

Workshop Agenda – Feb 25th 2015

| Time | Presenter | Title |
|-------|------------|---|
| 09:30 | T. König | Talk – bwHPC Concept & bwHPC-C5 - Federated User Support Activities |
| 09:45 | R. Walter | Talk – bwHPC architecture (bwUniCluster, bwForCluster JUSTUS, ForHLR Phase I) |
| 10:00 | A. Fuchs | Talk – Cluster: Access, Data Transfer and Storage, GUI |
| 10:30 | | <i>Break</i> |
| 10:45 | R. Barthel | Talk – File System, Software System (modulefiles), Batch System |
| 11:10 | A. Fuchs | Tutorial – bwUniCluster: Access, Data Transfer, Compiling, Modulefiles, Batch Job Scripting |
| 11:50 | | <i>Lunch Break</i> |
| 13:00 | R. Barthel | Talk – Advanced Bash Scripting |
| 13:30 | R. Barthel | Tutorial – Advanced (Batch) Job Scripting |
| 14:15 | | <i>Break</i> |
| 14:30 | A. Fuchs | Tutorial – Compiling, Makefile, Parallelising |
| 15:15 | | User Forum – Solving User Cases |
| 16:00 | | <i>End</i> |



bw|HPC – C5

bwHPC course: File System, Software and Batch System

Robert Barthel, Simon Raffener



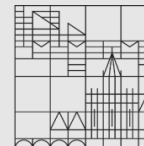
UNIVERSITÄT
HEIDELBERG
ZUKUNFT
SEIT 1386

Hochschule
für Technik
Stuttgart



Hochschule Esslingen
University of Applied Sciences

Universität
Konstanz



UNIVERSITÄT
MANNHEIM



Universität Stuttgart

EBERHARD KARLS
UNIVERSITÄT
TÜBINGEN



KIT
Karlsruher Institut für Technologie



ulm university universität
uulm



Reference: bwHPC-C5 Best Practices Repository

- Most information given by this talk can be found at <http://bwhpc-c5.de/wiki>:

The screenshot shows the bwHPC Wiki interface. At the top left is the logo, followed by the title 'bwHPC Wiki'. Below this is a search bar and navigation tabs for 'page', 'discussion', and 'view source'. The main content area is titled 'Main Page' and features a large yellow box with the text 'Knowledge Base Wiki of Baden-Württemberg's HPC services'. Below this is a welcome message: 'Welcome to the Knowledge Base Wiki of services and projects for high performance computing (HPC) and HPC data storage in the state of Baden-Württemberg, Germany. Hosted as a Best Practices Repository, the knowledge base contains user guides and best practice guides (BPG) and is maintained by members of Baden-Württemberg's federated HPC competence centers for clusters of tier 3 as well as by members of the HPC competence center for the ForHLR (tier 2). Federated HPC competence centers of tier 3 are an integral part of the project bwHPC-C5 which coordinates the federated user and science support for the HPC Infrastructure of tier 3 in the state of Baden-Württemberg.' Below the welcome message are two columns of service links. The left column, titled 'HPC Services', lists 'bwUniCluster', 'bwForCluster JUSTUS', and 'Best Practices Repository'. The right column, titled 'HPC Data Storage Services', lists 'bwFileStorage'. A sidebar on the left contains various navigation links such as 'Home', 'Best Practice Guides', 'Tools', and 'Personal tools'.

- Category:Hardware_and_Architecture
- Environment_Modules
- Batch_Jobs

Material: Slides & Scripts

- http://indico.scc.kit.edu/indico/event/workshop_2015-02_bwHPC
- @bwUniCluster/ForHLR/IC2/HC3
/pfs/data1/software_uc1/bwhpc/kit/workshop/2015-02-25

How to read the following slides

| Abbreviation/Colour code | Full meaning |
|---|--|
| <code>\$ command -option value</code> | <code>\$</code> = prompt of the interactive shell The full prompt may look like: user@machine:path\$ The command has been entered in the interactive shell session |
| <code><integer></code> <code><string></code> | <code><></code> = Placeholder for integer, string etc |
| <code>foo, bar</code> | Metasyntactic variables |

