

Deployment tools

This section looks at tools and other issues that may arise as part of a Grid system deployment. The information here is not specific to any particular Grid software: such package specific issues are covered in the [Appendix](#) ^[1].

Firewalls

Network controls can significantly reduce some of the most serious security vulnerabilities affecting Grid systems. For example, the serious incidents that have occurred when directories containing stored identity credentials were accidentally exported on a networked file system. This section suggests a number of ways that firewalls and Grid systems can be configured to work effectively together. The choice among these options will depend on the security requirements of individual sites and Grids, and the preferences of network and Grid operations managers.

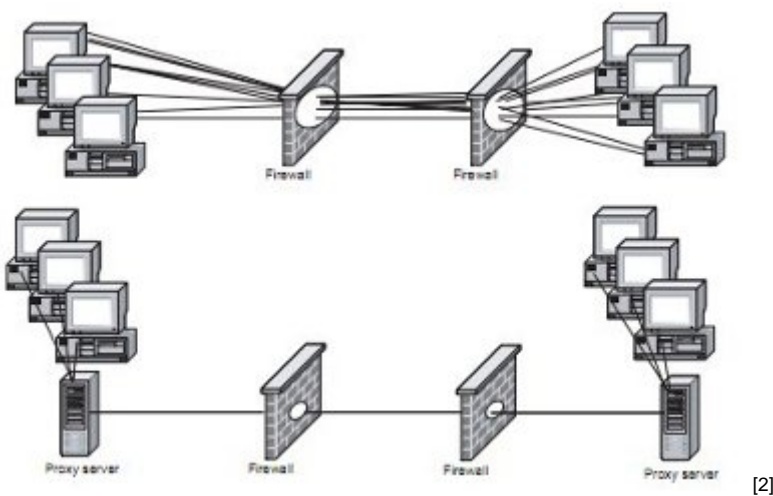
To be effective, a firewall needs to be able to distinguish between 'good' traffic which should be permitted to pass and 'hostile' traffic which should be blocked. Traditionally this was done using port numbers. A typical firewall might allow connections on ports 80 and 443 to the web server, and port 25 to the mail server. For simple protocols using a single TCP connection, this level of description by port and destination was adequate. However, Grid protocols tend to involve multiple connections between groups of machines, making them considerably more complex to describe.

For some complex protocols in common use, firewall and router vendors have developed software to enable their systems to interpret the protocols and support the use of ephemeral ports by those protocols. For example, it is common for firewalls to inspect the control channel of conventional FTP connections and make temporary changes to their rules to allow the related data connections to pass. This approach is attractive as it allows very specific alterations to be made to the firewall at particular times, thus minimising the additional exposure to threats from the network. Unfortunately, Grid communications are often encrypted between their endpoints, which means that intermediate network devices cannot see the information they would need to extract from the protocol exchanges to do such dynamic reconfiguration. At present, no standard firewalls or routers are known to provide dynamic support for Grid protocols.

Static firewall configurations can still be used with Grid systems, especially if they are located at points on the network where flows are simpler (for example, at the perimeter of, rather than within, a Condor® cluster). The configuration of such firewalls will be much easier if IP addresses are allocated to place Grid and non-Grid systems in separate blocks. Grid operators need to be prepared to co-operate with such addressing plans and also to use the controls provided by their software to make their communications easier for routers or firewalls to identify, for example using defined port ranges. Demanding that firewalls be thrown wide open exposes the Grid system and its surroundings to considerable and

unnecessary risks.

When using a static firewall to manage a dynamic protocol exchange it is inevitable that a greater range of ports will be left open to a greater range of clients and server than are actually required at any particular time. This means that Grid systems will be more exposed to attacks than others for which these ports are not opened. This need not be a serious security problem, but it does mean that all systems to which ports are opened need to be managed securely, as discussed in Section 4.3 System Management, and their users need to be aware of the greater exposure. If ports or addresses are opened to allow Grid traffic, care must be taken that this does not increase the threat to computers that never use Grid applications. One possible way to reduce the exposure of internal systems is to pass all Grid traffic through a dedicated proxy server at each site. As shown in Figure 3 overleaf, routing traffic via proxies means that a much smaller opening is needed in the firewall and reduces the number of systems that are exposed to direct attack. It appears that a Globus gatekeeper system could be used as a local proxy, but it is not known whether this configuration has ever been used. The proxy machine is still exposed to attack so needs to be designed and managed securely.



