

Hack the Gibson

Exploiting Supercomputers

44Con Edition - September 2013

John Fitzpatrick Luke Jennings



Public

EXTERNAL

Labs.mwrinfosecurity.com | © MWR Labs



- Introduction
- Important Security Considerations
- Job Schedulers
 - MOAB
 - Torque
- OpenMPI
- Distributed Filesystems
 - GPFS
 - Lustre
- Cluster Management
- Local Privilege Escalation



Introduction

- Important Security Considerations
- Job Schedulers
 - MOAB
 - Torque
- OpenMPI
- Distributed Filesystems
 - GPFS
 - Lustre
- Cluster Management
- Local Privilege Escalation



What is High Performance Computing?

- Computer hardware and software aimed at solving very heavy computational tasks
- Common Uses
 - Weather forecasting
 - Data Mining
 - Cryptanalysis
 - Nuclear weapons simulation
 - Molecular dynamics





Some Numbers

Tianhe-2

Top500 rank: #1 (June 2013)

- China's National University of Defense Technology
- 33.86 petaflops/s (Linpack)
- 16,000 nodes (2x Intel Xeon IvyBridge + 3x Xeon Phi processors)
- Total 3,120,000 cores
- Memory: 1,375 TiB ~1,500 TB
- Storage: 12.4 PB
- Cost: ~£256 Million

Blue Joule (UK Fastest)

Top500 rank: #18 (June 2013)

- Daresbury Laboratory
- IBM Blue Gene/Q
- 1.25 petaflops/s (Linpack)
- Total 114,688 cores
- Storage ~5 PB



Parallel Computing - Limits

 Not all tasks can be parallelised e.g. Task B requires the result of Task A



Public



Message Passing Interface

- Achieves Cluster Computing Goal
- Implements the concept of "Distributed Shared Memory"

```
main(int argc, char **argv)
{
    int node;
    MPI_Init(&argc,&argv);
    MPI_Comm_rank(MPI_COMM_WORLD, &node);
    printf("I am %d\n",node);
    MPI_Finalize();
}
```



Public



Supercomputer vs Cluster

- What really makes a true "supercomputer" is the use of extremely high performance networking kit to link the nodes
- "Distributed Shared Memory" with normal Ethernet network connections will be orders of magnitude slower than RAM
- Infiniband and Myrinet
 - Upto 300 Gbit/s with link aggregation!
 - 1 µs latency (that's micro not milli)



Supercomputer – what about file systems?

- Distributed file systems are used
 - Much higher performance
- Each node will generally have access to identical data
 - Though its boot drive may be unique
- Examples
 - GPFS
 - Lustre



Supercomputer – Job Scheduling

- Not all nodes should be useable by anyone all the time.... Wait your turn!
- Generally a user may write a parallel program using MPI and then submit it to a job scheduler
- This then manages all the submissions and decides what runs on what nodes and what time etc.



Summary Architecture



Public EXTERNAL



- Introduction
- Important Security Considerations
- Job Schedulers
 - MOAB
 - Torque
- OpenMPI
- Distributed Filesystems
 - GPFS
 - Lustre
- Cluster Management
- Local Privilege Escalation



Important Security Considerations

Considerations

- Users can typically connect to login nodes
- Users run code on compute nodes
- Trust relationships exist (passwordless key based SSH)
- Job scheduling happens across many nodes
- Nodes may need regular rebuilding

Implications

- Local priv-esc is a very big deal
- Users have access to inner components
- Remote imaging is pretty much essential
- Several critical OS mechanisms are now network accessible

Public



- Introduction
- Important Security Considerations
- Job Schedulers
 - MOAB
 - Torque
- OpenMPI
- Distributed Filesystems
 - GPFS
 - Lustre
- Cluster Management
- Local Privilege Escalation



Job Schedulers

- Resource Managers e.g. Torque, SLURM
 - Handle actual running of jobs on a single cluster
 - Basic logic for scheduling e.g. FIFO, round robin etc
- Workload Managers (Schedulers) e.g. MOAB
 - Sit on top of resource managers
 - Can handle multiple clusters potentially with different resource managers
 - Intelligent, flexible, scalable scheduling



