



Project no. IST-033576

# XtreemOS

Integrated Project

BUILDING AND PROMOTING A LINUX-BASED OPERATING SYSTEM TO SUPPORT VIRTUAL ORGANIZATIONS FOR NEXT GENERATION GRIDS

## Study of XtreemOS testbed extension

### D4.3.2

Due date of deliverable: November, 30, 2007

Actual submission date: December, 12, 2007

Start date of project: June 1<sup>st</sup> 2006

Type: Deliverable  
WP number: WP4.3

Responsible institution: INRIA  
Editor & and editor's address: Yvon Jégou  
IRISA/INRIA  
Campus de Beaulieu  
Rennes Cedex  
FRANCE

Version 1.0 / Last edited by Yvon Jégou / Dec, 12, 2007

Project co-funded by the European Commission within the Sixth Framework Programme		
Dissemination Level		
<b>PU</b>	Public	✓
<b>PP</b>	Restricted to other programme participants (including the Commission Services)	
<b>RE</b>	Restricted to a group specified by the consortium (including the Commission Services)	
<b>CO</b>	Confidential, only for members of the consortium (including the Commission Services)	

**Revision history:**

<b>Version</b>	<b>Date</b>	<b>Authors</b>	<b>Institution</b>	<b>Section affected, comments</b>
0.1	09/11/07	Yvon Jégou	INRIA	Initial document
0.2	07/12/07	Yvon Jégou	INRIA	Integrated reviewer's comments
1.0	12/12/07	Yvon Jégou	INRIA	Final version

**Reviewers:**

Paolo Costa (VUA) and Gregor Pipan (XLAB)

**Tasks related to this deliverable:**

<b>Task No.</b>	<b>Task description</b>	<b>Partners involved<sup>°</sup></b>
T4.3.2	Preparing the extension of XtremOS testbed to other grids	INRIA*, VUA, ICT

<sup>°</sup>This task list may not be equivalent to the list of partners contributing as authors to the deliverable

\*Task leader

## Executive summary

Operating system developments must be validated on real hardware. The behavior of a grid operating system depends on so many parameters such as the number of nodes, their heterogeneity (memory, cpu, devices), the structure of the Grid (small or large clusters), the interconnection network (structure, latency, bottlenecks), the dynamicity of the grid, the stability of the grid (node failures), the efficiency of grid services, etc. that there is no possible way to evaluate it through simulation. Grid operating systems such as XtremOS must be validated on a realistic testbed. This grid testbed must provide a significant number of computation nodes for scalability evaluation. The testbed nodes must allow full reconfiguration of the software stack for operating system experimentations from the low level communication layers up to the grid services.

XtremOS is a grid operating system targeting thousands of nodes. The Grid'5000 platform is the initial testbed for XtremOS and can provide up to 2000 physical nodes. A larger testbed is expected for the second half of the project through the integration of the DAS-3 and the CNGrid platforms.

This document analyzes the three experimental grid platforms, Grid'5000 in France, DAS-3 in The Netherlands and CNGrid in China under various aspects: hardware and software architecture, platform management, usage policy, Internet policy... In order to allow realistic operating system experiments, each platform must provide the possibility to run operating system experiments (on raw hardware or on virtual machines), the usage policies must allow the co-reservation of a significant number of resources for a few hours and the Internet policies must provide direct IPV6 communication between all computation nodes.

The main conclusion of this deliverable is that it is not possible to aggregate these three platforms for XtremOS experimentations because of their current configuration. It is necessary to negotiate various settings, mainly:

- all platforms must support IPV6;
- Grid'5000 has very restricted connectivity with the outside world;
- Only application and middleware experimentations are possible on DAS-3. Operating system experimentations must be possible for XtremOS;
- Operating system experimentations are possible only on part of CNGrid.

Further steps are necessary before DAS-3 and CNGrid can be integrated to the XtremOS testbed. Apart from some technical limitations such as IPV6, firewalling rules or reservation rules which can be easily solved, the major point to be negotiated is the capability for external users to make operating system experimentations on DAS-3 and CNGrid. This capability has a major impact on the security of the platforms.

This deliverable will serve as a negotiation base between XtremOS partners and the platform managers in order to provide a larger testbed for the second half of the XtremOS project.



