

Virtual Clusters for Grid Communities

I. Foster^{1,2}, T. Freeman¹, K. Keahey^{1,2}, D. Scheftner¹, B. Sotomayor¹, and X. Zhang¹

¹University of Chicago, Chicago, IL

²Argonne National Laboratory, Argonne, IL

{foster, tfreeman, keahey}@mcs.anl.gov {dscheftn, borja, hai}@cs.uchicago.edu

Abstract

A challenging issue facing Grid communities is that while Grids can provide access to many heterogeneous resources, the resources to which access is provided often do not match the needs of a specific application or service. In an environment in which both resource availability and software requirements evolve rapidly, this disconnect can lead to resource underutilization, user frustration, and much wasted effort spent on bridging the gap between applications and resources. We show here how these issues can be overcome by allowing authorized Grid clients to negotiate the creation of virtual clusters made up of virtual machines configured to suit client requirements for software environment and hardware allocation. We introduce descriptions and methods that allow us to deploy flexibly configured virtual cluster workspaces. We describe their configuration, implementation, and evaluate them in the context of a virtual cluster representing the environment in production use by the Open Science Grid. Our performance evaluation results show that virtual clusters representing current Grid production environments can be deployed and managed efficiently, and thus can provide an acceptable platform for Grid applications.

1. Introduction

Many significant Grid deployments today, such as Open Science Grid [1, 2] and TeraGrid [3], rely on clusters to provide a powerful computation platform for their user communities. Sharing such clusters between different virtual organizations (VOs) [4] is not always easy. Problems arise when the requirements of different VOs that want to use the same resources conflict or are incompatible with site policies. Furthermore, the software available on clusters today does not provide the capabilities needed to implement isolation of different communities and guarantee resource availability while also ensuring good utilization of site resources. Both issues can result in user dissatisfaction (as resources with suitable

hardware and software configuration cannot be guaranteed or even obtained on a best-effort basis), hardware under-utilization (as available resources do not match the current demand), and significant overhead to ensure software compliance and portability.

To overcome these problems, the use of virtual machines (VMs) in Grid computing has been proposed [5, 6]. Virtual machines offer the ability to instantiate a new, independently configured guest environment on a host resource; multiple, different such guest environments can be deployed on one resource at the same time. In addition to providing convenience to a resource's users, this arrangement can also increase resource utilization as more flexible, but strongly enforceable, sharing mechanisms can be put in place. Furthermore, the ability to serialize and migrate the state of a VM opens new opportunities for better load balancing and improved reliability that are not possible with traditional resources. Modern VM implementations, such as Xen [7] and VMware [8] also provide outstanding isolation and enforcement properties, as well as excellent performance making the use of virtual machines cost-effective.

While the combining of Grid and virtualization technology has generated much interest, relatively little attention has been paid to evaluating to what extent the VM technology can be applied to the most common Grid fabric – clusters. Deploying generic virtual clusters of diverse topologies requires the ability to deploy many VMs in a coordinated fashion so that sharing of infrastructure, such as disks and networking, can be properly orchestrated. Due to the need for deployment and configuration of many VMs, cluster deployments can be more costly than the deployment of single VMs [9]. Furthermore, many Grid applications are composed of interdependent tasks requiring frequent I/O calls, a workload that is typically hard for virtual machines to execute efficiently.

Our key contribution in this paper is a description of protocols allowing an authorized Grid client to deploy a virtual cluster and an evaluation of how such

a cluster can be used in terms of speed of deployment and management and suitability for Grid applications. We also discuss our implementation of virtual clusters.

We first provide our work context by describing the related work. We proceed to describe how generic virtual clusters of diverse topologies can be represented to ensure deployment with the required characteristics. We follow with a description of extensions to the Workspace Service [10] required to deploy such clusters, and the implementation that goes with it. We then describe the structure and operation of a typical OSG cluster. Finally, we evaluate our implementation in the context of a virtual cluster modeled on Grid clusters currently in production use by one of the leading Grid communities: the Open Science Grid (OSG) [2]. We complete our investigation by evaluating the performance of an OSG application on the virtual cluster and conclude.

