

The Grid-Ireland Deployment Architecture

Brian Coghlan, John Walsh, and David O'Callaghan

Department of Computer Science, Trinity College Dublin, Ireland.
coghlan@cs.tcd.ie, john.walsh@cs.tcd.ie, david.ocallaghan@cs.tcd.ie

Abstract. Grid-Ireland is unusual in its integrated infrastructure, and the stress that is laid on homogeneity of its core. The major benefit is the decoupling of site details from the core infrastructure, and the resulting freedom for heterogeneity of site resources. We describe the efforts to support this heterogeneity in a systematic way. We also describe the deployment architecture and a methodology to increase the availability of the core infrastructure.

1 Introduction

The HPC facilities in Ireland are very limited. There is an Origin 3800/40, a 64-way Xeon cluster and a 6TB disk farm at NUI, Galway, and a 20-CPU SGI Altix 3700 has just been installed. There is a 100-CPU cluster in the Boole Centre at UCC in Cork, and two 96-CPU clusters plus funding for three extra 96-CPU clusters in the NMRC at UCC. There are 80-CPU P3 and 130-CPU Xeon clusters and a 4TB disk farm at TCD in Dublin, and a fully immersive VR cave is to be installed. A 256-CPU cluster will be installed at UCD in Dublin by the end of 2004. There is a 84-CPU cluster at NUIM in Maynooth. There are several small-scale clusters. There is funding from Science Foundation Ireland (SFI) and the Higher Education Authority (HEA) for several medium-scale clusters (100-500 CPUs) and two medium-scale data farms in 2004/5.

Whilst the experimental and theoretical science paradigms remain strongly embedded in Irish science, there is strong growth in the hybrid paradigm, computational science. Most of this scientific computing is still done on local facilities. It involves a wide range of application areas, but few truly parallel applications. Most users develop their codes but use commercial libraries and tools. The reference architectures for these are a major factor in the choice of HPC architecture, i.e. most of the deployed architectures are mission-specific.

Currently there is no large-scale facility in Ireland. Until very recently there was no identified governmental intention to have one. In August, however, SFI announced that they wish to enhance the high-end computational capabilities of the overall Irish research community by the creation of a National Centre for High End Computing within the Republic of Ireland. It is very likely that the resulting centre will include a mix of mission-specific architectures.

These limited and mostly mission-specific resources should not be wasted by being inaccessible to the Grid simply because their architectures are not those of the reference ports.

1.1 Grid-Ireland

Grid-Ireland provides grid services above the Irish research network, allowing researchers to share Irish computing and storage resources using a common interface. It also provides for international collaborations by linking Irish sites into the European grid infrastructures being developed under such EU projects as EGEE, LCG and CrossGrid. The grid infrastructure is currently based on LCG2, the common foundation that ensures interoperability between participating scientific computing centres around the world. Internationally, members of Grid-Ireland are involved in the EU EGEE, CrossGrid, JetSet and COST 283 iAstro projects, and there are links to the UK GridPP and e-Science programs. The Grid-Ireland OpsCentre is the EGEE Regional Operations Centre (ROC) for Ireland.

Grid-Ireland currently encompasses six sites, at TCD, UCC, NUIG, DIAS, UCD and QUB, with an Operations Centre in TCD. It aims to make Grid services accessible to an additional eleven Irish third-level institutions in the near future as a result of a generous donation by Dell Ireland Limited. The Irish NREN (HEAnet) are substantively assisting in this initiative. Grid-Ireland will then encompass 17 sites, i.e. the majority of academic institutions in Ireland will be connected to the national grid.

There are three major national VOs. The first, CosmoGrid, is a collaborative project entitled Grid-enabled computational physics of natural phenomena, explicitly aimed at inter-institutional and interdisciplinary compute-intensive research in natural physics, with nine participating institutions, led by the Dublin Institute for Advanced Studies (DIAS). Astrophysics are a key and central element of the project centred on astrophysical objects ranging from supernova remnant (with strong collisionless shocks), forming stars (jets and outflows) to neutron stars (radiative processes) and the sun (the solar transition region). In addition to those areas, studies on gravitational waves, adaptive optics, meteorology (regional climate models), geophysics (full simulation of a digital rock) and atmospheric physics are being pursued.