## Contents

<b>1</b>	<b>Introduction</b>	<b>5</b>
<b>2</b>	<b>Presentation</b>	<b>5</b>
2.1	Presentation of Grid'5000 . . . . .	5
2.2	Presentation of DAS-3 . . . . .	5
2.3	Presentation of CNGrid . . . . .	6
<b>3</b>	<b>Platform architecture</b>	<b>8</b>
3.1	Grid'5000 . . . . .	8
3.1.1	Grid'5000 nodes . . . . .	8
3.1.2	Grid'5000 local networks . . . . .	8
3.1.3	Grid'5000 interconnection . . . . .	9
3.2	DAS-3 . . . . .	9
3.2.1	DAS-3 nodes . . . . .	9
3.2.2	DAS-3 local networks . . . . .	10
3.2.3	DAS-3 interconnection . . . . .	10
3.3	CNGrid . . . . .	12
3.3.1	CNGrid nodes . . . . .	12
3.3.2	CNGrid local networks . . . . .	12
3.3.3	CNGrid interconnection . . . . .	12
<b>4</b>	<b>Platform management</b>	<b>14</b>
4.1	Grid'5000 . . . . .	14
4.1.1	Grid'5000 usage policy . . . . .	14
4.1.2	Grid'5000 reservation system . . . . .	14
4.1.3	Operating system experimentations on Grid'5000 . . . . .	15
4.2	DAS-3 . . . . .	15
4.2.1	DAS-3 usage policy . . . . .	15
4.2.2	DAS-3 reservation system . . . . .	16
4.2.3	Operating system experimentations on DAS-3 . . . . .	16
4.3	CNGrid . . . . .	16
4.3.1	CNGrid reservation system . . . . .	16

4.3.2	Operating system experimentations . . . . .	16
<b>5</b>	<b>Internet policy</b>	<b>16</b>
5.1	Grid'5000 . . . . .	16
5.2	DAS-3 . . . . .	17
5.3	CNGrid . . . . .	17
<b>6</b>	<b>Issues</b>	<b>17</b>
6.1	General Issues . . . . .	17
6.2	Specific Issues: Grid'5000 . . . . .	18
6.3	Specific Issues: DAS-3 . . . . .	18
6.4	Specific Issues: CNGrid . . . . .	18
<b>7</b>	<b>Conclusion</b>	<b>19</b>

## 1 Introduction

During the first half of the project, the French Grid'5000 infrastructure is used by all partners as the XtremOS grid testbed. The XtremOS project plans to integrate DAS-3 and the China National Grid (CNGrid) platforms to XtremOS testbed in order to provide this large-scale international platform during the last two years of the project.

This deliverable compares these platforms from three main points of view: reservation system, capability to run operating system experimentations and testbed interconnection. The last section lists actions to be taken before these platforms can be integrated into the XtremOS testbed.

## 2 Presentation

### 2.1 Presentation of Grid'5000

This short presentation has been extracted mainly from the Grid'5000 home page [5]. The Grid'5000 platform is an experimental testbed for research on all layers of the software stack used in Grid computing: networking protocols (improving point to point and multi-points protocols in the Grid context, etc.), operating systems mechanisms (virtual machines, single system image, etc.), Grid middleware, application runtimes (object oriented, desktop oriented, etc.), applications (in many disciplines: life science, physics, engineering, etc.) and problem solving environments. Research in these layers concerns scalability (up to thousands of CPUs), performance, fault tolerance, QoS and security.

The Grid'5000 platform is distributed on 9 sites in France: Bordeaux, Grenoble, Lille, Lyon, Nancy, Orsay, Rennes, Sophia-Antipolis and Toulouse. Each site is equipped with one special node, the site front-end. A site front-end supports the management tools for all clusters of the site (reservation tools, deployment tools, monitoring tools). All nodes of the same site are interconnected through 1 gigabit Ethernet. Other network equipments can be present, for instance Myrinet-2000, Myrinet-10g or Infiniband on some clusters. Grid'5000 sites are interconnected by a dedicated 10 Gb/s optical link.

### 2.2 Presentation of DAS-3

This short presentation has been extracted from the DAS-3 home page [2]. The Distributed ASCI Supercomputer (DAS) is an experimental testbed for research on wide-area distributed and parallel applications. The system was built for the Advanced School for Computing and Imaging (ASCI), a Dutch research school in which several universities participate. The goal of DAS is to provide a common computational infrastructure for researchers within ASCI, who work on various aspects of parallel and distributed systems, including communication substrates, programming environments, and applications.

DAS-3 (The Distributed ASCI Supercomputer 3) is a five-cluster wide-area distributed system designed by ASCI. DAS-3 is funded by NWO/NCF (the Netherlands Organization for Scientific Research), the VL-e project, and the participating universities and organizations. As one of its distinguishing features, DAS-3 employs a novel internal wide-area interconnect based on light paths.