- Introduction
- Important Security Considerations
- Job Schedulers
 - MOAB
 - Torque
- OpenMPI
- Distributed Filesystems
 - GPFS
 - Lustre
- Cluster Management
- Local Privilege Escalation



- Submit jobs using the "msub" command
- This talks to a web server using XML

```
<job>
<user>jbloggs</user>
<group>jbloggs</group>
<command user="jbloggs" group="jbloggs">helloworld-proj</command>
<nodelist>node1,node2,node3</nodelist>
</job>
<hmac>2465ebc8cd6e412cdc1ab9fef40bcae6</hmac>
```

- HMAC verifies authenticity of user information submitted
- Default behaviour is insecure key to generate hmac is hardcoded in binary



- More secure option is to use "mauth" authentication
- Key is configured in a protected file

-r----- 1 root root 15 2007-04-05 03:47 /opt/moab/.moab.key

- When "msub" is called, it makes use of a SUID root binary called "mauth" to obtain the HMAC for the XML
- "mauth" can read the key (SUID root)
- If the user information submitted does not match the caller, it rejects the request
- Should be secure right?

Public



```
<job>
<user>jbloggs</user>
<group>jbloggs</group>
<command user="jbloggs" group="jbloggs">helloworld-proj</command>
<nodelist>node1,node2,node3</nodelist>
</job>
<hmac>2465ebc8cd6e412cdc1ab9fef40bcae6</hmac>
```





Public



- Web service and "mauth" check different user IDs ③
- By default, "root" jobs are not allowed
- However, web service allows "dynamic reconfiguration"
- Pretend to be "root" to enable root jobs, then submit after ③

Public



- Introduction
- Important Security Considerations
- Job Schedulers
 - MOAB
 - Torque
- OpenMPI
- Distributed Filesystems
 - GPFS
 - Lustre
- Cluster Management
- Local Privilege Escalation



Public EXTERNAL

Job Schedulers - Torque

- Resource manager
- Terascale Open-Source Resource and Queue Manager
- Can integrate with MAIU/MOAB
- Based on the original PBS project



2.4.X – support ended August 2012

2.5.X – Widely used, recommended, considered reliable 2.5.13 = Latest

3.0.X - has everything 2.5.x has plus it supports NUMA architectures

4.0.X - has significant enhancements to scalability for petaflop and large environments

4.1.X - has added significant enhancements for Cray systems

4.2.X – Includes support for Intel Xeon Phi 4.2.3.1 = Latest





3. Find node for job

Public



\$ qstat

Job id	Name	User	Time Use	S	Queue
				-	
85.server1	STDIN	hpc1	0	Q	batch
87.host1	STDIN	rd002	0	Q	batch
88.host2	STDIN	user1	0	Q	batch
89.host2	STDIN	user1	0	Q	batch
90.server1	STDIN	testuser	0	Q	batch
91.host1	STDIN	rd002	0	Q	batch





Public EXTERNAL



Public EXTERNAL

Job Schedulers - Torque



So let's hack the Gibson...

- unauthenticated remote root ;)



- Introduction
- Important Security Considerations
- Job Schedulers
 - MOAB
 - Torque
- OpenMPI
- Distributed Filesystems
 - GPFS
 - Lustre
- Cluster Management
- Local Privilege Escalation



Running standard program e.g. uname

root@bt5:~/mpi# mpirun -H node,localhost uname -a Linux bt5 3.2.6 #1 SMP Fri Feb 17 10:34:20 EST 2012 x86_64 GNU/Linux Linux bt 2.6.39.4 #1 SMP Thu Aug 18 13:38:02 NZST 2011 i686 GNU/Linux