The second VO, MarineGrid, is a data-intensive collaboration between Geological Survey of Ireland, Marine Institute and four Universities (NUIG, UCC, UCD and UL). The Irish National Sea-bed Survey is taking place at present through bathymetric mapping of the seabed. Ireland is an island with nine tenths of its area under water. The seabed survey spans 525,000km² and currently contains approximately 6TB of data. Detailed knowledge of the seabed topography with location resolutions of up to 2m will have significant economic implications on for example fishing and mineral resource management. An unexpected spin-off is exploitation of historically valuable wrecks. The Survey is a valuable government resource with associated security concerns.

The third VO, WebCom-G, is investigating an alternative to existing von Neumann grid execution models, which are not appropriate to their high-latency, loosely-coupled infrastructure. UCC, TCD, NUIG and QUB are creating a condensed graph grid engine that exploits laziness and speculation and is compatible with and uses traditional grids. There is a depth of interest in Irish computer

science circles about issues of languages, programming models[1] and execution models[2][3] for heterogenous environments, and this VO is a good example. Grid-Ireland specifically wishes to support these research directions.

1.2 Homogeneous Core Infrastructure, Heterogenous Resources

There are three further motivations:

- (a) To minimize the demand on human resources by minimizing the proportion of the software that needs to be ported. The simplest component of most grid software frameworks is that relating to the worker nodes.
- (b) To minimize the demand on human resources by maximizing the proportion of the software that does *not* need to be ported. Thus all the non-worker node components should use the reference port, i.e. the core infrastructure should be homogeneous.
- (c) To maximize the availability of the infrastructure. Grid-Ireland has designed a transactional deployment system to achieve this[4]. This requires that the core infrastructure be homogeneous, and also centrally managed.

Thus *Homogeneous Core Infrastructure, Heterogenous Resources* is a pervasive motto that encapsulates explicit and implicit principles:

- (a) *Explicit homogeneous core infrastructure*: this principle enables a uniform dedicated core national grid infrastructure, which supports a uniform architecture based on reference ports of the grid software, and thereby frees resources for maximum focus on the critical activities such as security and monitoring/information systems. Logically, it allows a uniform control of the grid infrastructure that guarantees uniform responses to management actions. It also assures a degree of deterministic grid management. Furthermore it substantially reduces the complexity of the release packaging and process. Minimizing this complexity implies maximizing the uniformity of the deployed release, i.e. the core infrastructure should be homogeneous.
- (b) *Implicit centralized control via remote management*: this principle enables simpler operations management of the infrastructure. Remote systems management enables low-level sub-actions to be remotely invoked if required. It also enables remote recovery actions, e.g. reboot, to be applied in the case of deadlocks, livelocks, or hung hardware or software. Realistically, all infrastructure hardware should be remotely manageable to the BIOS level.
- (c) *Implicit decoupling of grid infrastructure and site management*: this principle enables the infrastructure and sites to be independent. It can encompass policies, planning, design, deployment, management and administration. In particular it allows the infrastructure upgrade management to be independent of that of the site, and non-reference mission-specific architectures to be deployed at the site.

Grid-Ireland has been designed with this approach since mid-2001. Funding was sought and eventually granted, a senior Grid Manager appointed, a Grid Operations Centre established and staffed, infrastructure specified, purchased and

installed, and finally middleware and management tools deployed. We consider that the use of these principles has been highly beneficial.

The Operations Centre (see 1) hosts approximately 20 national servers, a certification TestGrid of approximately 40 machines and 4TB disk farm, a cluster of 64-CPU and 4TB disk farm for the various testbeds, plus support staff. The TestGrid serves multiple purposes: it implements a fully working replica of the national servers and sites; it permits experimentation without affecting the national services; it acts as the testing and validation platform; and it acts as a non-reference porting platform.