2. Related work

With superior isolation properties, fine-grained resource management, and the ability to instantiate independently configured guest environments on a host resource, virtual machines provide a good platform for Grid computing [5, 6]. The Xenoserver project [11] is building a distributed infrastructure as an extension of the Xen virtual machine effort [7]. The In-Vigo project [12, 13] proposed a distributed Grid infrastructure based on virtual machines, while the Virtuoso [14] and Violin [15] projects explore relevant networking issues.

The approach described here extends the workspace abstraction [10], which treats virtual machines as one of the potential ways of providing virtual resources in the Grid and focuses on using virtual machines to represent *virtual clusters*. Similar work has been done in the context of physical resources by the Cluster on Demand (COD) project [16] which enables a Grid client to obtain a physical cluster partition based on credentials. In this paper, we focus specifically on developing the abstraction of a cluster as well as ways of implementing such clusters using virtual machines.

3. Virtual clusters

Virtual workspaces (VWs) [10] provide an metadata-defined abstraction of an execution environment that can be made available dynamically to authorized clients by using well-defined protocols. We briefly summarize this abstraction here, and then describe how we extended the workspace metadata schema to support a virtual cluster and explain how the schema was used to support its deployment. We also describe changes to Workspace Service interfaces and

implementation. We provide a comprehensive description of both metadata schema and interfaces elsewhere [10].

3.1. The Virtual Workspace

A virtual workspace combines the notion of an execution environment and resources allocated to that environment.

The former is described by *workspace metadata*, which contains all the information needed for deployment. An *atomic workspace*, representing a single execution environment, specifies the data (e.g., VM images) that must be obtained and deployment information (such as networking setup) that must be configured on deployment. As part of a workspace deployment request, a client specifies a requested *resource allocation*, which describes how much resource (CPU, memory, etc.) should be assigned to the workspace.

Workspace deployment has been implemented in the *Workspace Service*, a Globus Toolkit 4 (GT4) [22] based service used to deploy workspaces. The Workspace Service implements a WSRF-based protocol allowing remote clients to start, stop, manage and inspect workspaces. All client actions are authorized using the Grid Security Infrastructure (GSI).

3.2. Virtual Cluster Representation

Atomic workspaces, described previously [9], can be combined to form virtual clusters. The key extension proposed in this work is thus an *aggregate workspace* that contains one or more *workspace sets*: sets of atomic workspaces with the same configuration. A combination of such sets can be used to define complex heterogeneous clusters. For example, a typical OSG cluster, described in Section IV, is composed of two workspace sets: a first set containing one “service node” (with service node configuration) and a second set of “worker nodes” (all with the same worker node configuration). This cluster specification can be easily redefined to include, for example, multiple service nodes of the same or different configurations or different sets of worker nodes.

A virtual cluster’s networking is described to reflect the potentially different configurations of its atomic components: each atomic workspace can have a number of differently configured network connections [10]. Continuing with our OSG cluster example, each worker node has one NIC element configured to obtain an IP on a private network – this could be done by DHCP, pre-arranged, obtained from an external service, or (as in our experiments) allocated by the

Workspace Service. The service node is described by two NIC elements, one of them sharing the subnet of the worker nodes, the other using the pre-arranged option, giving the service node a static, public IP address.

Further configuration options, such as information about shared disk partitions, are also recorded in the definition sections of the workspace metadata for all atomic descriptions. This information is extracted by the Workspace Service and passed as a kernel option when the corresponding node is propagated. When a virtual node is booted, an OS boot script is executed to customize the NFS sharing configuration and update the necessary configuration files.

3.3. Workspace Service Interfaces for Virtual Cluster

Performing operations on the virtual cluster requires allocating resources to a group of VMs as well as deploying and managing them as a group. Therefore, we extended the Workspace Service operations to handle aggregate workspaces. To match the resource needs of aggregate workspaces, we have defined a corresponding aggregate resource allocation type for allowing the user to specify appropriate resource allocations. The resource properties for an aggregate workspace reflect its structure and allow the client to operate on it.

