

Research Report

Preparing the Worldwide LHC Computing Grid for MPI applications



Research Report for RP2
University of Amsterdam
MSc in System and Network Engineering

Class of 2005-2006

Richard de Jong, Matthijs Koot
{rjong,mrkoot}@os3.nl

30th June 2006



Abstract

This document is the result of a four-week research project, performed in the course of our MSc in System and Network Engineering. We describe the current and desired state of support for running parallel applications using the MPI library on the LHC Computing Grid, that is currently being developed to support the experiments on CERNs Large Hadron Collider. We propose a solution to achieve the first phase of the desired state (single-site MPI) and we describe the key problem areas and possible solutions in achieving the second phase (cross-site MPI). We include a description and an evaluation of YAIM, a tool used for installing and configuring Grid middleware. We found that YAIM fulfills its goals, but only if it is used strictly in the environments it was designed for: small and/or non-complex sites.

Acknowledgement

We'd like to thank dr. David Groep of NIKHEF for being our supervisor, and moreover for being so helpful in provisioning our grid accounts and getting us to CERN. Without him, this project wouldn't have been possible. We hope our work approaches your expectations.

We'd like to thank Louis Poncet MSc of CERN for playing a key role in getting us enrolled in HR, providing us with the inside information that enabled us to get the project started immediately, and mostly for being an allround nice guy. We enjoyed the drinks at the Lake of Geneva, Louis. Happy wake-boarding to you!

We'd like to thank dr. Maarten Litmaath of CERN, dr. Charles Loomis of LAL, dr. Stephen Childs of Grid-Ireland, and our colleagues at the IT-GD team for their feedback and knowledge sharing.

We'd like to thank the University of Amsterdam, and especially dr. Karst Koymans, for providing us with the possibility of performing a Master of Science study in the very interesting and exciting field of System and Network Engineering.

Lastly but certainly not leastly, we'd like to thank our families for their support. You've been great!

Contents

1	Introduction	7
1.1	The Worldwide LHC Computing Grid project	7
1.2	Grid	8
1.3	Research goals	10
1.4	Acronyms	10
1.5	Copyrights and ownership	11
2	Enabling MPI on the WLCG	12
2.1	Introduction to parallel computing	12
2.2	MPI on the WLCG, current situation	15
2.3	MPI on the WLCG, desired situation	15
2.3.1	Single-site MPI (first phase)	16
2.3.1.1	Implications for the end-user	16
2.3.1.2	Implications for RBs	17
2.3.1.3	Implications for CEs	17
2.3.1.4	Implications for WNs	19
2.3.2	Cross-site MPI (second phase)	19
2.3.2.1	More implications	19
2.3.2.2	Prior art	22
2.3.3	Remarks	23
2.3.3.1	Unmet requirements	23
2.3.3.2	Efficiency of fabric usage	23
2.3.3.3	Integration of MPI and LRMS	24
2.4	Implementations of MPI and their fitness for the WLCG	25
2.4.1	LAM	26
2.4.1.1	LAM-6.5.9	26
2.4.1.2	LAM-6.5.9 on the WLCG	26
2.4.1.3	LAM-7.0.6	26
2.4.1.4	LAM-7.0.6 on the WLCG	26
2.4.1.5	LAM-7.1.2	27
2.4.1.6	LAM-7.1.2 on the WLCG	27
2.4.2	MPICH	27
2.4.2.1	MPICH-1.2.6 on the WLCG	27
2.4.2.2	MPICH-1.2.7 on the WLCG	27
2.4.3	MPICH2	27
2.4.3.1	MPICH2 on the WLCG	28
2.4.4	OpenMPI	28
2.4.5	OpenMPI on the WLCG	29
2.4.6	Mpiexec	29
2.4.7	Mpiexec on the WLCG	29
2.5	Modifications for deployment of MPI	30
2.6	Site Functional Test for MPI	30



2.7	Conclusion on MPI on the WLCG	31
2.8	Future work	31
3	YAIM	32
3.1	Structure	32
3.2	Scope	33
3.3	Evaluation	33
3.4	Conclusion on YAIM	34
3.5	Future work	34
	Bibliography	35
A	The BSD license for this project	39
B	Requirements	40
B.1	User requirements	40
B.2	Site requirements	41
C	Changes to YAIM	42
C.1	functions/config_mpi	42
C.2	examples/site-info.def	46
C.3	scripts/node-info.def	47
D	A Quick Guide to using MPI on the WLCG	48
D.1	mpi-example.jdl	48
D.2	MPItest.c	48
E	Wrapper scripts	50
E.1	lam-6-wrapper.sh without Torque/PBS	50
E.2	lam-7-wrapper.sh without Torque/PBS	52
E.3	lam-7-pbs-wrapper.sh with Torque/PBS	54
E.4	mpich-wrapper.sh without Torque/PBS	56
E.5	mpich-pbs-wrapper.sh with Torque/PBS	58
E.6	mpich2-wrapper.sh without Torque/PBS	59
E.7	mpich2-pbs-wrapper.sh with Torque/PBS	61
E.8	openmpi-wrapper.sh without Torque/PBS	63
E.9	openmpi-wrapper.sh with Torque/PBS	65
F	Site Functional Test for MPI	67
F.1	tests/sft-mpi.def	67
F.2	tests/sft-mpi	67

List of Figures

1.1	Data handling and computation for physics analysis	8
1.2	The layered grid model	9
2.1	Flynn-Johnson taxonomy of computer systems.	13
2.2	Serial computing: example of a SISD scenario.	13
2.3	Parallel computing: example of a MISD scenario.	14
2.4	Parallel computing: example of a SIMD scenario.	14
2.5	Parallel computing: example of a MIMD scenario.	14
2.6	Key (problem) areas of cross-site MPI.	20
2.7	Run-time incompatibility between MPI implementations.	21
2.8	Cross-site MPI with MPICH and MPICH-G2	23
2.9	Each MPI implementation needs to be integrated with a site's LRMS.	24
3.1	The dependability tree	34

Chapter 1

Introduction

As part of our Master of Science study in the field of System and Network Engineering at the University of Amsterdam, we have done research at CERN, Geneva [2], on enabling MPI-based parallel computing on the LHC Computing Grid (LCG). The research was performed on behalf of NIKHEF [3], Amsterdam, and was supervised by David Groep. The work was done in close cooperation with Louis Poncet from the IT Grid Deployment group at CERN.

1.1 The Worldwide LHC Computing Grid project

The Worldwide LHC Computing Grid, or *WLCG*, is the project to create a grid to support the physicists that will analyze the data coming from the Large Hadron Collider, a particle accelerator. The section below is taken from the LCG Project Overview¹ [4]:

“The Large Hadron Collider (LHC), currently being built at CERN near Geneva, is the largest scientific instrument on the planet. When it begins operations in 2007, it will produce roughly 15 Petabytes (15 million Gigabytes) of data annually, which thousands of scientists around the world will access and analyse. The mission of the LHC Computing Project is to build and maintain a data storage and analysis infrastructure for the entire high energy physics community that will use the LHC.

The data from the LHC experiments will be distributed around the globe, according to a four-tiered model. A primary backup will be recorded on tape at CERN, the “Tier-0” centre of LCG. After initial processing, this data will be distributed to a series of Tier-1 centres, large computer centres with sufficient storage capacity for a large fraction of the data, and with round-the-clock support for the Grid.

The Tier-1 centres will make data available to Tier-2 centres, each consisting of one or several collaborating computing facilities, which can store sufficient data and provide adequate computing power for specific analysis tasks. Individual scientists will access these facilities through Tier-3 computing resources, which can consist of local clusters in a University Department or even individual PCs, and which may be allocated to LCG on a regular basis.

Discovering new fundamental particles and analysing their properties with the LHC accelerator is possible only through statistical analysis of the massive amounts of data gathered by the LHC detectors ATLAS, CMS, ALICE and LHCb, and detailed comparison with compute-intensive theoretical simulations. The goals of the LCG project include:

- Developing different software components to support the physics application software in a Grid environment.

¹The W in WLCG is added later than the quoted text was written, to denote the difference between the grid middleware software stack which was used to be called LCG, and the worldwide grid using this software

- Developing and deploying computing services based on a distributed Grid model.
- Managing users and their rights in an international, heterogeneous and non-centralized Grid environment.
- Managing acquisition, installation, and capacity planning for the large number of commodity hardware components that form the physical platform for the LCG.”

The workflow for processing the LHC data is depicted in figure 1.1, taken from [5].

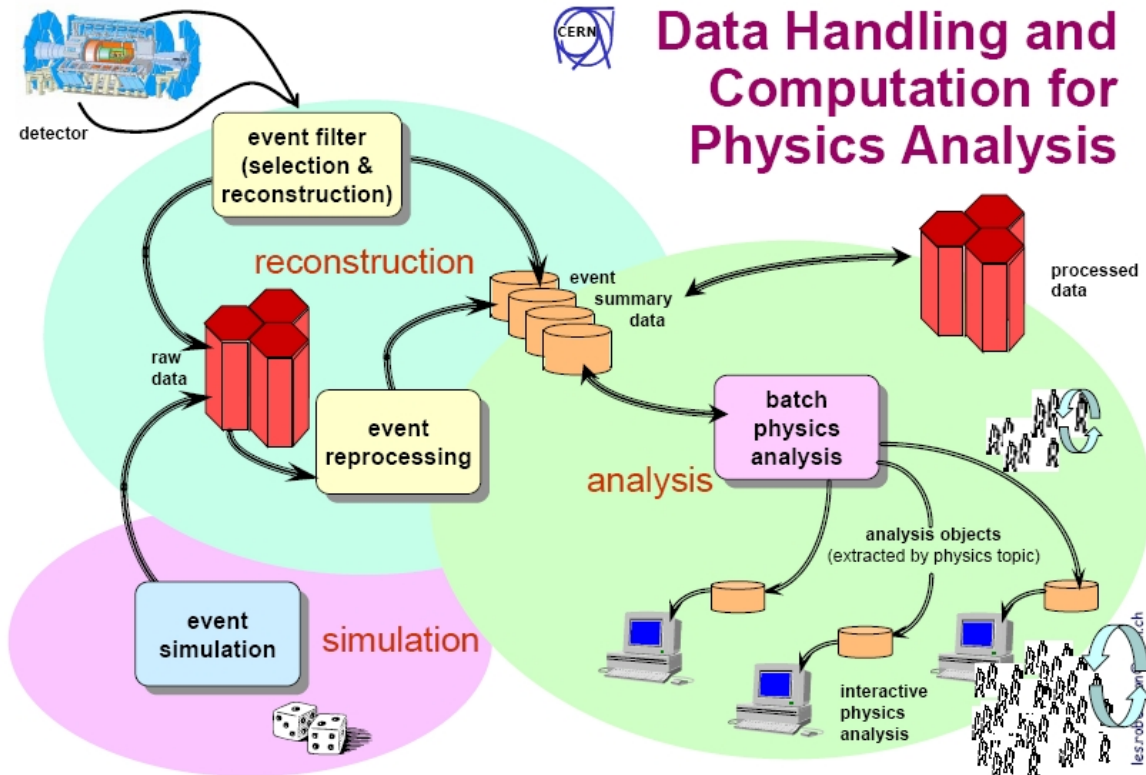


Figure 1.1: Data handling and computation for physics analysis

1.2 Grid

So, what exactly is a *grid*? According to Ian Foster [9], one of the ‘fathers of the grid’, a grid is a system that:

1. ...coordinates resources that are not subject to centralized control...
2. ...using standard, open, general-purpose protocols and interfaces...
3. ...to deliver nontrivial qualities of service.

Thus, it is more than simply a bunch of interconnected computational resources (which might rather be called a *cluster*); it includes social and political aspects as well. The most prevalent vision on grid computing is that of ‘service-oriented computing’, i.e. considering computational resources to be like water and electricity facilities [10, 11]. Exploitation of these means for scientific purposes is also denoted with the term *e-Science*, as coined by John Taylor [10]. Building on knowledge and code from Globus Toolkit (the US-born *de facto* grid middleware), the European DataGrid project, or *EDG*,

and the Enabling Grids for E-sciencE project, or *EGEE*, the LCG project at CERN aims to deliver a production-quality world-wide grid for scientific purposes (and perhaps commercial usage at a later stage). The WLCG actually comprises multiple grids, making it a ‘grid of grids’: besides the LCG, it also connects the US Open Science Grid, or *OSG*, and the Nordic Data Grid Facility, or *NDGF*. This implies a need for grid interoperability. It is believed that one day, like every system needs an IP-address to be connected to the Internet, every system to be connected to ‘the Grid’ shall need to support the Open Grid Services Architecture, or *OGSA*. The OGSA webservice interfaces encourage interoperability between heterogeneous grids. At this time, OGSA is not yet addressed in the gLite development; instead, the interoperability is endeavored by providing common interfaces in the Virtual Data Toolkit, or *VDT* (which provides some fundamental grid services).

During years of research on distributed computing, a general multilayer model has been developed for grid architectures, which is depicted in figure 1.2. This model is semantically comparable with the TCP/IP model used for Internet. The functions placed between the *User Applications* layer and the *Fabric* layer are provided by *Grid middleware*. Several middleware projects exist; in the case of the WLCG, the most prevalent middleware is called *gLite*. gLite, which was previously ambiguously named *LCG* (referring to both the operational grid *and* to the grid middleware itself), is a collection of software delivered by the JRA-1 group of EGEE. The gLite middleware is typically deployed through Yet Another Installation Method, or *YAIM*, an installation method developed in the course of Grid deployment. The OGSA specifies interfaces for all layers.

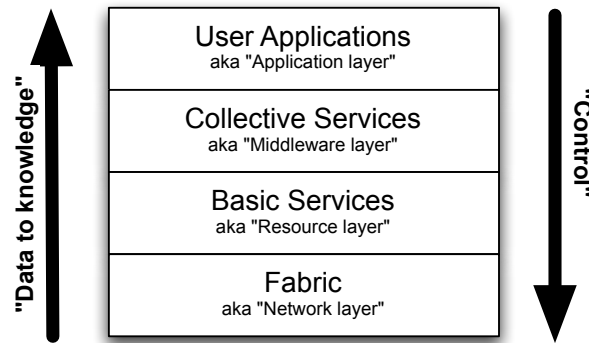


Figure 1.2: The layered grid model

In WLCG terminology, a Computing Element, or *CE*, is an abstraction of Worker Nodes, or *WNs*, and a single ‘head node’ providing a job queue. CEs represent an interface to computing resources. One CE can be seen as a single cluster, owned by an organization that controls this cluster, defines who has access to it and, for instance, is responsible for scheduling decisions. When an end-user submits a *job* described in the Job Description Language, or *JDL*, through a User Interface, or *UI*, the job is sent to a Resource Broker, or *RB*. The RB will notify the Logging & Bookkeeping component, or *LB*, and then query the Grid Index Information Service, or *GIIS*, to find CEs that are capable of executing the job. If any such CEs are found, the RB delegates the job to ‘the best one’ by accessing it through its *Gatekeeper*. The Gatekeeper then calls the CE Job Manager, which distributes the job to the WNs. The components relevant to our research are visualized in figure 2.6.



1.3 Research goals

The research goals for our research were as follows:

Single-site MPI The primary goal of our research was to allow MPI-jobs to be submitted through the grid interface and be executed within a single CE. This consists of the following subgoals:

- Define the requirements for integrating MPI into the gLite middleware and the YAIM deployment scheme, so that, if enabled during gLite deployment, MPI jobs may be submitted through the WLCG grid interface (this includes scalability and dependency issues).
- Integrate (and demonstrate) MPI into the gLite middleware and the YAIM deployment scheme.

YAIM MPI will be deployed and configured by YAIM. Our secondary goal is to assess YAIM on elementary software quality attributes, as well as on adherence to best practices and standards that have been defined by the Global Grid Forum, or *GGF*².

Cross-site MPI The tertiary goal for our research was to pinpoint and document the issues that arise when trying to implement Cross-site MPI.

1.4 Acronyms

The following acronyms are used in this report:

API Application Programming Interface

BDII Berkeley Database Information Index

CE Computing Element

CERN European Laboratory for Particle Physics

EDG European DataGrid

EGEE Enabling Grids for E-sciencE

GGF Global Grid Forum

GIIS Grid Index Information Server

GRAM Globus Resource Allocation Manager

GRIS Grid Resource Information Service

GSI Grid Security Infrastructure

IS Information Service

JDL Job Description Language

LB Logging and Bookkeeping Service

LDAP Lightweight Directory Access Protocol

LHC Large Hadron Collider

LCG LHC Computing Grid

LRMS Local Resource Management System

LSF Load Sharing Facility

MPI Message Passing Interface

OS Operating System

PBS Portable Batch System

RB Resource Broker

RFIO Remote File Input/Output

² GGF is often called ‘the IETF of Grid computing’.



SE Storage Element

SFT Site Functional Tests

SMP Symmetric Multi Processor

QWG Quattor Working Group

VDT Virtual Data Toolkit

VO Virtual Organization

WLCG Worldwide LHC Computing Grid

WMS Workload Management System

WN Worker Node

YAIM Yet Another Installation Mechanism

1.5 Copyrights and ownership

All documents which are created as a part of this project will be licensed under the Creative Commons 2.5 Attribute license [66]. All source and object code which is produced as a part of this project will be licensed under the revised BSD license [67].

Chapter 2

Enabling MPI on the WLCG

2.1 Introduction to parallel computing

Parallel programming is the art of using multiple processors to solve a single problem. The traditional paradigm of computer architecture is the exact opposite: to solve multiple problems with a single processor. This is also known as *serial computing*. The problems which are dealt with in scientific and R&D communities nowadays are often too big to be solved using scalar processors; using a single scalar processor to solve a large-scale problem, such as depicted in figure 2.2, is likely to exceed the acceptable amount of time to solve that problem (should be $< \Delta time$). With millions of scalar computers now connected through the Internet, the perspective on computing is slowly changing to that of ‘the network is the computer’; instead of using a single processor to solve a problem, a global network of processors is now available to solve (the same or larger) problems.

In the Flynn-Johnson¹ taxonomy of computer systems [18], as depicted in figure 2.1, this kind of (non-parallel) system is classified as a Single Instruction, Single Data machine, or *SISD*. To handle a problem with multiple processors, one may divide it into smaller problems through *decomposition* and orchestrate them with *parallel programming*.

In *functional decomposition*², each processor is assigned a different role/responsibility, such as depicted in figure 2.3; this type of system is classified as Multiple Instruction, Single Data, or *MISD*. In *domain decomposition*³, each processor is assigned a different part of the data, such as depicted in figure 2.4; this type of system is classified as Single Instruction, Multiple Data, or *SIMD*; in perfect domain decomposition, given a SIMD with N processors, the total execution time may be decreased to $\Delta time/N$ (best case scenario). The last type of system in Flynn’s original taxonomy, which dates back to 1972, is the Multiple Instruction, Multiple Data, or *MIMD* machine; an example is depicted in figure 2.5. In 1988, Johnson extended the taxonomy by differentiating the MIMD class on *Communication/Synchronization* and *Memory* attributes.