[2]

An alternative to a static firewall that leaves all necessary communications ports open all the time is a dynamic firewall that only opens ports as they are required. This requires:

- the Grid systems, or their authenticated users, to inform the firewall before they attempt to open and close a new connection
- the firewall to trust the Grid system that this is a properly authenticated request and not an intruder attempting to disable the firewall's protection.

This approach makes the system controlling the firewall a critical part of the site's, and the Grid's, security protection and so it needs to be designed, implemented and managed with great care. An intruder who can compromise any part of the control system can use it to modify firewall settings in any way they choose, making the firewall at least useless and potentially hostile. Particular care must therefore be given to the authentication of the user and the security of the computer system performing the authentication, since these are obvious points of attack. As with other Grid systems, the security of the user's client computer is also a potential weakness if an intruder can either steal a usable authentication credential or take over a session that has already been authenticated. It is a principle of good security design that critical systems should be simple, so a fully integrated Grid and firewall system

may be difficult to achieve. However, a number of simpler approaches could be used. One that has been used successfully uses an SSL (Secure Sockets Layer) connection to carry an initial, one time password authentication of the user. As long as the SSL connection remains active, the firewall will allow Grid protocols to and from the same client system [PPPL]. This apparently works well although there were initial problems when the SSL connection was judged to be idle and closed automatically, thereby shutting off the Grid connections as well. This system also requires the user to log in twice: once to gain access through the firewall and once to the Grid system. In principle, a single login process could be used to give access to both the network and the Grid systems, provided it provides sufficiently good proof of the user's identity, although it is not known if this has been attempted in practice.

Dynamic firewalls may also be implemented on the Grid server or gateway itself. In this case, the server will only accept calls on particular ports from a particular IP address once a user has authenticated from that IP address and been authorised to use services on those ports. This requires that at least the authentication port is left open, or can be opened by a particular combination of packets. An example of this is the Dyna-Fire project [DynaFire].

Since there are some ports that will never be used by Grid protocols, it is still recommended that static controls on network devices are used in addition to protect those ports, both as a backup precaution and in case any accident or malicious activity results in them being opened and vulnerable on the Grid server itself.

Any security design involves a trade-off between different risks, for example, the risk of permanently open ports against the risk of a more complex dynamic control system. The decision on which risks are the most significant, and therefore the appropriate way to manage a firewall, will depend on the nature of individual Grids, networks and organisations.

Tunnelling

Where it is not possible to allow native Grid protocols to flow across networks, it may instead be possible to communicate using tunnels constructed using a variety of standard and proprietary protocols. Tunnels that involve an overhead in both bandwidth and processing at the endpoints are unlikely to be suitable for extreme bandwidth applications but they may be a possibility for lower bandwidth applications, such as between clients and servers.

Tunnels use packets belonging to one protocol, for example SSH (Secure Shell), to carry complete packets making up a second protocol. Tunnels are established between two endpoints. At the first endpoint, the packets to be tunnelled are encapsulated within the tunnelling protocol. At the other endpoint the encapsulation is stripped off and the original packet re-created. Since tunnels generally use a single TCP connection between predictable endpoints, they require a smaller opening in a firewall so may be more acceptable to network managers. Conversely, the tunnel provides an unobstructed route through the firewall between the two endpoints, so any security problem at one endpoint is likely to spread rapidly to the other. Some networks may therefore regard tunnels as a greater risk than firewall holes. As with firewall configurations, tunnels require an appreciation and balancing of risks.

The Globus Toolkit™, having the most complex protocols, is the most common application for tunnels. A paper describing how connections between two GT2 (Globus Toolkit™ v2) systems were established through an SSH tunnel has been published [Graupner & Reimann] and the principles explained there should be applicable to any TCP-based protocol. A number of

commercial VPN products exist and it seems likely that they could be used in a similar way. These types of tunnel have no authentication: all traffic that enters one endpoint will be passed to the other. The Globus team are understood to be considering a tunnelling system that would require Globus authentication before allowing a connection. This might be more acceptable to some sites than an unrestricted tunnel, as it would provide nearly the same function as a fully Globus-aware firewall but on a less complex system.

System Management

Although Grid systems are likely to run specialised software, they are normally based on standard computer hardware and operating systems, and may well run standard applications as well. To protect Grid security it is essential to ensure that this software is secure and not concentrate effort only on the Grid specific parts. If a computer that forms part of the Grid (whether as a server or a client workstation) is compromised through a security weakness, the Grid itself is likely to be compromised soon afterwards.