Software System

File System

- bwUniCluster (**U**) & ForHLR I (**F**) @ Karlsruhe, bwForCluster JUSTUS (**J**) @ Ulm

| Property) | \$HOME | \$TMPDIR | workspaces | \$WORK | \$PROJECT |
|----------------|-----------------------------|--------------------------------|-----------------------------|-------------|------------------------------|
| Where | U+F+J (all) | all | all | U + F | F |
| Visability | global | local | global | global | global |
| Lifetime | permanent | batch job wallt. | 240 d (U+F) 90 d (J) | 28 d | permanent |
| Disk space: | 469 TB (U+F) 200 TB (J) | | 938 TB (U+F) 200 TB (J) | 938 (U+F) | 469 TB |
| Disk @ thin n. | | 2 TB (U+F) 0/1 TB (J) | | | |
| Disk @ fat n. | | 7 TB (U), 8 TB (F) 2 TB (J) | | | |
| Quotas | yes: 1GB (F), 100 GB (J) | no | if required (U+F) no (J) | if required | yes (defined by proposal) |
| Backup | yes | no | no | no | yes |

- \$HOME, \$WORK and workspaces: on the parallel file system Lustre

→ **BUT:** only \$HOME under backup

\$HOME = Home directory

■ \$HOME:

@ bwUniCluster/ForHLR:

■ Current quota: `$ lfs quota -u $(whoami) $HOME`

■ Diskusage: `$ grep -E "$(whoami)|Account" ~/../../diskusage`

@ KIT: \$HOME directories of bwUniCluster, ForHLR, IC2, HC3 are the same

■ But: different operational systems, hardware, libraries, queueing etc.

→ bwUniCluster & ForHLR & (OS = REHL) vs. IC2 & HC3 (OS = SLES)

→ generalise your scripts to work on all systems using **\$CLUSTER**

```
if [ _${CLUSTER} == "uc1" ]; then  
  <operations>  
fi
```

@ JUSTUS:

■ Diskusage: `$ less ~/../../diskusage/$(whoami)`

\$PROJECT = Project directory of ForHLR

- ONLY ForHLR:

- All features of \$HOME
- Access granted based on approved projects
 - assigned „name/acronym“
 - \$PROJECT_GROUP

- Access project home directory: `$ cd $PROJECT`

- **Do not use: \$HOME** → since it has very low quota for the project group

- Quota of Project: `$ lfs quota -g ${PROJECT_GROUP} ${PROJECT}`

\$WORK = Working directory

- bwUniCluster/ForHLR → **additional parallel file system** with **limited lifetime, no redundancy, quotas**
 - especially designed for parallel access and for a high throughput to large files
- 2 concepts of access via:
 - (A) → \$WORK
 - (B) → *workspaces*
- (A) **\$WORK**:
 - Change to it via: `$ cd $WORK`
 - Quota: `$ lfs quota -u $(whoami) $WORK`
 - But: files no longer needed should be removed
 - any file inside your \$WORK older than 28 days will be deleted

Workspaces = Working directory

■ (B) Workspaces: lifetime on allocated folder

■ HowTo:

→ http://www.bwhpc-c5.de/wiki/index.php/BwUniCluster_File_System#Workspaces

| | |
|-----------------------|--|
| \$ ws_allocate foo 10 | Allocate a workspace named <i>foo</i> for 10 days |
| \$ ws_list -a | List all your workspaces |
| \$ ws_find foo | Get absolute path of workspace <i>foo</i> |
| \$ ws_extend foo 5 | Extend lifetime of your workspace <i>foo</i> by 5 days from now. You can extend 3 times → max. lifetime of <i>foo</i> = 240 days (U+F) 90 days (J) |
| \$ ws_release foo | Manually erase your workspace <i>foo</i> |

Example:

```
$ ws_allocate scratch
$ SDIR=$(ws_find scratch)
$ echo $SDIR
/work/workspace/scratch/ab1234-scratch-0
```



Software System

Environment modules

- Default → manual setup of
 - compilers, libraries and software packages etc.
→ complicated if multiple versions of same software installed
- Solution:
 - dynamic modification of the session environment by
→ instruction sets stored in *modulefiles*
- HowTo?
 - *load* and *unload* instruction sets (= modulefiles)
- How to use modulefiles in general?

