# TORQUE Resource Manager 5.0.3 release notes

The release notes file contains the following sections:

# New Features

The following is a summary of key new features in TORQUE Resource Manager.

## 5.0.3

No new features.

## 5.0.2

No new features.

## 5.0.1

### Modification of the Output Location

TORQUE now allows for the modification of the output location based on the Mother Superior hostname. An environment variable ($HOSTNAME) has been added to the job's environment.

## 5.0.0

### CPU frequency control

TORQUE can now set the CPU frequency on requested nodes for submitted jobs. The request is made with the new cpuclock resource extension. A user can specify a desired clock frequency in megahertz, a Linux power governor policy name, or an ACPI performance state (P-state) number. For more information, see cpuclock in Requesting Resources in the TORQUE Resource Manager RPM Installation Guide.

The qalter -l command has been updated to allow modification of the requested CPU frequency on an already submitted job.

The pass_cpuclock server parameter was added allowing administrators to track, but not grant, the CPU frequency request portion of a job submission. For more information, see pass_cpuclock (Appendix B: Server Parameters) in the TORQUE Resource Manager RPM Installation Guide.

### qrerun all command

When you execute the `qrerun all` command, you will be prompted for confirmation. TORQUE will then place all running jobs in a queued state without contacting the MOMs. You should only use this when the entire cluster is down and you cannot contact it.

### Node power state control

TORQUE can now set the power state of a node. Depending on the hardware and software capabilities of the node, TORQUE can set the power state to Running, Standby, Suspend, Hibernate, or Shutdown. A new -m option was added to the pbsnodes command to make this work. For more information, see pbsnodes (Appendix A: Commands Overview) in the TORQUE Resource Manager RPM Installation Guide.

The syntax of the command is:

```
pbsnodes -m [running|standby|suspend|hibernate|shutdown] <space delimited list of
nodes to alter>
```

In order to wake a node from a low-power state, Wake-on-LAN must be supported and configured on the node. For more information, see Changing node power states in the TORQUE Resource Manager RPM Installation Guide.

# Differences

This section contains differences in previously existing features that require a change in configuration or routine.

### 5.0.3

### $prologalarm is Always Honored

$prologalarm was ignored on the prologue for a job. Also when the epilogue was run the $prologalarm value was ignored if it was more the 300. Now the $prologalarm value is always honored regardless of how large it is for both prologue and epilogue scripts. The default timeout is still 300 seconds.

### Log Milliseconds

Added milliseconds in TORQUE's log files.

### Abnormal Job Exit States Combined Into Single Exit Code

Added parameter "exit_code_canceled_job" to force all canceled jobs to have the same exit code, regardless of exit state. Previously, exit codes would be inconsistent depending on the type of cancellation or abnormal exiting of each job. The "exit_code_canceled_job" gets rid of extraneous exit codes and allows users to set their own exit code. For example, if a user sets "exit_code_canceled_job = 300" in qmgr, the exit code for any job that exits with a non-zero status becomes 300.

### pbsnodes -c Option Brings All Nodes Back Online

Users can now use the "pbsnodes -c all" or "pbsnodes -c ALL" to set offline nodes to a free state.

### Node Health Checks Now Run On Sister Nodes

Before, health checks would only run on the mother superior. Now, they run on sister nodes as well.

## 5.0.2

### Disable Automatic Requeueing of Jobs

Added the ability to disable the automatic requeuing of jobs due to transient failures.

### pbs_mom now sets environment variable for NVIDIA GPUs

A new mom config parameter, $cuda_visible_devices, was added to specify whether pbs_mom sets the CUDA_VISIBLE_DEVICES environment variable when it starts a job. The default is TRUE.

## 5.0.1

### qmgr Server Parameter "copy_on_rerun"

A new qmgr option: set server copy_on_rerun=[True|False] is available. When set to True, Torque will copy the OU, ER files over to the user-specified directory when the qrerun command is executed (i.e a job preemption). This setting requires a pbs_server restart for the new value to take in effect. Note that the MOMs and the pbs_server must be updated to this version before setting copy_on_rerun=True will behave as expected.

### qmgr Server Parameter "job_exclusive_on_use"

A new qmgr option: job_exclusive_on_use=[True|False] is available. When set to True, pbsnodes will report job-exclusive anytime 1 or more processors are in use. This resolves discrepancies between Moab and TORQUE node reports in cases where Moab is configured with a SINGLEJOB policy.

### TORQUE Accounting Improvements

Two new fields were added to the accounting file for completed jobs: total_execution_slots and unique_node_count. total_execution_slots should be 20 for a job that requests nodes=2:ppn=10. unique_node_count should be the number of unique hosts the job occupied.