An *aggregate resource allocation* is a set of homogeneous sets of atomic resource allocations (CPU, memory, etc.). For example, an aggregate resource allocation for an OSG cluster might variously (i) be a set of identical resource allocations (if the service node and all worker nodes have the same requirements), (ii) reflect a different resource allocation for the service node and identical ones for the worker nodes, or (iii) be yet another configuration. The structure of the aggregate workspace type and the aggregate resource allocation need not be the same – differently configured workspaces may require the same type of resource allocation, and vice versa. For the purpose of matching workspaces to resource allocations, an ordering is imposed on both sets. In addition, our current implementation assumes that the resource allocation will match the workspace exactly; in the future, we plan to extend this functionality to work with bulk allocations (such as represented by WS-Agreement [17]).

3.4. Workspace Service Implementation

In our implementation we assume that the Workspace Service executes on a service node of a physical cluster and provides a secure gateway to a set of resources that can support the deployment of virtual

machines. We further assume that all data necessary for deployment (such as VM images) has already been staged to a node in the trusted computing base (TCB). The service node of the physical cluster runs a GT4 container and the Workspace Service. The hosts are configured with the Xen hypervisor, Workspace Service back-end scripts and a means to invoke them such as SLURM [18], as well as some means of transferring image files and other data relevant to the workspace from within the TCB.

Our current implementation accepts workspace creation requests based on resource availability. The Workspace Service maintains a database of information about physical hosts available for workspace deployment. For each physical host it records availability, CPU type, total/available memory size, total/available disk size, and system information. When the Workspace Service receives the cluster workspace creation request, it searches the database for a set of resources matching the resource allocation request, defines a matching set, marks it as reserved, and maps the resource allocation onto it. When the workspace is terminated, the resources are reclaimed and the database is modified accordingly. Our implementation does not yet support advance reservations, but the database can support the time dimension. In addition to allocating resources, the Workspace Service also handles local IP address allocation.

Workspace creation results in the creation of a workspace WSRF resource as well as workspace deployment up to a specified state [10]. A workspace is deployed through the invocation of workspace back-end scripts via local schedulers; our current implementation works with SLURM [18] and PBS [19]. The first step of workspace deployment involves propagating the images to the target nodes: workspace scripts executing on each node download the images from a specified location. This step is separate to create an opportunity for pre-staging of the images without actually starting the VMs. To deploy a workspace, the back-end scripts work with the Xen hypervisor and complete the configuration of the workspace. Configuration information that cannot be processed by Xen (such as networking) is set up by calling an OS boot script preinstalled in the VM images. After a workspace is deployed, it can be managed by invoking start and stop operations with different parameters to pause/unpause or shut down a workspace. These operations are broadcast to all participating nodes

4. OSG clusters and infrastructure

The Open Science Grid (OSG) [2, 20] is a U.S.

production-quality Grid for large-scale science, enabling scientists to access shared resources using common Grid infrastructure tools. The OSG software stack is based on the NSF Middleware Initiative distribution, which includes Condor [21] and Globus [22] technologies, as well as additional utilities provided by the Virtual Data System (VDS) [23].

The OSG Grid infrastructure is structured so that users with work to do contact a *submit host* that organizes the work as a set of inter-dependent tasks and orchestrates their execution on OSG resources. To do this, the submit host uses tools provided by VDS to express task workflow, planners such as Pegasus [24] to translate it into an executable form, and Condor DAGMan [25] as well as the Condor-G task manager [21] to schedule them on OSG resources.

OSG resource are typically organized as clusters, consisting of one or more *service nodes* that provide secure access to a group of *worker nodes* on which application jobs are executed. Worker nodes are sometimes configured to be accessible on a site's private network only, while service nodes can always be accessed via both this private network and the public Internet. Services executing on the service nodes implement various management functions. For example, *compute elements* (CEs) provide access to the cluster's compute resources, while *storage elements* (SEs) provide access to its storage resources. In this paper, we focus on providing access to a cluster's computational resources.