The goal of DAS-3 is to provide a common computational infrastructure for researchers within ASCI, who work on various aspects of parallel, distributed, and grid computing, and large-scale multimedia

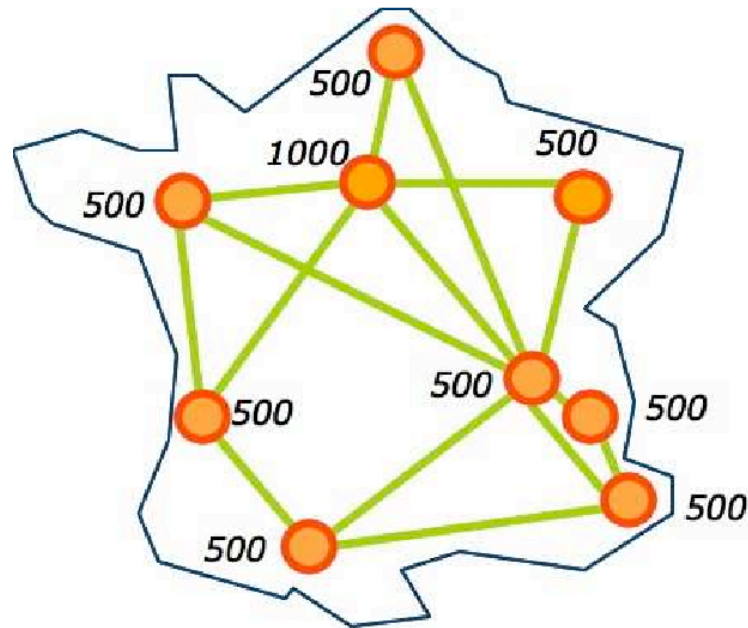


Figure 1: Overview of Grid'5000.

content analysis. The following institutes and organizations are directly involved in the realization and running of DAS-3:

- Vrije Universiteit, Amsterdam (VU)
- Leiden University (LU)
- University of Amsterdam (UvA)
- Delft University of Technology (TUD)
- The MultimediaN Consortium (UvA-MN)

As one of its distinguishing features, DAS-3 employs a very novel internal wide-area interconnect based on light-paths.

### 2.3 Presentation of CNGrid

The China National Grid Project is a key project supported by the National High-Tech R&D Program (the 863 program), which is a continuing effort to the HPCE project. The China National Grid (CNGrid) Project is a testbed for the new generation of information infrastructure by integrating high performance computing and process transaction capacity. It supports various applications including scientific research, resource and environment research, advanced manufacturing and information service by sharing resources, collaborating and service mechanism. It also propels the progress of national informatization and related industry through technology innovation.



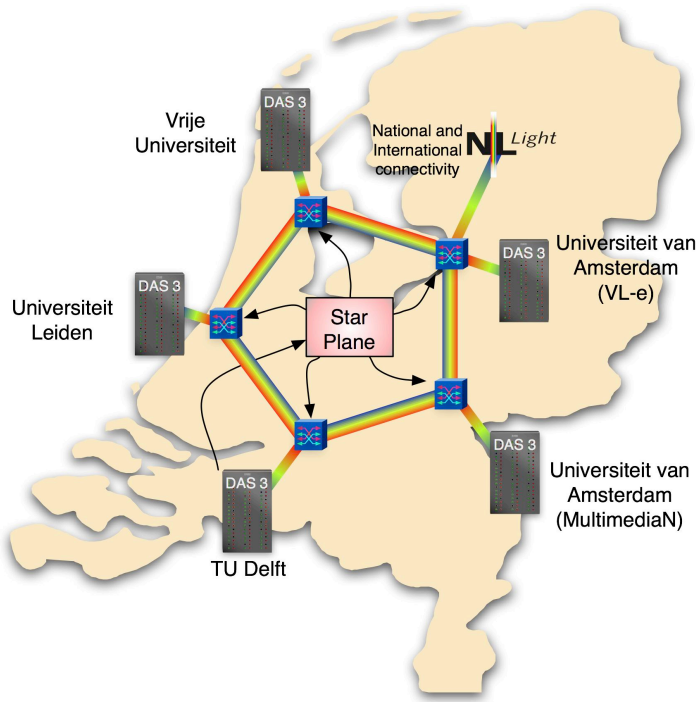


Figure 2: Overview of DAS-3



Figure 3: Overview of CNGrid.

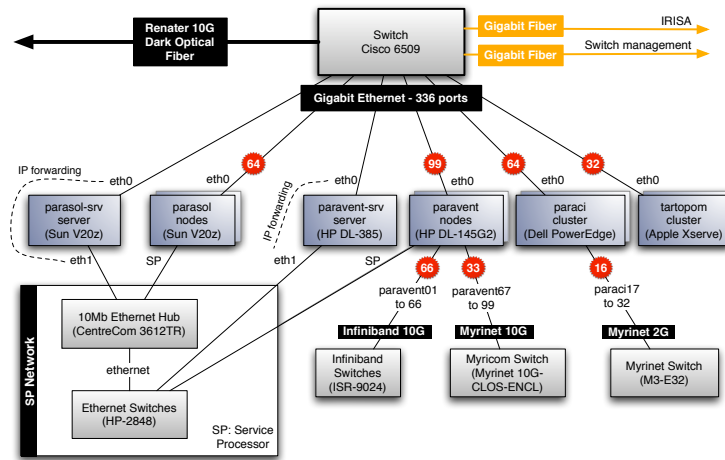


Figure 4: Grid'5000-Rennes local network

The CNGrid currently consists of 8 sites. The nodes are located at Beijing (Computer Network & Information Center of CAS, Tsinghua University, Applied Physics and Computing Institute), Shanghai (Shanghai Supercomputer Center), Hefei (University of Science and Technology of China), Changsha (University of National Defense Science and Technology), Xian (Xian Jiaotong University), and Hong Kong (Hong Kong University). The nodes are inter-connected by public networks such as CERNet and CSTNet. Each node is distinguished by its unique applications, for example, bio-information applications at Tsinghua University, weather forecasting at Shanghai Supercomputer Center, etc.