Historic methods of parallel programming included threading, Inter Process Communication (IPC) and Parallel Virtual Machines (PVM). None of them, however, are really suitable for a ‘heavily distributed’ environment like the globe-spanning WLCG [6]. Today, there are two main approaches to parallel programming [14]. In the *directives-based data-parallel language* approach, serial code is parallelized by adding directives that tell the machines how to distribute data and use multiple processors. Examples of such languages are OpenMP and High Performance Fortran. In the *message passing* approach, which was designed to succeed PVM and is suitable for large-scale distribution, the programmers themselves are responsible for dividing data and work across the processors (striving for optimal load balancing), and manage communication between them (striving for minimal and non-blocking communication) by explicitly calling specific library functions. The prime quality attributes of the *de facto* standard for message passing, the Message Passing Interface, or *MPI*, are source code portability — to support

¹ Other and more exhaustive approaches to classification of parallel systems exist, but the Flynn-Johnson appears most widely accepted and suffices for our purposes.

² Parallelism that arises through functional decomposition is also known as *task-parallelism*.

³ Parallelism that arises through domain decomposition is also known as *data-parallelism*.

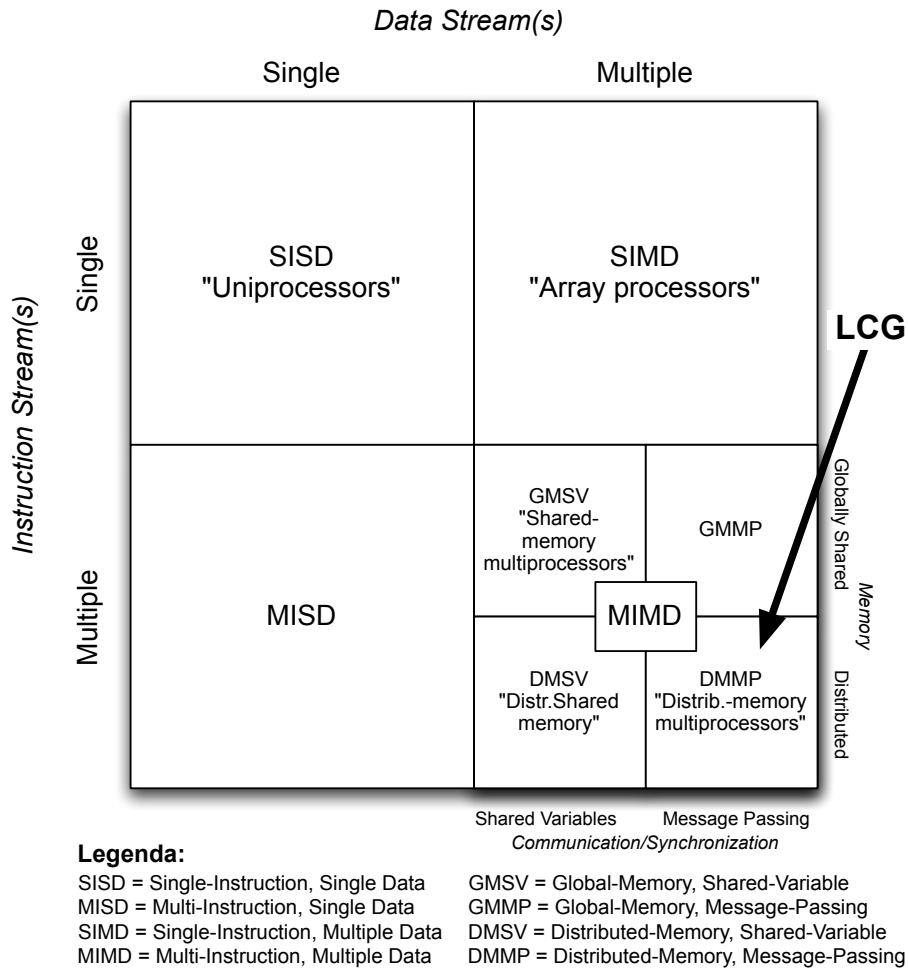


Figure 2.1: Flynn-Johnson taxonomy of computer systems.

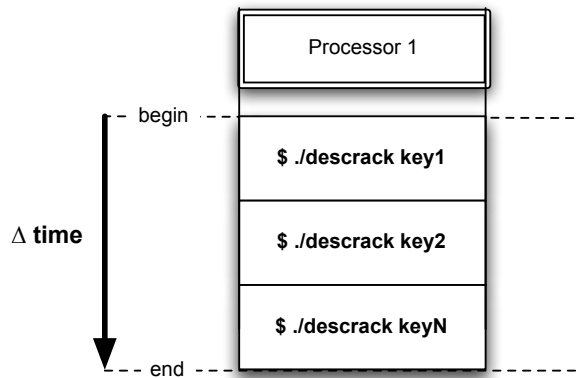


Figure 2.2: Serial computing: example of a SISD scenario.

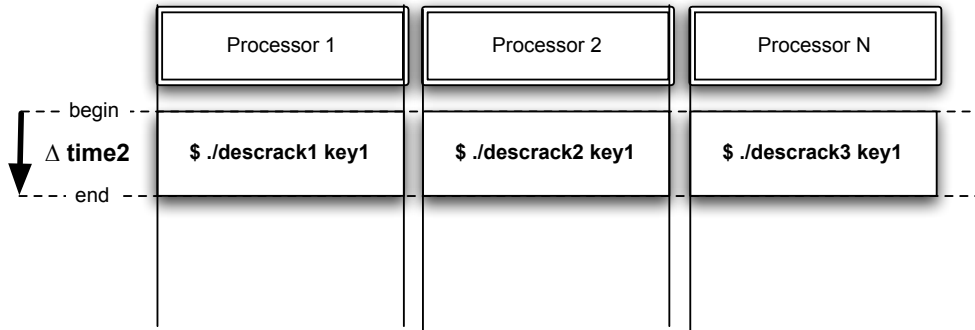


Figure 2.3: Parallel computing: example of a MISD scenario.

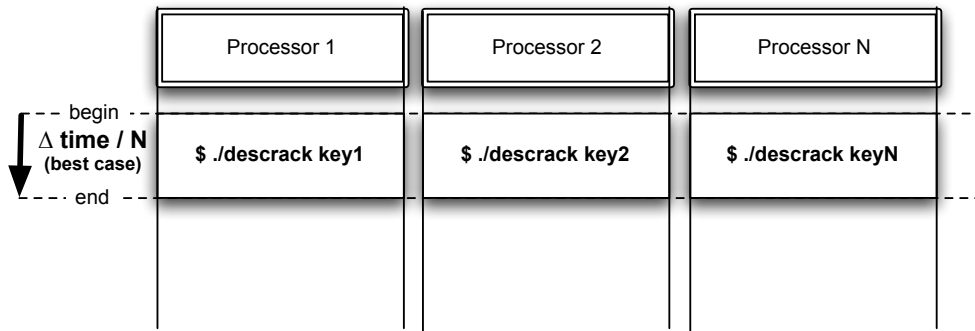


Figure 2.4: Parallel computing: example of a SIMD scenario.

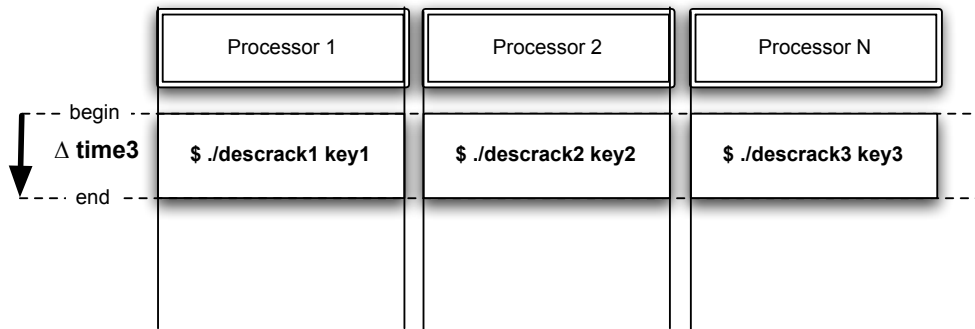


Figure 2.5: Parallel computing: example of a MIMD scenario.



heterogeneous parallel architectures — and to allow efficient implementation [7, 8]. Although MPI may be used in shared-memory architectures (SMP, NUMA), its original design was focused on distributed memory architectures. It is the latter type of environment in which MPI is used at the WLCG; many scalar processors with their own memory, mostly grouped into scientific computing clusters, connected over the Internet.

Message passing has been an established paradigm for parallel computing since the early nineties. Due to the proliferation of several incompatible libraries for parallel computing, the Message Passing Interface Forum, or *MPI Forum*, was established to provide for standardization. The group, consisting of some 60 people from 40 organizations, published version 1.0 of the Message Passing Interface standard in 1994 [15]. Version 1.1 was released in 1995, mainly to correct some errors and make clarifications in the first document. The most recent MPI is version 2.0, which was released in 1997 and is an extension (*not* a modification) to version 1.1. Today, MPI is the most widely-used standard for message passing on distributed-memory parallel computers.

2.2 MPI on the WLCG, current situation

When discussing parallel computing on the WLCG, one is firstly discussing distributed-memory, message-passing MIMD-multiprocessors, although a possible scenario would be to use distributed-memory for inter-CE computation and shared-memory for intra-CE or intra-WN computation (i.e. multiprocessor WNs or CEs with a high-speed inter-WN memory bus, such as Myrinet or Infiniband). Currently, preliminary support for running MPI jobs on the WLCG is already available. It is possible to run an MPI job, provided the following conditions are met:

- SSH host-based authentication is set up between the Worker Nodes abstracted by the CE;
- The attribute “MPICH” is advertised by the CE;

When installing a Worker Node, the mpich package is installed. Mpich is one of the Open Source implementations of MPI. The user can specify a Jobtype of “MPICH” in the job and request a number of nodes. This number of nodes is internally interpreted as the number of CPUs, and it is currently impossible to specify both the required number of nodes and the number of processors. The list of Resource Broker types that are capable of running MPI jobs is hard-coded in the gLite software. The user needs to submit his job as a precompiled binary, linked to the same version of mpich as is installed on the CE his job is sent to. This leaves the user little flexibility in choice of an MPI implementation and in variations on its version. The user’s application is wrapped inside a script that uses Secure Copy (SCP) to copy the executable to all the assigned hosts. Some sites prefer to use a shared filesystem and therefore do not need this. The wrapper script invokes the `mpirun` command, which does not integrate very well with the gLite software. Altogether, the current solution for running MPI jobs is not complete, nor is it very flexible.

2.3 MPI on the WLCG, desired situation

Based on the experiences with MPI on the WLCG, in Q1 2006, a new working group was formed within the EGEE Technical Coordination Group (EGEE-TCG), specifically focussing on MPI, called EGEE-TCG-MPI [16]. This working group has elicited several requirements for MPI on the WLCG. The requirements are cited in Appendix B and together, they define the desired situation. In the next sections, these requirements will be referred to using codes — please refer to the appendix for their meaning.

In the next subsections, the relationship between those requirements and the LCG will be discussed. Facilitating MPI is neutral in some respects, but not in others; this mostly depends on whether only single-site MPI is facilitated, or users want to have their MPI jobs span multiple sites.

Most of the requirements are met when single-site MPI is deployed throughout the WLCG. To also meet the remaining requirements, cross-site MPI has to be facilitated. We propose to split up the effort in two phases: first enable single-site MPI, then enable cross-site MPI.



2.3.1 Single-site MPI (first phase)

Single-site MPI refers to the ability to submit an MPI job through the grid interface, and have it executed within a single CE (i.e. a single cluster).

2.3.1.1 Implications for the end-user

End-users will be able to execute MPI jobs by submitting them from a UI. MPI is a standardized interface specification, of which a whole breed of different implementations is available (in the dozens [13]). The main difference between these implementations is their prioritization of certain quality attributes; whereas one implementation focusses on sheer performance, others focus on portability and supportability, et cetera. Since MPI merely is an interface specification, code that behaves correctly when linked against, say, OpenMPI, should behave correctly when linked against, for example, MPICH2. Theoretically, anyway. In reality, the differences in implementation *may* lead to an ‘unsafe’ program that behaves differently on one system than it does on another [14]:

“Several outcomes are consistent with the MPI specification and the actual outcome depends on the precise timing of events. In order to complete successfully, an ‘unsafe’ program may require resources that are not guaranteed by the MPI implementation of the system on which it is running. Because of the varying MPI implementations, it is important to be aware of the issues you should consider when writing code that you want to be portable. Some of these issues are:

- **Buffering Assumptions.** In standard mode, blocking sends and receives should not be assumed to be buffered. The reason for this is that buffer memory is finite and all computers will fail under sufficiently large communication loads. If you write a program using buffering assumptions, it will work under some conditions and fail under others. Hence, it is not portable.
- **Barrier Synchronization Assumptions for Collective Calls.** In MPI, collective communications are always assumed to be blocking. However, your program should not depend on whether collective communication routines, like broadcast commands, act as barrier synchronizations. An MPI implementation of collective communications may or may not have the effect of barrier synchronization. One obvious exception to this is the `MPLBARRIER` routine.
- **Communication Ambiguities.** When writing a program, you should make sure that messages are matched by the intended receive call. Ambiguities in the communication specification can lead to incorrect or non-deterministic programs if race conditions arise. Use the message tags and communicators provided by MPI to avoid these types of problems.”

The end-users of the WLCG are advised to write MPI-code that is suitable for (tested with) an MPI-implementation that is (will be) explicitly provided by sites. If requirement `u3` (see `B`) is adequately fulfilled, this issue is unlikely to be an issue for most end-users; a range of different implementations and versions would then be available to them. An extension to MPI was suggested, called *Dyn-MPI*, to improve performance of parallel applications in non-dedicated cluster environments in which resources are continuously competed for [32]; it may be worthwhile to consider.

Requirement `u1+u3+u5` implies that the user shall be responsible for calling `mpicc` on the WNs, as well as calling *the preferred implementation* of `mpicc`. This may be done by writing a Bash script in which the path is prepended with an environment variable `$MPI_<implementation>_PATH/bin`, in which `<implementation>` must refer to one of the CE-provided implementations of MPI (i.e. `$MPI_MPICH_PATH` for MPICH, as advertised by the CE). WNs shall be required to provide such variables if the CE advertises availability of MPI (one variable for each implementation), ref. section 2.3.1.3. The script is then assigned to the `Executable` parameter in the JDL-file, and the required source files are assigned to `InputSandbox`. Example usage:



```
#!/bin/sh
#
# MPIjob.sh
#

(...)

# Set up the paths
PATH=$MPI_MPICH_PATH/bin:$PATH
LD_LIBRARY_PATH=$MPI_MPICH_PATH/lib:$LD_LIBRARY_PATH

mpicc helloworld.c -o helloworld
mpiexec -n 2 ./helloworld
(...)

/* job.jdl */
(...)
Executable = "MPItest.sh";
InputSandbox = "helloworld.c";
(...)
```

Requirement u2+s2 implies that if the user needs a shared directory, he will explicitly refer to that directory by referring to the `$MPI_PATH_TO_SHARED_HOME` environment variable. WNs shall be required to provide this variable if the CE advertises availability of a shared directory, ref. section 2.3.1.3. Example usage:

```
#!/bin/sh
#
# MPIjob.sh
#

(...)
echo '/bin/hostname' >> $MPI_PATH_TO_SHARED_HOME/unsafe-example.txt
(...)
```

2.3.1.2 Implications for RBs

Requirement u7 implies that an end-user is able to pass two MPI-related parameters to the RB, which denote the number of nodes required and the number CPUs to use on each node. As is, JDL specification v0.8 provides a single such parameter called *NodeNumber*, which may be assigned a positive numerical value. **Depending on how a JDL-job is submitted, this number indicates either the number of required ‘nodes’ (submission through Network Server⁴ [20]) or the number of required CPUs (submission through WMPProxy⁵ [21]).** Due to interpretation issues, a *CPUNumber* parameter has been unofficially suggested, but its current status is unknown. Be it as it may, the RB will be required to explicitly process those parameter(s) while matching the job to available CEs. In gLite-3.0, this means that parts of these packages need to be modified: `org.glite.wms.helper`, `org.glite.wms.jdl`, `org.glite.wms.jdlj`, `org.glite.wms-ui.gui-java`.

2.3.1.3 Implications for CEs

The very fact that MPI needs to be facilitated has implications for the Local Resource Management Systems⁶, or *LRMSs*, that are used by a CE. The implications are learned from [22]:

⁴ “The *NodeNumber* attribute is an integer greater than 1 specifying the number of nodes needed for an MPI job.”

⁵ “The *NodeNumber* attribute is an integer greater than 1 specifying the number of CPUs needed for an MPI job.”

⁶ *Local Resource Management Systems* are sometimes referred to as *Job Management Systems*, or *JMSs*, and several other synonyms have been seen in the wild; examples of such systems are OpenPBS, its successor Torque, and LSF.



“Since in MPI, every process is a member of some communicator; when one allows MPI to create or destroy processes, all of the communicators that that process belongs to change. In order to keep collective operations on communicators meaningful (for example, what does a reduction mean when a process joins the reduction during the operation; for that matter, how is ‘during’ defined), all changes to communicators are collective operations.”

Thus, to correctly manage an MPI job (e.g. provide for correct signal delivery, free resources, provide accounting), the LRMS and MPI implementation will need some coupling [23, 24, 25]. This integration may, for example, be done by having the LRMS offer an external interface to job scheduling functions and instruct the MPI library to use this API when spawning new processes. This approach has been demonstrated in the integration of LAM/MPI with OpenPBS [23] through the PSCHED Task Management API [33], which was an attempt to standardize such an external API. While this integration is still available in OpenMPI and Torque (the successors of respectively LAM-MPI and OpenPBS), PSCHED does not appear to be used by other LRMSs [31]. CEs that employ an LRMS that does not explicitly support parallelism cannot service MPI jobs and should therefore *never* advertise themselves as being MPI-capable. This issue is elaborated in section 2.3.3.3.

Furthermore, scheduling a mixture of MPI and non-MPI jobs may pose a problem. It is up to the site administrator to evaluate the capabilities of a site’s LRMS to perform such scheduling. If mixed scheduling can not be provided economically, a site should not deploy MPI, unless the site would be *dedicated* to it (which in turn requires the WMS to somehow exclude it from matching it to non-MPI jobs).