A typical OSG cluster has the following configuration:

- There is at least one service node.
- The service node(s) and worker nodes run Debian Linux 3.1 (Sarge).
- A local batch scheduler, such as Condor-C [26], PBS [19], LSF [27], or SGE [28], is installed on the cluster and used to manage it.
- All worker and service nodes have access to a shared file system, such as NFS; thus, application executables, data and other files placed on service nodes are available to worker nodes.
- Grid infrastructure, in current deployments typically the Globus GRAM [29] and GridFTP [30] services, run on the service node to authenticate and authorize job requests, stage required files, and submit jobs to the local batch scheduler. An OSG user first authenticates to an OSG submit host. If the submission is authorized, executables and other data are staged to a selected OSG cluster followed by a GRAM submission (e.g., via Condor-G) of the

requested tasks to the cluster. The submission is orchestrated and managed through the cluster's service nodes.

5. Experimental results

We conducted two sets of experiment to evaluate the performance of an OSG virtual cluster. The first set measures the cost of deploying and managing the cluster itself, and the second the efficiency of the cluster for various application classes.

We ran all experiments on a partition of the Chiba City cluster at Argonne National Laboratory [31]. Each Chiba node has two 500 MHz Intel PIII CPUs (with a 512 KB cache per CPU), 512 MB main memory, and 9 GB of local disk. Nodes are connected by a 100 Mbps LAN.

To host the virtual cluster, Chiba City nodes were configured with Xen 2.0 testing distribution (both domain 0 and user domain run port of Linux 2.6.11) and rebooted with XenLinux. We chose domain 0 memory size to allow the best possible performance: neither too large (so that as much memory as possible can be given to the VM) nor too small (so that it does not interfere with VM performance by skimping on buffering space). We found that for I/O-intensive applications a memory size of 96 MB works well, and we used this size in the experiments described here. In the experiments described here, SLURM was used to schedule VMs on Chiba City nodes.

The configuration of the OSG virtual cluster replicates in detail the configuration of the OSG cluster described in section 4. (We find that even seemingly irrelevant configuration details may impact performance significantly.) Since the images represent inactive VMs, the primary constituent of the image is the VM's disks. The service node file system is made up of three disk partitions mounted from three disk images, including a 600 MB root file system image containing a Debian Sarge Linux installation, a 750 MB image required for the Virtual Data Toolkit (VDT) v1.3.6 OSG middleware stack, and a 1 GB image containing OSG application executables and data. The latter two partitions are exported as an NFS volume so that they can be shared by worker nodes. The worker node's root file system is built from a 600MB root file system image containing Debian Sarge Linux installation, similar to the head node. Worker nodes are configured to share the OSG middleware partition exported by the head node when it is booted.

5.1. Evaluation of Cluster Deployment and Management

We report first on experiments that measure the time it takes to (a) deploy and (b) manage the virtual cluster.

In our timings we assume that all the data to deploy the cluster has already been staged to a known place within the TCB. We split VM deployment into two major components: (1) a virtual cluster scheduling and image propagation phase, in which physical nodes are selected and VM images are copied to those nodes, and (2) customizing and deploying the virtual machines. The first phase takes by far the most time and, disk space and other considerations permitting, can be executed the actual VM deployment thus eliminating this factor from deployment time.

Figure 1 shows timing results for the *create* operation when using the “propagate only” option [10], as measured on the server from the moment of receiving the request to the request’s completion. The graph shows the request time as a sum of (1) processing time, including the time required to find physical resources for deployment, (2) broadcast time, and (3) image propagation time, i.e., the time for the Workspace Service backend to pull images to worker nodes. We see that the broadcast time is negligible, but the operation processing time is relatively high (on the order of seconds versus ~100ms for operations -- see Figure 2) as it includes all of the deployment overhead (i.e., the time required to map virtual machines to resources).

The deployment time is dominated by image propagation. Although the workspace service backend can be configured to use a variety of transport mechanisms for image propagation, the data we show here is obtained using the Linux *cp* command to copy files from a locally mounted NFS directory to a directory on a local disk on the shared file system (NFS v3) available on the physical cluster (using GridFTP in the same scenario yielded similar results). As expected, the propagation time increases linearly with the size of the physical resource allocation (i.e., number of nodes) to which the cluster is assigned. This is because the 100Mbps network link is easily saturated by the 600 MB image transfer and each connection adds a fixed amount of time to the transfer. Using Chiba City’s default NFS server, equipped with a Gigabit Ethernet connection to the switch, the same operation resulted in a significant speedup and flat transfer times for up to 10 transfers, after which the incoming connections to target nodes get saturated.