\$ module help
- More information:
 - http://www.bwhpc-c5.de/wiki/index.php/Environment_Modules

modulefiles: available / search

■ Display all modulefiles

```
$ module avail
```

```
----- /opt/bwhpc/kit/modulefiles -----  
cae/abaqus/6.13-5 cae/ansys/15.0 cae/comsol/4.4 system/d-default  
cae/adina/9.0 cae/ansys/15.0.7 cae/starccm+/9.4  
  
----- /opt/bwhpc/common/modulefiles -----  
bio/bismark/0.10.1 lib/boost/1.55.0  
bio/bowtie/1.0.1 lib/matplotlib/1.3.1  
bio/bowtie2/2.1.0 lib/netcdf/3.6.3-gnu-4.8  
bio/bowtie2/2.2.3 lib/netcdf/3.6.3-intel-13.1  
bio/cufflinks/2.2.0 lib/pnetcdf/1.4.1  
bio/qiime/1.8.0 math/R/3.0.2  
bio/samtools/0.1.19 math/matlab/R2013a  
bio/tophat/2.0.11 math/matlab/R2013b  
bio/trimmomatic/0.32 math/matlab/R2014a  
cae/ansys/15.0.7_bw mpi/impi/4.1.0-gnu-4.4  
cae/ansys/15.0_bw mpi/impi/4.1.0-gnu-4.5  
cae/openfoam/1.6-ext mpi/impi/4.1.0-intel-12.1
```

■ Display all modulefiles with prefix „compiler“

```
$ module avail compiler
```

```
----- /opt/bwhpc/common/modulefiles -----  
compiler/gnu/4.5 compiler/gnu/4.8 compiler/intel/12.1  
compiler/gnu/4.7(default) compiler/gnu/4.9 compiler/intel/13.1(default)
```

modulefiles: help / whatis

- Show help of modulefiles, e.g. `$ module help compiler/intel`

```
----- Module Specific Help for 'compiler/intel/13.1' -----  
  
This module provides the Intel(R) compiler suite version 13.1.3 via  
commands 'icc', 'icpc' and 'ifort', the debugger 'idb' as well as the Intel(R)  
Threading Building Blocks TBB and the Integrated Performance Primitives IPP  
libraries (for details see also 'http://software.intel.com/en-us/intel-compilers/').  
  
The related Math Kernel Library MKL module is 'numlib/mkl/11.0.5'.  
The related Intel MPI module is 'mpi/impi/4.1.1-intel-13.1'.  
The Intel icpc should work well with GNU compiler 4.7.  
  
Commands:  
icc          # Intel(R) C compiler  
icpc         # Intel(R) C++ compiler  
ifort        # Intel(R) Fortran compiler  
idb          # Intel(R) debugger in GUI mode  
idbc         # Intel(R) debugger in console mode  
  
Local documentation:  
Man pages: man icc; man icpc; man ifort  
firefox $INTEL_DOC_DIR/documentation_c.htm  
firefox $INTEL_DOC_DIR/documentation_f.htm
```

- Show short info modulefile

```
$ module whatis compiler/intel
```

```
compiler/intel      : Intel(R) compiler suite (icc, icpc, ifort), debugger (idb), IPP and TBB ver 13.1.3
```

modulefiles: show

- Show all instructions of modulefile `$ module show compiler/gnu/4.7`

```
/opt/bwhpc/common/modulefiles/compiler/gnu/4.7:
```

```
module-whatis  GNU compiler suite version 4.7.3 (gcc, g++, gfortran)
setenv        GNU_VERSION 4.7.3
setenv        GNU_HOME /opt/bwhpc/common/compiler/gnu/4.7.3/x86_64
setenv        GNU_BIN_DIR /opt/bwhpc/common/compiler/gnu/4.7.3/x86_64/bin
setenv        GNU_MAN_DIR /opt/bwhpc/common/compiler/gnu/4.7.3/x86_64/share/man
setenv        GNU_LIB_DIR /opt/bwhpc/common/compiler/gnu/4.7.3/x86_64/lib64
prepend-path  PATH /opt/bwhpc/common/compiler/gnu/4.7.3/x86_64/bin
prepend-path  MANPATH /opt/bwhpc/common/compiler/gnu/4.7.3/x86_64/share/man
prepend-path  LD_RUN_PATH /opt/bwhpc/common/compiler/gnu/4.7.3/x86_64/lib
prepend-path  LD_LIBRARY_PATH /opt/bwhpc/common/compiler/gnu/4.7.3/x86_64/lib
prepend-path  LD_RUN_PATH /opt/bwhpc/common/compiler/gnu/4.7.3/x86_64/lib64
prepend-path  LD_LIBRARY_PATH /opt/bwhpc/common/compiler/gnu/4.7.3/x86_64/lib64
setenv        CC gcc
setenv        CXX g++
setenv        F77 gfortran
setenv        FC gfortran
setenv        F90 gfortran
setenv        TEST_MODULE_SCRIPT /opt/bwhpc/common/compiler/gnu/4.7.3/install-doc/test-compiler-gnu.sh
setenv        TEST_MODULE_NAME compiler/gnu/4.7
conflict      compiler/gnu
```