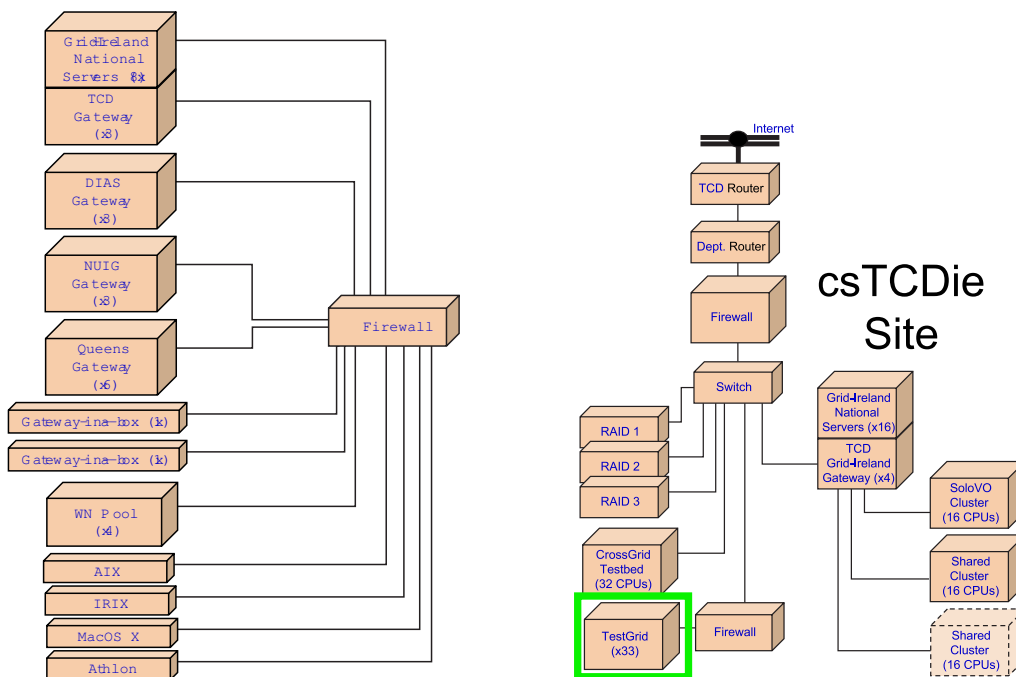


Fig. 1. (a) Test Grid (b) TCD Site

2 Heterogeneity

Grid-Ireland wished, in the first instance, that the porting of the LCG2 software to other platforms would focus on the ability to execute Globus and EDG jobs on worker nodes, and that replica management, R-GMA and VOMS would be supported. There was also a desire that MPI, replica management and the OpenPBS client be provided on each worker node. In some cases Torque might be required since newer versions of operating systems are not always provided for in OpenPBS. Also the R-GMA information system producer and consumer APIs and the VOMS client were required.

In summary we wished to port:

1. VDT
2. MPI
3. OpenPBS or Torque client
4. R-GMA producer and consumer APIs
5. VOMS client

There are a number of on-going issues, but we have successfully ported the functionality for job submission to Fedora Core 2, IRIX 6.5.14 and 6.5.17m, AIX 5.2L and Red Hat 9. We also plan to do this for Mac OS X v10.3 very soon, and a number of other platforms if the need arises within Grid-Ireland.

A number of CVS repositories are used to build all the necessary software for a worker node. The head version of VOMS is obtained from INFN's own repository. The whole of LCG2 is extracted using CVS checkouts directly from CERN's lcgware repository. The CrossGrid software is obtained by directly copying the CVS repository to a local repository. Nightly builds are then done from this local repository. The RAL repository of R-GMA will also need to be added soon, since LCG2 no longer maintain the most recent version of R-GMA.

Figure 2 shows the status of the build system in November 2004. The results change quite regularly as new ports are completed.

