

## PBS Pro – Documentation

### Introduction

Most jobs will require greater resources than are available on individual nodes. All jobs must be scheduled via the batch job system. The batch job system in use is PBS Pro. Jobs are submitted to PBS specifying required resources, including the queue, number of CPUs, the amount of memory, and the length of time needed. PBS will then run a job or jobs when the resources are available, subject to constraints on maximum resource usage.

### Basic PBS commands

Some basic PBS commands are:

Command	Description
<code>qsub <i>commandfile</i></code>	<p>Submit jobs to the queues. The simplest use of the <code>qsub</code> command is typified by the following example:</p> <pre>[username@hpc-login-prd-t1 ~]\$ qsub commandfile</pre> <p>or</p> <pre>[username@hpc-login-prd-t1 ~]\$ qsub -q default -l \ walltime=20:00:00,vmem=3000MB commandfile</pre> <p>where <i>commandfile</i> is an ascii file containing PBS commands (<b>not</b> the compiled executable which is a binary file).</p> <p>See <a href="#">qsub Options</a> below for more options.</p>
<code>qstat -u <i>username</i></code>	Displays the status of PBS jobs and queues for the user <i>username</i> . See <i>man qstat</i> for details of options
<code>qdel <i>jobid</i></code>	Delete your job from a queue. The <i>jobid</i> is returned by <code>qsub</code> at job submission time, and is also displayed in the <code>qstat</code> output.
<code>qhold <i>jobid</i></code>	Place a hold on your job in the queue and stops it from running.
<code>qrls -h u <i>jobid</i></code>	Release a user hold on your job and allows it to be run.
<code>qrerun <i>jobid</i></code>	Terminate an executing job and return it to a queue.

## PBS Pro – Documentation

<code>qmove <i>jobid</i></code>	Move a job to a different queue or server.
---------------------------------	--

### qsub Options

There are two methods of specifying qsub options:

1. Within a PBS commandfile, and
2. On the qsub command-line.

Below is a simply example showing both methods.

```
[username@hpc-login-prd-t1 ~]$ qsub commandfile
```

Where *commandfile* contains:

```
#!/bin/bash -l
#
#PBS -l select=1:ncpus=16
#PBS -l walltime=1:00:00
#PBS -q default

### Specify the executable...
./an_executable
```

or

```
[username@hpc-login-prd-t1 ~]$ qsub -q default -l select=1:ncpus=16,walltime=1:00:00 an_executable
```

Below are some commonly used qsub options.

qsub Option	Description
<b>#PBS -A <i>acct</i></b>	Causes the job time to be charged to " <i>acct</i> ".
<b>#PBS -N <i>myJob</i></b>	Assigns a job name. The default is the name of PBS job script.
<b>#PBS -l nodes=4:ppn=2</b>	The number of nodes and processors per node. (deprecated)
<b>#PBS -l select=1:ncpus=2</b>	The number of chunks or nodes and processors per.
<b>#PBS -l ngpus=2</b>	The number of gpus required.
<b>#PBS -l walltime=01:00:00</b>	Sets the maximum wall-clock time during which this job can run. ( <i>walltime=hh:mm:ss</i> )
<b>#PBS -l mem=<i>n</i>{mb gb}</b>	Sets the maximum amount of memory allocated to the job.
<b>#PBS -l vmem=<i>n</i>{mb gb}</b>	Sets the maximum amount of virtual memory allocated to the job. (deprecated)
<b>#PBS -q <i>queuename</i></b>	Assigns your job to a specific queue.
<b>#PBS -o <i>mypath/my.out</i></b>	The path and file name for standard output.
<b>#PBS -e <i>mypath/my.err</i></b>	The path and file name for standard error.
<b>#PBS -j oe</b>	Join option that merges the standard error stream with the standard output stream of the job.
<b>#PBS -M <i>email-address</i></b>	Sends email notifications to a specific user email address.

## PBS Pro – Documentation

<b>#PBS -m {a b e}</b>	Causes email to be sent to the user when: <ul style="list-style-type: none"> <li>• a - the job aborts</li> <li>• b - the job begins</li> <li>• e - the job ends</li> </ul>
<b>#PBS -P project</b>	Specifies what project the job belongs to.
<b>#PBS -r n</b>	Indicates that a job should not rerun if it fails.
<b>#PBS -S shell</b>	Sets the shell to use. Make sure the full path to the shell is correct.
<b>#PBS -V</b>	Exports all environment variables to the job.
<b>#PBS -W</b>	Used to set job dependencies between two or more jobs.
<b><u>NOTE</u></b>	PBS directives are all at the start of a script, that there are no blank lines between them, and there are no other non-PBS commands until after all the PBS directives.