Running custom MPI program

```
root@bt5:~/mpi# mpirun -H node,localhost helloworld
Hello from process 0 of 2 running on bt (0)
Processor 0 sending buf data
Hello from process 1 of 2 running on bt5 (0)
Processor 1 received buf data: Hello :)
Process 0 received reply: Wagwan!
```



/usr/bin/ssh -x node orted

- --daemonize
- -mca ess env
- -mca orte_ess_jobid 76873728
- -mca orte_ess_vpid 1
- -mca orte_ess_num_procs 2
- --hnp-uri "76873728.0;tcp://192.168.154.130:43622"

	00000000 00 00000010 2f	00 00 00 00 00 00 00 00 9c 00 00 00 00 00 00 00	2† 90 00 00	00 00 00 00 00 00 01 / 00 02 00 00 00 00 /	Public
LABS	000000000000000000000000000000000000000	00 00 00 00 00 00 00 00 2f 9c 00 00 00 00 00 00	0 00 Z	f 9c 00 00 00 00 00 00 /	EXTERNAL
Open MPI	00000020 00000024 2f 00000034 2f 00000044 00 00000054 00 00000064 37 00000074 2f 00000084 31	00 00 00 00 9c 00 00 00 00 01 9c 00 00 00 00 00 01 9c 00 00 00 00 00 00 00 00 00 0a 2f 9c 00 00 0a 00 00 00 00 00 00 0a 39 39 38 37 35 32 37 36 2f 31 39 32 2e 31 36 3a 34 38 38 35 00	2f 9c 00 00 00 00 38 2e 38 2e 00 00	00 00 00 00 00 01 / / 00 04 00 00 00 4d / /	
	00000094 00 00000024 00000034 00000054 00000064 00000064 00000004 00000004 00000004 00000004 000000	$\begin{array}{cccccccccccccccccccccccccccccccccccc$	0 00 2 0 01 0 0 00 0 0 00 0 0 00 0 0 00 0 0 00 0 0 00 0	f 9c 00 00 00 00 00 / 00 00 00 04 00 00 02 11 / 00 00 00 02 11 / 00 00 00 2f 9c 00 00 00 00 00 01 09 00 00 00 00 00 01 00 00 00 00 44 00 00 01 00 00 00 04 62 74 35 00 00 00 04 62 74 35 00 00 00 00 05 6e 6f 64 65 00	 e. 8. 8. 8. 8. 8.
	00000194 000001A4 000001B4 000001C4 000001D4 000001E4 00000214 00000224 00000234 00000244	6c 6c 6f 77 6f 72 6c 48 4f 4d 50 49 5f 4d 72 65 63 6f 6e 64 69 73 70 6f 72 74 73 3d 62 62 32 65 66 39 62 66 35 66 61 33 34 38 6f 74 2f 6d 70 69 00 00 00 00 01 00 00 00 61 6c 68 6f 73 74 00 00 00 00 2b 00 00 00 00 00 00 02 00 00 00 00 00 00 00 00 00 00	c 64 0 d 43 4 9 74 6 d 31 3 2 2d 3 8 31 0 0 00 0 0 0f 6 0 00 0 0 01 0 0 00 0 0 00 0	00 00 00 01 00 00 01 100 00 00 110world 1 5f 6f 72 74 65 5f 70 HOMPI_MC A_orte 9 6f 6e 5f 74 72 61 6e recondit ion_tr 66 1 39 64 35 34 37 64 sports=1 oasd54 64 38 30 35 30 34 36 64 bb2ef9b- 480504 00 00 00 00 00 00 ot/mpi. 00 00 00 00 00 00 00 00 00 00 00	 7d 6d ro oc







- Vulnerable to a MITM...but most MITM attacks require root privileges e.g. ARP spoofing
- If we have root on a node, we have root on everything anyway
- The "Super-fast" network is only going to be physically accessible from within the data centre
- ...so is this vulnerability actually applicable to the environment?