Requirement u2 implies CEs need to explicitly advertise the (un)availability of a shared home directory. There is not yet an accepted way of advertising *availability of a shared home*, but this could be done by having the CE advertise “MPI_HOME_SHARED” (see instructions below). The exact path to the shared directory may not be known to the CE; therefore, as is, the LRMS wraps the job executable inside a preprocessing script that is run at the WNs and interprets such variables (i.e., `cd` to the right directory) and finalizes by initiating the actual job. *Unavailability of a shared home* is currently advertised as “MPI_HOME_NOTSHARED”. This is done by modifying the `GlueHostApplicationSoftwareRunTimeEnvironment` value in the `/opt/lcg/var/gip/ldif/lcg-info-static-cluster.ldif` file (or rather, the `$CE_RUNTIMEENV` variable in your `site-info.def` and run `configure_node` again):

```
(...)
GlueHostApplicationSoftwareRunTimeEnvironment: MPI_HOME_NOTSHARED
(...)
```

The exact strings that may be advertised are *arbitrary*, and as is, coordination of naming conventions is up to the end-users; such responsibility belongs to the User Applications layer, of which the gLite middleware is no part.

Requirement u3+u5 implies CEs need to explicitly advertise the various MPI implementations they provide, i.e. by tagging themselves with {“OPENMPI-1.0.2”, “MPICH-1.2.7”, “MPICH-G2-1.2.7”, ...}. This is done, again, by modifying the `GlueHostApplicationSoftwareRunTimeEnvironment` value in the `/opt/lcg/var/gip/ldif/lcg-info-static-cluster.ldif` file:

```
(...)
GlueHostApplicationSoftwareRunTimeEnvironment: OPENMPI-1.0.2
GlueHostApplicationSoftwareRunTimeEnvironment: MPICH-1.2.7
GlueHostApplicationSoftwareRunTimeEnvironment: MPICH-G2-1.2.7
(...)
```

Requirement u7 implies that an end-user is able to pass two parameters to the LRMS, instructing it to assign a certain number of nodes and the number of CPUs per node for a job. As is, JDL specification v0.8 provides a parameter called *NodeNumber*, which may be assigned a positive numerical value designating that required number of nodes [49]. Due to interpretation issues, a *CPUNumber* parameter has been unofficially suggested, but its current status is unknown. The current JDL specification does not provide for a way to fulfill requirement u7 as cited. Be it as it may, presuming that



the JDL will provide for such parameters, the LRMS will be required to process those parameter(s). In gLite-3.0, this means that parts of these packages need to be modified: `org.glite.ce.blahp`, `org.glite.ce.cream`, `org.glite.cream-api-java`, `org.glite.ce.cream-cli`.

Requirement s1 implies that the administrator of a site should check if his LRMS is able to run MPI jobs, before advertising support for MPI. To our best knowledge, Torque, PBS, LSF and lcgpbs have this ability. We were only able to test with Torque.

2.3.1.4 Implications for WNs

Requirement u1+u3+u5 implies that WNs shall provide working `mpicc` compiler and run-time environments for multiple implementations of both MPIv1 and MPIv2. The paths to these environments shall be available for user scripts by having the WN provide a `$MPI_<implementation>_PATH` environment variable for each MPI-implementation, in which `<implementation>` is the name of the preferred implementation. For example, when OpenMPI-1.0.2 and MPICH-1.2.6 are supported, the WN shall execute:

```
#!/bin/sh
# some WN/bootscript

(...)
export MPI_MPICH_PATH=/opt/mpich-1.2.6
export MPI_OPENMPI_PATH=/opt/openmpi-1.0.2
(...)
```

Requirement u2 implies that WNs shall provide an environment variable called `$MPI_PATH_TO_SHARED_HOME` which contains the path of the shared directory. This is not required if the CE does not advertise the availability of a shared directory.

2.3.2 Cross-site MPI (second phase)

Cross-site MPI refers to the ability of submitting an MPI job through the grid interface, and have it executed *across multiple CEs* (i.e. spanning a single MPI job over multiple clusters). Requirement `u4+s3+s4+s5+s9` states that MPI jobs may indeed span multiple CEs. The implications of such requirements on the WLCG are discussed in the next sections.

2.3.2.1 More implications

Cross-site MPI adds several implications to those discussed for single-site MPI, and is a non-trivial requirement. Figure 2.6 depicts the key problem areas. The problems have been numbered by priority ('common sense'-style).

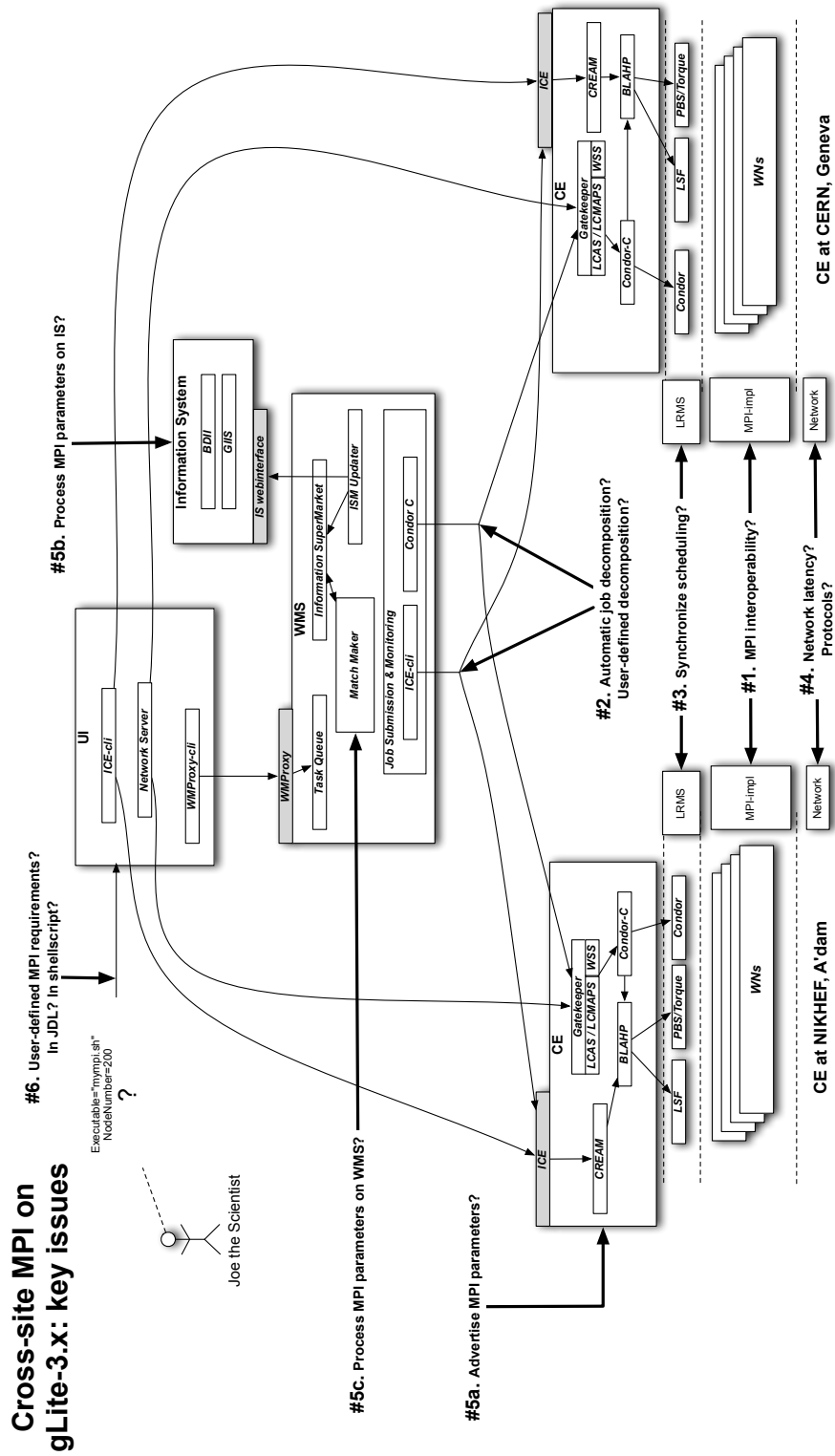


Figure 2.6: Key (problem) areas of cross-site MPI.

Problem #1 MPI interoperability?

Description By default, MPI implementations of different vendors are *not interoperable* (this also goes for the open-source ones). It is not possible to span a job over CEs that do not provide a common MPI implementation, as depicted in more detail in figure 2.7 [35]. An open, low-level protocol called *IMPI* has been suggested to solve the interoperability problem, which has (only) been implemented in LAM/MPI [34], WMPI II, HP MPI [36] and GridMPI [37]. More recently, HeteroMPI was suggested [38].

Solution 1 Require different sites to run the same implementation(s) of MPI.

Solution 2 Require different sites to run an IMPI-supporting implementation of MPI.

Solution 3 Initiate a working group to build yet another implementation of MPI that fulfills the grid-specific needs.

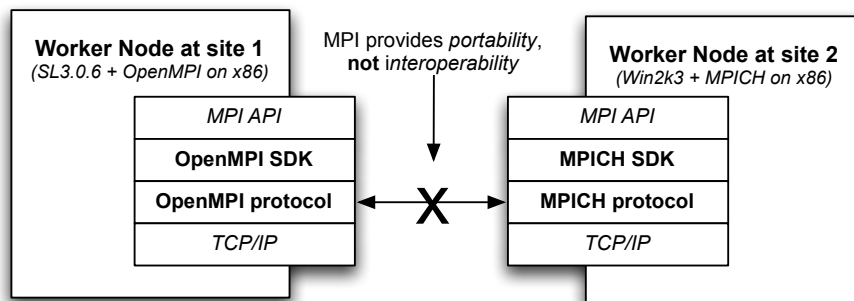


Figure 2.7: Run-time incompatibility between MPI implementations.

Problem #2 Automatic job decomposition? User-defined decomposition?

Description To span a single job over multiple CEs, the problem needs to be decomposed in a way that allows efficient fabric usage and is forgiving to network latency; i.e. independently schedulable jobs. Although it may currently be done for some classes of problems, manually, it is unclear how to do so in general — let alone how to automate it.

Solution Future work; this is area of current academic research. It seems that any solution to this problem will involve both end-users (i.e. MPI programmers) and grid middleware engineers. One relevant source is the MSc-thesis of Sean Philip Peisert in 2000, in which he proposed a programming model for automated decomposition on heterogeneous clusters [39].

Problem #3 Synchronize scheduling?

Description Jobs which cannot be decomposed in a way that completely avoids the need for inter-CE communication, require the LRMSs (or their abstractions) of the participating CEs to provide some form of synchronization and/or reservation mechanisms.

Solution Future work. It yields new requirements to the CE and LRMS components. A solution may come from analyzing the Globus DUROC scheduling component as used by MPICH-G2 [30] and the KOALA grid scheduling system developed at TU Delft [40], which also introduces an economic co-allocation scheduling policy ⁷. Furthermore, the Korean K*Grid Middleware Initiative has forked MPICH-G2 into MPICH-GX to allow co-allocation using the GRASP scheduler of their MoreDream middleware [41]. Lastly, some work has been done on ‘hybrid preemptive scheduling’ of MPI applications [42].

⁷This policy is named *Incremental Claiming Policy* and “optimistically postpones the claiming of the processors for the job to a time close to the estimated job start time”; this has been demonstrated in the DAS-2 system.

**Problem #4** Network latency? Protocols?

Description Latency is an important issue for MPI jobs. In [43], it is shown most MPI benchmarks only scale fine up to 20ms RTT; in real life, however, the latency might be up to hundreds of ms, depending on topological and geographical distance between sites participating in the job. Walking up the stack of the OSI networking model, another problem appears at the transport layer: it has been suggested that TCP is rather unsuitable as a carrier of MPI communication (TCP is designed for contiguous data, whereas MPI is typically non-contiguous).

Solution 1 It has been suggested that some three small modifications to TCP may improve its performance for MPI purposes. This includes the use of a *pacifier* at start-up of an MPI job to prevent TCP from going into ‘Slow Start’ mode [44, 45], and switching TCP parameters at different stages of the MPI execution flow [45]. Applying these modifications may improve performance by 10% to 30% [45, 46];

Solution 2 It has been suggested that SCTP is more suitable for MPI than TCP [47]. Depending on message length and acceptable packet loss, using SCTP instead of TCP may increase MPI communication performance between 250% and 1100%;

Solution 3 Design latency-tolerance into parallel jobs (this is area of current academic research).

Problem #5a Advertise MPI parameters?**Problem #5b** Process MPI parameters on IS?**Problem #5c** Process MPI parameters on WMS?

Description The availability of MPI implementations and several MPI-specific parameters need to be processed throughout several grid components, including the WMS, the IS and the CE. The CE shall need to advertise the MPI implementations it supports and the (un)availability of a shared home directory to the IS. The WMS shall need to accept such parameters into its ISM, and the MatchMaker component shall need to be able to perform matchmaking on those MPI-specific requirements.

Solution See section 2.3.1; this is a topic for the JRA-1 group.

Problem #6 User-defined MPI requirements? In JDL? In shellsript?

Description The end-user shall need to specify the requirements to the MPI environment in the job description or the job executable.

Solution See section 2.3.1 and appendix D.

2.3.2.2 Prior art

Cross-site MPI has already been demonstrated in 2001 using the now-proprietary Legion MPI [27, 28], and in 2005 on TeraGrid using MPICH-G2 [29]. Similar environments are probably out there. MPICH-G2, a ‘grid-enabling’ extension to MPICH [30], in essence allows MPICH to walk the grid stack to the level needed to coordinate an MPI job over multiple sites, as depicted in figure 2.8. Its implementation is tightly coupled to the Globus Toolkit.

We strongly advise the WLCG team to consider MPICH-G2 as a pilot for facilitating cross-site MPI (assuming cross-site MPI indeed needs to be facilitated). MPICH-G2 has already been demonstrated to work with the Globus Toolkit 2.x grid service architecture, and full integration with the future webservices architecture is on its current roadmap [26]:

“COMING SOON: We have a version of MPICH-G2 that is fully integrated with the web services distributed in version v4.x of the Globus Toolkit. We are currently working on that version to complete its integration with the new MPICH-2 library and expect to have our first release in the first quarter of 2006”

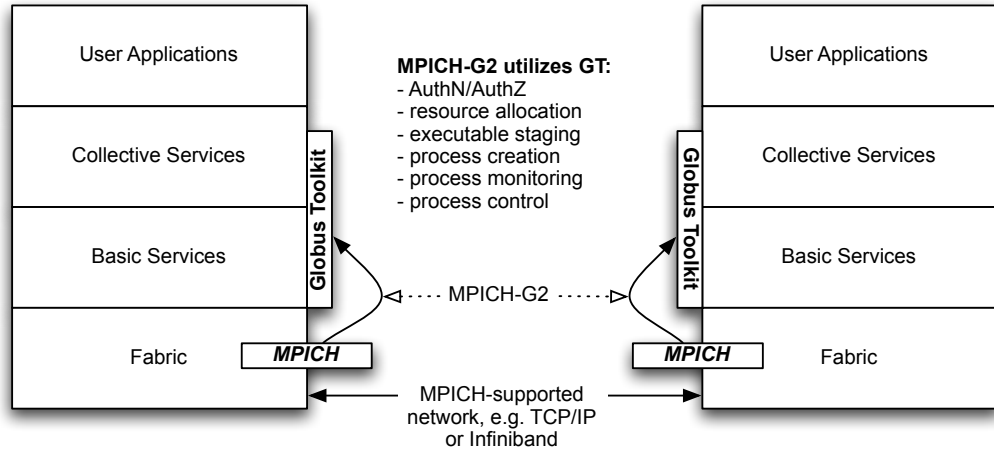


Figure 2.8: Cross-site MPI with MPICH and MPICH-G2

2.3.3 Remarks

Some of the implications discussed in the previous sections are reflected in YAIM code. Please refer to appendix C for a listing of the exact changes. The following section contains a listing of comments on the requirements which have not been (fully) addressed by us.

2.3.3.1 Unmet requirements

Requirement u2

Remark Scalability, reliability and availability of the shared directory are a responsibility of the site, and will not be addressed in our work.

Requirement u5

Remark Commercial/proprietary compilers are beyond the scope of our research due to anticipated mind-numbing software licensing issues.

Requirement s1+s7

Remark Due to the required integration (albeit loose or tight) between LRMS and MPI libraries and a lack of standardization in that respect, it is currently not feasible to fulfill this requirement. Instead, it is chosen to provide support for one major LRMS, namely Torque.

Requirement s6

Remark This requirement is being addressed through Savannah (bug-reporting) tickets and is the responsibility of the EGEE JRA-1 group.

Requirement s8

Remark This requirement is considered a ‘recommendation’ and has been used as such.

2.3.3.2 Efficiency of fabric usage

It is hard to estimate the performance of a message passing job and thus to know *ex ante* if the grid fabric will be unnecessarily stressed. This is actually true for *all* parallel programs, since the datasets they process might contain inefficiencies as well. It is the responsibility of the end-users to provide efficient code (i.e., do not use parallelism when you do not need to); global deployment of MPI will not by itself affect the efficiency of fabric usage on the WLCG.

2.3.3.3 Integration of MPI and LRMS

One of the problems that needs to be solved in the course of enabling single-site or cross-site MPI is fairly fundamental: the MPI run-time environment needs to collaborate with the LRMS environment. This collaboration is important, because it enables accounting of resource use and correct termination of processes: without collaboration, one host is assigned by the LRMS to start the application on the assigned hosts, and only this host is reported in the LRMS statistics. When there is collaboration, the LRMS takes care of starting the application on each of the hosts and therefore all hosts will be reported in the statistics.

Some implementations of MPI do not provide *any* such integration, some provide only the TM interface that comes with Torque/OpenPBS, some support SLURM and/or SGE. Few, if any, implementations support all possible LRMSs. This is probably due to the fact that there is no standardized interface for applications to interact with LRMSs. There has been an effort to provide such a standardized interface, the PSCHED Task Management API (TM) [33], but it has not really caught on: only PBS with its derivatives support it with the TM interface.

If there is no collaboration between the MPI implementation and the LRMS, it is often still possible to submit and run jobs. Most LRMSs provide a list of hosts that are assigned to the job through an environment variable. This list can be used by the user's wrapper script, but there is no guarantee that he will not pick other hosts to use as well. This requires a lot of trust in the user who has to provide a wrapper script to read out this environment variable and start his application. The areas of interest are depicted in figure 2.9.

