



1967-30

Advanced School in High Performance and GRID Computing

3 - 14 November 2008

How to set-up a Queue system on your cluster

CALUCCI Piero
S.I.S.S.A.
International School for Advanced Studies
Via Beirut 2-4
34014 Trieste
ITALY

Piero Calucci

Obtaining and compiling TORQUE and Maui

Configuration

Diagnostics & Troubleshooting

Setting up Queue Systems with TORQUE & Maui

Piero Calucci

Scuola Internazionale Superiore di Studi Avanzati
Trieste

November 2008
Advanced School
in High Performance and Grid Computing



Piero Calucci

Obtaining and compiling TORQUE and Maui

Configuration

Diagnostics & Troubleshooting

Outline

1 Obtaining and compiling TORQUE and Maui

2 Configuration



3 Diagnostics & Troubleshooting DEMOCRITOS/SISSA



Piero Calucci

Obtaining and compiling TORQUE and Maui

Obtaining and compiling TORQUE

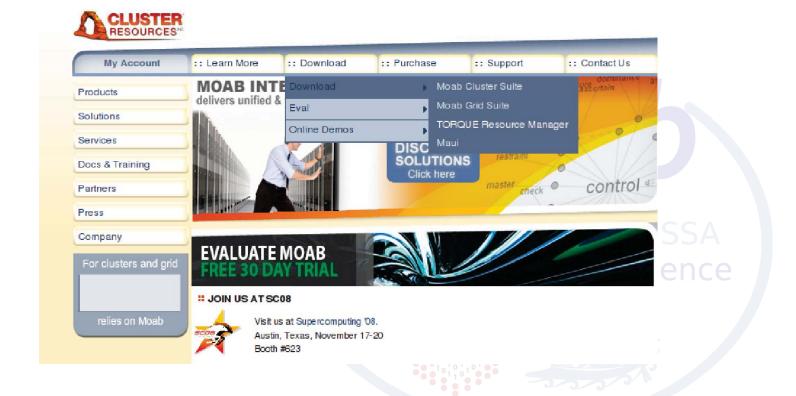
Obtaining and compiling Maui

Configuration

Diagnostics & Troubleshooting

TORQUE Source Code

TORQUE is available from www.clusterresources.com



Piero Calucci

Obtaining and compiling TORQUE and Maui

Obtaining and compiling TORQUE

Obtaining and compiling Maui

Configuration

Diagnostics & Troubleshooting

Building TORQUE

- configure --prefix=/whatever/you/like make
 make install
- not very clean, actually:
 quite a lot of important files go into /var/spool
 including configuration files!

You can build only the server or MOM components, just tell --disable-mom or --disable-server

My favorite install uses a directory that is shared among the masternode and the computing nodes, so that I need to build only once.

Piero Calucci

Obtaining and compiling TORQUE and Maui

Obtaining and compiling TORQUE

Obtaining and compiling Maui

Configuration

Diagnostics & Troubleshooting

Maui Source Code

Maui too is available from www.clusterresources.com You need to register to their site to download the code, and they may contact you later and ask what are you going to do with their software (and offer commercial support for it)

MAUI CLUSTER SCHE	EDULER DOWNL	OADS		
manageability and efficiency of m configurable tool capable of supp	aui Scheduler) is our first g nachines ranging from clus porting a large array of sch nost advanced scheduler in	generation cluster scheduler. Maui is an adv sters of a few processors to multi-teraflop su eduling policies, dynamic priorities, extensi in the world', and is currently in use at hundn		
** AVAILABLE DOWNLOAD The patch releases below are the				
Maui 3.2.6 - Patch 21 Maui 3.2.6 - Patch 20 Maui 3.2.6 - Patch 19 Maui 3.2.6 - Patch 18 Maui 3.2.6 - Patch 17	This new system will str may experience at this: Username: Password: Remember Me: Login If you have forgotten you	menting a new login system. During the roll of reamline much of the download process and time. calucci	will make things easier in the future	ence

◆□▶ ◆□▶ ◆■▶ ● り ()

Piero Calucci

Obtaining and compiling TORQUE and Maui

Obtaining and compiling TORQUE

Obtaining and compiling Maui

Configuration

Diagnostics & Troubleshooting

Building Maui

- same "configure; make; make install"
- maui build system need to know where TORQUE has been installed
- again, important files go into /var/spool

Joint DEMOCRITOS/SISSA Laboratory for @-Science



Piero Calucci

Obtaining and compiling TORQUE and Maui

Configuration

TORQUE Configuration

Maui Configuration
Prologue and
Epilogue

Diagnostics & Troubleshooting

TORQUE Common Configuration Files

 pbs_environment contains the environment variables for TORQUE; any minimal set will do e.g.

 server_name contains the «official» name of the machine where pbs_server runs (this is usually your master node)

The server name must be identical to the FQDN e.g.