### 3 Platform architecture

#### 3.1 Grid'5000

##### 3.1.1 Grid'5000 nodes

Grid'5000 currently hosts 22 clusters on 9 sites for a total of 1454 dual-processors computation nodes. All processors are of type 64 bits x86. Each site provides a shared home file-system for all its nodes. The Linux distribution running on the nodes is site-dependent: debian, ubuntu or fedora.

##### 3.1.2 Grid'5000 local networks

Each site provides a basic Gigabit Ethernet interconnection of all its nodes. Some sites provide extra high performance local networks such as Myrinet and Infiniband. Figure 4 shows a typical local network installation in the site of Rennes including Myrinet-10G and Infiniband on some clusters. Each Grid'5000 site provides one 10 Gb/s connection to the Grid'5000 dedicated network and one 1 Gb/s connection to the local LAN.

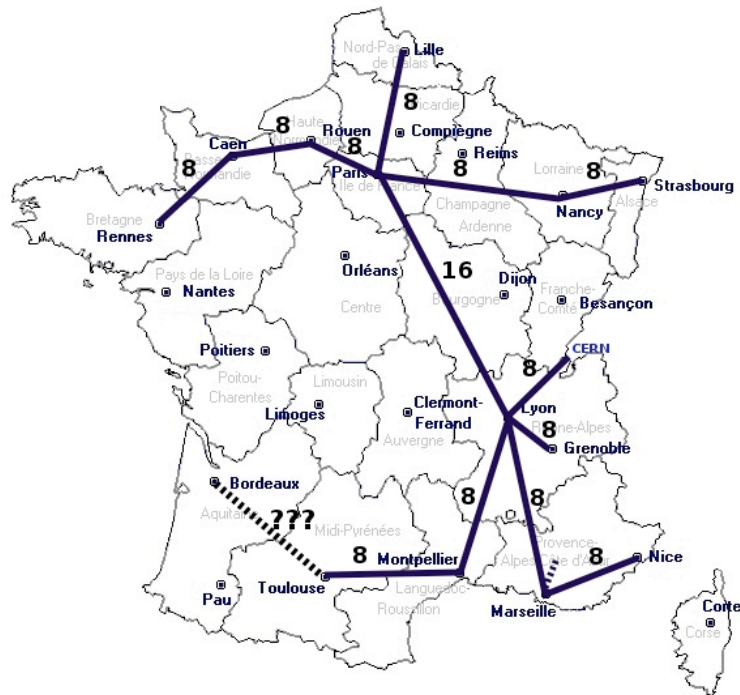


Figure 5: Grid'5000 dedicated network

### 3.1.3 Grid'5000 interconnection

Renater [10], the French National Telecommunication Network for Technology, Education and Research provides a private network interconnecting all Grid'5000 sites through dedicated 10 Gb/s optical links. This network provides level 2 routing between the sites. Level 2 routing allows experimentations on Internet protocols between Grid'5000 sites. Figure 5 shows the Grid'5000 dedicated network. Each site provides controlled connectivity to Internet for some or all of its local nodes.

## 3.2 DAS-3

### 3.2.1 DAS-3 nodes

DAS-3 consists of 272 dual AMD Opteron compute nodes, spread out over five clusters, located at the four universities. The system has been built by ClusterVision. Unlike its predecessor, DAS-2, DAS-3 is rather heterogeneous in design:

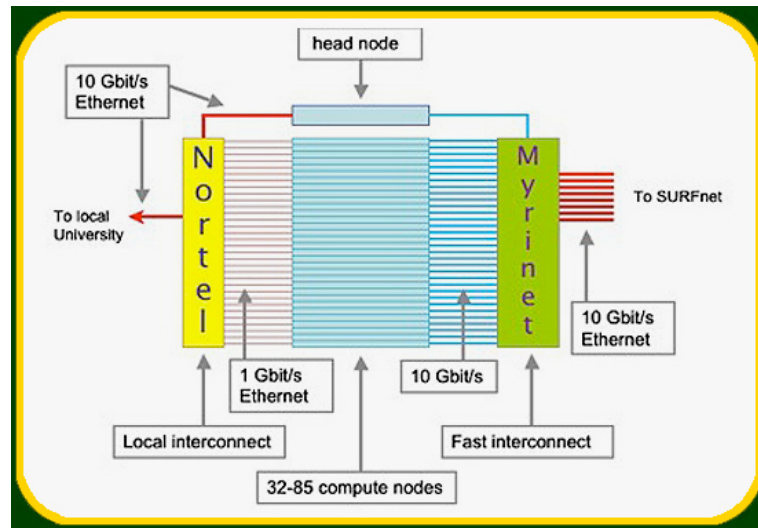


Figure 6: DAS-3 local network

Cluster	Nodes	Type	Speed	Memory	Storage	Node HDDs	Network
VU	85 dual	dual-core	2.4 GHz	4 GB	10 TB	85 x 250 GB	Myri-10G and GbE
LU	32 dual	single-core	2.6 GHz	4 GB	10 TB	32 x 400 GB	Myri-10G and GbE
UvA	41 dual	dual-core	2.2 GHz	4 GB	5 TB	41 x 250 GB	Myri-10G and GbE
TUD	68 dual	single-core	2.4 GHz	4 GB	5 TB	68 x 250 GB	GbE (no Myri-10G)
UvA-MN	46 dual	single-core	2.4 GHz	4 GB	3 TB	46 x 1.5 TB	Myri-10G and GbE

### 3.2.2 DAS-3 local networks

The connectivity scheme of each local cluster is shown in figure 6. There is a LAN connection to the local university and a WAN connection to SURFnet. The head node and the cluster nodes connect to the LAN with Ethernet interfaces: 1 GE for the cluster nodes and 10GE for the head node. The Ethernet switches used for this connection are stackable Nortel 5530s and 5510s.

### 3.2.3 DAS-3 interconnection

Myri-10G is used both as an internal high-speed interconnect as well as to directly communicate with extremely low latency and jitter at up to 10 Gbit/s (or even multiples of that) to remote DAS-3 computing resources via a fully optical (DWDM) backbone, SURFnet6. The WAN connections, as shown in the drawing, go via a Myrinet switch. Both the head node and the cluster nodes use 10 Gbit/s connections on NIC set to operate in Myrinet mode. The Myrinet switch itself has 8 10GE connections to SURFnet,

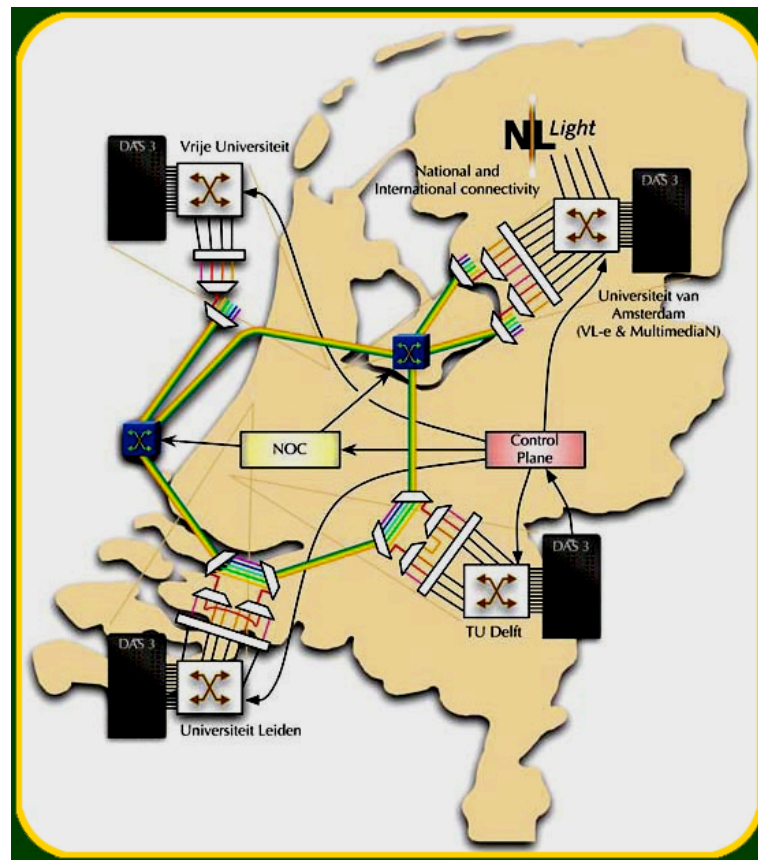


Figure 7: Interconnection DAS-3 sites

and more specifically to a Nortel OME 6500.

Applications running on the cluster will contact the management plane to request a suitable network path through SURFnet6. The StarPlane control plane will configure both the network portion in the DAS-3 domain, such as the Myrinet switches, and negotiate the path with the CPL components on SURFnet6.

### **3.3 CNGrid**

CNGrid currently hosts 10 clusters, more than 3620 processors and 4000 cores on 8 sites.