Most security incidents on networked computers involve weaknesses in standard software or operating system components, and there is no reason to expect that Grid computers will be any different. Indeed the short history of Grid incidents already confirms this: attackers will not bother to learn complex new software if they can gain complete control of the systems by exploiting known weaknesses. Insecure file sharing systems, web servers or other software can allow intruders to gain full control of computers, thus completely bypassing the authentication procedures intended to protect Grid systems from unauthorised use.

Grid systems that are exposed to the Internet must therefore be managed at least as well as any other type of Internet server, and follow the same basic principles. Software must not be installed if it is not required, and both applications and operating systems must be kept up to date by installing security patches provided by their authors (for more detail, see the CERT@-CC (Computer Emergency Response Team – Co-ordination Centre) guide to securing networked servers [CERT@-CC SI]). Where Grid packages incorporate standard software, every effort must be made to install and distribute patches for that software as soon as possible. Once a patch is released by the author, intruders know of the problem and will start to try to exploit it, so any delay in installing the patch represents a considerably increased risk. If security patches are available but not yet installed, consideration should be given to increasing other security measures such as firewall controls and monitoring to protect the vulnerable system.

Most computers are managed remotely across the network, and care is required to ensure that the channels used for system management do not themselves present a security risk. Software used to provide remote access must be kept up to date, and access to it should be restricted if possible using IP address or other access controls. Any remote access system that carries passwords across a shared network, especially those of system administrators, should use encryption to prevent the passwords being read as they pass across the network. Systems based on SSH or other VPN protocols should provide encryption as well as allowing the administrator to confirm that they are connected to the correct computer before entering their password.

Measurement and Monitoring

One of the aims of Grids is to achieve the maximum effect from the available resources.

Measuring performance is therefore important to identify bottlenecks or faults that may restrict what users can achieve, as well as to plan future use and supply of resources. Some Grid systems can also distribute work dynamically to take account of changing load and so require real-time information about the resources available. The requirement to collect and distribute such information may need to be considered as part of any Grid deployment.

Grid applications may be affected by different aspects of performance, some of which may not matter to other, general purpose uses of the same systems and networks. It is likely that Grids will need their own specific performance measurements in addition to those provided by normal systems and network management tools. Parameters of interest are likely to include total capacity and current load as well as details of performance (e.g. packet loss and jitter) at all levels of the system including hardware (e.g. disk and CPU), networks and software applications. These values are likely to require both examining the performance of real jobs on the Grid (referred to here as passive monitoring) and generating additional jobs or traffic solely for the purpose of obtaining information (referred to here as active measurement). Measurement allows specific types of packet or job to be used to gather particular types of information, but it also itself consumes some of the resources of the Grid, thus reducing the resources available for real jobs.

Passive monitoring can usually be done locally within the Grid systems themselves, so is less likely to require specific provision when systems are deployed. Active measurement will more often use dedicated systems and network flows so is more likely to affect the deployment. For example packet loss between two Grid data centres is likely to be measured by sending a known stream of packets between the sites, rather than waiting for Grid users to generate traffic. The packets will usually be sent and received by dedicated measurement systems rather than the Grid systems themselves, to reduce the impact of other jobs running at the same time. Clearly, the measuring systems need to be able to send and receive external traffic so routers and firewalls need to be configured to allow this. The ports and destinations needed will depend on the measurement tools used. One set of network measurement tools is described in [GRIDMON]. Allowing external traffic to reach the measurement systems may expose them to possible threats from the external network, so these systems must be secured. It may also be appropriate to restrict the access from the measurement systems to the rest of the internal network to reduce the impact of any security problems that may occur.

Most Grids will involve a number of different sites, so obtaining an overall view of the performance of the Grid will require information to be collected from local measurement and monitoring systems at each site. Various projects have proposed hierarchical schemes for the publication and collation of this information: these are likely to require further data flows between local measurement systems and a central collation point. Where monitoring of live Grid jobs is used to provide information, it is recommended that this should not be made directly available for external queries, since this would represent an additional, unpredictable, load on Grid systems, but that the information should be collected on a regular basis by a local monitoring system which can then publish it.

Intrusion Detection Systems

IDSs (Intrusion Detection Systems) are a useful technique for detecting possible security incidents. An IDS is a software or hardware system that inspects some aspect of a computer or network and attempts to identify abnormal activity that may indicate a security incident. It should be noted that an IDS can only raise an alarm after the abnormal activity has begun,

and for many attacks this will happen after the initial attack has succeeded. Even in this case, prompt action in response to an IDS alarm may still significantly reduce the attack's impact by preventing its spread either within the attacked system or to other systems. For this reason IDSs are most useful when deployed in combination with firewalls and other preventive measures, to indicate when those measures may be insufficient for a new type of attack. Using IDSs in combination with Grids involves some new challenges that are the subject of current research.