OS Type	Version	VDT	Basic	VOMS	RGMA	RM	Colour	Meaning
Redhat	7.3	RPMS	RPMS	RPMS	RPMS	RPMS		
Redhat	9.0	RPMS	RPMS	RPMS	RPMS	RPMS		
Fedora Core	2	RPMS	RPMS	RPMS	tarball	tarball		
					RPMS	RPMS		
SGI	6.5.14	tarball	tarball	tarball	tarball	tarball	Red	To be started
AIX	5.2L	tarball	tarball	tarball	tarball	tarball	Yellow	Started
Darwin	10	tarball	tarball	tarball	tarball	tarball	Green	Done

Fig. 2. Auto-build Results for Worker Nodes

3 Deployment

As stated above, Grid-Ireland has installed a grid computing infrastructure that is fully homogeneous at its core. Each of the sites connects via a grid gateway. This infrastructure is centrally managed from Trinity College Dublin. These gateways have at their core a set of seven machines: a firewall, a LCFGng install server, a compute element (CE), a storage element (SE), a user interface machine (UI), and a worker node that is used for gateway tests only. All the sites are identically configured. The grid software is initially based on LCG2, but later will follow the EGEE releases.

As can be seen from Figure 3, the site resources (shown as a cluster of worker nodes) are outside the domain of the gateway; these resources belong to the site, and the site is in charge of their own resources. One of the key departures from the structures for deployment commonly used in Europe is to allow those resources to be heterogeneous with respect to the gateways. As explained in Section 2, Grid-Ireland is attempting to provide ported code for any potential platform that will be used for worker nodes (WNs). The rationale behind this is described in the early sections of this paper.

The only requirement on the site is that the worker nodes be set up to cater for data- and computation-intensive tasks by installing the worker-node software outlined in Section 2. This includes both the Replica Management software from LCG2 and the revised versions of MPICH-G2 from CrossGrid.

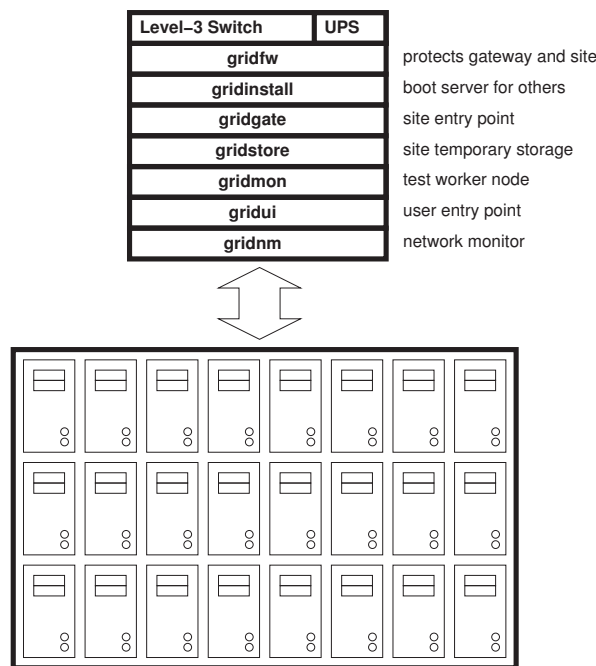


Fig. 3. Generic Grid-Ireland Site

Grid-Ireland has specified its gateways to ensure minimal divergence from standard site configuration, minimal hardware and space costs per site, and minimal personnel costs per site. The basic technology for this is the use of virtual machines. Currently there are two physical realisations of this architecture. At minimum, a generic Grid-Ireland gateway comprises a single physical machine, a switch, and a UPS unit. Each machine will run its own OS plus a number of virtual machines that appear to be a normal machine both to the host OS and to external users. The Linux OS and grid services will be remotely installed,

upgraded and managed by the Operations Centre, without needing any attention at the site. The firewall and LCFG server run concurrently on the host operating system. All other servers are hosted as virtual machines. Eleven such gateways are presently being prepared for deployment.

For more demanding sites the functionality is spread over four physical machines, with the firewall, LCFG server, CE and SE running on the host operating systems. The other servers are hosted as virtual machines: the test WN on the CE, and the UI and NM on the SE. Six such gateways are already deployed.

Apart from the firewall, all other servers on the gateways are installed from the LCFG server. The LCFG server itself is manually installed, but thereafter it is updated with new releases from the central Grid-Ireland CVS repository, see Figure 4. It is usual that this is a manual process involving CVS checkouts. Whilst this takes place the site is essentially in an inconsistent state. Grid-Ireland, however, have designed an automatic process that is specifically intended to reduce inconsistency to a minimum.