Load modulefiles (3)

- Modulefiles are sorted in categories, software name and versions:

```
$ module load <category>/<software_name>/<version>
```



- Load a default software:

```
$ module load <category>/<software_name>
```

- e.g. Intel compiler

```
$ module load compiler/intel mpi/impi
```


→ loads currently Intel compiler suite 13.1

→ loads currently Intel-MPI 4.1.1 for Intel compiler 13.1

```
$ module list
```

- Display all loaded modules

```
Currently Loaded Modulefiles:  
 1) compiler/intel/13.1(default)  2) mpi/impi/4.1.1-intel-13.1(default)
```



modulefiles: categories & dependencies

- Module names already implicate dependencies:

→ **Category/softwarename/version_attributes-dependencies**

e.g. **numlib/fftw/3.3.3-impi-4.1.1-intel-13.1**

→ fftw package version 3.3.3, compiled with Intel 13.1 and Intel-MPI 4.1.1

- Categories:

| | |
|-----------|--|
| compiler/ | for compiler, e.g. intel, gnu, pgi, open64 |
| devel/ | for debugger, e.g. ddt, and development tools, e.g. cmake, itrac |
| mpi/ | for MPI libraries, e.g. impi, openmpi, mvapich(2) |
| numlib/ | for numerical libraries, e.g. Intel MKL, ACML, nag, gsl, fftw |
| lib/ | for other libraries, e.g. netcdf, global array |
| bio/ | for biology software, e.g. bowtie, abyss, mrbayes |
| cae/ | for CAE software, e.g. ansys, abaqus, fluent |
| chem/ | for chemistry software, e.g. gromacs, dacapo, turbomole |
| math/ | for mathematics software, e.g. matlab, R |
| phys/ | for physics software, e.g. geant4 |
| vis/ | for visualisation software, e.g. vmd, tigervnc |



modulefiles: conflicts

■ Conflicts:

- a) load different software version in the same session, e.g. Intel:

```
$ module load compiler/intel/12.1  
$ module load compiler/intel/13.1
```

```
compiler/intel/13.1(394):ERROR:150: Module 'compiler/intel/13.1' conflicts  
with the currently loaded module(s) 'compiler/intel/12.1'
```

- b) load module with dependencies on other modules

```
$ module load mpi/openmpi/1.6.5-intel-13.1
```

```
Loading module dependency 'compiler/intel/13.1'.  
compiler/intel/13.1(394):ERROR:150: Module 'compiler/intel/13.1' conflicts  
with the currently loaded module(s) 'compiler/intel/12.1'
```

modulefiles: unload/swap

- To remove module *foo*:

```
$ module unload foo
```

```
$ module remove foo
```

be aware that you might create **inconsistencies**,

e.g. you can remove

compiler/intel/13.1 while *mpi/openmpi/1.6.5-intel-13.1* is still loaded

- Swap = remove + load

e.g.:

```
$ module swap compiler/intel/12.1 compiler/intel/13.1
```

Private modulefiles

- Each user can create own modulefiles:

e.g. modulefiles that adds path of own programs, `$HOME/special`, to `$PATH`

→ content of this modulefile „*mybin*“

```
#%Module1.0  
  
Append-path    PATH    "$env(HOME)/special"
```

→ place „*mybin*“ under `$HOME/privatemodules`

→ to make all own modules visible to “module avail” command, enter:

```
$ module load use.own    or    $ module use $HOME/privatemodules
```

→ former: own modules have lower priority than system ones if equally named

→ latter: own module have higher priority

- Remove own modules:

```
$ module unload use.own  or  $ module unuse $HOME/privatemodules
```

Batch System

Resource management

- Components of management system (Batch System)
 - **resource manager**
 - control over jobs and distributed compute nodes
 - SLURM (bwUniCluster, ForHLR)
 - TORQUE (bwForCluster JUSTUS)
 - **workload manager (scheduler)**
 - scheduling, managing, monitoring, reporting
 - MOAB