#### **3.3.1 CNGrid nodes**

Most sites of CNGrid host one or more clusters; several sites are also equipped with MPP/SMP machines (e.g. IBM SP2). Clusters contain dual processor Intel Itanium or AMD Opteron. Some nodes are dual cores. Several experimental nodes are quad cores.

The two major CNGrid clusters are Dawning 4000A and Lenovo DeepComp 6800. Dawning 4000A, [3] located in Shanghai Supercomputer Center, runs 2560 AMD Opteron 2200 MHz processors and is ranked #121 on the TOP500 list (June 2007). Lenovo DeepComp 6800, [7] located at Chinese Academy of Science, runs 1024 Itanium2 1.3 GHz processors and is ranked #473 on the TOP5000 list (June 2007).

#### **3.3.2 CNGrid local networks**

The nodes in clusters are connected by gigabit Ethernet or Myrinet. New clusters like Dawning 5000 will be equipped with Infiniband in the future (available after 2008).

#### **3.3.3 CNGrid interconnection**

All sites of CNGrid are connected through CERNET and CSTNET.

With approximately 1,300 connected organizations and 15 million end users, CERNET is currently China's largest national Research and Education Internet backbone, and the second largest network backbone in China. It provides Internet services for Chinese universities, institutes, schools and other non-profit organizations. Managed by the Computer Network Information Center of the Chinese Academy of Sciences (CNIC, CAS), the Chinese Scientific and Technology Network (CSTNET) (figure 9) connects with international counterparts, including the United States at 2.5G, Russia at 155M, Korea at 10G and Japan at 1G, as well as with Hong Kong district at 2.5G and Taiwan district at 1G. For a detailed description of CERNET and CSTNET and the possibilities of Europe-China network interconnection, please refer to EuChinaGRID deliverable [4].

All sites provide external IPs (public ipv4 addresses) for access gateways. Access to these external access gateway is very limited under strict firewall regulations.



Figure 8: CERNET topology



Figure 9: CSTNET topology

## 4 Platform management

### 4.1 Grid'5000

The default Linux installation on each Grid'5000 node provides a minimal set of basic software for users: Java, MPI, compilers, etc.

#### 4.1.1 Grid'5000 usage policy

The Grid'5000 platform is intended to support research in all areas of computer science related to large scale distributed processing and networking. Users should use Grid'5000 in the perspective of large scale experiments (at least 3 sites and 1000 CPUs).

It is a shared tool, used by many people with different and varying needs. The administrators pursue the following objectives, the main one being the first of this list:

1. Make the tool available to experiments involving a significant number of nodes (in the 1000's). To make this possible, reservation fragmentation must be avoided as much as possible.
2. Keep the platform available for the development of experiments during the day. Therefore, reservations using all the nodes available on one site during work hours (in France) should be avoided in general.

The platform can be used with two different modes: submissions and reservations.

- Submission: the user submits an experiment when he lets the scheduler decide when to run it.
- Reservation: when the user makes a reservation, he gain usage of the platform at the time he explicitly specified. He will then need to launch his experiment interactively.

#### 4.1.2 Grid'5000 reservation system

Each Grid'5000 site hosts a batch OAR [9] queuing system. The OAR system reserves the requested number of nodes for the duration of a program run. It is also possible to reserve a number of hosts in advance, terminate running jobs, or query the status of current jobs. OAR handle background jobs as well as interactive jobs. OAR also distinguishes best effort jobs (which run as long as no other reservation requests the nodes), normal jobs and deployable jobs (which allow users to experiment their own operating system).

Each OAR system is in charge of all clusters of one site. The heterogeneity of the clusters (type of nodes, disks, local networks) is managed through OAR properties: it is, for instance, possible to request for nodes connected to the same Ethernet switch.

Multi-site reservations are also possible using the oargid [8] wrapper.



### 4.1.3 Operating system experimentations on Grid'5000

Operating system experimentation are possible using the OAR *deploy* job type. During a deploy job, a user gets the capability to reformat disk partitions, to populate a disk partition and to reboot the nodes. Grid'5000 provides a specific tool, kdeploy [6] for the management of operating system images: kdeploy integrates taktuk [13], an optimized disk copying system and a set of hooks which are run before and after operating installation (local configuration file setup).

In the current setup of Grid'5000, reservation requests with operating system deployment capabilities reserve whole nodes: nodes cannot be shared in deploy mode. A user can deploy a virtualization infrastructure (Xen, OpenVZ) and install his own virtual machines. Each Grid'5000 site manages a pool of IP addresses, routable across all sites, which can be allocated to the virtual machines. The possibility to deploy virtual machines should be provided in the future.

## 4.2 DAS-3

The operating system the DAS-3 runs is Scientific Linux [11]. In addition, software from many sources is available to support research on DAS-3: the Grid Engine resource management system, various MPI implementations (e.g., MPICH and OpenMPI), the Globus Grid toolkit, optimizing compilers, visualization packages, performance analysis tools, etc.