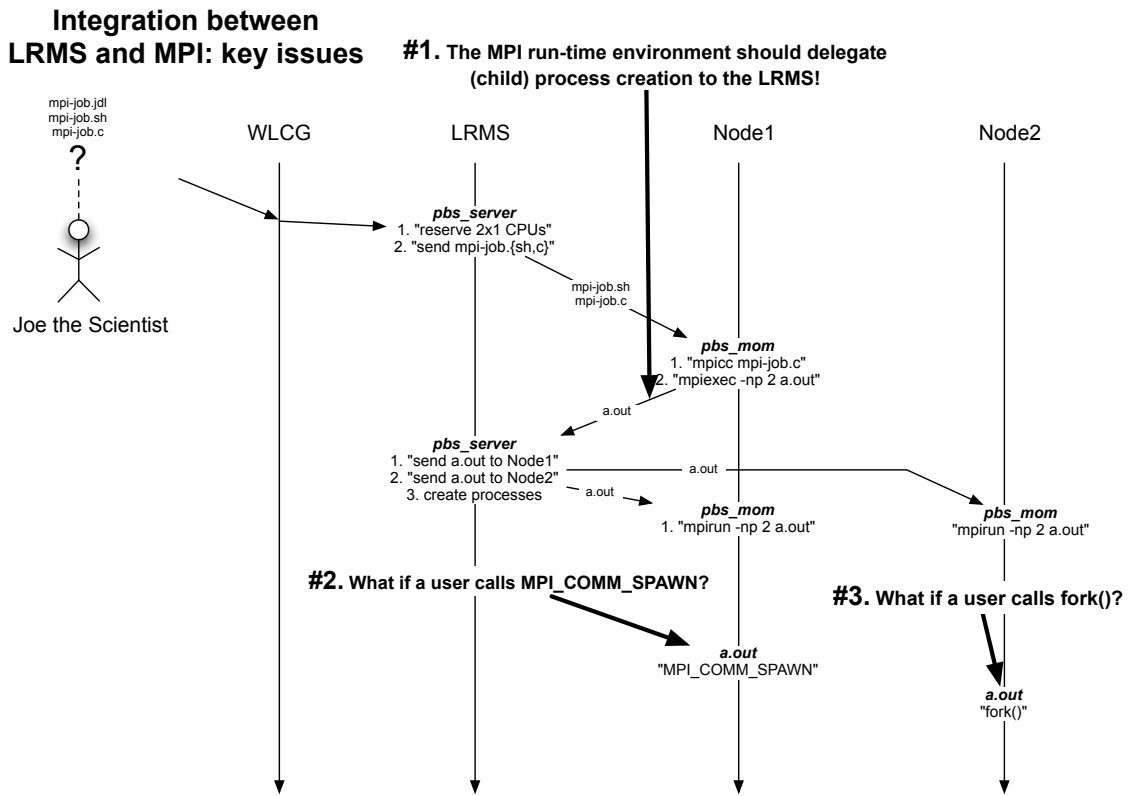


Figure 2.9: Each MPI implementation needs to be integrated with a site's LRMS.

Problem #1 The MPI run-time environment should delegate (child) process creation to the LRMS.

Description Once mpi-job.jdl has been interpreted by the CE gatekeeper and the LRMS has reserved the requested resources, mpi-job.sh is executed on Node 1. This script compiles the MPI source



code (`mpi-job.c`) into object code (`a.out`), which is then ready to be distributed and executed (on Node 1 and 2). But who creates the new processes? As explained, process creation should only be performed by the LRMS, not by plain forking or `ssh`'ing. Each (WLCG-supported) MPI implementation will thus need to support each (WLCG-supported) LRMS implementation. In table 2.1, the currently available (vendor-claimed) integration is depicted.

LRMS-MPI integration						
	Torque/PBS	LSF	SGE	SLURM	XGrid	POE
MPICH	yes	unknown	unknown	unknown	unknown	unknown
MPICH2	yes	unknown	unknown	unknown	unknown	unknown
LAMv6	yes	unknown	unknown	unknown	unknown	unknown
LAMv7	yes	unknown	unknown	unknown	unknown	unknown
OpenMPI	yes	unknown	unknown	yes	yes	yes

Table 2.1: Current state of LRMS-MPI integration.

Solution 1 Execute the object code through a special, implementation-specific MPI wrapper; this is the case for MPICH and MPICH2, which use a specialized `mpirexec`, as well as OpenMPI, which uses the Open Run-Time Environment (`orte`) for such purposes, as well as LAMv6/7, which offer a specialized boot sequence for Torque/PBS;

Solution 2 Compile the source code with an LRMS-aware MPI compiler;

Solution 3 Provide LRMS-integration in the source code (not a realistic option for grid environments).

Problem #2 What if a user calls `MPI_COMM_SPAWN`?

Description From the official MPI-2 standard, we cite [7]:

“The `MPI_COMM_SPAWN` and `MPI_COMM_SPAWN_MULTIPLE` routines provide an interface between MPI and the runtime environment of an MPI application.”

“Complex interaction of an MPI application with its runtime environment should be done through an environment-specific API.”

Solution 1 Kindly ask end-users not to call `MPI_COMM_SPAWN` or `MPI_COMM_SPAWN_MULTIPLE`;

Solution 2 Compile the source code with an LRMS-aware MPI compiler;

Problem #3 What if a user calls `fork()`?

Description Calling `fork()` results in the creation of a new process, by default subverting the LRMS. Evidently, this is the case for *any* application calling `fork()`; not just MPI-applications.

Solution Future work; we consider this beyond our scope.

2.4 Implementations of MPI and their fitness for the WLCG

A large number of implementations of MPI exist. They vary in a number of degrees. They may implement MPI-1, MPI-1 and MPI-2, or MPI-1 and a subset of MPI-2. It appears to be very difficult to implement MPI-2 completely, since only one of the implementations currently available offers full MPI-2 support.

There are different mechanisms for distributing the executable and starting the executable on the assigned machines. Most of the implementations rely on some form of remote shell access, for instance through the `RSH` or `SSH` commands. Most implementations use a daemon framework which can be set up before starting the job and can often remain active so that the next job can be run quickly.

We continue to describe a number of available Open Source implementations of MPI, that may or may not fulfill all the requirements, and describe how well they work in the WLCG.



2.4.1 LAM

LAM, which stands for Local Area Multicomputer, comes standard with Red Hat Enterprise Linux, and is therefore included in Scientific Linux, the Linux distribution created by CERN and Fermilab and other labs and universities to provide one standard base Operating System for various researchers [50]. It implements MPI-1.2 and much of MPI-2.

The project is in maintenance mode, because LAM is one of the contributing implementations to OpenMPI which is discussed in section 2.4.4. OpenMPI will obsolete LAM.

2.4.1.1 LAM-6.5.9

Version 6.5.9 of LAM is shipped with Red Hat ES 3. This version does not support any LRMS and is therefore not the best candidate for use within the WLCG. It offers the `lamboot` command to start up the LAM environment. It requires a list of hosts to which it logs in using SSH by default. It uses this to start a daemon under the current user account on all the hosts that listens on 2 ports, 32870 and 32871 for requests.

If this has succeeded, the `mpirun` command can be issued to start the MPI application on some or all of the hosts that have been set up. The `mpirun` command has the option of uploading (*staging*) the executable to all the required hosts before starting it. It does not have a solution for uploading the data to process.

When the application has finished running, the `lamhalt` command can be issued to stop all the daemons that have been started. For this SSH logins to all hosts are performed.

2.4.1.2 LAM-6.5.9 on the WLCG

Although not the best solution, it is possible to use LAM on the WLCG. LAM is installed easily by running something like:

```
apt-get install lam
```

on a Scientific Linux 3 machine. Running jobs on LAM requires the user to create a wrapper script to boot the LAM environment using the hosts provided by the LRMS. An example script can be found in Appendix E.1. It takes the name of the program to run on the commandline. It boots the LAM environment if necessary, compiles the `.c` file, runs the executable, and breaks down the LAM environment if required.

Setting the correct paths is not much of an issue in this case, because LAM is installed in the default system locations, such as `/usr/bin`, that are already in the path.

The script assumes SSH authentication has been set up between the Worker Nodes.

2.4.1.3 LAM-7.0.6

Version 7.0.6 is provided with Redhat ES 4, so it is also included in Scientific Linux 4. As of release 7, LAM supports TM (the PSCHED API), which makes it an interesting solution because now it integrates neatly with Torque/PBS, the default LRMS in the gLite Grid middleware.

2.4.1.4 LAM-7.0.6 on the WLCG

If a worker nodes runs Red Hat ES 4 or one of it's derivatives, it can simply be installed with the usual means, for instance by running

```
apt-get install lam
```

The integration with Torque/PBS should make it straightforward to start a LAM/MPI job, just run `lamboot` without any arguments. It will detect it is running from within Torque/PBS and will use `tm_spawn` function calls to set up the required daemons. Then an invocation of `mpirun` will use the hosts provided by Torque/PBS.

Setting the correct paths is not much of an issue in this case, because LAM is installed in the default system locations, such as `/usr/bin` that are already in the path.



2.4.1.5 LAM-7.1.2

The latest version of LAM is 7.1.2. On its website [57] the RPM and SRPM (Source RPM) can be downloaded. The SRPM builds nicely on SLC3 (Scientific Linux 3, CERN edition) and can be installed with no trouble. A reminder here is that Grid middleware is not supposed to just override Operating System provided packages unless necessary. Therefore it would be a good idea to modify the spec file⁸ in the SRPM, to build and install into a different directory, for instance in `/opt/lam-7.1.2`. Unfortunately, we had some trouble building this RPM due to some peculiarities in the spec file.

2.4.1.6 LAM-7.1.2 on the WLCG

As mentioned before, take care to not just upgrade components of the Operating System. If one manages to install LAM-7.1.2, the use of it is just as easy as described in section 2.4.1.4. Care must be taken however to determine the correct path. Wrapper scripts can be found in Appendix E.2 and E.3.

2.4.2 MPICH

MPICH is a freely available, portable implementation of MPI-1, created by Argonne National Lab. It uses SSH to connect to remote hosts and launch the requested application. It offers multiple ways to communicate between processes. For a distributed system with TCP/IP communication between the nodes, `ch_p4` should be chosen.

MPICH has no notion of LRMS's, so it cannot cooperate with any of them. To work around this problem, `Mpiexec` was developed and should be used in Torque/PBS environments, see section 2.4.6.

2.4.2.1 MPICH-1.2.6 on the WLCG

MPICH-1.2.6 has been in the LCG middleware stack for quite some time, and is still included in the `gLite 3.0` release. There is rudimentary support for MPI jobs available based on this software. It has been preconfigured to use the `ch_p4` device for communication.

One problem with this version of MPICH as provided in `gLite`, is that it installs in the normal system directories, just as LAM does. Therefore MPICH and LAM cannot co-exist without modifying one or the other. Since LAM is a default system component, it makes sense to modify the MPICH RPM so that it installs in another location. This has already been done for MPICH-1.2.7, see section 2.4.2.2. Therefore we propose to upgrade MPICH from 1.2.6 to 1.2.7, which will solve this problem.

Another problem with the current MPI support in `gLite` is that the invocation of `/usr/bin/mpirun` is hard coded in the jobwrapper template. Use of `mpirun` is not the preferred way of starting MPI applications. This should be fixed by removing the invocation of `mpirun` and instead let the user provide a script that executes `mpiexec`.

2.4.2.2 MPICH-1.2.7 on the WLCG

The latest version of MPICH is 1.2.7p1. An RPM has already been created by Charles Loomis of the Laboratoire de l'Accélérateur Linéaire (LAL) [51] of the Institut national de physique nucléaire et de physique des particules (IN2P3) [52] in France and can be found in this directory [53]. This RPM installs the software in `/opt/mpich-1.2.7`, so it does not clash with a LAM installation. Therefore the path needs to be modified if a user wants to start an MPICH MPI job, by prepending the `$PATH` variable accordingly. A wrapper script can be found in Appendix E.4. A wrapper script using `mpiexec` can be found in Appendix E.5.

2.4.3 MPICH2

From the MPICH site of Argonne National Lab [54]:

⁸A spec file describes how to create an RPM from the provided tarball with the source code



“MPICH2 is an all-new implementation of MPI, designed to support research into high-performance implementations of MPI-1 and MPI-2 functionality. In addition to the features in MPICH, MPICH2 includes support for one-side communication, dynamic processes, intercommunicator collective operations, and expanded MPI-IO functionality. Clusters consisting of both single-processor and SMP nodes are supported. With the exception of users requiring the communication of heterogeneous data, we strongly encourage everyone to consider switching to MPICH2. Researchers interested in using using MPICH as a base for their research into MPI implementations should definitely use MPICH2.”

MPICH2 uses a different mechanism to start processes from that of MPICH, which is called MPD. This mechanism looks a lot like that of LAM. First the `mpdboot` command is used to start an `mpd` daemon on the hosts specified. This can be checked using the `mpdtrace` or `mpdringtest` commands. When the `MPD-ring` is set up, the `mpiexec` command can be given to run the command on all or some of the hosts. When the task is finished, the `mpdallexit` command is used to shut down all the daemons.

2.4.3.1 MPICH2 on the WLCG

The current version of MPICH2 is 1.0.3. It is distributed in `.tar.gz` format, but an RPM has been provided by Charles Loomis and can be found at [53]. This RPM installs itself in the `/opt/mpich2-1.0.3` directory, so it does not conflict with the other MPI implementations.

To be able to run MPICH2 jobs on the WLCG, wrapper scripts are provided in Appendix E.6 and E.7.

2.4.4 OpenMPI

From the OpenMPI website [55]:

“Open MPI is a project combining technologies and resources from several other projects (FT-MPI, LA-MPI, LAM/MPI, and PACX-MPI) in order to build the best MPI library available. A completely new MPI-2 compliant implementation, Open MPI offers advantages for system and software vendors, application developers and computer science researchers.”

The implementations above have migrated into the OpenMPI project, which will make this an important implementation for the future. It offers full MPI-2 standards conformance. One of the things they have already created, is yet another way of starting an MPI job. It is called OpenRTE, which stands for The Open Run-Time Environment [56]. This is a spin-off from the Open MPI initiative. Its overarching goals are to create a portable, high-performance run-time environment for both serial and parallel applications in a wide variety of hardware and software environments, and to transparently support both single and multiple applications operating across a variety of environments.

OpenRTE supports a number of “launchers”, listed below:

- rsh / ssh
- Recent versions of BProc (e.g., Clustermatic)
- PBS Pro, Torque, and Open PBS (the TM system)
- Yod (Red Storm)
- POE
- SLURM
- XGrid



2.4.5 OpenMPI on the WLCG

Since OpenMPI includes OpenRTE, and OpenRTE provides support for Torque, it is an excellent candidate for use on the WLCG. The `orterun` command is used to start an MPI application. It automatically recognizes if it is running under Torque/PBS and then uses the TM interface.

Just as for the other implementations, Charles Loomis has provided RPM's for OpenMPI as well. Unfortunately they are a little outdated, since the RPM's are for version 1.0.1, while 1.1 is the current release. It installs nicely into a Worker Node, because it is installed in `/opt/openmpi-1.0.1`. Sample wrapper scripts can be found in Appendix E.8 and E.9.

2.4.6 Mpiexec

Mpiexec is not an MPI implementation, but a helper program, to ease integration with Torque/PBS. Is created by OSC, the Ohio Supercomputer Center. To quote their website [58]:

“Mpiexec is a replacement program for the script `mpirun`, which is part of the `mpich` package. It is used to initialize a parallel job from within a PBS batch or interactive environment. Mpiexec uses the task manager library of PBS to spawn copies of the executable on the nodes in a PBS allocation.”

The following MPI implementations are supported by Mpiexec:

MPICH/p4 Basic implementation of MPI over ethernet devices.

MPICH/gm MPI using GM or MX message-passing software on Myrinet hardware.

MVAPICH An MPI implementation for InfiniBand networks, based on MPICH.

MPICH/rai An MPI implementation for the IB-like network on the Cray XD1 products.

MPICH/shmem MPI using only the local memory bus on a single-node multi-processor machine.

MPICH2/PMI Rewrite of the venerable MPICH, with a different job launch interface (called PMI) that is intended to be common across all interconnects. Start `mpiexec` with `-comm=pmi` to specify this MPI.

MVAPICH2 An InfiniBand implementation based on MPICH2. Uses the same PMI interface for startup (`-comm=pmi`).

Intel MPI Based on `mpich2` and uses the same PMI interface for startup.

EMP Experimental programmable gigabit ethernet MPI implementation.

none A way to start up many processes which are not related to each other. Handy for startup and shutdown activities or multiple single-processor tasks in a single script.

2.4.7 Mpiexec on the WLCG

Mpiexec is already deployed on all Worker Nodes (WN's), but in a rather old version, 0.77. This version had no support for MPICH2 yet. The current stable version is 0.80 and does have support for MPICH as well as MPICH2 as listed above. An RPM is provided by Charles Loomis that installs in `/opt/mpiexec-0.80/`. This means the path to `mpiexec` has to be provided for the user.



2.5 Modifications for deployment of MPI

We have modified YAIM to provide for configuration of MPI. This modification consisted of three steps:

- adding configuration parameters to `examples/site-info.def`;
- creating a new function `functions/config_mpi`;
- calling it in `scripts/node-info.def`;

In the course of doing this, we also came across an inconsistency in YAIM. The PBS configuration (`/var/spool/pbs/server_priv/nodes`) should contain a list of WNs, with their number of processors as an argument for every WN. In a pre-gLite CE this was done correctly by YAIM, but in a gLite-CE, the number of processors was hardcoded in the `gLiteContainers.xml` file and set to 1. This was fixed by adding the line

```
CE_SMPSIZE:torque-wn.number.processors
```

to `gLite.def`. This file creates mappings between environment variables (from `site-info.def` and Python variables).

With this change, the `NodeNumber` parameter in a JDL file is actually interpreted by the LRMS as the number of CPUs (processors) to use for a job.

The exact changes are listed in appendix C, and have been committed to the gLite CVS repository. Charles Loomis has provided RPM packages for OpenMPI, MPICH and MPICH2, which will be included in the appropriate gLite ‘metapackage’ by the gLite package and release managers. LAM-6 is included in the default Scientific Linux 3.0.6 distribution. Deployment of MPI has the following implications to grid components:

- on WNs, a file will be created at `/etc/profile.d/mpi.sh` — this file provides end-users with site-specific information of MPI (e.g. paths, availability of shared home), and is supposed to be ‘sourced’ in the job they submit;
- on CEs, new attributes will be advertised: one for each deployed implementation of MPI (currently MPICH, MPICH2, OPENMPI, LAM_6 and/or LAM_7), and one stating the (un)availability of a shared home (`MPI_HOME_SHARED` or `MPI_HOME_NOTSHARED`). The latter is required to be advertised if any MPI implementation is advertised, and NOT to be advertised if a site doesn’t provide MPI.

There is no impact on other components.

2.6 Site Functional Test for MPI

Site Functional Tests, or SFTs are part of a framework to test the operation of the LCG grid. They consist of small scripts that each test one or more functionalities. These scripts are wrapped inside one job and submitted to the CE that is to be tested. To read more about this framework, please see [59].

We have created a simple script that tests all the available implementations of MPI that are installed when using our suggestions. These can be found in Appendix F. To use this test, copy the SFT scripts, the wrapper scripts and the file `MPItest.c` from the Appendices or from [1] and add `sft-mpi` to the list of tests to be executed. Remember to add the required extra variables for `NodeNumber` and `JobType` to the JDL file.