Figure 2 shows timing results for cluster

management (start and stop) operations for different node numbers. Virtual cluster operation requests are broadcast using cluster job submission tools. All request information is broadcast to all the nodes; each node then selects information relevant to it from the broadcast data. This implementation results in flat broadcast time as shown in the graph. The time spent in each operation is dominated by the time it takes to perform the requested actions, and thus varies little with the cluster size as the operations on each node are independent from each other. The start/running operation is slightly slower than start/unpause since in the latter case the image is already in memory. Note that operation processing time for start and shutdown operations (100 ms) is significantly less than that for the “propagate only” *create* operation seen in figure 1 (a few seconds). As explained previously, this is because the processing time for that operation includes matching VMs to resources as well as resource setup, while the processing time for start/stop operations requires no such setup.

5.2. OSG Applications on Virtual Clusters

To assess the practicality and trade-offs involved in using a virtual cluster, we ran experiments to evaluate the impact of running on a virtual cluster for various OSG applications.

In these tests, we compared the execution time of an application instance on clusters made up of physical Chiba nodes to the execution time of the same application instance on a virtual cluster running on physical Chiba nodes. Unlike Barham et al. [7], we did not replicate the exact hardware conditions for virtual and physical nodes. Thus, the virtual machines deployed on the cluster had access to fewer resources (memory and CPU), as some resources had to be

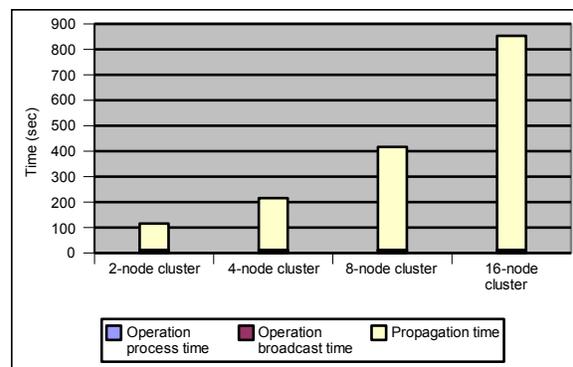


Figure 1. Propagation operation timing results on four clusters.

assigned to hypervisor functions (domain 0). We believe that these tests give a better measure of how well an application running in a VM can perform in practice on a given resource, because we thus include the overhead of running a virtual machine.

For our evaluation we use the Fast Ocean Atmosphere Model (FOAM) [32], a climate science application that uses coupled climate modeling techniques to address questions that require many simulated years of interaction. FOAM is an MPI-based data-parallel application organized into three components: atmosphere, land/sea-ice, and ocean. The first two are collocated data-parallel programs that use MPI for intra-component communication and shared memory for inter-component communication on each node. The third component executes on a distinct set of processors and uses MPI to communicate with the other two. Altogether, the communication makes up about 10% of run time. Other I/O requests issued by the application consist of using GridFTP to transfer the input data files from a storage node (in our experiment, located on the same LAN as the execution cluster) at the beginning of the computation. The size of the data files is less than 22MB.

Figure 3 presents the execution time of the FOAM 1.5 workflow running on both physical and virtual clusters. The time measured is from the submission of the *mpirun* command to its completion. Table 1 shows the performance difference between them.

5.3. Analysis

While much faster than node re-imaging, virtual cluster deployment can be costly—taking on the order of minutes—due to the necessity of copying images to target nodes. Where NFS is present, in principle nodes could be loaded directly from the sources. However, we find that NFS is not always capable of supporting the network traffic that results as some images need to

be frequently written to. A more practical optimization is to ensure that images are staged to advantageous locations (e.g., an NFS server with high speed location). We are also investigating schemes that would optimize propagation by pre-staging the most commonly used images, or their parts, to the worker nodes or by placing less-frequently used read-only OS libraries on separate partitions that could be mounted via NFS with acceptable performance. Another possibility is using broadcast techniques to propagate all the common blocks similar images to several nodes, and then sending each node the blocks that are specific to its VM image.