### A Job Script Example

A working job submission script takes the following form:

---

```
#!/bin/bash -l
#PBS -N Example_Job
#PBS -q default
#PBS -l select=2:ncpus=16
#PBS -l walltime=<hh:mm:ss>
#PBS -o <output-file>
#PBS -e <error-file>

module load matlab/r2016b

matlab -nodisplay -nosplash -r example_job.m
```

---

Where the line "**-l select=2:ncpus=16**" is the number of processors required for the job. **select** specifies the number of nodes (or chunks of resource) required; **ncpus** indicates the number of CPUs per chunk required.

As this is not the most intuitive command, the following table is provided as guide to how this command works:

## PBS Pro – Documentation

select	ncpus	Description
<b>2</b>	16	32 Processor job, using 2 nodes and 16 processors per node
<b>4</b>	8	32 Processor job, using 4 nodes and 8 processors per node
<b>16</b>	1	16 Processor job, using 16 nodes and 1 processor per node
<b>8</b>	16	128 Processor job, using 8 nodes and 16 processors per node

The line "**-l walltime=<hh:mm:ss>**" is the time limit for the job. If your job exceeds this time the scheduler will terminate the job. It is recommended to find a usual runtime for the job and add some more (say 20%) to it. For example, if a job took approximately 10 hours, the walltime limit could be set to 12 hours, e.g. "**-l walltime=12:00:00**". By setting the walltime the scheduler can perform job scheduling more efficiently and also reduces occasions where errors can leave the job stalled but still taking up resource for the default much longer walltime limit (for queue walltime defaults run "**qstat -q**" command).

### Job management

The **qstat** command displays the status of the PBS scheduler and queues. Using the flags **-Qa** shows the queue partitions available. If no queue is defined, it will use the queue called *default*. The following table shows the commonly using queues:

express:

- all nodes available
- low priority
- 8 hours of run time available

serial:

- all nodes available
- high priority
- 168 hours of run time available

short:

- all nodes available
- standard priority
- 24 hours of run time available

## PBS Pro – Documentation

medium:

- all nodes available
- standard priority
- 72 hours of run time available

long:

- all nodes available
- standard priority
- 168 hours of run time available

### PBS Job States

The table below describes the different job states through the life cycle of a job. There are some attributes that are only applicable when submitting jobs to an Enterprise PBS Professional complex.

Job State	Description
<b>B</b>	Job arrays only: job array has Begun.
<b>E</b>	Job is Exiting after having run.
<b>F</b>	Job has Finished exiting and execution. The job was completed successfully and had no application errors.
	Job has Finished exiting and execution; however, the job experienced application errors.
<b>H</b>	Job is Held. A job is put into a held state by the server or by a user or administrator. A job stays in a held state until it is released by a user or administrator.
<b>Q</b>	Job is Queued, eligible to run or be routed.
<b>R</b>	Job is Running.
<b>S</b>	Job is Suspended by server. A job is put into the suspended state when a higher priority job needs the resources.
<b>T</b>	Job is in Transition (being moved to a new location).
<b>U</b>	Job is User-suspended.
<b>W</b>	Job is Waiting for its requested execution time to be reached or job specified a staging request which failed for some reason.
<b>X</b>	Sub jobs only; sub job is finished (expired).

## PBS Pro – Documentation

### Queue Limits

	express	serial	short	medium	long
<b>Priority</b>	161	140	160	160	160
<b>Max CPU per job</b>	500	1	200	100	40
<b>Max Node</b>	29	29	29	29	29
<b>Min Walltime (hr)</b>	1	1	8	24	72
<b>Max Walltime (hr)</b>	8	168	24	72	168
<b>Default Walltime (hr)</b>	1	24	24	24	72
<b>Default Memory (gb)</b>	2	2	2	2	2
<b>Max Running Jobs</b>	500	450	450	400	200
<b>Max Queued Jobs</b>	20000	10000	20000	5000	2000

### Queue Scheduling Issues

The scheduling algorithm used on the HPC aims to:

- promote large scale parallel use of the HPC
- allow equal access to resources for all users
- provide good turnaround for all users
- minimize the impact of jobs on one another

Some of the scheduler features to achieve these aims are:

- resources are strictly allocated so jobs will not start unless there is sufficient free resources (e.g. cpus and memory).
- queued jobs are shuffled so that jobs from different users are "interleaved". This means your first job should appear near the top of the queue even if there are many jobs in the queue.

From a user's perspective, it is very important that you minimize your requests for resources (e.g. walltime, memory and cpus). Otherwise your job may be queued or suspended longer than necessary. Of course, make sure you ask for sufficient resources - a little experimentation might help.