Both these files reside in the spool directory (/var/spool/torque)

Piero Calucci

Obtaining and compiling TORQUE and Maui

Configuration

TORQUE Configuration

Maui Configuration
Prologue and
Epilogue

Diagnostics & Troubleshooting

pbs_server configuration

The nodes file

server_priv/nodes contains the list of available computing nodes and a list of attributes for each node.

node name	# of CPUs	«features» (list of arbitrary strings, can be used later to select a node type)
node01 node02	np=2 np=2	opteron myri opteron myri Laboratory for e-Science
node51 node52	np=4 np=4	opteron IB opteron IB

Piero Calucci

Obtaining and compiling TORQUE and Maui

Configuration

TORQUE Configuration

Maui Configuration
Prologue and
Epilogue

Diagnostics & Troubleshooting

pbs_server configuration

Creating the Configuration Database

The bulk of pbs_server configuration is written in a (binary) database. You first need to create the empty database with

pbs_server -t create

This will destroy any existing configuration, create the empty database and start a pbs_server.

Configuration can then be edited using the qmgr tool.

Configuration data are written to server_priv/serverdb as well as in various other files.

Piero Calucci

Obtaining and compiling TORQUE and Maui

Configuration

TORQUE Configuration

Maui Configuration
Prologue and
Epilogue

Diagnostics & Troubleshooting

pbs_server configuration Security Note

qmgr doesn't actually edit the configuration database. It only sends configuration commands to pbs_server which in turn writes the configuration.

This means that:

- you need a running pbs_server to use qmgr (no big issue)
- the pbs_server process needs write access to its own configuration files this is usually considered very bad in any security-conscious environment – unfortunately no easy workarounds are available

Piero Calucci

Obtaining and compiling TORQUE and Maui

Configuration

TORQUE Configuration

Maui Configuration
Prologue and
Epilogue

Diagnostics & Troubleshooting

pbs_server configuration

Sample Configuration

[root@borg]# qmgr

```
Qmqr:
      create queue batch
       set queue batch queue_type = Execution
Qmgr:
       set queue batch resources_max.walltime = 01:00:00
Qmgr:
       set queue batch resources default.nodes = 1
Qmgr:
       set queue batch resources default.walltime = 00:01:00
Qmqr:
Omar:
       set queue batch enabled = True
       set queue batch started = True
Qmgr:
       set server managers = maui@borg.cluster
Qmqr:
       set server managers += root@borg.cluster
Qmqr:
       set server operators = maui@borq.cluster
Qmgr:
       set server operators += root@borq.cluster
Omar:
```

One of the most common configuration issues, that prevents the batch system from running any job, involves missing or incorrect set server managers and/or set server operators lines.

Piero Calucci

Obtaining and compiling TORQUE and Maui

Configuration

TORQUE Configuration

Maui Configuration
Prologue and
Epilogue

Diagnostics & Troubleshooting

pbs_mom configuration

pbs_mom configuration can be fairly minimal, the only thing the Mom needs to know is the hostname where pbs_server is running on.

Useful additions include log configuration, how to handle user file copy and which filesystem to monitor for available space.

```
mom_priv/config:
```

```
$clienthost master.hpc
$logevent 0x7f
$usecp *:/home /home
size[fs=/local_scratch]
```

Piero Calucci

Obtaining and compiling TORQUE and Maui

Configuration

TORQUE Configuration

Maui Configuration

Prologue and Epilogue

Diagnostics & Troubleshooting

Maui Configuration

How to Connect to Resource Manager

- simpler approach: a single configuration file (maui.cfg)
- Maui needs to know what RM to connect to and how

SERVERHOST

RMCFG[BORG.CLUSTER]

RMPOLLINTERVAL

SERVERPORT

SERVERMODE

ADMIN1

borg.cluster

TYPE=PBS

00:00:30

42559

NORMAL

root



Piero Calucci

Obtaining and compiling TORQUE and Maui

Configuration

TORQUE Configuration

Maui Configuration

Prologue and Epilogue

Diagnostics & Troubleshooting

Maui Configuration

Job Prioritization

Job priority is recomputed at each scheduler iteration, according to site-defined parameters. If no parameters are set only queue time is taken into account, i.e. the scheduling is strictly FIFO.