### 4.2.1 DAS-3 usage policy

Each DAS-3 cluster consists of a file/compile server and a number of compute nodes. The file/compile server is called fsX (X is 0, 1, 2, 3 or 4 according to the cluster base); the worker nodes are named nodeX[0-9][0-9].

The DAS-3 must be used in the following way:

- Program development (editing, compiling) is done on the file/compile servers.
- Program execution must be done on the compute nodes, never on a file/compile server. The file/compile servers are usually heavily loaded just carrying out the tasks they are intended for.
- Jobs must be run on the worker nodes via the DAS-3 cluster scheduling system Sun Grid Engine (SGE), which offers sufficient abstraction and flexibility for nearly all needs.
- The user interface of the original DAS scheduler *prun* has also been ported to the DAS-3 scheduler. For many users this will be the most convenient way to start jobs on one of the clusters.
- The default run time for jobs scheduled on DAS-3 is 15 minutes, which is also the maximum during working hours, from 08:00 to 20:00. Users should not monopolize the clusters for a longer time, since that makes interactively running short jobs on a large number of nodes practically impossible.
- During day time, DAS-3 is specifically not a cluster for doing production work. It is meant for people doing experimental research on parallel and distributed programming. Only during the night and in the weekend, when DAS-3 is regularly idle, it is allowed to run long jobs.

#### **4.2.2 DAS-3 reservation system**

Cluster management is done using ClusterVision's ClusterVisionOS [1]. Programs are started on the DAS-3 compute nodes using the Sun Grid Engine (SGE) [12] batch queueing system. The SGE system reserves the requested number of nodes for the duration of a program run. It is also possible to reserve a number of hosts in advance, terminate running jobs, or query the status of current jobs.

There is a single GridEngine server per DAS-3 site, which runs autonomously. Co-allocation of nodes on different sites is done via other means: e.g. by additional tools like Koala (TUDelft), or simply by hand.

Additionally, the tool "prun" is supported which provides a convenient command-line user interface to SGE.

#### **4.2.3 Operating system experimentations on DAS-3**

In the current configuration, no operating system experimentation is possible on the DAS-3 platform.

### **4.3 CNGrid**

#### **4.3.1 CNGrid reservation system**

The main job manager system on big clusters is Platform LSF. Small clusters are managed by OpenPBS. The reservation system depends on the cluster vendor.

#### **4.3.2 Operating system experimentations**

The major parts of the CNGrid clusters are production clusters and cannot be interrupted for operating system experimentations. Several experimental clusters present on some sites of CNGrid allow users to experiment their own OS/distribution. However it is recommended to experiment OS/distributions in virtual machines.

## **5 Internet policy**

### **5.1 Grid'5000**

Because users are allowed to experiment their own implementations of operating systems, the security of grid'5000 sites is based on strong firewalling rules:

- no restriction between Grid'5000 sites,
- all communications from a Grid'5000 node to a non Grid'5000 node are blocked, except for a strict controlled list of external services (debian mirror for instance)

- all communications from a non Grid'5000 node to a Grid'5000 node are filtered, except for a strict controlled list of services such as the site front-end.

Each site provides one front-end for its local user. The front-end accepts connections from local user origin domains (a small range of IP addresses). Once a user is connected to his Grid'5000 home site, he can access all Grid'5000 sites. Some Grid'5000 sites provide routable IPV4 addresses for their experimentation nodes. Routable addresses allow experimentations with the external world provided that the firewall can be configured. Some Grid'5000 sites recognize satellite sites. A satellite site is an external experimental cluster. The site firewall can be configured to allow experimentations with satellite sites.

## 5.2 DAS-3

The basic Internet policy of DAS-3 is: allow communication between a site and the external world in both directions to facilitate cross-site experiments, coupling with remote Grids and datasets, etc. All DAS-3 nodes have routable IPV4 addresses. All DAS-3 head- and compute nodes have access lists for standard system services like ssh to allow access from a limited number of sites only.

## 5.3 CNGrid

Except for gateway (front-end) machines or experimental clusters, there is no direct Internet connection between sites and the external world.

- routable IP addresses: the basic policy is applied through explicit firewall rules.
- private IP addresses: nodes cannot be directly addressed from outside. A site can provide access to Internet (a limited list of addresses) for the nodes through NAT and/or proxies (for instance to some open source repositories). Experimentations with external world are possible (only) through tunnelling.

A limited number of experimental clusters have full access to the Internet using public IPV4 addresses. For the other sites, experimentations with the external is (still) possible (only) through tunnelling.

# 6 Issues

## 6.1 General Issues

The three platforms, Grid'5000, DAS-3 and CNGrid are currently configured using IPV4 IP addressing. But XtreamOS requires IPV6 addressing (XtreamOS deliverable D4.2.1, Requirements Capture and Use Case Scenarios) and the implementation of the scalable and fault tolerant layer of XtreamOS is based on mobile IPV6.