The potentially relatively high cost of image propagation makes it imperative that a client, or a scheduler, be given the opportunity to execute this operation before the virtual machine is started. We incorporated this observation into our interface design by allowing for the workspace to progress through the deployment states only up to a specific state. This approach gives a client the opportunity to halt the workspace deployment at well-defined points, image propagation being one of them. Even given this flexibility, when using medium to large sized images virtual cluster deployment is potentially expensive; it is therefore most suitable for scenarios that can offset this relatively high deployment cost, such as either hosting several applications, hosting a long-running application or simply contributing a cluster to OSG.

OSG application results are promising. For our evaluation we chose FOAM because its data-parallel nature makes it potentially hard for virtual machines. Furthermore, although our benchmark study of MPI behavior under Xen [33] indicates that certain patterns of MPI communication over Ethernet may be potentially expensive in Xen, we observe the FOAM workflow has close execution time performance on

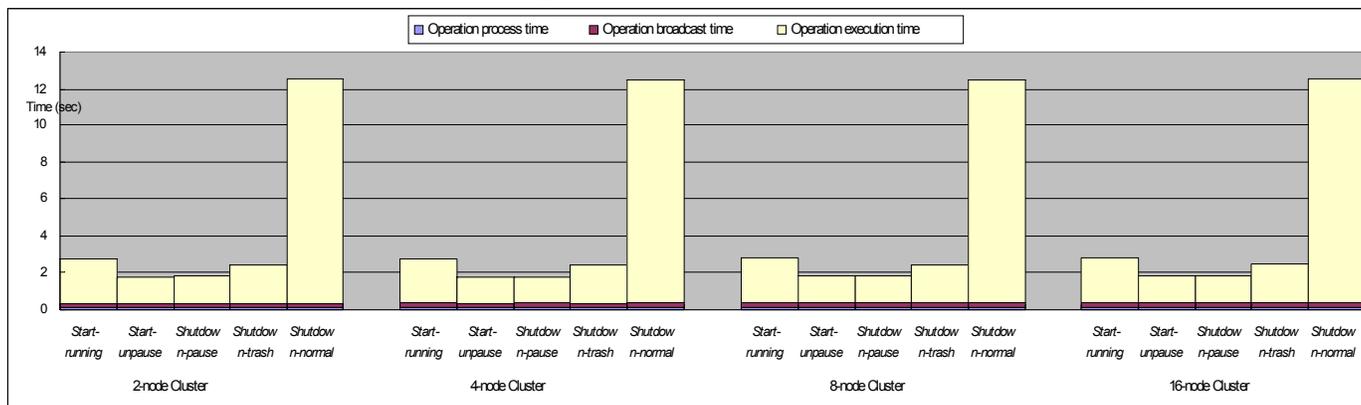


Figure 2. Start and shutdown time as a function of virtual cluster size.

both physical and virtual clusters.

The 5% performance degradation is due to the fact that the overall proportion of time spent on communication within the application is small. In practice, applications that are heavily communication-bound tend to scale poorly and are not a good match for clusters thus making it likely that applications of this type will respond favorably to virtual clusters. These results show that, when communication constitutes a small part of an application, the performance issues observed in [34] are not significant. In fact, evaluating other OSG applications we observed that some tasks appear to execute faster on a virtual than on a physical machine—likely due to double caching of I/O operations. We are currently investigating these effects and their trade-offs in realistic settings.

6. Conclusions

We have described *virtual clusters*, groups of virtual machines designed to execute, and share infrastructure, within a trusted computing base. We showed how to define cluster descriptions so that atomic workspaces can be composed flexibly into more complex constructs, while organizing infrastructure sharing among the virtual nodes. Cluster deployment allows us to specify different resource allocations for different members of aggregates defined in this way. We also described how such clusters can be deployed, evaluated their deployment, and integrated its results into our design.

Our application execution results are promising. The slowdown suffered by the FOAM application from virtual machine impact on execution as well as from the resource overhead of using virtual machines was less than expected: within 5%. Considering that virtual machines offer unprecedented flexibility in terms of matching clients to available resources, this performance impact can be viewed as an acceptable trade-off. Preliminary results from investigating other OSG applications of more complex dependency patterns are equally promising and lead us to believe that virtual clusters have the potential to be a popular solution in production settings.