Resource and workload manager

```
#!/bin/bash
#MSUB -l nodes=1:ppn=1
#MSUB -l walltime=00:01:00
#MSUB -l pmem=50mb

echo "Hello from job"
exit 0
```

(2) MOAB parses the job script:
→ where & when to run job

(1) User creates a job script and submits it to MOAB via the “msub” command

MOAB

(3) Job execution:
delegated to resource manager on the node

compute resources:
TORQUE/SLURM

(4) The resource manager (TORQUE/SLURM) executes the job and communicates status information to MOAB

Job's life circle

- Setup job script:

```
#!/bin/bash
#MSUB -l nodes=1:ppn=1
#MSUB -l walltime=00:01:00
#MSUB -l pmem=50mb

echo "Hello from job"
exit 0
```

- Submit job to workload manager **ONLY** with “msub”

```
$ msub <resource_options> <job_script>
<job_ID>
```

- Job waits for free resources in queue

```
$ showq
<job_ID> state “Idle” → “Running”
```

- Job is finished → check output (default job name)

```
bwUniCluster/ForHLR:   job_<uc1,fh1>_<job_ID>.out
JUSTUS:                STDIN.o<job_ID> or STDIN.e<job_ID>
```


msub options

■ http://www.bwhpc-c5.de/wiki/index.php/Batch_Jobs#msub_Command

■ msub options: command line or in your job script

| Command line | Script | Purpose |
|---------------------------|---------------------------------|--|
| <code>-l resources</code> | <code>#MSUB -l resources</code> | Defines the resources that are required by the job. See the description below for this important flag. |
| <code>-N name</code> | <code>#MSUB -N name</code> | Gives a user specified name to the job. |
| <code>-q queue</code> | <code>#MSUB -q queue</code> | Defines the queue class |
| <code>-m bea</code> | <code>#MSUB -m bea</code> | Send email when job begins (b), ends (e) or aborts (a). |

→ command line option overwrites script option

`msub -l resource_list`

■ http://www.bwhpc-c5.de/wiki/index.php/Batch_Jobs#msub_-l_resource_list

| Resource | Purpose |
|------------------------------------|--|
| <code>-l nodes=2:ppn=16</code> | Number of nodes and number of processes per node |
| <code>-l walltime=600</code> | Wall-clock time (seconds) |
| <code>-l walltime=01:30:00</code> | HH:MM:SS format |
| <code>-l pmem=1000mb</code> | Max. amount of physical memory used by one process of the job (kb,mb,gb) |
| <code>-l mem=1000mb</code> | Max. total physical memory used by the job |

→ for workshop: `-l advres=workshop.54`

→ resources can be combined, but must be separated by comma:

```
$ msub -l nodes=1:ppn=1,walltime=00:01:00,pmem=1gb <job_script>
```

msub -q *queues* (bwUniCluster)

■ www.bwhpc-c5.de/wiki/index.php/Batch_Jobs_-_bwUniCluster_Features#msub_-q_queues

| <i>queue</i> | <i>default resources</i> | <i>MIN resources</i> | <i>MAX resources</i> |
|----------------------------------|-----------------------------|-------------------------------------|--|
| automatic queue choosing | | | |
| develop | <i>procs=1, mem=4000mb</i> | <i>nodes=1</i> | <i>walltime=00:30:00, nodes=1:ppn=16</i> |
| singlenode | <i>procs=1, mem=4000mb</i> | <i>walltime=00:30:01, nodes=1</i> | <i>walltime=3:00:00:00, nodes=1:ppn=16</i> |
| multinode | <i>procs=1, mem=4000mb</i> | <i>nodes=2</i> | <i>walltime=2:00:00:00, nodes=16:ppn=16</i> |
| explicit queue definition | | | |
| verylong | <i>procs=1, mem=4000mb</i> | <i>walltime=3:00:00:01, nodes=1</i> | <i>walltime=6:00:00:00, nodes=1:ppn=16</i> |
| fat (fat nodes) | <i>procs=1, mem=32000mb</i> | <i>nodes=1</i> | <i>walltime=3:00:00:00, nodes=1:ppn=32</i> |