Priority components include:

- Queue Time: how long the job has been idle in the queue
- Credentials: a static priority can be assigned on a user, group, queue basis
- Fair Share: historical usage data
- Resources requested for the job



Piero Calucci

Obtaining and compiling TORQUE and Maui

Configuration

TORQUE Configuration

Maui Configuration

Prologue and Epilogue

Diagnostics & Troubleshooting

Maui Configuration

Job Prioritization: Queue Time and Credentials

QUEUETIMEWEIGHT XFACTORWEIGHT

CLASSCFG[batch]

CLASSCFG[fast]

GROUPCFG[quests]

GROUPCFG[users]

GROUPCFG[devel]

USERCFG [DEFAULT]

USERCFG[luser1]

10

PRIORITY=1

PRIORITY=1000

PRIORITY=1

PRIORITY=1000

PRIORITY=10000

PRIORITY=2000

PRIORITY=0



Piero Calucci

Obtaining and compiling TORQUE and Maui

Configuration

TORQUE Configuration

Maui Configuration

Prologue and Epilogue

Diagnostics & Troubleshooting

Maui Configuration

Job Prioritization: Fair Share

The FS priority component must be explicitly enabled by setting its weight to a non-zero value.

FSINTERVAL	8640
FSDEPTH	30
FSDECAY	0.90
FSWEIGHT	1
FSGROUPWEIGHT	240
FSUSERWEIGHT	10

duration of each FS window number of FS windows decay factor applied to older FS windows





Piero Calucci

Obtaining and compiling TORQUE and Maui

Configuration

TORQUE Configuration

Maui Configuration

Prologue and Epilogue

Diagnostics & Troubleshooting

Maui Configuration

Job Prioritization: Fair Share

Usage targets can be set on a per-user, per-group and per-queue basis.

USERCFG[DEFAULT] FSTARGET=1

GROUPCFG[users] FSTARGET=30

GROUPCFG[devel] FSTARGET=40

USERCFG[master]

You can set also FS floors or caps so that priority is affected only when usage drops below the floor or goes above the cap:

GROUPCFG[guests] FSTARGET=5-

give a negative priority

component if usage is

above 5%

FSTARGET=20+

give a priority boost if

usage is below 20%



Piero Calucci

Obtaining and compiling TORQUE and Maui

Configuration TORQUE

Configuration

Maui Configuration

Prologue and Epilogue

Diagnostics & Troubleshooting

Prologue & Epilogue scripts

pbs_mom looks for scripts in its configuration directory mom_priv. If found, the prologue script is executed just before job start and the epilogue script at job termination. The prologue script performs any initialization that is requered on the node for the job to run, while the epilogue undoes the modifications.

/etc/security/access.conf

before prologue

-: ALL EXCEPT

root:ALL

disallows login to everybody except root, from anywhere after prologue

ightarrow -:ALL EXCEPT root

someuser:ALL

now allows someuser to

login



Piero Calucci

Obtaining and compiling TORQUE and Maui

Configuration

Diagnostics & Troubleshooting

TORQUE Diagnostics

Maui Diagnostics

momctl

Query and control remote pbs_mom:

momctl -d3 -h i602

diagnostics complete

```
master.hpc Version: 1.2.0p6
      i602/i602.hpc Server:
Host:
 HomeDirectory:
                           /var/spool/PBS/mom priv
                           6907718 seconds
 MOM active:
 Last Msq From Server:
                           213582 seconds (DeleteJob)
                           1 seconds
 Last Msq To Server:
                          45 seconds
 Server Update Interval:
                           10 hellos/2 cluster-addrs
 Init Msgs Received:
 Init Msgs Sent:
                           190 hellos
                           0 (use SIGUSR1/SIGUSR2 to adjust)
 LOGIEVEL:
 Communication Model:
                           RPP
                           20 seconds
 TCP Timeout:
 Prolog Alarm Time:
                           300 seconds
                           0 of 10 seconds
 Alarm Time:
 Trusted Client List:
 JobList:
                           NONE
```

Piero Calucci

Obtaining and compiling TORQUE and Maui

Configuration

Diagnostics & Troubleshooting

TORQUE Diagnostics

Maui Diagnostics

checknode

Check who is doing what on a node and show node capabilities

```
# checknode a034
```

```
checking node a034
State: Busy (in current state for 1:13:38:12)
Configured Resources: PROCS: 2 MEM: 3949M SWAP: 7242M DISK:
59G
Utilized Resources: PROCS: 2 DISK: 10G
Dedicated Resources: PROCS: 2
Opsys: DEFAULT Arch: [NONE]
Speed: 1.00 Load: 2.000 (ProcSpeed:
                                     2600)
Network: [DEFAULT]
          [myri][opteron][opteron-sc]...
Attributes: [Batch]
Classes: [smp2 2:2][smp4 2:2][mpi4 0:2][mpi8 2:2]...
Total Time: 25:14:33:36 Active: 25:04:53:26 (98.43%)
Reservations:
Job '30069' (x2) -1:13:38:44 -> 2:10:20:16 (3:23:59:00)
JobList: 30069
```

Piero Calucci

Obtaining and compiling TORQUE and Maui

Configuration

Diagnostics & Troubleshooting

TORQUE Diagnostics

Maui Diagnostics



Elab

calucci@sissa.it>
Laboratory for e -Science