### PBS Variables

PBS sets multiple environment variables at submission time. The following PBS variables are commonly used in command files:

## PBS Pro – Documentation

Variable Name	Description
<b>PBS_ARRAYID</b>	Array ID numbers for jobs submitted with the -t flag. For example a job submitted with #PBS -t 1-8 will run eight identical copies of the shell script. The value of the PBS_ARRAYID will be an integer between 1 and 8.
<b>PBS_ENVIRONMENT</b>	Set to PBS_BATCH to indicate that the job is a batch job; otherwise, set to PBS_INTERACTIVE to indicate that the job is a PBS interactive job.
<b>PBS_JOBID</b>	Full jobid assigned to this job. Often used to uniquely name output files for this job, for example: <i>mpirun -np 16 ./a.out &gt;output.\${PBS_JOBID}</i>
<b>PBS_JOBNAME</b>	Name of the job. This can be set using the -N option in the PBS script (or from the command line). The default job name is the name of the PBS script.
<b>PBS_NODEFILE</b>	Contains a list of the nodes assigned to the job. If multiple CPUs on a node have been assigned, the node will be listed in the file more than once. By default, <i>mpirun</i> assigns jobs to nodes in the order they are listed in this file
<b>PBS_O_HOME</b>	The value of the HOME variable in the environment in which qsub was executed.
<b>PBS_O_HOST</b>	The name of the host upon which the qsub command is running.
<b>PBS_O_PATH</b>	Original PBS path. Used with pbsdsh.
<b>PBS_O_QUEUE</b>	Queue job was submitted to.
<b>PBS_O_WORKDIR</b>	PBS sets the environment variable <i>PBS_O_WORKDIR</i> to the directory from which the batch job was submitted
<b>PBS_QUEUE</b>	Queue job is running in (typically this is the same as PBS_O_QUEUE).

### Interactive PBS Jobs

Use of PBS is not limited to batch jobs only. It also allows users to use the compute nodes interactively, when needed. For example, users can work with the developer environments provided by Matlab or R on compute nodes, and run their jobs (until the walltime expires). Instead of preparing a submission script, users pass the job requirements directly to the qsub command. For instance, the following PBS script:

```
#PBS -l nodes=7:ppn=4
#PBS -l mem=2gb
#PBS -l walltime=15:00:00
#PBS -q default
```

This corresponds to:

```
qsub -l -X -q default -l select=7:ncpus=4,walltime=15:00:00,mem=2gb
```



## PBS Pro – Documentation

Hence, the PBS scheduler will allocate  $7*4=28$  cores to the user as soon as nodes with given specifications become available, then automatically log the user into one of the compute nodes. From now on, the user can work interactively using these cores until the walltime expires. Note that there should be no space between the parameters being passed to `-l` (as in `'L'ima`) flag, only commas!

Here, `-l` (as in `'l'ndia`) stands for `'interactive'` and `-X` allows for GUI applications.

### PBS Job Dependencies

In some situations a job or jobs will be dependent on the output of another job in order to run. To add a job dependency, the option `-W [additional attributes]` is used when submitting a job. In the example below the `afterok` rule will be used, but there are several other rules that may be useful. In this example two PBS command files will be used:

- `number.pbs` - generates a list of numbers in the file `number.list`
- `order.pbs` - sorts the list of numbers generated by `number.pbs`

If both jobs were submitted as:

```
[username@hpc-login-prd-t1 ~]$ qsub number.pbs ; qsub order.pbs
```

the error output from `order.pbs` will be `order: open failed: number.list: No such file or directory` If `order.pbs` was submitted with a dependency on `number.pbs` as in:

```
[username@hpc-login-prd-t1 ~]$ qsub number.pbs
4674.hpc-admin-prd-t1
[username@hpc-login-prd-t1 ~]$ qsub -W depend=afterok:4674 order.pbs
4675.hpc-admin-prd-t1
[username@hpc-login-prd-t1 ~]$ qstat -u $USER
```

```
hpc-admin-prd-t1.usq.edu.au:
```

Job ID	Username	Queue	Jobname	SessID	NDS	TSK	Req'd Memory	Req'd Time	Elap S	Time
4674.hpc-admi	username	short	number.pbs	18029	1	1	-- 48:00 R	00:00		
4675.hpc-admi	username	short	order.pbs		1	1	-- 48:00 H	--		

Notice the `order.pbs` is in a hold state however once the dependent job completes the order job run as:



## PBS Pro – Documentation

```
[username@hpc-login-prd-t1 ~]$ qstat -u $USER
```

```
hpc-admin-prd-t1.usq.edu.au:
```

Job ID	Username	Queue	Jobname	SessID	NDS	TSK	Req'd Memory	Req'd Time	Elap S	Time
4675.hpc-admi	username	short	order.pbs		1	1	--	48:00	R	--

Other options to -W include:

- `afterany:jobid[:jobid...]` implies that the job may be scheduled for execution after jobs `jobid` have terminated, with or without errors.
- `afterok:jobid[:jobid...]` implies that job may be scheduled for execution only after jobs `jobid` have terminated with no errors.
- `afternotok:jobid[:jobid...]` implies that job may be scheduled for execution only after jobs `jobid` have terminated with errors.

### References:

1. [PBS Professional 14 User Guide](#)
2. [PBS Professional 14 Administrator's Guide](#)
3. [PBS Professional - HPC Cluster Workload Manager](#)