



Enabling Grids for E-science

Torque testing and extensions to other batch systems

Nikos Voutsinas, Dimitris Apostolou

UoA

Ioannis Liabotis

GRNET

*SA3 All Hands meeting
Barcelona, 21-24 May 2007*

www.eu-egee.org
www.glite.org



- It all started from the certification of some Torque/Maui related patches
- We needed a methodology to streamline these tests and be able to thoroughly test the batch system and be able to perform the same methodology when needed
- The methodology should be useful for other batch systems tests
- **Outline**
 - Test Planning
 - Test Framework
 - Performed Test
 - Results

- **Two Types of tests**
 - Test batch system as part of the full gLite stack
 - Using gLite style job submission
 - Test batch system as a standalone component
 - Using batch system specific commands
- **Installation and Configuration Certification**
 - 1. The installation procedure doesn't fail (Procedure Completeness Tests)
 - 2. The installation and configuration procedure brings the system to an acceptable initial state (Install Integrity Tests)
 - 3. The upgrade from the previous version works as expected, ie. without breaking existing functionality and system configuration (Upgrade Completeness Tests)

- **System Operation Certification:**
 - 1. The system operates in accordance with the functional requirements (Functionality Tests)
 - 2. The system fails at established peak load conditions. Subject to extreme data and event traffic. (Stress Tests)
 - 3. The system scales at an established rate. (Scalability Tests)
 - 4. The system is resistant against compromise attempts. Security Tests)
 - 5. The system meets its performance requirements under various circumstances. Performance tests)
 - 6. The system uses an established amount of resources (memory, disk, network, ..).(Resource usage tests)
 - 7. The system works properly with other applications. Compatibility Tests)
 - 8. The system is available over an extended period or number of requests Reliability/Availability Tests)
 - 9. The system after failure recovers to a previous working state with minimum loss of information (Resilience Tests)

- **JDLs**
 - Various types of submitting jobs have been used to simulate different types of loads
- **Glite Job submission process**
 - Multiple UIs and multiple WMS are used to submit jobs to the LRMS that is being tested.
- **Batch System job submission process**
 - Jobs are submitted directly to the LRMS that is being tested.
- **Monitoring process and Results Archiving**
 - LRMS processes and system nodes are monitored using simple scripts to facilitate results validation and verification
- **Certification documents and report**
 - The certification method, the utilities used and the certification results.

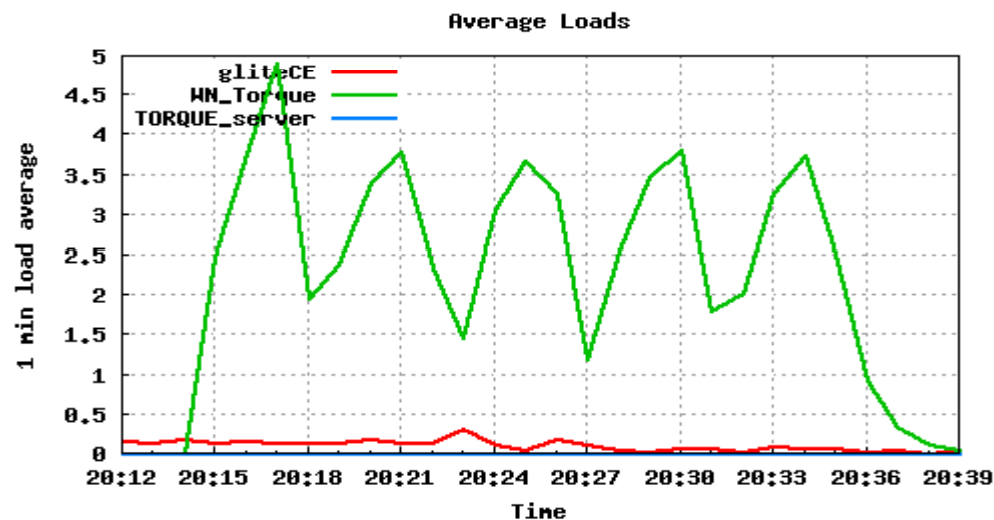
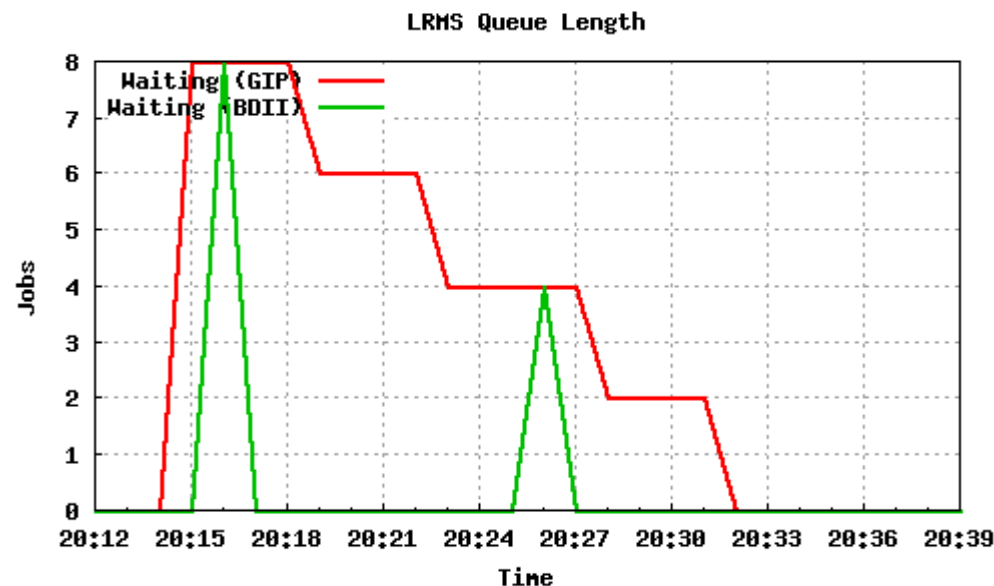
- **Limitations**
 - Number of available worker nodes (1 in our case)
 - Number of queues (1 in our case)
 - User Certificates (1 in our case)
- **Tools**
 - Bash shell scripts
 - Testing hardware (Certification Site and submission node)
- **Presentation and archival**
 - Bash scripts
 - Cron Jobs
 - Text files/Database
- **Monitored Data**