TABLE 1
FOAM Performance loss on Virtual Cluster relative to Physical Cluster.

4 Nodes	8 Nodes	16 Nodes
4.55%	4.05%	4.10%

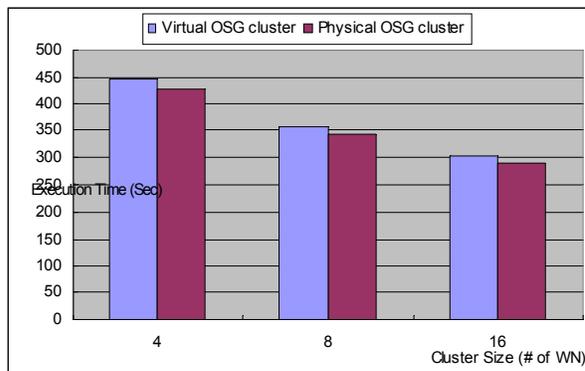


Figure 3: Comparison of FOAM execution time on virtual and physical clusters.

Acknowledgements

We acknowledge the help of Jens Voeckler, Mike Wilde, and Robert Jacob with explanation of the configuration of an OSG cluster as well as application setup. We also thank Rick Bradshaw for assistance with the Chiba City cluster testbed at Argonne. This work was supported in part by the Mathematical, Information, and Computational Sciences Division.

References

- [1] Foster, I. and others. The Grid2003 Production Grid: Principles and Practice. in IEEE International Symposium on High Performance Distributed Computing. 2004: IEEE Computer Science Press.
- [2] Open Science Grid (OSG). 2004: www.opensciencegrid.org.
- [3] TeraGrid: www.teragrid.org.
- [4] Foster, I., C. Kesselman, and S. Tuecke, The Anatomy of the Grid: Enabling Scalable Virtual Organizations. International Journal of Supercomputer Applications, 2001. 15(3): p. 200-222.
- [5] Figueiredo, R., P. Dinda, and J. Fortes. A Case for Grid Computing on Virtual Machines. in 23rd International Conference on Distributed Computing Systems. 2003.
- [6] Keahey, K., K. Doering, and I. Foster. From Sandbox to Playground: Dynamic Virtual Environments in the Grid. in 5th International Workshop in Grid Computing. 2004.
- [7] Barham, P., B. Dragovic, K. Fraser, S. Hand, T. Harris, A. Ho, R. Neugebar, I. Pratt, and A. Warfield. Xen and the Art of Virtualization. in ACM Symposium on Operating Systems Principles (SOSP).
- [8] VMware: http://www.vmware.com/.
- [9] Keahey, K., I. Foster, T. Freeman, X. Zhang, and D. Galron. Virtual Workspaces in the Grid. in Europar. 2005. Lisbon, Portugal.
- [10] Keahey, K., I. Foster, T. Freeman, and X. Zhang, Virtual Workspaces: Achieving Quality of Service and Quality of Life in the Grid. accepted for publication in the Scientific Programming Journal, 2005.