### scan_for_terminated

Improved performance by moving scan_for_terminated to its own thread.

### Port Using the Munge API Instead of Forking

TORQUE now uses the Munge API, rather than forking, when configured with the --enable-munge-auth option.

### pbsdsh -o Option Captures sderr

The pbsdsh -o option is modified to add sderr capturing.

<u>5.0.0</u>

### *The job_stat_rate Parameter Has a New Default and Function*

Before this release, pbs_server asked the mother superior of every job for an update on the job every **job_stat_rate** seconds. The mother superior now sends updates on every job with its regular status, so there is no need for pbs_server to regularly poll. Instead, this parameter sets a timeout.

### *Two threadpools and New Default for max_threads*

The **max_threads** parameter has a new default: The value of min_threads ((2 * the number of procs listed in /proc/cpuinfo) + 1) * 20. In previous versions, 20 was 10.

Additionally, threadpools are now split. One-fourth of the threads are allocated for background tasks and three-fourths of the threads are allocated for incoming requests from MOMs and through the API (client commands, Moab, and so forth). Additionally, incoming requests no longer build up indefinitely. If a new request comes in, pbs_server evaluates whether it is too busy to service the request. For managers, the request is serviced as long as there are at least two threads available in the threadpool. For non-managers, the request is serviced as long as at least 5% of the threadpool is available. When pbs_server is too busy, it returns the error code PBSE_SERVER_BUSY with the message: "Pbs Server is currently too busy to service this request. Please retry this request."

### *Job Status Polling Removed from TORQUE*

Pbs_server now polls a mom for a job's information only if it hasn't received the information in 5 minutes. Otherwise, the information is communicated with the mom's status information.

### *TORQUE no longer searches linearly for the nodes in a node list*

TORQUE now recognizes when a request to run a job specifies a node list. It directly accesses those nodes instead of searching for them linearly.

### *The exec_host list has one entry per node*

The exec_host list has been condensed to contain one entry per node instead of one entry per execution slot. The node entry contains a string specifying each execution slot index.

TORQUE no longer displays the value of exec_port in qstat.

### *The qstat -f output for CPUs has been reduced*

The output of qstat -f has been condensed to reduce clutter in the CPU section.

# Known Issues

The following are known issues in TORQUE Resource Manager. Following each issue description is an associated issue number in parentheses.

Known issues are aggregated and grouped by the release version for which they first occurred or where reported.

### 5.0.3

No known issues.

### 5.0.2

No known issues.

### 5.0.1

- Some limitations exist in the way that pbsdsh can be used. Please note the following situations are not currently supported:

  - Running multiple instances of pbsdsh concurrently within a single job. (TRQ-2851)

  - Launching a large number of processes in succession (causes pbdsdsh to hang). (TRQ-2890)

### 5.0.0

No known issues.

# Resolved issues

The following is a list of some key bugs fixed in TORQUE Resource Manager. Following each issue description is an associated issue number in parentheses.

Resolved issues are aggregated and grouped by the release version in which they were resolved.

### 5.0.3

- **qstat Performance Suffered When Submitting 10k-20k Node Jobs.** When large jobs were submitted and launch, Torque performance would decline, with qstat becoming unresponsive when a large job was submitted. Performance improvements have been implemented. (TRQ-3110)

- **With kill_delay and $exec_with_exec set, a job would be set to a completed state after running qrerun instead of getting set back to queued.** This issue is fixed. (TRQ-2993)

- **Occasionally a random group name would show up for a user who did not belong in the group.** A race condition was fixed by changing to thread safe calls to get group and user ids. (TRQ-3190)

- **Threadpool in pbs_mom was not starting properly.**When $thread_unlink_calls is set to true in /var/spool/torque/mom_priv/config, job files were not being deleted at job end in the mom. This caused a holdup when starting the pbs_mom threadpool. This issue is fixed. (TRQ-3232)

- **Read timeouts were being retried indefinitely by pbs_server.** This issue is fixed. (TRQ-3306)

- **Reporter mom did not correctly handle UNKNOWN role.** This issue is fixed. (TRQ-3245)

- **A hang in pbsdsh occured if the pbs_mom daemon was started with a -q or -r option.** This issue is fixed. (TRQ-3308)

- **X11 very slow with qsub -I -X submit options.** The port forwarding buffer size has been increased to improve performance when enabling X11 forwarding from qsub (-X switch). (TRQ-2740)

- **End-of-job times were not consistent.** End-of-job times are now recorded closer to when nodes are free for use. (TRQ-2840)