2.7 Conclusion on MPI on the WLCG

There is a need to run MPI-parallelized jobs on the WLCG. It seems that support for MPI has been beyond the scope of the EGEE working groups, at least until Q1 2006 when the EGEE-TCG-MPI group was formed. Preliminary support has been available since then - but establishing a working MPI environment still required site administrators to perform a fairly large amount of (poorly documented) customization. Even after doing so, the resulting MPI environment was inflexible in that it contained several hard-coded elements (at the RB) and was not integrated with the LRMSs. Our work provides support for multiple concurrent implementations of MPI, and allows for better integration with LRMSs. End users of the grid are now able to submit MPI jobs for any of these implementations. It is up to the site administrators to determine if their site is able to provide MPI functionality; primarily, this means that they need to evaluate the capacity of the site's LRMS to schedule both MPI and non-MPI jobs. If not, they may choose to dedicate it for MPI, or probably shouldn't expose MPI capacities to the WLCG. Although we only provide single-site MPI at this time, we have described the problems that need to be solved for establishing cross-site MPI, and believe our report provides a small pointer towards the right direction.

2.8 Future work

Due to time constraints, we have not yet been able to determine possible issues related to mixed scheduling of MPI and non-MPI jobs. It is not uncommon to see clusters dedicated for MPI, but this is probably not an acceptable requirement to most WLCG sites, nor for the future. This is an issue for both single-site and cross-site MPI, and suggest it to be researched in the near future. It appears that the *Maui* scheduler plugin for Torque has some capabilities in this respect.

We observed a lack of integration between MPI implementations and various LRMSs. Support for Torque/PBS seems to be available in some common MPI implementations, but integration with other LRMSs may be important to the success of grid-based MPI (that is, assuming that Torque/PBS is not the sole LRMS used in WLCG sites). Again, this is an issue for both single-site and cross-site MPI. A solution may be found through collaboration of the MPI and LRMS communities.

With regards to cross-site MPI, a mechanism for co-allocation of computing resources needs to be defined. It may be sensible to build on knowledge from MPICH-G2, MPICH-GX and KOALA (as mentioned).

Reliability and performance tests should be performed on the work we provided, and the Site Functional Test (SFT) we created should be checked and then used for monitoring purposes.

Lastly, in the context of e-Science, it may be interesting to review support for Cactus Code, a well-known problem solving environment for scientists and engineers.

Chapter 3

YAIM

Yet Another Installation Method, or *YAIM*, is a tool for installation and configuration of grid middleware. It was originally created by Louis Poncet¹, Oliver Keeble and Laurence Field. It started as a straightforward, simple and easy to use tool, to help system administrators configure their machines. Now, it is the default and supported deployment system. Before YAIM was created, system administrators had to work through a 300 page installation and configuration manual. This manual still exists, but now more as a reference guide. The installation and configuration manual is now reduced to some 30 pages.

3.1 Structure

YAIM is entirely written in Bash, a UNIX scripting language, except for a few scripts in Python that are used to create gLite configuration files. The gLite middleware itself contains Python scripts for configuration, which can be controlled and run from YAIM, but not after the configuration files have been created.

YAIM is a very modular system, the components are configuration files, scripts, functions and utilities.

Configuration files are what guide the configuration of the machine.

- The `site-info.def` file looks like a shell script, consisting only of attribute-value assignments. It is not to be regarded as a shell script, but solely as a configuration file. Functionality beyond setting a variable should be done in function files. In principle, one `site-info.def` should suffice for an entire site.
- The `node-info.def` file defines what steps need to be taken to configure a node to take on a certain role. This means that for every role, there is a list of functions that need to be called.

Scripts are the glue of YAIM.

- The `install_node` script manages the installation of grid middleware components. This boils down to installing a meta-package, which is an RPM file named after a grid middleware node type, that has dependencies on all software necessary to install this node type. Node types can be CE, WN, Torque-server, etcetera. It is possible for a node to have multiple types.
- The `configure_node` script is run after the `install_node` script and takes care of configuring the components based on the configuration files listed above. It does this by scanning the command line for arguments that match node types. These node types are then looked up in `node-info.def` and the resulting functions are run.

¹Who happens to be our supervisor at CERN, indeed



Functions are the basic elements of configuration. For every element that needs to be configured, a file is created named after the function, and it only contains this function².

Sometimes a site needs to override the default version of a function. This is possible by creating a file named after the function in the `local` directory below the `functions` directory.

Utilities are small scripts that contain a routine that may be called by various functions.

3.2 Scope

As mentioned above, YAIM originated as a solution to help reduce the burden on system administrators. It was not built as a full featured configuration management system for large sites. Actually such a system has been built, as a result of Work Package 4 of EDG. Its task definition was [60]:

“WP4 Installation system provides a mean to install operating system and applications through the network, to configure system parameters and thereby to apply site policies. It makes available all the necessary applications and user environments on the node. It runs automatically or on request to examine the system and eventually to bring it to the externally defined state stored in the ‘configuration database’.”

The resulting system from WP4 is Quattor [61]. This is in fact in use at a number of the bigger WLCG sites. But the scope of Quattor (and WP4), was installation and configuration of the Operating System. Not grid middleware.

There is a working group, called the Quattor Working Group (QWG), whose aim is to install, configure and maintain the grid middleware using Quattor [62]. This effectively extends the use of Quattor so that YAIM is no longer required³.

3.3 Evaluation

At the time we wrote our research plan, we assumed that CMU’s Software Quality Attributes would be a suitable (albeit somewhat general) framework for evaluating YAIM. After analyzing YAIM, this seemed a bit overdone; the CMU Quality Attributes comprise *performance*, *dependability*, *security* and *safety*, whereas the nature of YAIM — a tool for grid middleware deployment —, suggests only *dependability* to have any relevance; the other attributes are probably more relevant to the *grid middleware* and *middleware configuration*. The goal of YAIM is defined as follows:

“Yaim’s aim is to be able to configure the most common site architectures quickly and simply. The code should be easy to read and should be regarded as unambiguous documentation as much as a set of scripts. Yaim does not need to cover the more obscure requirements of highly specialized sites.”

The so-called *dependability tree*, as depicted in figure 3.1, visualizes the attributes of dependability, the means to achieving it, and the impairments [63, 64, 65]. *A posteriori*, we note the following with regards to the dependability of YAIM:

Not atomic: lots errors or warnings may be produced when running `configure_node`, but still the end result may be: “Configuration complete”.

Unstructured output: when configuring multiple node types or services, it is not obvious from observing the output which part is being configured.

Reconfiguring: when things have to be changed after a system has been configured, one would want to be able to rerun YAIM with a new `site-info.def`, and be sure that everything will be updated accordingly. This is not always the case.

²A helper function may be defined in a function file, but only if it is not possible to make it so generic that it should be a utility

³Or perhaps only for the QWG to do reverse engineering

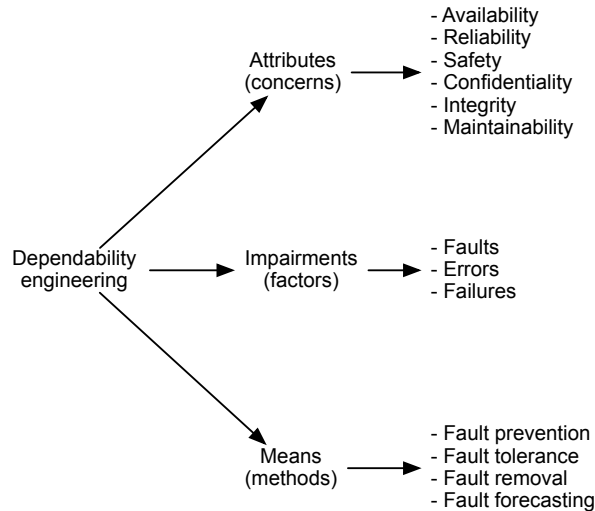


Figure 3.1: The dependability tree

Unconfiguring/removing a service: how to uninstall a service? This needs some scripting to list the requirements of the meta-packages and uninstalling all of them.

Perhaps the worst part is that site administrators may get used to these problems. Two anecdotal stories: We watched someone use YAIM and ignore some error messages, his comment was: “It’s no problem, it’s only GridICE”. Which happened to be a component that was not installed, but YAIM *did* try to configure. Another example was a presentation from someone from the QWG who said that when using YAIM, “some errors you have to ignore, and some you do not.”

We acknowledge the suggestively low-profile goal of YAIM, but would like to point out that excluding it from fulfilling ‘obscure requirements’ does not exclude it from having to be dependable — we believe YAIM should provide the aforementioned fundamentals, *especially* for ‘the most common’ sites who may not have the expertise to deal with deployment peculiarities. That being said, if a YAIM deployment run completes successfully and no rigorous reconfiguration is required, it does exactly what it was built for; deploying grid components in an presumed-default environment.

3.4 Conclusion on YAIM

YAIM is designed to deploy grid components on a default OS environment, and it does a fine job as such. We believe, however, that YAIM lacks several fundamental properties that may be essential to deployment by unexperienced site administrators, or to deployment on non-default environments. Excluding it from fulfilling requirements of ‘highly specialized’ sites does not exclude it from being dependable; it may be worthwhile to elicit possible new requirements that site administrators — the end-users of YAIM, after all — may have, and hear their opinion in this.

3.5 Future work

We suggest that a(nother) round of requirements gathering is performed amongst current and future site administrators, and that its output is used to improve YAIM (if necessary). One of the possible to-do’s we learned from site administrators is to allow them to (automatically) map VO-members that are assigned an ‘ladmin’-role to separate pool accounts, instead of having all members with this role map to the same account — this may aid in auditing and responsibility issues. In more general terms, they want YAIM to be more flexible in user configuration.

Bibliography

- [1] Jong de, Richard and Koot, Matthijs: *Wiki with our scripts and other online information* <https://twiki.cern.ch/twiki/bin/view/LCG/ImplementationOfMpi>
- [2] CERN: European Organization for Nuclear Research, <http://www.cern.ch>
- [3] NIKHEF: Dutch Institute for High Energy and Nuclear Physics, <http://www.nikhef.nl>
- [4] LCG Project Overview, <http://lcg.web.cern.ch/LCG/overview.html>
- [5] Shiers, Sheets of SC4/WLCG Workshop, 2006, <http://indico.cern.ch/materialDisplay.py?materialId=slides&confId=1148>
- [6] Gropp, William and Lusk, Ewing: *"PVM and MPI are completely different"*, 1997, <http://citeseer.ist.psu.edu/573977.html>
- [7] MPI Forum: Message Passing Interface (MPI) Forum homepage, 1998, <http://www.mpi-forum.org/docs/docs.html>
- [8] Fagg, Graham and London, Kevin: *"MPI Inter-connection and Control"*, 1998, <http://citeseer.ist.psu.edu/400213.html>
- [9] Foster, Ian: *"What is the Grid? A Three Point Checklist."*, 2002, <http://www-fp.mcs.anl.gov/~foster/Articles/WhatIsTheGrid.pdf>
- [10] Foster, Ian: *"Service-Oriented Science."*, 2005, <http://www.sciencemag.org/cgi/content/short/308/5723/814>
- [11] Papazoglou, M. and Georgakopoulos, D.: *"Service-oriented computing: Introduction"*, 2003, <http://portal.acm.org/citation.cfm?doid=944217.944233>
- [12] Gropp, William and Lusk, Ewing: *"Goals Guiding Design: PVM and MPI"*, 2002, <http://citeseer.ist.psu.edu/568858.html>
- [13] Trustees of Indiana University: *"MPI Implementation List"*, 2005, <http://www.lam-mpi.org/mpi/implementations/>
- [14] NCSA / PACS Training Group: *"Introduction to MPI"*, 2001, <http://webct.ncsa.uiuc.edu:8900/public/MPI/>
- [15] MPI-Forum: *"MPI: A Message-Passing Interface Standard"*, 1995, <http://www.mpi-forum.org/docs/mpi-11-html/mpi-report.html>
- [16] EGEE TCG: *"Message Passing (MPI) WG"*, <http://egee-intranet.web.cern.ch/egee-intranet/NA1/TCG/wgs/mpi.htm>
- [17] EGEE MPI Working Group: Minutes of the 2006-02-23 meeting, <http://agenda.cern.ch/fullAgenda.php?ida=a061204>
- [18] Johnson, Mike: *"Superscalar Microprocessor Design"*, Prentice Hall, 1991

- [19] IETF (S. Bradner): “*RFC 2119: Key words for use in RFCs to Indicate Requirement Level*”, 1997, <http://www.ietf.org/rfc/rfc2119.txt>
- [20] EGEE JRA-1: “*Job Description Language Attributes Specification (submission through Network Server)*”, 2006, <https://edms.cern.ch/document/555796/1>
- [21] EGEE JRA-1: “*Job Description Language Attributes Specification (submission through WM-Proxy)*”, 2006, <https://edms.cern.ch/document/590869/1>
- [22] Gropp, William and Lusk, Ewing: “*Dynamic Process Management in an MPI Setting*”, 1995, <http://citeseer.ist.psu.edu/gropp95dynamic.html>
- [23] Barret, Brian et al: “*Integration of the LAM/MPI environment and the PBS scheduling system*”, 2003, <http://citeseer.ist.psu.edu/633226.html>
- [24] Saphir, William et al: “*Job Management Requirements for NAS Parallel Systems and Clusters*”, 1995, <http://citeseer.ist.psu.edu/saphir95job.html>
- [25] Sistare, Steve et al: “*An Architecture for Integrated Resource Management of MPI Jobs*”, 2002, <http://doi.ieeecomputersociety.org/10.1109/CLUSTR.2002.1137769>
- [26] Northern Illinois University: “*MPICH-G2*”, 2005, <http://www3.niu.edu/mpi/>
- [27] Humphrey, Marty et al: “*Legion MPI: High Performance in Secure, Cross-Site, Cross-Architecture MPI Applications*”, 2001, <http://legion.virginia.edu/presentations/LegionMPI-NAVO-Jun2001/>
- [28] Natrajan, Anand et al: “*Capacity and Capability Computing using Legion*”, 2001, <http://citeseer.ist.psu.edu/natrajan01capacity.html>
- [29] Dong, Suchuan et al: “*Cross-Site Computations on the TeraGrid*”, 2005, Computing in Science and Engineering 7(5):14 – 23, ISBN: 1521-9615
- [30] Karonis, N. et al: “*MPICH-G2: A Grid-enabled implementation of the Message Passing Interface*”, 2003, Journal of Parallel and Distributed Computing, 63(5):551 – 563
- [31] Open MPI: “*Development Mailing List Archives*”, 2006, <http://www.openmpi.org/community/lists/devel/2006/06/0939.php>
- [32] Weatherly, D. Brent et al: “*Dyn-MPI: Supporting MPI on Non Dedicated Clusters*”, 2003, <http://citeseer.ifi.unizh.ch/weatherly03dynmpi.html>
- [33] The PSCHED API Working Group: “*PSCHED: An API for Parallel Job/Resource Management*”, 1996, <http://www.pbspro.com/docs/psched-api-report.ps>
- [34] Squyres, Jeffrey et al: “*The interoperable message passing interface (IMPI) extensions to LAM/MPI*”, 2000, <http://citeseer.ist.psu.edu/squyres00interoperable.html>
- [35] Skjellum, Anthony et al: “*Interoperability of Message-Passing Interface (MPI) Implementations: A Position Paper*”, 1997, <http://citeseer.ist.psu.edu/68559.html>
- [36] HP: “*HP MPI User’s Guide - Chapter 3*”, 2003, <http://docs.hp.com/en/B6060-96013/ch03s03.html>
- [37] GridMPI: “*GridMPI version 1.1*”, 2006, <http://www.gridmpi.org/gridmpi-1-1/>
- [38] Lastovetsky, Alexey et al: “*HeteroMPI: Towards a message-passing library for heterogeneous networks of computers*”, 2006, Journal of Parallel and Distributed Computing, 66(6):197 – 220
- [39] Peisert, Sean Philip: “*A Programming Model for Automated Decomposition on Heterogeneous Clusters of Multiprocessors*”



- [40] Mohamed, H.H. et al: “*The Design and Implementation of the KOALA Co-Allocating Grid Scheduler*”, 2005, <http://www.st.ewi.tudelft.nl/koala/papers/egc-final00.ps>
- [41] K*Grid Middleware Initiative (KMI): “*MPICH-GX (Grid?eXtension of MPICH)*”, 2005, <http://kmi.moredream.org/doc/MPICH-GX-Manual.pdf>
- [42] Bouteiller, Aurelien et al: “*Hybrid Preemptive Scheduling of MPI Application on the Grids*”, 2004, Proceedings of the Fifth IEEE/ACM International Workshop on Grid Computing (GRID’04), p.130 – 137
- [43] Matsuda, Motohiko et al: “*Evaluation of MPI Implementations on Grid-connected Clusters using an Emulated WAN Environment*”, 2003, Proceedings of the 3st International Symposium on Cluster Computing and the Grid (page 10), ISBN:0-7695-1919-9
- [44] IETF (R. Braden): “*RFC 1122: Requirements for Internet Hosts – Communication Layers*”, 1989, <http://www.ietf.org/rfc/rfc1122.txt>
- [45] Matsuda, Motohiko et al: “*TCP Adaptation for MPI on Long-and-Fat Networks*”, 2005, Proceedings of the IEEE International Conference on Cluster Computing (Cluster 2005)
- [46] Matsuda, Motohiko et al: “*The Design and Implementation of an Asynchronous Communication Mechanism for the MPI Communication Model*”, 2004, Proceedings of the IEEE International Conference on Cluster Computing (Cluster 2004)
- [47] Kamal, Humaira et al: “*SCTP versus TCP for MPI*”, 2005, Proceedings of the ACM/IEEE SC 2005 Conference (SC’05), p.30 – ?
- [48] CERN: “*LCG Computing Grid - Technical Design Report v1.04*”, 2005, http://lcg.web.cern.ch/LCG/tdr/LCG-TDR_v1.04.pdf
- [49] JRA-1: “*Job Description Language - Attributes Specification v0.8*”, 2006, <https://edms.cern.ch/file/555796/1/EGEE-JRA1-TEC-555796-JDL-Attributes-v0-8.pdf>
- [50] Scientific Linux, website, <https://www.scientificlinux.org/>
- [51] Laboratoire de l’Accélérateur Linéaire (LAL), website, <http://www.lal.in2p3.fr/>
- [52] Institut national de physique nucléaire et de physique des particules (IN2P3), website, <http://www.in2p3.fr/>
- [53] Index of MPI Packages, <http://quattor.web.lal.in2p3.fr/packages/mpi/>
- [54] Argonne National Laboratory, website, <http://www.anl.gov/>
- [55] OpenMPI, A High Performance Message Passing Library, website, <http://www.open-mpi.org/>
- [56] The Open Run-Time Environment (ORTE) Project, website, <http://www.open-rte.org/>
- [57] LAM/MPI Parallel Computing, website, <http://www.lam-mpi.org/>
- [58] Ohio Supercomputer Center, website, <http://www.osc.edu/~pw/mpiexec/index.php>
- [59] Site Functional Tests, http://goc.grid.sinica.edu.tw/gocwiki/Site_Functional_Tests
- [60] Polok Jaroslaw, Iven Jan: “*WP4 Task definition: Installation*”, 2001, <http://iven.home.cern.ch/iven/wp4/wp4-inst-task.pdf>
- [61] QUATTOR: QUattor is an Administration ToolkiT for Optimizing Resources, website, <http://quattor.web.cern.ch/quattor/>
- [62] Quattor Working Group, <https://svn.lal.in2p3.fr/LCG/QWG/web/index.html>



- [63] CMU-SEI: “*Technical Report CMU/SEI-95-TR-021 - Quality Attributes*”, 1995, <http://www.sei.cmu.edu/pub/documents/95.reports/pdf/tr021.95.pdf>
- [64] Baracci, Mario et al: “*Principles for Evaluating the Quality Attributes of a Software Architecture*”, 1997, <http://citeseer.ist.psu.edu/barbacci97principles.html>
- [65] IFIP WG10.4: “*WG 10.4 on DEPENDABLE COMPUTING AND FAULT TOLERANCE*”, 1980-now, <http://www.dependability.org/wg10.4/>
- [66] Creative Commons: Creative Commons Attribution 2.5 license, www.creativecommons.org
- [67] Open Source Initiative, The BSD License, www.opensource.org

Appendix A

The BSD license for this project

Copyright (c) 2006, Richard de Jong and Matthijs Koot All rights reserved.