Nodes	data types							
	Load Average	%cpu,size				Jobs Status		
		pbs_server	pbs_mom	maui	BLParserPBS	GIP	BDII	qstat
gliteCE/SiteBDII	X					X	X	
Worker	X		X					
Torque/Maui	X	X		X	X			X

- **Site Used**
 - **EGEE-SEE-CERT**
- **Hosts**
 - **ctb03.gridctb.uoa.gr**
 - **site BDII**
 - **gliteCE**
 - **ctb07.gridctb.uoa.gr**
 - **Torque head node**

- **check Batch System information published through BDII**
 - simple Idapsearch commands and observation of correct results
- **check Batch System configuration**
 - View configuration files and batch system configuration using various commands
 - `qmgr -c 'print server'`
- **check network ports and services**
 - Check which Batch System related services run on which ports, on related nodes (gliteCE, TORQUE_server, WN_torque)
 - `netstat`
- **check logging**
 - Check all relevant log files during installation, startup and operation of services
- **checking General Information Providers (GIPs)**
 - Watch Information Providers' results while running long lived, CPU intensive jobs. Compare with results received directly from the LRMS.
 - `qstat`, `showstate`, `bdii` information

- job submission of few long lived, cpu intensive jobs
 - In particular submission of 10 jobs through the following path:
 - UI -> WMS -> gliteCE -> Torque_server -> WN
- **Results:**
 - All jobs done
- **Why BDII is updated both from rgma and edguser in different time intervals?**