- **Ability to have the procct values displayed on the queues where it was set was lost.** This issue is fixed. (TRQ-3135)

- **Jobs were stuck in EXITING state on mother superior.** This issue is fixed. (TRQ-3165)

- **Completed jobs were not cleaned up properly.** Fixed various issues relating to restarting dependency jobs. (TRQ-3175)

- **Jobs were run in the wrong order and canceled prematurely.** Running jobs are no longer deleted due to dependencies. (TRQ-3189)

- **Some threads would cause Torque to hang and become completely unresponsive.** Count_proc() now returns 0 on a NULL argument. (TRQ-3196)

- **Server entered deadlock around completed jobs map.** This issue is fixed. (TRQ-3226)

- **Kill_delay sent multiple signals.** Signals sent to job pids are now tracked to eliminate redundant sending. (TRQ-3239)

- **Improperly escaped xml caused job files and dependencies to fail upon reload attempt.** This issue is fixed. (TRQ-3284)

- **pbs_submit() method changed for stable job submissions.** The qsub command currently utilizes the pbs_submit_hash() method to submit jobs and is stable. Submitting jobs via the pbs_submit() method would occasionally render pbs_server unresponsive. Now pbs_submit() is functioning properly. (TRQ-3314)

- **Flooding server with client commands disabled some functionality.** Resolved issue where server gets stuck due to large volume of client requests. (TRQ-3296)

- **Server timed out to MOM when running a lot of small jobs.** Added load balancing to login nodes when login nodes begin to get busy. (TRQ-3367)

- **Init.d scripts failed if process id was owned by another process.** This issue is fixed. (TRQ-3034)

- **tm_poll() did not wait for event to complete, causing pbsdsh to hang.** This issue is fixed. (TRQ-3077)

- **Could not add remote user to operator or manager list.** This issue is fixed. (TRQ-3200)

- **"qdel all" command would frequently return "nonexistent job id: all" error.** This issue is fixed. (TRQ-3301)

- **Jobs would be created incorrectly if the path name contained spaces.** This issue is fixed. (TRQ-3347)

## 5.0.2

- **max_user_queuable issues.** Fixed some bugs causing incorrect max_user_queuable values. (TRQ-2841)

- **Delays reported with preemption.** Speed up process data collection to help diminish time to rerun a job, including correctly sum job memory and cpu time resources for processes created under setsid() calls within a job, and avoid an infinite loop in pbs_mom when running under a Linux Container (lxc). (TRQ-3026)

- **Problem reported with interactive jobs not staying on the node from which they were submitted.** This issue is fixed. (TRQ-3122)

- **When primary job is deleted, dependent job immediately disappears.** When a job is deleted because its dependency can no longer be satisfied, that job will follow the keep_completed parameter set - if any - for pbs_server. Previously, jobs were purged immediately. (TRQ-2949)

- **PBS_NUM_NODES on multi-req job wrong when using specific hostnames.** This issue is fixed. (TRQ-1949)

- **Jobs were getting stuck in a running state when an asynchronous run failed.** This issue is fixed. (TRQ-3114)

- **qalter was not working for running jobs, when disable_job_server_suffix was set to true.** This issue is fixed. (TRQ-3078)

- **Crash reported on start up when reading empty array file.** Fixed start up bug related to empty job array (.AR) files. (TRQ-2980)

- **Torque will not compile with --enable-debug on configure.** This issue is fixed. (TRQ-2969)

- **Unable to submit from acl_hosts login nodes.** Removed a code change that blocked the use of host.equiv. (TRQ-3014)

- **pbs_server would count completed jobs against queue limits when pbs_server was restarted.** This issue is fixed. (TRQ-3087)

## 5.0.1

- **Torque job can kill processes not owned by the job owner.** Fixed issue around unauthorized termination of processes. (CVE-2014-3684, TRQ-2885)

- **qstat -Q <bad_queue_name> prints queue name twice.** Fixed bug where giving a bad queue name to qstat -Q results in duplicate output. (TRQ-2025)

- **Mail output for qsub -m options was failing to output new lines.** Fixed bug in qsub -m when TORQUE is configured --with-sendmail. Some missing newlines were added. (TRQ-2937)

- **mppnodes hostlist was being re-ordered.** This bug is fixed. (TRQ-2112)

- **Some tasks were incorrectly listed as 0 in qstat -a when requested specific nodes.** This has been resolved. (TRQ-2292)

- **TORQUE accounting problems - jobs without accounting records.** Fixed bug related to accounting records on large systems. (TRQ-2367)

- **qstat wouldn't parse anything after a bad job ID.** Improved qstat behavior in cases where bad job IDs were referenced in the command. (TRQ-2410)