■ **Automatic queue choosing** - walltime, nodes, processes

msub -q queues (ForHLR)

 http://www.bwhpc-c5.de/wiki/index.php/Batch_Jobs_-_ForHLR_Phase_I_Features

| <i>queue</i> | <i>default resources</i> | <i>MIN resources</i> | <i>MAX resources</i> |
|--------------------------------|--|-----------------------|---|
| explicit queue choosing | | | |
| develop | <i>procs=1, mem=3200mb, walltime=00:10:10</i> | <i>nodes=1</i> | <i>walltime=00:30:00, nodes=1:ppn=20</i> |
| singlenode | <i>procs=1, mem=3200mb, walltime=00:10:10</i> | <i>nodes=1</i> | <i>walltime=3:00:00:00, nodes=1:ppn=20</i> |
| multinode | <i>procs=1, mem=3200mb, walltime=00:10:10</i> | <i>nodes=2</i> | <i>walltime=3:00:00:00, nodes=128:ppn=20</i> |
| fat (fat nodes) | <i>procs=1, mem=160000mb, walltime=00:10:10</i> | <i>nodes=1</i> | <i>walltime=3:00:00:00, nodes=1:ppn=32</i> |

msub -q *queues* (JUSTUS)

■ www.bwhpc-c5.de/wiki/index.php/Batch_Jobs_-_bwForCluster_Chemistry_Features#Queues

| Queue name | Walltime MIN | Walltime MAX | MAX nodes (total per user) | MAX run/idle jobs (total per user) |
|------------|--------------|----------------|-------------------------------|---------------------------------------|
| quick | 00:00:01 | 00:05:00 | 2 | 1/1 |
| short | 00:05:01 | 2d. 48:00:00 | 64 | |
| normal | 48:00:01 | 7d. 168:00:00 | 16 | |
| long | 168:00:01 | 14d. 336:00:00 | 4 | |

■ Automatic queue choosing

- all queues
- based on requested **walltime, nodes**

■ DO NOT use “-q *queue*” or “#MSUB -q *queue*” by job-submitting

Environment variables

- www.bwhpc-c5.de/wiki/index.php/Batch_Jobs#Environment_Variables_for_Batch_Jobs
- **bwUniCluster + ForHLR + JUSTUS**

queue =

| Environment variables | Description |
|-----------------------|---------------------------------------|
| MOAB_CLASS | Class name |
| MOAB_GROUP | Group name |
| MOAB_JOBID | Job ID |
| MOAB_JOBNAME | Job name |
| MOAB_NODECOUNT | Number of nodes allocated to job |
| MOAB_PARTITION | Partition name the job is running in |
| MOAB_PROCCOUNT | Number of processors allocated to job |
| MOAB_SUBMITDIR | Directory of job submission |
| MOAB_USER | User name |

```
$ printenv | grep MOAB
```

- Using in scripts:

```
## add suffix to job output file  
./program > $program_`${MOAB_JOBID}`.log
```

Interactive jobs

■ Common

- Access to compute nodes
→ start your application direct there
- Specify resources what you need
- Auto logout when job is finished
- Submit job via “`msub -I -V`”

```
$ msub -I -V -l nodes=1:ppn=1,walltime=02:00:00
```

- `-I` = interactive
- `-V` = all environment variables are exported to the compute node

■ bwUniCluster

- www.bwhpc-c5.de/wiki/index.php/Batch_Jobs_-_bwUniCluster_Features#Interactive_Jobs

■ JUSTUS

- www.bwhpc-c5.de/wiki/index.php/Batch_Jobs_-_bwForCluster_Chemistry_Features#Interactive_jobs

Check/change status of your jobs (1)

- after submission → msub returns <job-ID>

```
$ msub job.sh
```

```
659562
```

- **commands:**

| | |
|--------------------------|--|
| \$ showq -r | All your active (running) jobs |
| \$ showq -i | eligible(idle) jobs |
| \$ showq -b | blocked jobs |
| \$ showq -c | completed jobs |
| \$ showstart <job-ID> | Get information about start time of job with <job-ID> |
| \$ showstart 16@12:00:00 | Get information about start time of 16 procs with run time of 12 hours |
| \$ checkjob <job-ID> | Get detailed information of your job → explains why your job is pending |
| \$ canceljob <job-ID> | Cancel the job with <job-ID> |

Check status of your jobs (2)

■ Command “showq”:

```
$ showq
```

active jobs-----

| JOBID | USERNAME | STATE | PROCS | REMAINING | STARTTIME |
|-------|----------|----------------|-------|-----------|---------------------|
| 12345 | /// | Running | 1 | 00:04:58 | Thu Jan 22 19:21:56 |