- Exploitable as a low privileged user?
 - What if I connect to the port first, before the remote node?
 - Will the master node give up listening, so I can hijack the port?
 - Race condition?
- Can I win the race?
 - Legitimate workflow involves an SSH connection then a TCP connect back
 - My code can all run on the master node
 - No network comms = I have the speed advantage



• The port we could brute force, but we need that ID too...

--hnp-uri "76873728.0;tcp://192.168.154.130:43622"

 ...but Linux normally gives away command line arguments via /proc file system e.g. "ps aux" below

```
daemon
           833
                0.0
                     0.0
                          18932
                                   392 ?
                                                Ss
                                                     Nov05
                                                              0:00 atd
           834
                0.0
                     0.0 21124
                                  1024 ?
                                                     Nov05
                                                              0:15 cron
root.
                                                Ss
                                   640 ?
                                                             0:12 /usr/sbin/irqbalance
           835
                     0.0 11332
                                                     Nov05
                0.0
                                                Ss
root
           844
                0.0
                     0.1 120676
                                  3644 ?
                                                     Nov05
                                                              0:00 /usr/sbin/console-kit-daemon --no-daemon
root.
                                                Sl
           965
                0.0
                     0.1 49356
                                  4080 ?
                                                S
                                                     Nov05
                                                              0:00 /opt/metasploit/postgresgl/bin/postgres -D
postgres
/opt/metasploit/postgresgl/data -p 7337
postgres 1278
                0.0
                     0.0
                          49356
                                  1244 ?
                                                Ss
                                                     Nov05
                                                              0:20 postgres: writer process
postgres 1279
                0.0
                     0.0
                          49356
                                   976 ?
                                                     Nov05
                                                              0:17 postgres: wal writer process
                                                Ss
                                                              0:19 postgres: autovacuum launcher process
         1280
postgres
                0.0
                     0.0
                          49492
                                  1180 ?
                                                Ss
                                                     Nov05
postgres 1281
                0.0
                     0.0
                          20764
                                  1008 ?
                                                              0:03 postgres: stats collector process
                                                Ss
                                                     Nov05
root.
          1303
                0.0
                    0.0 59480
                                   716 ?
                                                Ssl
                                                     Nov05
                                                              0:00 /usr/sbin/vmware-vmblock-fuse -o
subtype=vmware-vmblock, default permissions, allow other /var/run/vmblock-fuse
          1359
                          85000
                                  4008 ?
                                                     Nov05
                                                              1:28 /usr/sbin/vmtoolsd
root
                0.0 0.1
                                                S
          1423
                0.0 0.0
                                     0 ?
                                                S
                                                     Nov05
                                                              0:05 [flush-8:0]
root.
                               0
                                                              0:15 dhclient3 -e IF METRIC=100 -pf
          1439
               0.0 0.0
                                   664 ?
                                                     Nov05
root
                           6604
                                                Ss
/var/run/dhclient.eth1.pid -lf /var/lib/dhcp3/dhclient.eth1.leases eth1
          1560 0.0 0.0 27912 1580 ?
                                                     Nov05
                                                              0:06 /usr/sbin/ntpd -p /var/run/ntpd.pid -g -u
                                                Ss
ntp
113:121
```



The Exploit

- 1. Monitor /proc file system
- 2. When the ssh connection is detected, read the port number and application specific ID
- 3. Connect to port and "talk" MPI to it along with the correct ID so the master node is happy execution has started
- 4. Master node will then release the port, so we can listen on it
- 5. If we are quick enough, we receive the connect back from the remote node and "talk" MPI and issue our own "command"



Public EXTERNAL

Open MPI

The Exploit

DEMO



- Introduction
- Important Security Considerations
- Job Schedulers
 - MOAB
 - Torque
- OpenMPI
- Distributed Filesystems
 - GPFS
 - Lustre
- Cluster Management
- Local Privilege Escalation



Public EXTERNAL

Distributed Filesystems

- Gbit/sec
- Infiniband
- Petabytes
- Tens of thousands of disks



Public EXTERNAL

Distributed Filesystems

- Commonly used:
 - GPFS
 - NFS
 - Lustre
- Performance is Key



- Introduction
- Important Security Considerations
- Job Schedulers
 - MOAB
 - Torque
- OpenMPI
- Distributed Filesystems
 - GPFS
 - Lustre
- Cluster Management
- Local Privilege Escalation