- **Separate headers for multiple jobs IDs provided to qstat.** Fixed output format bug in cases where multiple job IDs are passed into qstat. (TRQ-2411)

- **qsub did not process arguments correctly when a submitfilter is used.** Fixed bug where qsub did not process args correctly when using a submit filter. (TRQ-2646)

- **Parsing bug found when using hostlist ranges in qsub.** This bug is fixed. (TRQ-2652)

- **Build bug reported with MIC libraries.** Fixed build bug related to newer Intel MIC libraries installing in different locations. (TRQ-2653)

- **XML job log error.** Corrected mismatched <Job_Id> XML tags in the job log. (TRQ-2692)

- **TORQUE was not dividing GPUS amongst NUMA nodes.** Fixed problem where GPUs were not split between NUMA nodes. You now need to specify which gpus belong to each node board in the mom.layout file. (TRQ-2730)

  A sample mom.layout file might look like

  ```
  nodes=0 gpu=0
  nodes=1 gpu=1
  ```

  > ℹ This only works if you use nvml. The nvidia-smi command is not supported.

- **TORQUE was leaving behind error and out files when a job was preempted or requeued.** Fixed bug where OU files were being left in spool when job was preempted or requeued. (TRQ-2732)

- **Reported cput was incorrect.** Fixed bug where reported cput was incorrect. (TRQ-2759)

- **pbsnodes -l offline -n unexpected error.** Fixed unexpected error when running `pbsnodes -l offline -n`. (TRQ-2760)

- **max_user_queuable limit reached, however, there were no jobs in the queue**. Fixed bug where jobs rejected due to max_user_queuable limit reached, yet no jobs in the queue. (TRQ-2795)

- **momctl -q clearmsg didn't seem to clear error messages permanently.** Fixed bug where `momctl -q clearmsg` didn't properly clear error messages. (TRQ-2828)

- **Some bugs found that caused TORQUE core crashes.** These bugs are fixed.

- **TORQUE was crashing.** Fixed crashing.

- **pbs_server segfault after large array deletion.** Fixed segmentation fault. (TRQ-2835)

- **GPU nodes where not passed to sister nodes.** This bug is fixed. (TRQ-2837)

- **pbs_server did not write resource_default units to serverdb file.** This bug is fixed. (TRQ-2852)

- **pbs_mom filling up the logs in a HA environment.** Reduced verbosity in error logging in HA environments. (TRQ-2863)

- **Make trqauthd error messages more meaningful and non-repetitive.** Improved trqauthd error messages to be more meaningful and less redundant. (TRQ-2882)

- **Remote job submissions were being rejected, even when host is in submit_hosts list.** Fixed problem with remote client job submission during ruserok() calls. (TRQ-2918)

- **pbsdsh did not support running multiple instances concurrently.** Improved pbsdsh to better handle simultaneous use of -o and -s options. Also fixed some problems where -o output was sometimes getting truncated. (TRQ-2890)

- **TORQUE was not notifying Moab of completed jobs.** Fixed bug where TORQUE was not honoring KeepCompleted server parameter when job_nanny was set to true. (TRQ-2904)

- **HOST_NAME_SUFFIX was no longer adding suffix to job names.** This bug is fixed. (TRQ-2956)

- **Deadlock when running 'qdel -p' as non-root user.** Fixed deadlock issue. (TRQ-2919)

## 5.0.0

- **When the queue contained more than 50,000 jobs, TORQUE slowed down substantially.** TORQUE no longer slows down with a large number of jobs in the queue. (TRQ-2345)

- **When pbs_server had a high load, it would get stuck polling.** pbs_server no longer gets stuck polling under these conditions. (TRQ-2620)

- **When a job with a dependent job was deleted with qdel, TORQUE did not clean up the dependent job.** TORQUE now removes the dependent job when you delete its parent. (TRQ-2621)

- **When a lot of jobs were run at once, one job would get stuck in an exiting state.** This error no longer occurs. (TRQ-2622)

- **The afterok dependency did not work as expected.** afterok now works correctly. (TRQ-2626)

- **When running a large amount of jobs, the server would crash during job recycle.** The crash no longer occurs. (TRQ-2628)

- **pbsdsh requires FQDN even if other elements don't.** pbsdsh no longer requires FQDN. (TRQ-2632)

- **A deadlock would occur on job_save failure.** This error no longer occurs. (TRQ-2645)

- **Asynchronous job starts queued in TORQUE but not yet serviced caused Moab to reschedule jobs that would eventually run.** A new job sub-state in TORQUE prevents this from occurring. (TRQ-2715)