 2008.

[13] I. Rüngeler, M. Tüxen, and E. Rathgeb. Integration of SCTP in the OMNeT++ Simulation Environment. *International developer's Workshop on OMNeT++ (OMNeT++ 2008)*, March 2008.

[14] R. Stewart. Stream Control Transmission Protocol. *RFC 4960*, September 2007.