1 active job

eligible jobs-----

| JOBID | USERNAME | STATE | PROCS | REMAINING | STARTTIME |
|-------|----------|-------------|-------|-----------|---------------------|
| 12346 | /// | Idle | 1 | 00:04:58 | Thu Jan 22 19:21:56 |

1 eligible job

blocked jobs-----

| JOBID | USERNAME | STATE | PROCS | WCLIMIT | QUEUETIME |
|-------|----------|-------------|-------|----------|---------------------|
| 12347 | /// | Idle | 1 | 00:05:00 | Thu Jan 22 19:21:47 |

1 blocked job

Check status of your jobs (3)

■ STATE:

- Running OK, job is running
- Idle Job is waiting for free resources

- Deferred Buffer-state.
Job can not run (no free resources
or wrong resources)

- BatchHold Job is blocked by scheduler.
End-state.
Reasons: no resources,limits,failure

```
Idle → Running → Canceling == OK
```

```
Idle → Deferred → Idle → Deferred → ... → BatchHold → Canceling
```

Check status of your jobs (4)

- Check, why job can not start:

- `checkjob <job_ID>` get information of your job
- `checkjob -v -v -v <job_ID>` detailed information

Check status of your jobs (5)

example: MAXNODE limit

■ Submitted job (bwUniCluster)

```
$ msub -l nodes=1:ppn=8 -q fat <jobscrip>  
12345
```

showq:

```
blocked jobs-----  
JOBID          USERNAME      STATE PROCS    WCLIMIT          QUEUETIME  
12345          ///          Idle    5    00:05:00  Fri Jan 23 15:31:05
```

checkjob 12345:

```
State: Idle  
class:fat  
...  
NodeCount: 1  
...
```

```
BLOCK MSG: job 12345 violates active  
HARD MAXNODE limit of 2 for class fat user partition ALL  
(Req: 8 InUse: 64) (recorded at last scheduling iteration)
```

Check status of your jobs (6)

example: organisation limits

■ Submitted job (bwUniCluster)

```
$ msub -l nodes=1:ppn=1 <jobscript>
```

```
55555
```

showq:

```
blocked jobs-----  
JOBID          USERNAME      STATE PROCS    WCLIMIT        QUEUETIME  
55555          ///          Idle    1    00:10:00  Fri Jan 21 15:31:05
```

checkjob -v -v -v 55555:

```
State: Idle  
class:develop  
...
```

```
BLOCK MSG: job 55555 violates active SOFT MAXPROC limit of 1000  
for acct university X partition ALL (Req: 1 InUse: 1010) ...
```

* limits for **university_X**
* TODO: only wait!

Change status of your jobs

■ Change commands

- `canceljob <job_ID>` cancel the job with <job_ID>
- `mjobctl -c <job_ID>` cancel the job (new command)
- `mjobctl -c -w state=Idle` cancel ALL idle jobs
- `mjobctl -c -w state=Running` cancel ALL running jobs
- `mjobctl -c -w state=BatchHold` cancel ALL hold jobs
- `mjobctl -c -w user=$USER` **cancel ALL your jobs!**

```
$ showq
active jobs-----
JOBID   USERNAME  STATE PROCS  REMAINING   STARTTIME
31172   ///       Running  1           00:04:58   Thu Jan 22 19:21:56
...
blocked jobs-----
JOBID   USERNAME  STATE PROCS  WCLIMIT     QUEUE TIME
31173   ///       Idle    1           00:05:00   Thu Jan 22 19:21:47
31174   ///       BatchHold 1           00:05:00   Thu Jan 22 19:21:48
```

Example

```
#!/bin/bash
#MSUB -l nodes=2:ppn=16
#MSUB -l walltime=01:00:00
#MSUB -l pmem=2gb
#MSUB -N serial-test

mpirun ./hello
```

→ Is equal to:

```
$ msub -l nodes=2:ppn=16,walltime=01:00:00,pmem=2gb -N serial-test
<job_script>
```

Common problems

- Wrong „ppn“ setting:

```
$ msub -l nodes=3:ppn=38,walltime=00:01:00,pmem=1gb <job_script>
```

- „mem“ instead of „pmem“:

```
$ msub -l nodes=4:ppn=16,walltime=00:01:00,mem=1gb <job_script>
```

- Wrong queue

- Data in \$HOME instead of \$WORK/\$TMPDIR

- # MSUB instead of #MSUB (note the space...)