Redistribution and use in source and binary forms, with or without modification, are permitted provided that the following conditions are met:

- Redistributions of source code must retain the above copyright notice, this list of conditions and the following disclaimer.
- Redistributions in binary form must reproduce the above copyright notice, this list of conditions and the following disclaimer in the documentation and/or other materials provided with the distribution.
- Neither the name of the University of Amsterdam nor the names of its contributors may be used to endorse or promote products derived from this software without specific prior written permission.

THIS SOFTWARE IS PROVIDED BY THE COPYRIGHT HOLDERS AND CONTRIBUTORS "AS IS" AND ANY EXPRESS OR IMPLIED WARRANTIES, INCLUDING, BUT NOT LIMITED TO, THE IMPLIED WARRANTIES OF MERCHANTABILITY AND FITNESS FOR A PARTICULAR PURPOSE ARE DISCLAIMED. IN NO EVENT SHALL THE COPYRIGHT OWNER OR CONTRIBUTORS BE LIABLE FOR ANY DIRECT, INDIRECT, INCIDENTAL, SPECIAL, EXEMPLARY, OR CONSEQUENTIAL DAMAGES (INCLUDING, BUT NOT LIMITED TO, PROCUREMENT OF SUBSTITUTE GOODS OR SERVICES; LOSS OF USE, DATA, OR PROFITS; OR BUSINESS INTERRUPTION) HOWEVER CAUSED AND ON ANY THEORY OF LIABILITY, WHETHER IN CONTRACT, STRICT LIABILITY, OR TORT (INCLUDING NEGLIGENCE OR OTHERWISE) ARISING IN ANY WAY OUT OF THE USE OF THIS SOFTWARE, EVEN IF ADVISED OF THE POSSIBILITY OF SUCH DAMAGE.

Appendix B

Requirements

In Q1/2006, a new working group was formed within EGEE Technical Coordination Group (EGEE-TCG), specifically focussing on MPI [16]. In their early meetings, preliminary end-user and (grid) site requirements were elicited. In the next subsections, the explicit requirements are listed and tagged for referencing purposes (i.e. traceability) and prioritized **for our own purposes** using the conventional {MUST,SHOULD,MAY} requirement levels described in RFC 2119 [19]. The prioritization has been refined during a phone conference with Charles Loomis, who is leader of the EGEE working group.

B.1 User requirements

The end-user requirements are as follows [17]:

- u1 MUST *“Users need to prepare the job before the MPI executable is launched. This typically involves compilation of the executable and preparation of the job inputs.”*
- u2 MUST *“Some jobs require a shared directory. During the execution of the job, some of the files in the shared directory are updated. These changes need to be visible on the child nodes.”*
- u3 MUST *“Multiple versions of MPI need to be supported. There is a mix of MPI-1 and MPI-2 jobs which need to be executed. Moreover, some users prefer a particular implementation. We should look at MPICH-1, MPICH-2, and LAM at a minimum.”*
- u4 MAY *“At least the computational chemistry application will need to have cross-site MPI available. This means looking at MPICH-G2 and other similar solutions. The experiences from CrossGrid will be helpful here.”*
- u5 MAY *“Many of the MPI codes are compiled with commercial compilers. We need to see how this can be accommodated on the production service.”*
- u6 SHOULD *“Some grid-enabled applications call grid services from the child processes. This means that these processes need access to the user proxy. This access is usually not possible because the proxy is put into a local space (/tmp). Currently the applications handle the distribution of the proxies directly.”*
- u7 MAY *“Some applications want to have detailed control over the distribution of jobs, specifying both the number of nodes and number of CPUs on each node.”*



B.2 Site requirements

The site requirements are as follows [17]:

- s1 MAY *“Need to ensure that MPI jobs are compatible with all of the job managers (lcgpbs, pbs, etc.).”*
- s2 MUST *“The details of the ‘shared home directory’ requirement needs to be better understood. First, is it necessary that this be the home directory or can another directory be provided? Second, is this required for all jobs; could some jobs survive without? Third, if jobs can survive without, what files need to be distributed automatically to the child nodes? The executable is the minimum, but proxies, etc. may be needed.”*
- s3 MAY *“It looks like most sites supporting MPI on the production service will support a mix of MPI and non-MPI jobs. Efficient scheduling can only be done in this situation if job backfilling is possible. This means that the local scheduler *must* have access to the job requirements. This requirement has come up in other contexts. It is critical for MPI jobs to avoid inefficient use of the fabric or job starvation. For the case of starvation, Fokke [Dijkstra] pointed out that the new ERT algorithm may help avoid some of the bad scheduling decisions.”*
- s4 MAY *“Need to look at efficient scheduling of the MPI jobs at the grid-level as well. For the CrossGrid modified resource broker some additional information was needed, such as the number of CPUs really free, number of CPUs available for MPI, and the like. May also need schema changes for indicating whether a shared home directory is available.”*
- s5 MAY *“CrossGrid also found that there were some problems with the standard LRMS information providers. For example for PBS, it would advertise a CPU as free even if it was loaded and could not be scheduled. Need to look in detail to ensure that the information is correct.”*
- s6 SHOULD *“Need to avoid hardcoding of things at the level of the resource broker. This makes it extremely difficult to customize MPI support at the site. For example, there is a desire for site running Torque/PBS to use mpiexec which permits correct accounting and control of the child processes. The job wrapper should not assume that this is mpirun, nor call it automatically at the start of the job.”*
- s7 SHOULD *“A variety of batch systems need to be supported. At the meeting there were people wanting to have SGE and Torque working with MPI. There are also many sites running LSF.”*
- s8 SHOULD *“Having a standard user environment came up in several of the points above. For example, it might be possible to provide an environment variable which indicates the location of a shared disk (not necessarily the home area). Also when supporting many different MPI implementations, there will be a need to indicate the locations of those implementations. Environment variables may also be a natural way of implementing site customizations, e.g. indicating the executable to start MPI.”*
- s9 MAY *“If cross-site MPI is supported, then the IP requirements for doing so must be determined.”*

Appendix C

Changes to YAIM

C.1 functions/config_mpi

We have added a file to YAIM, namely function/config_mpi:

```
#
# CONFIG_MPI
#
# Description: MPI-related configuration
#
# History:
#   2006-06-21, mkoot: File created.
#   2006-06-23, mkoot: Replaced "rpm -i" with "rpm -Uvh".
#                       Added $MPI_MPIEXEC.
#   2006-06-26, mkoot: Moved if-then clauses from site-info.def to this file;
#                       config_mpi() now changes CE_RUNTIMEENV.
#   2006-06-27, mkoot: Added timestamp to comments header of generated file (mpi.sh).
#   2006-06-28, mkoot: Complied the generated comment header of "mpi.sh" to the YAIM DevGuide
#
# Remarks:
# - if a site wants to allow 'scp' for job distribution, they must MANUALLY configure
#   their ssh_config and sshd_config files to allow it (i.e. "HostbasedAuthentication yes")
#
# TODO:
# - add MPI rpms to gliteWN metapackage
#
#
# The location of the MPI implementation is found by querying if the RPM is installed,
# searching for the mpicc command in the file list, extracting the path from it and
# testing if it exists in real life. You are free to use a fixed value for the
# $MPI_<foobar>_PATH and $MPI_<foobar>_VERSION, however, and arbitrary MPI implementations
# may be added (inform your users on the variable names). REMEMBER that you may want
# to have the CE advertise your custom MPI, and if so, the changes you make here
# need to be reflected in the CEs GlueSoftwareRunTimeEnvironment!
#
find_dir() {
    dir='rpm -ql $1 | awk -F/bin/mpicc '/mpicc$/ {print $1}''
    if [ "x$dir" != "x" ]
    then
        if [ -d $dir -a -f $dir/bin/mpicc ]
        then
            echo $dir
        fi
    fi
}
```



```

        fi
    fi
}

find_version() {
    version='rpm -q $1 | grep -v "is not installed" | awk -F- '{print $2}'
    if [ "x$version" != "x" ]
    then
        echo $version
    fi
}

config_mpi() {

#-----
# WN-related configuration for MPI
#-----

# testing for "WN" will also include "glite-WN" (Good Thing [tm])
if ( echo "${NODE_TYPE_LIST}" | egrep -q 'WN' > /dev/null ); then
    TIMESTAMP='date'
    cat <<EOF > /etc/profile.d/mpi.sh
#
# Script dynamically generated by YAIM function ${FUNCNAME[0]}
#
# MPI.SH (generated at $TIMESTAMP)
#
# Description:
# --
# This script is generated by YAIM when running "configure_node <foobar.def> glite-WN".
# It should be executed as a prelude to a MPI job and sets MPI environment variables,
# for example by having the end-user 'source' it in a shellsript he/she provides
# in the "Executable" parameter of foobar-job.jdl. If you don't provide MPI, this file
# merely acts as a place holder and should be empty (besides this comment header).
#
# WARNING: if you customize this file, please make sure to make a backup before
#          running configure_node again: this file will be overwritten!
#
EOF

    SOME_MPI_IS_SELECTED="false"

#----
# Handle MPICH
#----
if [ "x$WN_ENABLE_MPI_MPICH" == "xyes" ]
then
    SOME_MPI_IS_SELECTED="true"
    if [ "x$WN_ENABLE_MPI_MPICH_PATH" != "x" ]
    then
        echo "MPI_MPICH_PATH='$WN_ENABLE_MPI_MPICH_PATH'>> /etc/profile.d/mpi.sh
        echo "MPI_MPICH_VERSION='$WN_ENABLE_MPI_MPICH_VERSION'>> /etc/profile.d/mpi.sh
    else
        # rpm -Uvh "http://quattor.web.lal.in2p3.fr/packages/mpe/mpich-1.2.7p1-1.sl3.cl.1.i386.rpm"
        echo "MPI_MPICH_PATH='find_dir mpich'>> /etc/profile.d/mpi.sh
        echo "MPI_MPICH_VERSION='find_version mpich'>> /etc/profile.d/mpi.sh
    fi
    echo "export MPI_MPICH_PATH MPI_MPICH_VERSION">> /etc/profile.d/mpi.sh

```



```
fi

#----
# Handle MPICH2
#----
if [ "x$WN_ENABLE_MPI_MPICH2" == "xyes" ]
then
    SOME_MPI_IS_SELECTED="true"
    if [ "x$WN_ENABLE_MPI_MPICH2_PATH" != "x" ]
    then
        echo "MPI_MPICH2_PATH='$WN_ENABLE_MPI_MPICH2_PATH'">> /etc/profile.d/mpi.sh
        echo "MPI_MPICH2_VERSION='$WN_ENABLE_MPI_MPICH2_VERSION'">> /etc/profile.d/mpi.sh
    else
        # rpm -Uvh "http://quattor.web.lal.in2p3.fr/packages/mpi/mpich2-1.0.3-1.sl3.cl.1.i386.rpm"
        echo "MPI_MPICH2_PATH='find_dir mpich2'">> /etc/profile.d/mpi.sh
        echo "MPI_MPICH2_VERSION='find_version mpich2'">> /etc/profile.d/mpi.sh
    fi
    echo "export MPI_MPICH2_PATH MPI_MPICH2_VERSION">> /etc/profile.d/mpi.sh
fi

#----
# Handle OpenMPI
#----
if [ "x$WN_ENABLE_MPI_OPENMPI" == "xyes" ]
then
    SOME_MPI_IS_SELECTED="true"
    if [ "x$WN_ENABLE_MPI_OPENMPI_PATH" != "x" ]
    then
        echo "MPI_OPENMPI_PATH='$WN_ENABLE_MPI_OPENMPI_PATH'">> /etc/profile.d/mpi.sh
        echo "MPI_OPENMPI_VERSION='$WN_ENABLE_MPI_OPENMPI_VERSION'">> /etc/profile.d/mpi.sh
    else
        # rpm -Uvh "http://quattor.web.lal.in2p3.fr/packages/mpi/openmpi-1.0.1-1.sl3.cl.1.i386.rpm"
        echo "MPI_OPENMPI_PATH='find_dir openmpi'">> /etc/profile.d/mpi.sh
        echo "MPI_OPENMPI_VERSION='find_version openmpi'">> /etc/profile.d/mpi.sh
    fi
    echo "export MPI_OPENMPI_PATH MPI_OPENMPI_VERSION">> /etc/profile.d/mpi.sh
fi

#----
# Handle LAMv6
#----
if [ "x$WN_ENABLE_MPI_LAM_6" == "xyes" ]
then
    SOME_MPI_IS_SELECTED="true"
    if [ "x$WN_ENABLE_MPI_LAM_6_PATH" != "x" ]
    then
        echo "MPI_LAM_6_PATH='$WN_ENABLE_MPI_LAM_6_PATH'">> /etc/profile.d/mpi.sh
        echo "MPI_LAM_6_VERSION='$WN_ENABLE_MPI_LAM_6_VERSION'">> /etc/profile.d/mpi.sh
    else
        # TBD: rpm -Uvh SOMETHING?
        # apt-get install lam
        # echo "Note: there isn't a gLite-integrated .rpm for LAM-6 yet so it won't be usable..."
        echo "MPI_LAM_6_PATH='find_dir lam-6*'">> /etc/profile.d/mpi.sh
        echo "MPI_LAM_6_VERSION='find_version lam-6*'">> /etc/profile.d/mpi.sh
    fi
    echo "export MPI_LAM_6_PATH MPI_LAM_6_VERSION">> /etc/profile.d/mpi.sh
fi

#----
```



```

# Handle LAMv7
#----
if [ "x$WN_ENABLE_MPI_LAM_7" == "xyes" ]
then
  SOME_MPI_IS_SELECTED="true"
  if [ "x$WN_ENABLE_MPI_LAM_7_PATH" != "x" ]
  then
    echo "MPI_LAM_7_PATH='$WN_ENABLE_MPI_LAM_7_PATH'>> /etc/profile.d/mpi.sh
    echo "MPI_LAM_7_VERSION='$WN_ENABLE_MPI_LAM_7_VERSION'>> /etc/profile.d/mpi.sh
  else
    # TBD: rpm -Uvh SOMETHING?
    # echo "Note: there isn't a gLite-integrated .rpm for LAM-7 yet so it won't be usable..."
    echo "MPI_LAM_7_PATH='find_dir lam-7*'>> /etc/profile.d/mpi.sh
    echo "MPI_LAM_7_VERSION='find_version lam-7*'>> /etc/profile.d/mpi.sh
  fi
  echo "export MPI_LAM_7_PATH MPI_LAM_7_VERSION">> /etc/profile.d/mpi.sh
fi

#----
# Handle shared home and SSH host-based authentication parameters, but only
# if at least ONE implementation of MPI is being installed; otherwise it's
# just confusing to end-users and WN-sysadmins.
#-----
if [ "x$SOME_MPI_IS_SELECTED" = "true" ]; then
  echo "MPI_PATH_TO_SHARED_HOME='$WN_ENABLE_MPI_SHARED_HOME_PATH'>> /etc/profile.d/mpi.sh
  echo "MPI_SSH_HOST_BASED_AUTH='$WN_ENABLE_MPI_SSH_HOST_BASED_AUTH'>> /etc/profile.d/mpi.sh
  echo "export MPI_PATH_TO_SHARED_HOME MPI_SSH_HOST_BASED_AUTH">> /etc/profile.d/mpi.sh

  if [ "x$WN_ENABLE_MPIEXEC_EXECUTABLE" != "x" ]; then
    echo "MPI_MPIEXEC='$WN_ENABLE_MPIEXEC_EXECUTABLE'>> /etc/profile.d/mpi.sh
    echo "export MPI_MPIEXEC">> /etc/profile.d/mpi.sh
  fi
fi

#-----
# CE-related configuration for MPI
#-----

if ( echo "${NODE_TYPE_LIST}" | egrep -q 'CE' > /dev/null ); then

  SOME_MPI_IS_SELECTED="false"

  if [ "x$WN_ENABLE_MPI_MPICH" = "xyes" ]; then
    CE_RUNTIMEENV="${CE_RUNTIMEENV}
      MPICH"
    SOME_MPI_IS_SELECTED="true"
  fi

  if [ "x$WN_ENABLE_MPI_MPICH2" = "xyes" ]; then
    CE_RUNTIMEENV="${CE_RUNTIMEENV}
      MPICH2"
    SOME_MPI_IS_SELECTED="true"
  fi

  if [ "x$WN_ENABLE_MPI_OPENMPI" = "xyes" ]; then
    CE_RUNTIMEENV="${CE_RUNTIMEENV}
      OPENMPI"
  fi

```



```

    SOME_MPI_IS_SELECTED="true"
  fi

  if [ "x$WN_ENABLE_MPI_LAM_6" = "xyes" ]; then
    CE_RUNTIMEENV="{CE_RUNTIMEENV}
                  LAM_6"
    SOME_MPI_IS_SELECTED="true"
  fi

  if [ "x$WN_ENABLE_MPI_LAM_7" = "xyes" ]; then
    CE_RUNTIMEENV="{CE_RUNTIMEENV}
                  LAM_7"
    SOME_MPI_IS_SELECTED="true"
  fi

  #-----
  # Only advertise (un)availability of shared home if at least ONE implementation
  # of MPI is being installed; otherwise it's just confusing to end-users.
  #-----
  if [ "x$SOME_MPI_IS_SELECTED" = "true" ]; then
    if [ "x$WN_ENABLE_MPI_SHARED_HOME" = "xyes" ]; then
      if [ "x$WN_ENABLE_MPI_SHARED_HOME_PATH" = "x" ]; then
        echo "Error: WN_ENABLE_MPI_SHARED_HOME is set to 'yes' in your site-info.def file,"
        echo "but no WN_ENABLE_MPI_SHARED_HOME_PATH was given!"
        exit 1
      else
        CE_RUNTIMEENV="{CE_RUNTIMEENV}
                      MPI_HOME_SHARED"

        fi
      else
        CE_RUNTIMEENV="{CE_RUNTIMEENV}
                      MPI_HOME_NOTSHARED"

        fi
    fi
  fi

  return 0
}