Distributed Filesystems - GPFS

IBM – file system

"It's quite expensive - we bought an instance of GPFS from IBM and they gave us a free supercomputer with it"

Good tech overview here:

https://computing.llnl.gov/code/sio/ucrl-web-148419-pt-1.pdf





Distributed Filesystems - GPFS

- application node node running user app that accesses mounted GPFS
- VSD server node Virtual Shared Disk node (or I/O node) with disks attached
- RIO Remote I/O that allows for a connection between the server node and the disks
- SSA loop Serial Storage Architecture disks may be connected in a loop between nodes, allowing for failover:



(This is in contrast to a twin-tailed approach which is used for SCSI [Small Computer System Interface]disks to connect two nodes to a single disk.)





Bill Loewe: https://computing.llnl.gov/code/sio/ucrl-web-148419-pt-1.pdf



- Introduction
- Important Security Considerations
- Job Schedulers
 - MOAB
 - Torque
- OpenMPI
- Distributed Filesystems
 - GPFS
 - Lustre
- Cluster Management
- Local Privilege Escalation



Distributed Filesystems - Lustre

- Widely used
- Over 60 of the top 100 supercomputers use Lustre
- Resold in appliances



Distributed Filesystems - Lustre





Public EXTERNAL

Distributed Filesystems - Lustre

The Exploit (found by reading the documentation!)

DEMO



- Introduction
- Important Security Considerations
- Job Schedulers
 - MOAB
 - Torque
- OpenMPI
- Distributed Filesystems
 - GPFS
 - Lustre
- Cluster Management
- Local Privilege Escalation



Cluster Management

• Tens of thousands of nodes = remote imaging is the norm

• Used to this sort of technology from corporate environments e.g. PXE booting rebuilds etc.

- We've seen a couple of vendor implementations
 - Both used PXE + TFTP (as is common)
 - One used custom protocols for the second stage
 - The other used common protocols like rsync

Public



Cluster Management

- The nature of this type of technology makes it very difficult to secure directly
- But two things are key:
 - Read-only images -> otherwise suffer backdoors
 - No sensitive information in the images -> work to the assumption of no confidentiality
- Sensitive information
 - E.g. root's private SSH key or MOAB key files
 - This information can be added securely as a post imaging step within scripts



- Introduction
- Important Security Considerations
- Job Schedulers
 - MOAB
 - Torque
- OpenMPI
- Distributed Filesystems
 - GPFS
 - Lustre
- Cluster Management
- Local Privilege Escalation



Local Privilege Escalation

- When all your users run code on "servers", this is important
- Besides patch management and usual UNIX file permissions etc., pay attention to manufacturer added attack surface, e.g. SUID binaries
- On one manufacturer default install we found two non-standard SUID root binaries
- One looked particularly interesting as it had an "extract to directory" option



Local Privilege Escalation

When run normally it lists something related to the local directory contents
 0930292390394920testfile1.txt

0923902393902320testfile3.txt 2982983493898498testfile2.txt

• After much investigation, it turns out the cryptic numbers encode the ownership info and file permissions

• After further investigation, when supplied to the "extract" option, it modifies the local directory contents to match



Local Privilege Escalation

• ...Oh dear! So if we encode something like this on our own system...



- ...and then extract the output on the target system
- We get a SUID root shell. Nice ③
- Neither SUID binary even appeared to be needed, let alone require SUID permissions!
 - Remove unnecessary attack surface

Public



Summary

- Interesting Challenges
 - Performance and security
 - Authenticating users across multiple nodes securely



Questions?

http://labs.mwrinfosecurity.com | @mwrlabs

John.fitzpatrick@mwrinfosecurity.com | @j0hn__f

Luke.jennings@mwrinfosecurity.com