IDSs can work at many different levels to attempt to identify abnormal activity. Different levels may be combined in a single product, or the outputs of different products may be combined by central IDS management systems.

Host based IDS A host based IDS inspects the files on the computer's disk, looking for unauthorised changes to their content, permissions, ownership etc. This is usually done by taking a snapshot of a known good system and defining rules that indicate which aspects of each file are expected to change and which are not. For example, the content of a log file should change over time, but the permissions that prevent it being modified by an unauthorised user should not. The IDS program will then periodically inspect the file system confirming that no changes prohibited by the rules have occurred. Host based IDS packages such as Tripwire™ [Tripwire] and AIDE [AIDE] come with pre-defined rules for many standard operating systems. These rules may need to be further tailored for Grid systems to give the best protection.

Networked based IDS

A network based IDS examines packets on a network rather than files on a disk, looking for individual packets or combinations of packets that may indicate abnormal activity. At the simplest level, a network based IDS may have rules that compare packets against those generated by well-known exploit tools, and generate alarms when these are detected. More complex sets of rules can be defined to detect abnormal combinations of packets, for example, login attempts that are not followed by further activity which may be a sign of an attempted break in. A network-based IDS may be run on computers dedicated to the task, though on a switched network some additional engineering may be needed to ensure such a computer can see packets addressed to others, or on individual computers that it is particularly important to protect. However, an IDS can require significant CPU and network resources, so rulesets should be kept simple if the computer is not dedicated to running the IDS. The Grid environment presents some challenges for network-based IDSs as many Grid flows are encrypted, thus preventing the content of the packets from being inspected. However, work has been done with the BRO [BRO] network-based IDS to enable it to interpret Globus authentication traffic and further work is planned to allow it to inspect even Globus encrypted traffic [Chan] and [NERSC]. The other common open-source network IDS program, Snort® [Snort®] may also be made Globus-aware by the same team.

Flow based IDS

A flow based IDS also examines traffic on networks, but looks at the volumes of traffic exchanged between particular ports and systems. On traditional networks these can be highly effective as patterns of flows are quite predictable. If a desktop workstation starts to generate very large outbound traffic flows then there is a strong likelihood that it is behaving abnormally. However, many Grid projects are based around very large network flows, so the same pattern of activity from a Grid workstation could not so easily be classed as anomalous.

Once the traffic is examined more closely it is likely that genuine Grid activity would consist of a small number of long lived, high traffic connections, whereas an intrusion would more often result in a large number of short lived connections. It seems likely that Grids and flow based Intrusion Detection Systems will be developed to work together, but at present this is still a research topic.

Although IDSs can be very useful, there are still a number of problems with current products. All types of IDS rely on having a way to distinguish normal from abnormal activity. This may be done either by defining abnormal activity and assuming that everything else is normal (the most common approach in network and flow based IDS), or by defining normal activity and assuming that everything else is abnormal (most commonly used by host based IDS). Network based IDSs using definitions of normal traffic are often referred to as performing 'anomaly detection'. This is more often seen in research papers than in products, and is likely to be even harder to define and implement in a Grid environment. Most network and flow based IDSs should more properly be called Attempt Detection Systems, since they will rarely record whether an attack was successful or not, only that the attempt was made. In particular, unless an intruder immediately makes use of a compromised computer, an IDS alone will seldom be able to give any certainty that an attack failed.

Finally, all current IDSs generate a number of false alarms, where some change to the normal pattern of activity is reported as being abnormal (known as a false positive). Even in a state-of-the-art system that collected information from many different sensors, five per cent of alarms were found to be false positives [Hwang]. This is inevitable in systems that work by heuristics that are designed with the priority on detecting attacks wherever possible. However, the high rate does raise problems in determining the best way to respond to an alarm from an IDS. Given the speed with which incidents can propagate around a Grid, the ideal response would be to block any attack automatically. However, repeatedly closing down a Grid node in response to false positives generated by an IDS may be unacceptable to its users. The alternative is to have a human Incident Response Team ready to receive the alarms and perform further analysis before taking appropriate action. This should give fewer shutdowns for false positives, but at the cost of greater disruption and remedial costs when a rapidly spreading attack does occur and is not blocked until manual action is taken.