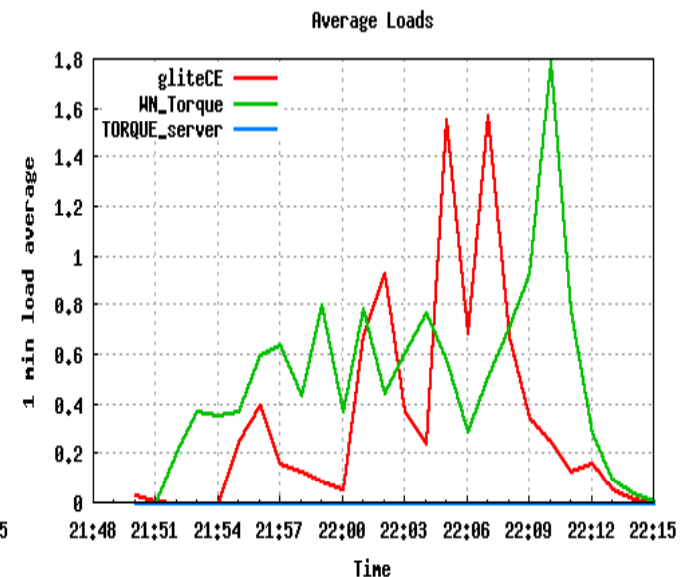
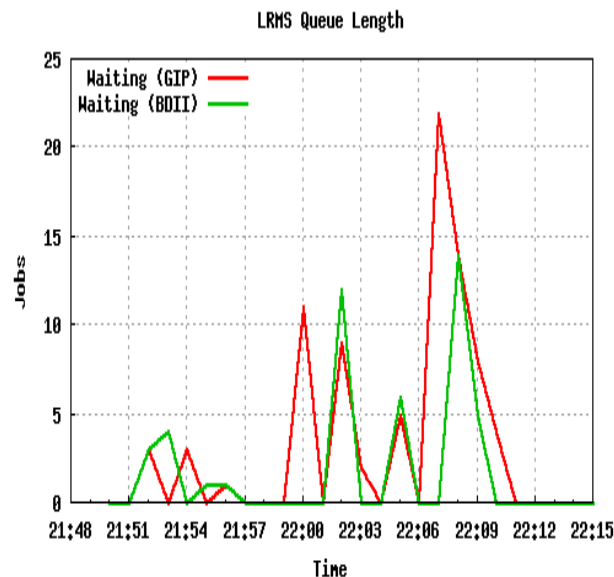
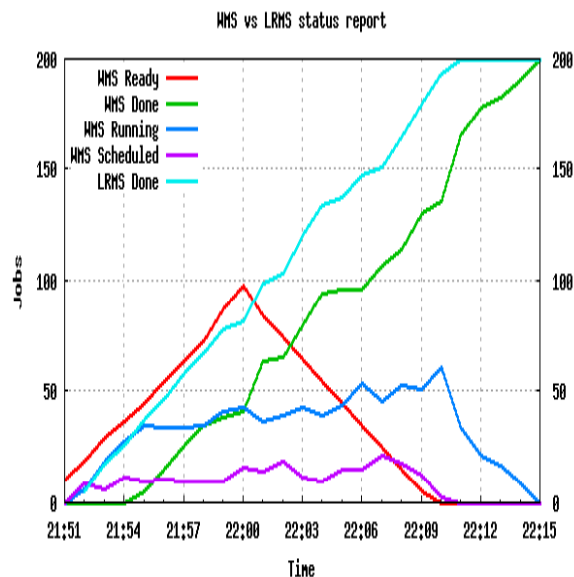


- **200 job submissions using 1 WMS**

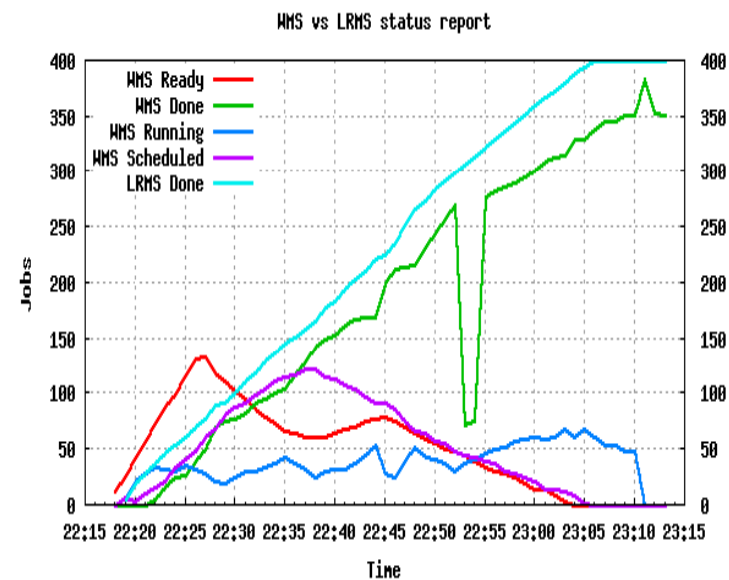
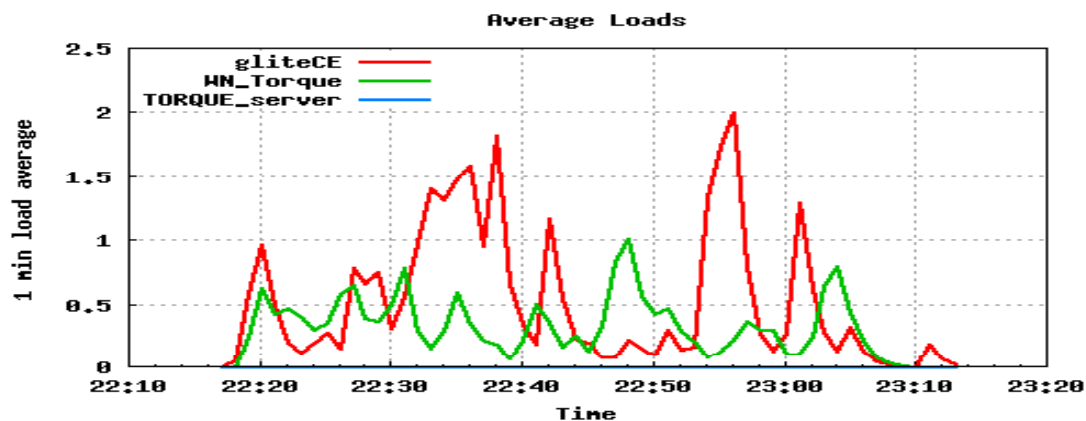
- Job submission of many short lived (instantaneous) simple jobs. In particular 100 job submissions using each one of the following paths:

- Path 1: UI-1 -> WMS-1 -> gliteCE -> Torque_server -> WN
- Path 2: UI-2 -> WMS-1 -> gliteCE -> Torque_server -> WN

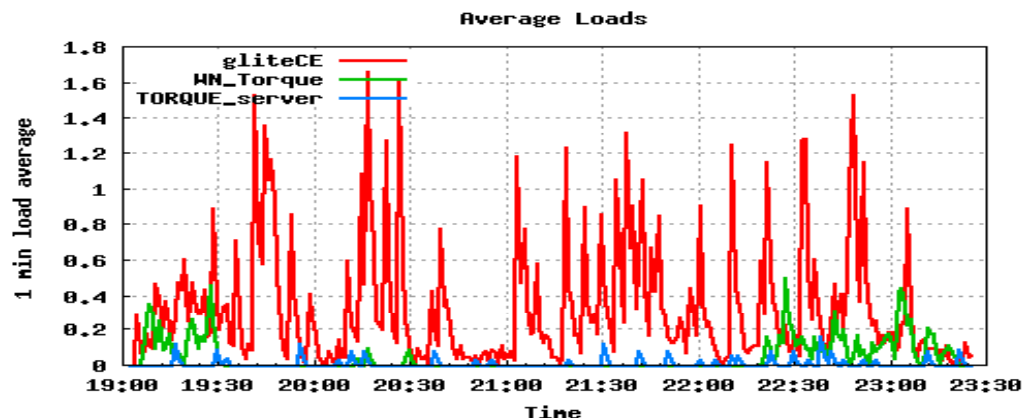
- Results: all Jobs done



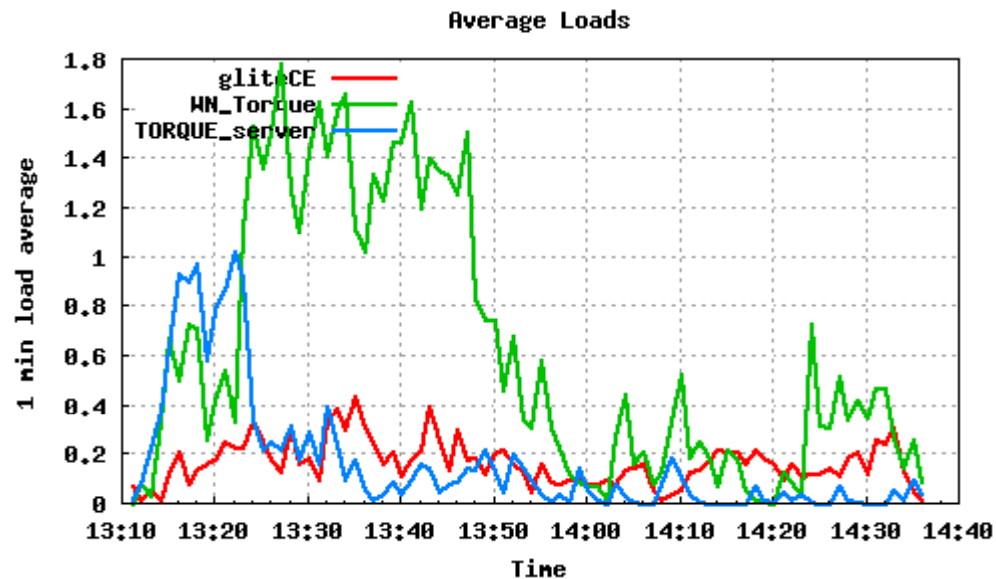
- **400 job submissions using 2 WMS**
 - Job submission of many short lived (instantaneous) simple jobs. In particular 100 job submissions using each one of the following paths:
 - Path 1: UI-1 -> WMS-1 -> gliteCE -> Torque_server -> WN
 - Path 2: UI-2 -> WMS-1 -> gliteCE -> Torque_server -> WN
 - Path 3: UI-1 -> WMS-3 -> gliteCE -> Torque_server -> WN
 - Path 2: UI-2 -> WMS-3 -> gliteCE -> Torque_server -> WN
 - Results: Test Results: **351 Jobs Done**
 - *PATH 1: Done=100*
 - *PATH 2: Done=100*
 - *PATH 3: Done=74 Aborted=26*
 - *PATH 4: Done=77 Aborted=23*



- **600 job submissions using 3 WMS**