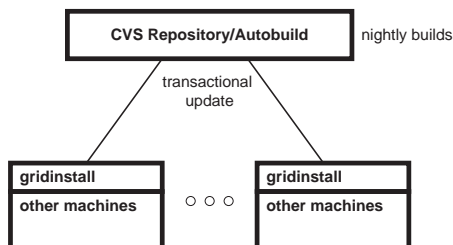


Fig. 4. Deployment Process

3.1 Consistency

Let us consider for instance that a new release of some infrastructural grid software is incompatible with the previous release, as is often the case. Once a certain proportion of the sites in a grid infrastructure are no longer consistent with the new release then the infrastructure as a whole can be considered inconsistent. Each grid infrastructure will have its own measures, but in all cases there is a threshold below which proper operation is no longer considered to exist. The infrastructure is no longer available. Thus availability is directly related to consistency. An inconsistent infrastructure is unavailable.

Assume N identical sites are being evaluated in independent experiments, and that at time t , $C(t)$ are consistent (have upgraded) and $I(t)$ are inconsistent (have yet to upgrade). The probabilities of consistency and inconsistency are:

$$P_C(t) = C(t)/N \qquad P_I(t) = I(t)/N \qquad (1)$$

The average time a site waits before it becomes consistent (the mean time to consistency MTTC) is:

$$\text{MTTC} = \int_0^{\infty} P_I(t) dt \qquad (2)$$

If the infrastructure contains M types of sites, then the probability of an inconsistent multi-site infrastructure is the probability of at least one site being inconsistent, i.e. the probability of NOT(all upgraded):

$$P_{\text{infra}}(t) = 1 - \prod_{m=1}^M (1 - P_{Im}(t)) \quad (3)$$

The more sites, the greater the probability. Clearly it is easier to understand the situation if all M sites are identical, i.e. the infrastructure is homogeneous, as is the case for Grid-Ireland. The MTTC is:

$$MTTC_{\text{infra}} = \int_0^{\infty} P_{\text{infra}}(t) dt \quad (4)$$

The interval between releases MTBR is quite independent of the MTTC. The MTTC is a deployment delay determined by the behaviour of the deployers, whilst the MTBR is dependent upon the behaviour of developers. Otherwise similar considerations apply.

3.2 The need for Transactionality

Maximizing availability means maximizing the proportion of time that the infrastructure is entirely consistent. This requires either the the time between releases MTBR to be maximized or the MTTC to be minimized. The MTBR is beyond the control of those who manage the infrastructure. On the other hand, if the MTTC can be minimized to a single, short action across the entire infrastructure then the availability will indeed be maximized.

However, an upgrade to a new release may or may not be a short operation. To enable the upgrade to become a short event the upgrade process must be split into a variable-duration prepare phase and a short-duration upgrade phase, that is, a two-phase commit. Later in this paper we will describe how this can be achieved.

If the entire infrastructure is to be consistent after the upgrade, then the two-phase commit must succeed at all sites. Even if it fails at just one site, the upgrade must be aborted. Of course this may be done in a variety of ways, but from the infrastructure managers' viewpoint the most ideal scenario would be that if the upgrade is aborted the infrastructure should be in the same state as it was before the upgrade was attempted, that is, the upgrade process should appear to be an atomic action that either succeeds or fails.

Very few upgrades will comprise single actions. Most will be composed from multiple subactions. For such an upgrade to appear as an atomic action requires that it exhibits transactional behaviour, that all subactions succeed or all fail, so that the infrastructure is never left in an undefined state.

Thus we can see that to maximize availability requires that an upgrade be implemented as a two-phase transaction.

3.3 Transactional Deployment System

We have implemented transactional deployment using three components. The first is a repository server, which hosts both the software to be deployed and the Transactional Deployment Service (TDS) logic. Secondly, there is a user interface, which has been implemented as a PHP page residing on an Apache web server. Finally there are the install servers at the sites that we are deploying to. These servers hold configuration data and a local copy of software (RPMs) that are used by LCFGng to maintain the configuration of the client nodes at the site. It is the state of these managed nodes that we are trying to keep consistent.