```

C.2 examples/site-info.def

We have appended several MPI-related variables to the site-info.def file:

```

#-----
# MPI-related configuration:
#-----
# Several MPI implementations are available in the package repository.
# If you do NOT want an implementation to be installed, set its variable
# to "no". Else, set it to "yes" (default). If you want to use an
# already installed version of an implementation, set its "_PATH" and
# "_VERSION" variables to match your setup (examples below).
#
# NOTE 1: the CE_RUNTIMEENV will be automatically updated in the file
# functions/config_mpi, so that the CE advertises the MPI implementations
# you choose here - you do NOT have to change it manually in this file.
# It will become something like this:
#
# CE_RUNTIMEENV="$CE_RUNTIMEENV

```



```
#           MPICH
#           MPICH2
#           OPENMPI
#           LAM_6
#           LAM_7"
#
# NOTE 2: the above requires yaim-3.0.0-20 or higher to work; otherwise you
# need to change CE_RUNTIMEENV manually for glite-CE nodes.
#
# NOTE 3: it is currently NOT possible to configure multiple concurrent
# versions of the same implementations (e.g. MPICH-1.2.3 and MPICH-1.2.7)
# using YAIM. Customize "/opt/glite/yaim/functions/config_mpi" file
# to do so.

WN_ENABLE_MPI_MPICH="yes"
WN_ENABLE_MPI_MPICH2="yes"
WN_ENABLE_MPI_OPENMPI="yes"
WN_ENABLE_MPI_LAM_6="yes"
WN_ENABLE_MPI_LAM_7="no"

#---
# Example for using an already installed version of MPI.
# Setting "_PATH" and "_VERSION" variables will prevent YAIM
# from downloading and installing the glite-provided packages.
# Just fill in the path to its current installation (e.g. "/usr")
# and which version it is (e.g. "6.5.9").
#---
#WN_ENABLE_MPI_LAM_6_PATH="/usr"
#WN_ENABLE_MPI_LAM_6_VERSION="6.5.9"

# If you do NOT provide a shared home, set $WN_ENABLE_MPI_SHARED_HOME to "no" (default).
# Else, set it to "yes" and assign its path to $WN_ENABLE_MPI_SHARED_HOME_PATH="/shared").
#
WN_ENABLE_MPI_SHARED_HOME="no"
WN_ENABLE_MPI_SHARED_HOME_PATH=""

#
# If you do NOT have SSH Hostbased Authentication between your WNs, set the below
# variable to "no" (default). Else, set it to "yes".
#
WN_ENABLE_MPI_SSH_HOST_BASED_AUTH="no"

#
# If you provide a 'mpiexec' for MPICH, please state the full path to
# that file here (http://www.osc.edu/~pw/mpiexec/index.php). Else, leave empty.
#
#WN_ENABLE_MPIEXEC_EXECUTABLE="/usr/bin/mpiexec"
WN_ENABLE_MPIEXEC_EXECUTABLE="/usr/bin/mpiexec"
```

C.3 scripts/node-info.def

We have modified scripts/node-info.def so that the functions/config_mpi script is called during WN and CE deployment.

Appendix D

A Quick Guide to using MPI on the WLCG

The following sections contain listings of two files, together comprising an example MPI job. Use the `mpich-wrapper.sh` from Appendix E.5 or E.4.

The expected output is:

```
Hello world! from processor 3 out of 6
Hello world! from processor 2 out of 6
Hello world! from processor 4 out of 6
Hello world! from processor 5 out of 6
Hello world! from processor 1 out of 6
Hello world! from processor 0 out of 6
```

D.1 `mpi-example.jdl`

```
Type = "Job";
JobType = "MPICH";
NodeNumber = 6;
Executable = "mpich-wrapper.sh";
Arguments = "MPItest";
StdOutput = "test.out";
StdError = "test.err";
InputSandbox = {"mpich-wrapper.sh", "MPItest.c"};
OutputSandbox = {"test.err", "test.out", "mpiexec.out"};
Requirements = member("MPICH", other.GlueHostApplicationSoftwareRunTimeEnvironment);
```

D.2 `MPItest.c`

```
#include "mpi.h"
#include <stdio.h>
int main(int argc, char *argv[])
{
    int numprocs; /* Number of processors */
    int procnum; /* Processor number */
    /* Initialize MPI */
    MPI_Init(&argc, &argv);
    /* Find this processor number */
    MPI_Comm_rank(MPI_COMM_WORLD, &procnum);
    /* Find the number of processors */
    MPI_Comm_size(MPI_COMM_WORLD, &numprocs);
```




```
printf ("Hello world! from processor %d out of %d\n", procnum, numprocs);
/* Shut down MPI */
MPI_Finalize();
return 0;
}
```

Appendix E

Wrapper scripts

E.1 lam-6-wrapper.sh without Torque/PBS

```
# Script to start an MPI job under lam-6
# Richard de Jong @ CERN 2006-06-29

echo "Using LAM-6"

if [ -z $MPI_LAM_6_VERSION ]
then
    echo "LAM-6, no such version available"
    exit
fi

PATH=$MPI_LAM_6_PATH/bin:$PATH
LD_LIBRARY_PATH=$MPI_LAM_6_PATH/lib:$LD_LIBRARY_PATH
SHUTDOWN="yes"
EXE=$1

if [ "x${PBS_NODEFILE}" != "x" ] ; then
    echo "PBS Nodefile: ${PBS_NODEFILE}"
    HOST_NODEFILE=${PBS_NODEFILE}
fi

if [ "x${LSB_HOSTS}" != "x" ] ; then
    echo "LSF Hosts: ${LSB_HOSTS}"
    HOST_NODEFILE='pwd'/lsf_nodefile.$$
    for host in ${LSB_HOSTS}
    do
        echo ${host} >> ${HOST_NODEFILE}
    done
fi

if [ "x${HOST_NODEFILE}" = "x" ] ; then
    echo "No hosts file defined. Exiting..."
    exit
fi

echo "*****"
HOST_NEEDED='cat ${HOST_NODEFILE} | wc -l | awk '{print $NF}''
echo "Node count: $HOST_NEEDED"
```



```
echo "Hosts in $HOST_NODEFILE:"
cat $HOST_NODEFILE
echo "*****"

echo "compiling program ${EXE}"
mpicc -o $EXE $EXE.c
if [ ! $? -eq 0 ]; then
    echo "Error compiling program. Exiting..."
    exit 1
fi

if [ "x${MPI_PATH_TO_SHARED_HOME}" != "x" ]
then
    if [ ! -d ${MPI_PATH_TO_SHARED_HOME} ]
    then
        echo "Shared home is set, but it is not accessible"
        exit 1
    else
        echo "Shared home exists and is ${MPI_PATH_TO_SHARED_HOME}"
    fi
else
    echo "No shared home is accessible"
    if [ "x${MPI_SSH_HOST_BASED_AUTH}" == "xyes" ]
    then
        echo "Using SSH to copy the files to all nodes"

        for i in `cat $HOST_NODEFILE | grep -v ${HOSTNAME} | sort -u`; do
            echo "Copying files to $i"
            ssh ${i} mkdir -p 'pwd'
            scp -rp . "${i}:'pwd'"
            # Copy the proxy certificate
            scp -p ${X509_USER_PROXY} ${i}:${X509_USER_PROXY}
        done
    else
        echo "Uh oh, cannot use SSH to copy the files to all WN's"
        exit 1
    fi
fi

# Check if there is already a running LAM environment set up:
if ( lamnodes > /dev/null 2>&1 )
then
    echo "lam is set up already"
else
    echo "lam not ready yet, booting lam environment using hosts from ${HOST_NODEFILE}"
    if ( lamboot ${HOST_NODEFILE} )
    then
        echo "lam booted succesfully"
    else
        echo "Error booting lam. Exiting..."
        exit 1
    fi
fi

# In case the lam environment had already been started, make sure we only use the hosts
```



```

# that are assigned to us by the LRMS. This provides at least a lightweight coupling.
AVAILABLE_NODES='lamnodes | grep -f ${HOST_NODEFILE} |awk '{print $1}' | tr '\n' ' ' '
echo "available nodes: ${AVAILABLE_NODES}"

MY_NODE_NUMBER='lamnodes h -c|awk '{print $1}' '
if [ "x${MY_NODE_NUMBER}" == "x" ]
then
    echo "Localhost is not one of the running lam nodes."
    echo "Can not continue because only this host has the executable.  Exiting..."
    exit 1
fi

# Run the command on the available nodes, using -s helps to get the executable to all hosts
mpirun -s ${MY_NODE_NUMBER} ${AVAILABLE_NODES} ${EXE}

#close down the lam environment?
if [ "${SHUTDOWN}" == "yes" ]
then
    lamhalt
fi

```

E.2 lam-7-wrapper.sh without Torque/PBS

```

# Script to start an MPI job under lam-7
# Richard de Jong @ CERN 2006-06-29

echo "Using LAM-7"

if [ -z $MPI_LAM_7_VERSION ]
then
    echo "LAM-7, no such version available"
    exit
fi

PATH=$MPI_LAM_7_PATH/bin:$PATH
LD_LIBRARY_PATH=$MPI_LAM_7_PATH/lib:$LD_LIBRARY_PATH
SHUTDOWN="yes"
EXE=$1

if [ "x${PBS_NODEFILE}" != "x" ] ; then
    echo "PBS Nodefile: ${PBS_NODEFILE}"
    HOST_NODEFILE=${PBS_NODEFILE}
fi

if [ "x${LSB_HOSTS}" != "x" ] ; then
    echo "LSF Hosts: ${LSB_HOSTS}"
    HOST_NODEFILE='pwd'/lsf_nodefile.$$
    for host in ${LSB_HOSTS}
    do
        echo ${host} >> ${HOST_NODEFILE}
    done
fi

if [ "x${HOST_NODEFILE}" = "x" ]; then
    echo "No hosts file defined. Exiting..."

```



```

        exit
    fi

    echo "*****"
    HOST_NEEDED='cat ${HOST_NODEFILE} | wc -l | awk '{print $NF}'
    echo "Node count: $HOST_NEEDED"
    echo "Hosts in $HOST_NODEFILE:"
    cat $HOST_NODEFILE
    echo "*****"

    echo "compiling program ${EXE}"
    mpicc -o $EXE $EXE.c
    if [ ! $? -eq 0 ]; then
        echo "Error compiling program. Exiting..."
        exit 1
    fi

    if [ "x${MPI_PATH_TO_SHARED_HOME}" != "x" ]
    then
        if [ ! -d ${MPI_PATH_TO_SHARED_HOME} ]
        then
            echo "Shared home is set, but it is not accessible"
            exit 1
        else
            echo "Shared home exists and is ${MPI_PATH_TO_SHARED_HOME}"
        fi
    else
        echo "No shared home is accessible"
        if [ "x${MPI_SSH_HOST_BASED_AUTH}" == "xyes" ]
        then
            echo "Using SSH to copy the files to all nodes"

            for i in `cat $HOST_NODEFILE | grep -v ${HOSTNAME} | sort -u`; do
                echo "Copying files to $i"
                ssh ${i} mkdir -p 'pwd'
                scp -rp . "${i}:'pwd'"
                # Copy the proxy certificate
                scp -p ${X509_USER_PROXY} ${i}:${X509_USER_PROXY}
            done
        else
            echo "Uh oh, cannot use SSH to copy the files to all WN's"
            exit 1
        fi
    fi

    # Check if there is already a running LAM environment set up:
    if ( lamnodes > /dev/null 2>&1 )
    then
        echo "lam is set up already"
    else
        echo "lam not ready yet, booting lam environment using hosts from ${HOST_NODEFILE}"
        if ( lamboot -ssi boot rsh -ssi boot_rsh_agent ssh ${HOST_NODEFILE} )
        then
            echo "lam booted succesfully"
        else

```



```

        echo "Error booting lam. Exiting..."
        exit 1
    fi
fi

# In case the lam environment had already been started, make sure we only use the hosts
# that are assigned to us by the LRMS. This provides at least a lightweigt coupling.
AVAILABLE_NODES='lamnodes | grep -f ${HOST_NODEFILE} |awk '{print $1}' | tr '\n' ' ' '
echo "available nodes: ${AVAILABLE_NODES}"

MY_NODE_NUMBER='lamnodes h -c|awk '{print $1}''
if [ "x${MY_NODE_NUMBER}" == "x" ]
then
    echo "Localhost is not one of the running lam nodes."
    echo "Can not continue because only this host has the executable.  Exiting..."
    exit 1
fi

echo "Executing application"
mpirun -v -np $HOST_NEEDED ${EXE}

#close down the lam environment?
if [ "${SHUTDOWN}" == "yes" ]
then
    echo "Shutting down the LAM environment"
    lamhalt
fi

```

E.3 lam-7-pbs-wrapper.sh with Torque/PBS

```

# Script to start an MPI job under lam-7
# Richard de Jong @ CERN 2006-06-29

echo "Using LAM-7"

if [ -z $MPI_LAM_7_VERSION ]
then
    echo "LAM-7, no such version available"
    exit
fi

PATH=$MPI_LAM_7_PATH/bin:$PATH
LD_LIBRARY_PATH=$MPI_LAM_7_PATH/lib:$LD_LIBRARY_PATH
SHUTDOWN="yes"
EXE=$1

if [ "x${PBS_NODEFILE}" != "x" ] ; then
    echo "PBS Nodefile: ${PBS_NODEFILE}"
    HOST_NODEFILE=${PBS_NODEFILE}
fi

if [ "x${LSB_HOSTS}" != "x" ] ; then
    echo "LSF Hosts: ${LSB_HOSTS}"
    HOST_NODEFILE='pwd`/lsf_nodefile.$$
    for host in ${LSB_HOSTS}

```



```

do
    echo ${host} >> ${HOST_NODEFILE}
done
fi

if [ "x${HOST_NODEFILE}" = "x" ]; then
    echo "No hosts file defined. Exiting..."
    exit
fi

echo "*****"
HOST_NEEDED=$(cat ${HOST_NODEFILE} | wc -l | awk '{print $NF}')
echo "Node count: $HOST_NEEDED"
echo "Hosts in $HOST_NODEFILE:"
cat $HOST_NODEFILE
echo "*****"

echo "compiling program ${EXE}"
mpicc -o $EXE $EXE.c
if [ ! $? -eq 0 ]; then
    echo "Error compiling program. Exiting..."
    exit 1
fi

if [ "x${MPI_PATH_TO_SHARED_HOME}" != "x" ]
then
    if [ ! -d ${MPI_PATH_TO_SHARED_HOME} ]
    then
        echo "Shared home is set, but it is not accessible"
        exit 1
    else
        echo "Shared home exists and is ${MPI_PATH_TO_SHARED_HOME}"
    fi
else
    echo "No shared home is accessible"
    if [ "x${MPI_SSH_HOST_BASED_AUTH}" == "xyes" ]
    then
        echo "Using SSH to copy the files to all nodes"

        for i in $(cat $HOST_NODEFILE | grep -v ${HOSTNAME} | sort -u); do
            echo "Copying files to $i"
            ssh ${i} mkdir -p 'pwd'
            scp -rp . "${i}:'pwd'"
            # Copy the proxy certificate
            scp -p ${X509_USER_PROXY} ${i}:${X509_USER_PROXY}
        done
    else
        echo "Uh oh, cannot use SSH to copy the files to all WN's"
        exit 1
    fi
fi

# Check if there is already a running LAM environment set up:
if ( lamnodes > /dev/null 2>&1 )
then

```



```

    echo "lam is set up already"
else
    echo "lam not ready yet, booting lam environment using hosts from ${HOST_NODEFILE}"
    if ( lamboot -ssi boot tm )
    then
        echo "lam booted succesfully"
    else
        echo "Error booting lam. Exiting..."
        exit 1
    fi
fi

# In case the lam environment had already been started, make sure we only use the hosts
# that are assigned to us by the LRMS. This provides at least a lightweigt coupling.
AVAILABLE_NODES='lamnodes | grep -f ${HOST_NODEFILE} |awk '{print $1}' | tr '\n' ' '
echo "available nodes: ${AVAILABLE_NODES}"

MY_NODE_NUMBER='lamnodes h -c|awk '{print $1}'
if [ "x${MY_NODE_NUMBER}" == "x" ]
then
    echo "Localhost is not one of the running lam nodes."
    echo "Can not continue because only this host has the executable.  Exiting..."
    exit 1
fi

# Run the command on the available nodes, using -s helps to get the executable to all hosts
echo "Executing application"
mpirun -s ${MY_NODE_NUMBER} ${AVAILABLE_NODES} ${EXE}

#close down the lam environment?
if [ "${SHUTDOWN}" == "yes" ]
then
    echo "Shutting down the LAM environment"
    lamhalt
fi

```

E.4 mpich-wrapper.sh without Torque/PBS

```

# Script to start an MPI job under MPICH
# Richard de Jong @ CERN 2006-06-29

echo "Using MPICH"

if [ -z $MPI_MPICH_VERSION ]
then
    echo "MPICH, no such version available"
    exit
fi

# Set up the paths
PATH=$MPI_MPICH_PATH/bin:$PATH
LD_LIBRARY_PATH=$MPI_MPICH_PATH/lib:$LD_LIBRARY_PATH

# The binary to execute
EXE=$1

echo "*****"

```



```

echo "Running on: $HOSTNAME"
echo "As:      " 'whoami'
echo "*****"

echo "*****"
echo "Compiling binary: $EXE"
echo mpicc -o ${EXE} ${EXE}.c
mpicc -o ${EXE} ${EXE}.c
echo "*****"

if [ "x$PBS_NODEFILE" != "x" ] ; then
    echo "PBS Nodefile: $PBS_NODEFILE"
    HOST_NODEFILE=$PBS_NODEFILE
fi

if [ "x$LSB_HOSTS" != "x" ] ; then
    echo "LSF Hosts: $LSB_HOSTS"
    HOST_NODEFILE='pwd'/lsf_nodefile.$$
    for host in ${LSB_HOSTS}
    do
        echo $host >> ${HOST_NODEFILE}
    done
fi

if [ "x${HOST_NODEFILE}" = "x" ]; then
    echo "No hosts file defined.  Exiting..."
    exit
fi

echo "*****"
HOST_NEEDED='cat ${HOST_NODEFILE} | wc -l | awk '{print $NF}''
echo "Node count: $HOST_NEEDED"
echo "Hosts in $HOST_NODEFILE:"
cat $HOST_NODEFILE
echo "*****"

if [ "x${MPI_PATH_TO_SHARED_HOME}" != "x" ]
then
    if [ ! -d ${MPI_PATH_TO_SHARED_HOME} ]
    then
        echo "Shared home is set, but it is not accessible"
        exit 1
    else
        echo "Shared home exists and is ${MPI_PATH_TO_SHARED_HOME}"
    fi
else
    echo "No shared home is accessible"
    if [ "x${MPI_SSH_HOST_BASED_AUTH}" == "xyes" ]
    then
        echo "Using SSH to copy the files to all nodes"

        for i in `cat $HOST_NODEFILE | grep -v ${HOSTNAME} | sort -u`; do
            echo "Copying files to $i"
            ssh ${i} mkdir -p 'pwd'
            scp -rp . "${i}:'pwd'"
            # Copy the proxy certificate
            scp -p ${X509_USER_PROXY} ${i}:${X509_USER_PROXY}
        done
    fi
fi