Using IPV6 on Grid'5000 should not introduce insolvable problems: the inter site interconnection is level 2 and all switching equipments are IPV6-compatible.

The possibility to use IPV6 IP addressing for XtremOS experimentations needs to be negotiated for DAS3 and CNGrid. Alternative solutions such as tunnelling are also possible in the case where the switching equipments do not support IPV6.

## **6.2 Specific Issues: Grid'5000**

Grid'5000 is the initial XtremOS testbed. XtremOS developers can request for a Grid'5000 account from the XtremOS office. Each external user must provide a small IP address range along with the account request in order to update the firewall rules for the Grid'5000 front-end in Rennes. Using this account, a user can log in the Grid'5000 front-end in Rennes and then has access to the whole Grid'5000.

Grid'5000 nodes are isolated from the external world. However, it is possible to request a firewall reconfiguration of some sites in order to run experimentations with nodes external from Grid'5000. This is the case of the site of Rennes which is partly managed by XtremOS partners.

Some Grid'5000 sites use private IPV4 IP addresses which cannot be routed outside Grid'5000. It is possible to overcome this limitation using tunnelling through one node from a Grid'5000 site using public addresses.

## **6.3 Specific Issues: DAS-3**

The current setup of DAS-3 does not allow operating system experimentations. This is the major point to be negotiated with DAS-3 staff in order to enable XtremOS experimentations.

Nodes from DAS-3 can communicate with Grid'5000 nodes and with CNGrid nodes through the public Internet. An experimental private 10 Gb/s optical link is also being set up between the Grid'5000 network and one DAS-3 site. This extra link should allow more efficient communications between both platforms once the routing tables are finalized by end 2007.

A typical reservation duration on DAS-3 is limited to 15 minutes during working hours. This duration is too short for experiments with multiple grid platform. Longer reservation durations need to be negotiated for XtremOS.

## **6.4 Specific Issues: CNGrid**

It is possible to connect some of the CNGrid experimental clusters to the XtremOS testbed via public Internet. The tunnelling approach between EU sites and CNGrid sites is also possible.

Access to other production clusters from CNGrid for XtremOS experimentations seems more difficult: very restricted connection capabilities to the external world, no general procedure for operating system experimentations.

In order to integrate parts of CNGrid in the XtremOS testbed, it is necessary to negotiate a strategy to reserve and to deploy experimental operating systems on the CNGrid experimental platforms.

## 7 Conclusion

The three Grid platforms, Grid'5000 in France, DAS-3 in The Netherlands and CNGrid in China target different goals:

- Grid'5000 targets experimentations on the whole software stack, from the operating system to the applications,
- DAS-3 targets experimentations on grid middleware and communication layers,
- CNGrid targets mainly experimentations on grid applications.

Grid'5000 is the initial testbed for XtremOS: it is well adapted for operating system experimentations. Extending this initial testbed to DAS-3 and to CNGrid necessitates some negotiation with the platform administrators:

- Grid'5000 has very limited connection capabilities with the external world. It is necessary to open routes to the other platforms through firewall reconfiguration when possible or through tunnelling.
- Operating system experimentations are not possible currently on DAS-3. We need to negotiate some procedure allowing XtremOS experimentations using DAS-3 nodes. Moreover the current reservation policy (15 minutes max for each experimentation during working day) does not seem realistic and need also to be negotiated.
- The major part of CNGrid is a production grid. Only some clusters can be used for operating system experimentations. We need further investigations to check if enough resources can be gathered for XtremOS experimentations.

Finally, XtremOS requires IPV6 addressing. It seems that there is very little experience on using IPV6 on the three platforms. This point must be investigated before building the final XtremOS testbed.

## References

- [1] Cluster vision os. <<http://www.clustervision.com/clustervisionos.html>>.
- [2] Das 3 home page. <<http://www.starplane.org/das3/>>.
- [3] Dawning 4000a. <<http://www.top500.org/system/details/7036>>.
- [4] Euchinagrid. <<http://www.euchinagrid.org/docs/EUChinaGRID-Del2.1-v6.pdf>>.
- [5] Grid'5000 home page. <<http://www.grid5000.org>>.
- [6] Kadeploy web site. <<http://www-id.imag.fr/Logiciels/kadeploy/index.html>>.
- [7] Lenovo deepcomp 6800. <<http://www.top500.org/system/6559>>.

- [8] Oar grid web site. <[https://www.grid5000.fr/mediawiki/index.php/OAR\\_Grid](https://www.grid5000.fr/mediawiki/index.php/OAR_Grid)>.
- [9] Oar web site. <<http://oar.imag.fr/>>.
- [10] Renater web site. <<http://www.renater.fr/>>.
- [11] Scientific linux. <<https://www.scientificlinux.org/>>.
- [12] Sun grid engine. <<http://gridengine.sunsource.net/>>.
- [13] Taktuk web site. <<http://taktuk.gforge.inria.fr/>>.