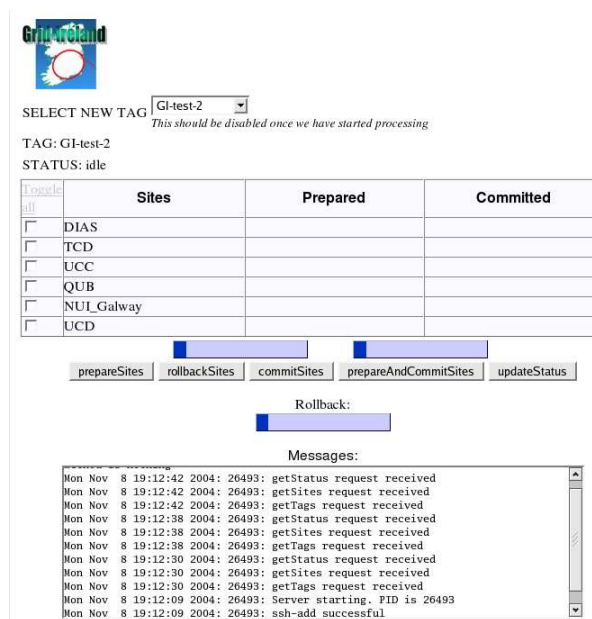


Fig. 5. Transactional Deployment System GUI

It will be a while before statistically significant data is available for transactional deployment to a real grid infrastructure such as Grid-Ireland. The MTBR, i.e. the time between releases, is the same with or without transactional deployment. Let us assume the CrossGrid production testbed MTBR by way of example. We have estimated the worst-case MTTC as 17.5 minutes. The resulting worst-case infrastructure availability is shown in Table 1.

If the LCFG client nodes could be signalled to update immediately, i.e. $T_{signal} = 0$, then the availability would be increased 99.92%, and in fact this is likely to be substantially better for realistic release updates.

Estimated MTBR	163 hours
Estimated MTTC	17.5 minutes
Estimated availability	99.82%

Table 1. Example MTBR, MTTC and availability for transactional deployment

4 Conclusions

Grid-Ireland is unusual in its integrated infrastructure, and the stress that is laid on homogeneity of its core. The principles behind this have been described above. The major benefit is the decoupling of site details from the core infrastructure, and the resulting freedom for heterogeneity of site resources.

We have described our efforts to support this heterogeneity by porting to non-reference platforms in a systematic way, with nightly autobuilding. This has allowed us to begin a most interesting set of benchmarking and heterogeneity experiments involving all of these platforms.

It is clear that the infrastructure has greatly enhanced availability with transactional deployment, improved from 87-93% to 99.8%, with the potential for 99.9% availability with a small amount of extra work. However, the most important benefit of the transactional deployment system is the ease it brings to deployment. Transactional deployment allows for an automated totally-repeatable push-button upgrade process, with no possibility of operator error. This is a major bonus when employing inexperienced staff to maintain the grid infrastructure.

5 Acknowledgements

We would like to thank Enterprise Ireland, the Higher Education Authority, Science Foundation Ireland and the EU for funding this effort. We gratefully thank DIAS for the SGI machine they have loaned to us for the IRIX port, and IBM and Dell for sponsoring us with machines to perform ports to their platforms. Most of all we would like to thank those at INFN and CERN for all their help in porting to each platform.

References

1. Lastovetsky, A.: Adaptive parallel computing on heterogeneous networks with mpc. *Parallel Comput.* **28** (2002) 1369–1407
2. Ryan, J.P., Coghlan, B.A.: Smg: Shared memory for grids. In: *The 16 th IASTED International Conference on Parallel and Distributed Computing and Systems*, Cambridge, MA, USA (2004)
3. Morrison, J.P.: *Condensed Graphs: Unifying Availability-Driven, Coercion-Driven, and Control-Driven Computing*. PhD thesis, Technische Universiteit Eindhoven (1996)
4. Coghlan, B., Walsh, J., Quigley, G., O’Callaghan, D., Childs, S., Kenny, E.: Principles of transactional grid deployment. Submitted to EGC 2005 (2004)