 - Job submission of many short lived (instantaneous) simple jobs. In particular 100 job submissions using each one of the following paths:
 - Path 1: UI-1 -> WMS-1 -> gliteCE -> Torque_server -> WN
 - Path 2: UI-2 -> WMS-1 -> gliteCE -> Torque_server -> WN
 - Path 3: UI-1 -> WMS-2 -> gliteCE -> Torque_server -> WN
 - Path 4: UI-2 -> WMS-2 -> gliteCE -> Torque_server -> WN
 - Path 5: UI-1 -> WMS-3 -> gliteCE -> Torque_server -> WN
 - Path 6: UI-2 -> WMS-3 -> gliteCE -> Torque_server -> WN
 - Results:
 - *PATH 1: Done=99 Aborted=1*
 - *PATH 2: Done=99 Aborted=1, PATH 3: Done=86 Aborted=13 Waiting=1, PATH 4: Done=78 Aborted=22, PATH 5: Done=13 Aborted=59. PATH 6: Done=7 Aborted=65 Waiting=1*
 - **WMS problems detected**

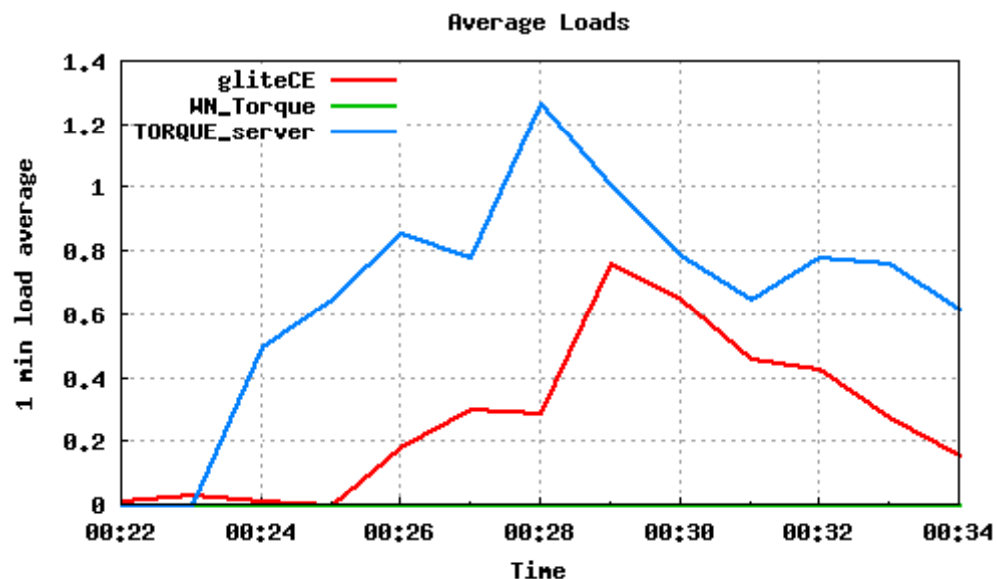


- **job submissions directly to Torque**
 - 1000 simple jobs were submitted, followed by 2000, 4000 and finally 5000 jobs.
 - E.g. 5000 Jobs directly to torque