- [11] Reed, D., I. Pratt, P. Menage, S. Early, and N. Stratford. Xenoservers: Accountable Execution of Untrusted Programs. in 7th Workshop on Hot Topics in Operating Systems. 1999. Rio Rico, AZ: IEEE Computer Society Press.
- [12] Adabala, S., V. Chadha, P. Chawla, R. Figueiredo, J. Fortes, I. Krsul, A. Matsunaga, M. Tsugawa, J. Zhang, M. Zhao, L. Zhu, and X. Zhu, From Virtualized Resources to Virtual Computing Grids: The In-VIGO System. *Future Generation Computer Systems*, 2004.
- [13] Krsul, I., A. Ganguly, J. Zhang, J. Fortes, and R. Figueiredo. VMPlants: Providing and Managing Virtual Machine Execution Environments for Grid Computing. in SC04. 2004. Pittsburgh, PA.
- [14] Sundararaj, A. and P. Dinda. Towards Virtual Networks for Virtual Machine Grid Computing. in 3rd USENIX Conference on Virtual Machine Technology. 2004.
- [15] Ruth, P., X. Jiang, D. Xu, and S. Goasguen, Towards Virtual Distributed Environments in a Shared Infrastructure. *IEEE Computer, Special Issue on Virtualization Technologies*, 2005.
- [16] Chase, J., L. Grit, D. Irwin, J. Moore, and S. Sprenkle, Dynamic Virtual Clusters in a Grid Site Manager. accepted to the 12th International Symposium on High Performance Distributed Computing (HPDC-12), 2003.
- [17] Andrieux, A., K. Czajkowski, A. Dan, K. Keahey, H. Ludwig, J. Pruyne, J. Rofrano, S. Tuecke, and M. Xu, Web Services Agreement Specification (WS-Agreement) Draft 20. 2004: <https://forge.gridforum.org/projects/graap-wg/>.
- [18] Yoo, A.B., M.A. Jette, and M. Grondona, SLURM: Simple Linux Utility for Resource Management, in *Job Scheduling Strategies for Parallel Processing*, L. Rudolph and U. Schwiegelshohn, Editors. 2003, SpringerVerlag. p. 44-60.
- [19] Portable Batch System. 2003: <http://www.openpbs.org>.
- [20] Pordes, R., The Open Science Grid. *Proceedings of the CHEP'04 Conference*, 2004.
- [21] Frey, J., T. Tannenbaum, I. Foster, M. Livny, and S. Tuecke. Condor-G: A Computation Management Agent for Multi-Institutional Grids. in 10th IEEE International Symposium on High Performance Distributed Computing. 2001: IEEE Computer Society Press.
- [22] Foster, I., Globus Toolkit version 4: Software for Service-Oriented Systems. *IFIP International Conference on Network and Parallel Computing*, 2005.
- [23] Foster, I., J. Voekler, M. Wilde, and Y. Zhao. Chimera: A Virtual Data System for Representing, Querying, and Automating Data Derivation. in 14th Intl. Conf. on Scientific and Statistical Database Management. 2002. Edinburgh, Scotland.
- [24] Deelman, E., J. Blythe, Y. Gil, C. Kesselman, G. Metha, S. Patil, M.-H. Su, K. Vahi, and M. Livny, Pegasus: Mapping Scientific Workflows onto the Grid. *Proceedings of the 2nd European Across Grids Conference*, 2004.
- [25] Frey, J., Condor DAGMan: Handling Inter-Job Dependencies: <http://www.cs.wisc.edu/condor/dagman/>.
- [26] Litzkow, M.J., M. Livny, and M.W. Mutka, Condor - A Hunter of Idle Workstations, in 8th International Conference on Distributed Computing Systems. 1988. p. 104-111.
- [27] LSF Web Site: <http://www.platform.com/products/wm/LSF/index.asp>.
- [28] Gentsch, W., Sun Grid Engine: Towards Creating a Compute Power Grid. *Proceedings of 1st International Symposium on Cluster Computing and the Grid*, 2001.
- [29] Czajkowski, K., I. Foster, and C. Kesselman, Resource and Service Management, in *The Grid: Blueprint for a New Computing Infrastructure (2nd Edition)*. 2004.
- [30] Allcock, W., GridFTP: Protocol Extensions to FTP for the Grid. 2003, Global Grid Forum.
- [31] Chiba City Homepage: <http://www.mcs.anl.gov/chiba>.
- [32] Jacob, R., C. Schafer, I. Foster, M. Tobis, and J. Anderson. Computational Design and Performance of the Fast Ocean Atmosphere Model, Version One. in *International Conference on Computational Science*. 2001: Springer-Verlag.
- [33] Zhang, X., The Effect of DomO Memory Size on the Performance of DomU Applications: <http://people.cs.uchicago.edu/~hai/vcluster/find-dom0-sweetpoint.pdf>.
- [34] Zhang, X., Evaluation of a Virtual Xen Cluster Using the Pallas MPI Benchmarks Suite: <http://people.cs.uchicago.edu/~hai/vcluster/PMB/>.