```



```

else
    echo "Uh oh, cannot use SSH to copy the files to all WN's"
    exit 1
fi
fi

echo "*****"
echo "Executing $EXE with mpirun"
chmod 755 $EXE
mpirun -np $HOST_NEEDED -machinefile $HOST_NODEFILE 'pwd'/$EXE
echo "*****"

```

E.5 mpich-pbs-wrapper.sh with Torque/PBS

```

# Script to start an MPI job under MPICH
# Richard de Jong @ CERN 2006-06-29

echo "Using MPICH"

if [ -z $MPI_MPICH_VERSION ]
then
    echo "MPICH, no such version available"
    exit
fi

# Set up the paths
PATH=$MPI_MPICH_PATH/bin:$PATH
LD_LIBRARY_PATH=$MPI_MPICH_PATH/lib:$LD_LIBRARY_PATH

# The binary to execute
EXE=$1

echo "*****"
echo "Running on: $HOSTNAME"
echo "As:      " 'whoami'
echo "*****"

echo "*****"
echo "Compiling binary: $EXE"
echo mpicc -o ${EXE} ${EXE}.c
mpicc -o ${EXE} ${EXE}.c
echo "*****"

if [ "x$PBS_NODEFILE" != "x" ] ; then
    echo "PBS Nodefile: $PBS_NODEFILE"
    HOST_NODEFILE=$PBS_NODEFILE
fi

if [ "x$LSB_HOSTS" != "x" ] ; then
    echo "LSF Hosts: $LSB_HOSTS"
    HOST_NODEFILE='pwd'/lsf_nodefile.$$
    for host in ${LSB_HOSTS}
    do
        echo $host >> ${HOST_NODEFILE}
    done
fi

```



```

if [ "x${HOST_NODEFILE}" = "x" ]; then
    echo "No hosts file defined.  Exiting..."
    exit
fi

echo "*****"
HOST_NEEDED=$(cat ${HOST_NODEFILE} | wc -l | awk '{print $NF}')
```

Node count: \$HOST_NEEDED
 Hosts in \$HOST_NODEFILE:
 cat \$HOST_NODEFILE

```

echo "*****"

if [ "x${MPI_PATH_TO_SHARED_HOME}" != "x" ]
then
    if [ ! -d ${MPI_PATH_TO_SHARED_HOME} ]
    then
        echo "Shared home is set, but it is not accessible"
        exit 1
    else
        echo "Shared home exists and is ${MPI_PATH_TO_SHARED_HOME}"
    fi
else
    echo "No shared home is accessible"
    if [ "x${MPI_SSH_HOST_BASED_AUTH}" == "xyes" ]
    then
        echo "Using SSH to copy the files to all nodes"

        for i in $(cat $HOST_NODEFILE | grep -v ${HOSTNAME} | sort -u); do
            echo "Copying files to $i"
            ssh ${i} mkdir -p 'pwd'
            scp -rp . "${i}:'pwd'"
            # Copy the proxy certificate
            scp -p ${X509_USER_PROXY} ${i}:${X509_USER_PROXY}
        done
    else
        echo "Uh oh, cannot use SSH to copy the files to all WN's"
        exit 1
    fi
fi

echo "*****"
echo "Executing $EXE with $MPI_MPIEXEC"
$MPI_MPIEXEC 'pwd' /$EXE
echo "*****"
```

E.6 mpich2-wrapper.sh without Torque/PBS

This script sets up the mpd-ring because there is no support in MPICH2 for other LRMS's than Torque/PBS

```

# Script to start an MPI job under MPICH2
# Richard de Jong @ CERN 2006-06-29
```

```

echo "Using MPICH2"
```

```

if [ -z $MPI_MPICH2_VERSION ]
then
```



```

        echo "MPICH2, no such version available"
    exit
fi

# Set up the paths
PATH=$MPI_MPICH2_PATH/bin:$PATH
LD_LIBRARY_PATH=$MPI_MPICH2_PATH/lib:$LD_LIBRARY_PATH

# The binary to execute
EXE=$1

echo "*****"
echo "Running on: $HOSTNAME"
echo "As:      " 'whoami'
echo "*****"

echo "*****"
echo "Compiling binary: $EXE"
echo mpicc -o ${EXE} ${EXE}.c
mpicc -o ${EXE} ${EXE}.c
echo "*****"

# Create the .mpd.conf file.
# This is required to boot the mpd daemons on all hosts
# $RANDOM creates a random integer between 0 and 32767.
echo "MPD_SECRETWORD=PASS$RANDOM" > ~/.mpd.conf
chmod 0600 ~/.mpd.conf

if [ "x$PBS_NODEFILE" != "x" ] ; then
    echo "PBS Nodefile: $PBS_NODEFILE"
    HOST_NODEFILE=$PBS_NODEFILE
fi

if [ "x$LSB_HOSTS" != "x" ] ; then
    echo "LSF Hosts: $LSB_HOSTS"
    HOST_NODEFILE='pwd'/lsf_nodefile.$$
    for host in ${LSB_HOSTS}
    do
        echo $host >> ${HOST_NODEFILE}
    done
fi

if [ "x${HOST_NODEFILE}" = "x" ] ; then
    echo "No hosts file defined. Exiting..."
    exit
fi

echo "*****"
HOST_NEEDED='cat ${HOST_NODEFILE} | wc -l | awk '{print $NF}''
echo "Node count: $HOST_NEEDED"
echo "Hosts in $HOST_NODEFILE:"
cat $HOST_NODEFILE
echo "*****"

if [ "x${MPI_PATH_TO_SHARED_HOME}" != "x" ]
then

```



```

if [ ! -d ${MPI_PATH_TO_SHARED_HOME} ]
then
  echo "Shared home is set, but it is not accessible"
  exit 1
else
  echo "Shared home exists and is ${MPI_PATH_TO_SHARED_HOME}"
fi
else
  echo "No shared home is accessible"
  if [ "x${MPI_SSH_HOST_BASED_AUTH}" == "xyes" ]
  then
    echo "Using SSH to copy the files to all nodes"

    for i in `cat $HOST_NODEFILE | grep -v ${HOSTNAME} | sort -u`; do
      echo "Copying files to $i"
      ssh $i mkdir -p `pwd`
      scp -rp . "$i:pwd"
      # Copy the proxy certificate
      scp -p ${X509_USER_PROXY} $i:${X509_USER_PROXY}
      # Copy the MPD secret file
      scp -p ~/.mpd.conf $i:

      done
    else
      echo "Uh oh, cannot use SSH to copy the files to all WN's"
      exit 1
    fi
  fi

# Boot the MPD ring
mpdboot -f ${HOST_NODEFILE}

echo "*****"
echo "Testing MPD ring with mpdtrace"
mpdtrace
echo "Testing MPD ring with mpdringtest"
mpdringtest
echo "*****"

echo "*****"
echo "Executing $EXE with mpirun"
mpirun -n ${HOST_NEEDED} `pwd`/${EXE}
echo "*****"

# Stop the MPD daemons
mpdallexit

```

E.7 mpich2-pbs-wrapper.sh with Torque/PBS

This wrapper script does not use the mpd commands, because the mpiexec that integrates with TM is used.

```

# Script to start an MPI job under MPICH2
# Richard de Jong @ CERN 2006-06-29

echo "Using MPICH2"

```



```

if [ -z $MPI_MPICH2_VERSION ]
then
    echo "MPICH2, no such version available"
    exit
fi

# Set up the paths
PATH=$MPI_MPICH2_PATH/bin:$PATH
LD_LIBRARY_PATH=$MPI_MPICH2_PATH/lib:$LD_LIBRARY_PATH

# The binary to execute
EXE=$1

echo "*****"
echo "Running on: $HOSTNAME"
echo "As:      " 'whoami'
echo "*****"

echo "*****"
echo "Compiling binary: $EXE"
echo mpicc -o ${EXE} ${EXE}.c
mpicc -o ${EXE} ${EXE}.c
echo "*****"

if [ "x$PBS_NODEFILE" != "x" ] ; then
    echo "PBS Nodefile: $PBS_NODEFILE"
    HOST_NODEFILE=$PBS_NODEFILE
fi

if [ "x$LSB_HOSTS" != "x" ] ; then
    echo "LSF Hosts: $LSB_HOSTS"
    HOST_NODEFILE='pwd`/lsf_nodefile.$$
    for host in ${LSB_HOSTS}
    do
        echo $host >> ${HOST_NODEFILE}
    done
fi

if [ "x${HOST_NODEFILE}" = "x" ] ; then
    echo "No hosts file defined.  Exiting..."
    exit
fi

echo "*****"
HOST_NEEDED='cat ${HOST_NODEFILE} | wc -l | awk '{print $NF}'`
echo "Node count: $HOST_NEEDED"
echo "Hosts in $HOST_NODEFILE:"
cat $HOST_NODEFILE
echo "*****"

if [ "x${MPI_PATH_TO_SHARED_HOME}" != "x" ]
then
    if [ ! -d ${MPI_PATH_TO_SHARED_HOME} ]
    then
        echo "Shared home is set, but it is not accessible"
    fi
fi

```



```

        exit 1
    else
        echo "Shared home exists and is ${MPI_PATH_TO_SHARED_HOME}"
    fi
else
    echo "No shared home is accessible"
    if [ "x${MPI_SSH_HOST_BASED_AUTH}" == "xyes" ]
    then
        echo "Using SSH to copy the files to all nodes"

        for i in `cat $HOST_NODEFILE | grep -v ${HOSTNAME} | sort -u`; do
            echo "Copying files to $i"
            ssh ${i} mkdir -p `pwd`
            scp -rp . "${i}:"`pwd`"
            # Copy the proxy certificate
            scp -p ${X509_USER_PROXY} ${i}:${X509_USER_PROXY}

        done
    else
        echo "Uh oh, cannot use SSH to copy the files to all WN's"
        exit 1
    fi
fi

echo "*****"
echo "Executing $EXE with mpiexec"
$MPI_MPIEXEC --comm=mpich2 `pwd`/${EXE}
echo "*****"

```

E.8 openmpi-wrapper.sh without Torque/PBS

```

# Script to start an MPI job under OpenMPI
# Richard de Jong @ CERN 2006-06-29

echo "Using OpenMPI"

if [ -z $MPI_OPENMPI_VERSION ]
then
    echo "OPENMPI, no such version available"
    exit
fi

# Set up the paths
PATH=$MPI_OPENMPI_PATH/bin:$PATH
LD_LIBRARY_PATH=$MPI_OPENMPI_PATH/lib:$LD_LIBRARY_PATH

# The binary to execute
EXE=$1

echo "*****"
echo "Running on: $HOSTNAME"
echo "As:      " `whoami`
echo "*****"

echo "*****"
echo "Compiling binary: $EXE"
echo mpicc -o ${EXE} ${EXE}.c
mpicc -o ${EXE} ${EXE}.c

```



```

echo "*****"

if [ "x$PBS_NODEFILE" != "x" ] ; then
    echo "PBS Nodefile: $PBS_NODEFILE"
    HOST_NODEFILE=$PBS_NODEFILE
fi

if [ "x$LSB_HOSTS" != "x" ] ; then
    echo "LSF Hosts: $LSB_HOSTS"
    HOST_NODEFILE='pwd'/lsf_nodefile.$$
    for host in ${LSB_HOSTS}
    do
        echo $host >> ${HOST_NODEFILE}
    done
fi

if [ "x${HOST_NODEFILE}" = "x" ]; then
    echo "No hosts file defined. Exiting..."
    exit
fi

echo "*****"
HOST_NEEDED='cat ${HOST_NODEFILE} | wc -l | awk '{print $NF}''
echo "Node count: $HOST_NEEDED"
echo "Hosts in $HOST_NODEFILE:"
cat $HOST_NODEFILE
echo "*****"

if [ "x${MPI_PATH_TO_SHARED_HOME}" != "x" ]
then
    if [ ! -d ${MPI_PATH_TO_SHARED_HOME} ]
    then
        echo "Shared home is set, but it is not accessible"
        exit 1
    else
        echo "Shared home exists and is ${MPI_PATH_TO_SHARED_HOME}"
    fi
else
    echo "No shared home is accessible"
    if [ "x${MPI_SSH_HOST_BASED_AUTH}" == "xyes" ]
    then
        echo "Using SSH to copy the files to all nodes"

        for i in `cat $HOST_NODEFILE | grep -v ${HOSTNAME} | sort -u`; do
            echo "Copying files to $i"
            ssh $i mkdir -p `pwd`
            scp -rp . "$i":`pwd`
            # Copy the proxy certificate
            scp -p ${X509_USER_PROXY} $i:${X509_USER_PROXY}
        done
    else
        echo "Uh oh, cannot use SSH to copy the files to all WN's"
        exit 1
    fi
fi

echo "*****"

```




```
echo "Executing $EXE with orterun, without tm integration"
chmod 755 $EXE
orterun --prefix $MPI_OPENMPI_PATH -np $HOST_NEEDED -machinefile $HOST_NODEFILE `pwd`/$EXE
echo "*****"
```

E.9 openmpi-wrapper.sh with Torque/PBS

```
# Script to start an MPI job under OpenMPI
# Richard de Jong @ CERN 2006-06-29

echo "Using OpenMPI"

if [ -z $MPI_OPENMPI_VERSION ]
then
    echo "OPENMPI, no such version available"
    exit
fi

# Set up the paths
PATH=$MPI_OPENMPI_PATH/bin:$PATH
LD_LIBRARY_PATH=$MPI_OPENMPI_PATH/lib:$LD_LIBRARY_PATH

# The binary to execute
EXE=$1

echo "*****"
echo "Running on: $HOSTNAME"
echo "As:      " `whoami`
echo "*****"

echo "*****"
echo "Compiling binary: $EXE"
echo mpicc -o ${EXE} ${EXE}.c
mpicc -o ${EXE} ${EXE}.c
echo "*****"

if [ "x$PBS_NODEFILE" != "x" ] ; then
    echo "PBS Nodefile: $PBS_NODEFILE"
    HOST_NODEFILE=$PBS_NODEFILE
fi

if [ "x$LSB_HOSTS" != "x" ] ; then
    echo "LSF Hosts: $LSB_HOSTS"
    HOST_NODEFILE=`pwd`/lsf_nodefile.$$
    for host in ${LSB_HOSTS}
    do
        echo $host >> ${HOST_NODEFILE}
    done
fi

if [ "x${HOST_NODEFILE}" = "x" ] ; then
    echo "No hosts file defined. Exiting..."
    exit
fi

echo "*****"
HOST_NEEDED=`cat ${HOST_NODEFILE} | wc -l | awk '{print $NF}'`
echo "Node count: $HOST_NEEDED"
```



```
echo "Hosts in $HOST_NODEFILE:"
cat $HOST_NODEFILE
echo "*****"

if [ "x${MPI_PATH_TO_SHARED_HOME}" != "x" ]
then
  if [ ! -d ${MPI_PATH_TO_SHARED_HOME} ]
  then
    echo "Shared home is set, but it is not accessible"
    exit 1
  else
    echo "Shared home exists and is ${MPI_PATH_TO_SHARED_HOME}"
  fi
else
  echo "No shared home is accessible"
  if [ "x${MPI_SSH_HOST_BASED_AUTH}" == "xyes" ]
  then
    echo "Using SSH to copy the files to all nodes"

    for i in `cat $HOST_NODEFILE | grep -v ${HOSTNAME} | sort -u`; do
      echo "Copying files to $i"
      ssh ${i} mkdir -p 'pwd'
      scp -rp . "${i}:'pwd'"
      # Copy the proxy certificate
      scp -p ${X509_USER_PROXY} ${i}:${X509_USER_PROXY}
    done
  else
    echo "Uh oh, cannot use SSH to copy the files to all WN's"
    exit 1
  fi
fi

echo "*****"
echo "Executing $EXE with orterun, tm way"
chmod 755 $EXE
orterun --prefix $MPI_OPENMPI_PATH -np $HOST_NEEDED 'pwd'/$EXE
echo "*****"
```

Appendix F

Site Functional Test for MPI

F.1 tests/sft-mpi.def

```
testName: sft-mpi
testTitle: SFT MPI test for multiple implementations
friendlyName: mpi
testHelp: https://twiki.cern.ch/twiki/bin/view/LCG/ImplementationOfMpi
```

F.2 tests/sft-mpi

```
#!/bin/bash

STATUS=

tests="lam-6-wrapper.sh lam-7-pbs-wrapper.sh lam-7-wrapper.sh mpich2-pbs-wrapper.sh
mpich2-wrapper.sh mpich-pbs-wrapper.sh mpich-wrapper.sh openmpi-pbs-wrapper.sh openmpi-wrapper.sh"

result=0
failed=
fcode=OK

function statusCode
{
    case $1 in
        0) echo "OK";;
        *) echo "ERROR" ;;
    esac
}

echo "<h2>SFT MPICH test</h2>"
echo "<p>Testing compilation and execution of an MPI application using various implementations of MP
I.<br />"
Note: this is a super-test that contains a number of tests!</p>"

log=sft-mpi.log
rm -f $log

echo "<table border=1><th>Test</th><th>Result</th>"

export PBS_NODEFILE=pbs_nodefile
echo $HOSTNAME > $PBS_NODEFILE

for test in $tests ; do
    echo "<h3>Now running: $test</h3>" >> $log
```



```
echo "<pre>" >> $log
data/$test data/MPItest >> $log 2>&1
res=$?
echo "</pre>" >> $log
code='statusCode $res'
if [ $res -gt 0 ] ; then
    result=$res
    fcode="$code:"
    failed="$failed$test "
fi
echo "<tr><td>$test</td><td>$code</td></tr>"
done

echo "</table>"

rm -f $PBS_NODEFILE

echo "<p>Overall summary: "
statusCode $result
echo "</p>"

echo "<p>LOG from sub-tests execution:</p>"
echo "<pre>"
cat $log
rm -f $log
echo "</pre>"

echo "result is $result"
if [ $result -gt 0 ]
then
    result=50
else
    result=10
fi

echo "summary: $fcode $failed"
exit $result
```