Incident Response

No matter how carefully preventive measures are designed, implemented and used, Grids are sufficiently complex systems that security vulnerabilities will arise, and will be misused. Preparing to detect and respond to security incidents is therefore an essential part of deploying Grid systems. Many aspects of Grid incident response are the same as for more conventional networked systems, dealt with by CSIRTs (Computer Security Incident Response Teams). Grid incident response has some of the characteristics of local site incident response, as administrators are likely to have direct control of their computers, but there is also likely to be a need to co-ordinate incident response across a number of different networks and organisations (see, for example [OSG Incident]). Where Grid projects are producing software with widespread distribution, there is a similarity to the work done by vendor CSIRTs in identifying, correcting and distributing solutions to software vulnerabilities. These different types of CSIRT work are described in [CERT-CC SoP].

Grid incident response must also work with existing CSIRTs, for example, those responsible for the local and national networks that Grids use, to ensure that incidents are identified and

handled appropriately. For example a network CSIRT may regard a novel high-bandwidth network flow as indicating an incident, whereas a Grid project may consider it normal behaviour. Conversely Grids may be much more concerned about 'minor' incidents that have the potential for identity compromise. Unless these differences of viewpoint and response are resolved before incidents occur, incident response teams may find their actions actually harming legitimate users. The issues of co-ordinating the response to incidents for standard networked computers are already well-documented (for example, in [NIST] and [CERT-CC HB]) and, from a Grid perspective, in [Demchenko]. This section therefore concentrates on the differences that arise from the particular characteristics of Grids.

Grid Incidents

The purpose of Grids is to allow resource intensive collaborations across computers that may be on different networks in different parts of the world. However, the same facilities that allow the easy transfer of data and jobs can also be used to propagate the effects of a security breach just as quickly and easily. Security incidents are therefore likely to spread very widely and very fast within groups of computers and networks that are designed to collaborate. Grid systems inherently trust each other and a malicious attacker can easily take advantage of that trust. Past experience with Internet worms demonstrates how hard it is to eliminate a security problem once it becomes widespread. Unless the problem can be eliminated from all systems at once, there is a strong likelihood that recovered systems will be re-infected by those that have not yet been dealt with. Incident response on the Grid must therefore be able to detect incidents early, and react rapidly to contain their spread. Unfortunately, the most effective way to contain the spread of an incident – disabling all trust relations and isolating compromised systems and accounts – is also the most disruptive for legitimate users. In some cases it will indeed be necessary to 'turn off' a Grid until all its components can once again be trusted, but if an incident can be detected and its potential impact assessed quickly then it may be possible to avoid this extreme measure by establishing an effective security perimeter between trusted and compromised systems. Containment measures are likely to involve some disruption: the best that can be hoped for is to make this less than the potential damage caused by an uncontained security breach. Incident plans and playbooks that outline the steps to be taken when an incident occurs are a good way to achieve this and are being developed by a number of Grid projects.

Identity Incidents

One of the unique features of the Grid is the ability to access a widespread collection of resources using a single set of credentials. However, this same feature means that the compromise of an identity may be a much more serious issue for a Grid than it is for conventional systems where multiple credentials are still common. An intruder who obtains a copy of a Grid identity certificate may be able to access hundreds of systems worldwide, apparently quite legitimately. Even the loss of a proxy certificate may be serious, though these are normally time limited to make them harder to exploit. The compromise of an identity certificate may therefore constitute a serious Grid incident as described, for example, in [Skow] and reported in [Post].

Unfortunately, it will often be hard to be certain that a credential has been compromised. Many incidents will give the intruder the potential to compromise a credential, which need involve no more than accessing a file, but there will rarely be clear evidence that the file either

has, or has not, been read. In this, identity compromises differ from system compromises, where there will normally be some definite change to a disk file or memory location that shows clearly that the system has been compromised.

Given this uncertainty, a balance must be struck between the risk of continuing to allow use of a credential that may have been compromised and the disruption to users if credentials have to be replaced. If there is a possibility that a credential capable of causing serious harm has been compromised and could be used by an intruder, then that credential should be revoked and replaced. Assessing the risk of harm may be difficult, as the Grid's trust relations may allow harm to be caused at a location far away from the original incident, to an organisation that may have only a slight connection with the individual and systems on which the problem was detected. Developing guidelines to ensure the prompt and appropriate response to identity incidents is an urgent and challenging task for Grid incident response.

Source URL: <https://community-stg.jisc.ac.uk/library/janet-services-documentation/deployment-tools>

Links

[1] <http://community.ja.net/library/janet-services-documentation/appendix-specific-package-issues-and-solutions>

[2] <http://community.ja.net/system/files/images/tg-deployinggrids-03.jpg>