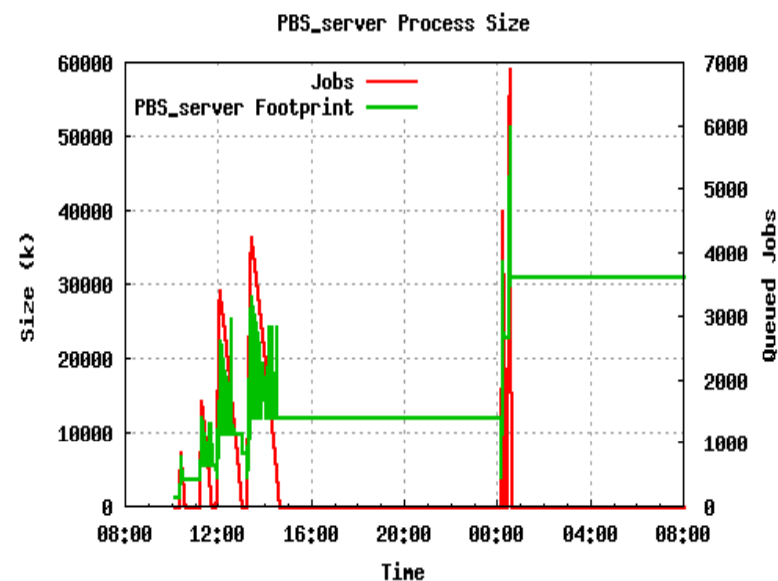
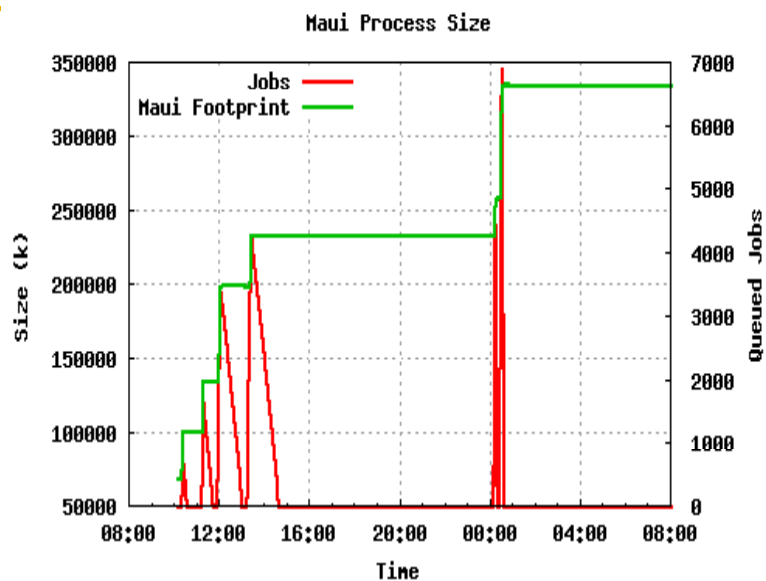


- **parallel job submissions directly to Torque**
 - As soon as the jobs are queued for execution the test script automatically requests their deletion from the queue. The objective is to stress the LRMS by simulating many requests via many paths.

- **E.g. 70x100 parallel jobs sent to torque**



- **stressing the LRMS memory management**
 - 1000-7000 instantaneous jobs were serially submitted, or 70 threads of 100 iterations each were in parallel performing qsub and qdel, in various combinations.
 - Results: memory is not freed properly



- **batch System Resilience Tests**
 - Submit jobs to batch system (using batch system commands) and shut down pbs_server and torque_server
 - Results: No jobs that have already started were lost after restart for the batch systems

- **Thorough testing of Torque/Maui Batch system has been done and various issues have been identified**
- **Defined a process for testing a batch system**
- **Can be used for other batch systems.**
- **This process has to be applied to different site configurations**
 - gCE/lcgCE, torque on same/separate machine
- **Detailed results can be found in:**
 - http://master.gridctb.uoa.gr/torque-report/Torque-Maui_testplan.html
- **Test Scripts can be found in:**
 - <http://email.uoa.gr/cgi-bin/viewvc.cgi/egee/ctb/>