

Using the RHUL-Psychology scheduling system

Tibor Auer

RHUL Department of Psychology

Research Fellow in MRI

Courtesy to Russel Thompson (MRC-CBSU)

Overview

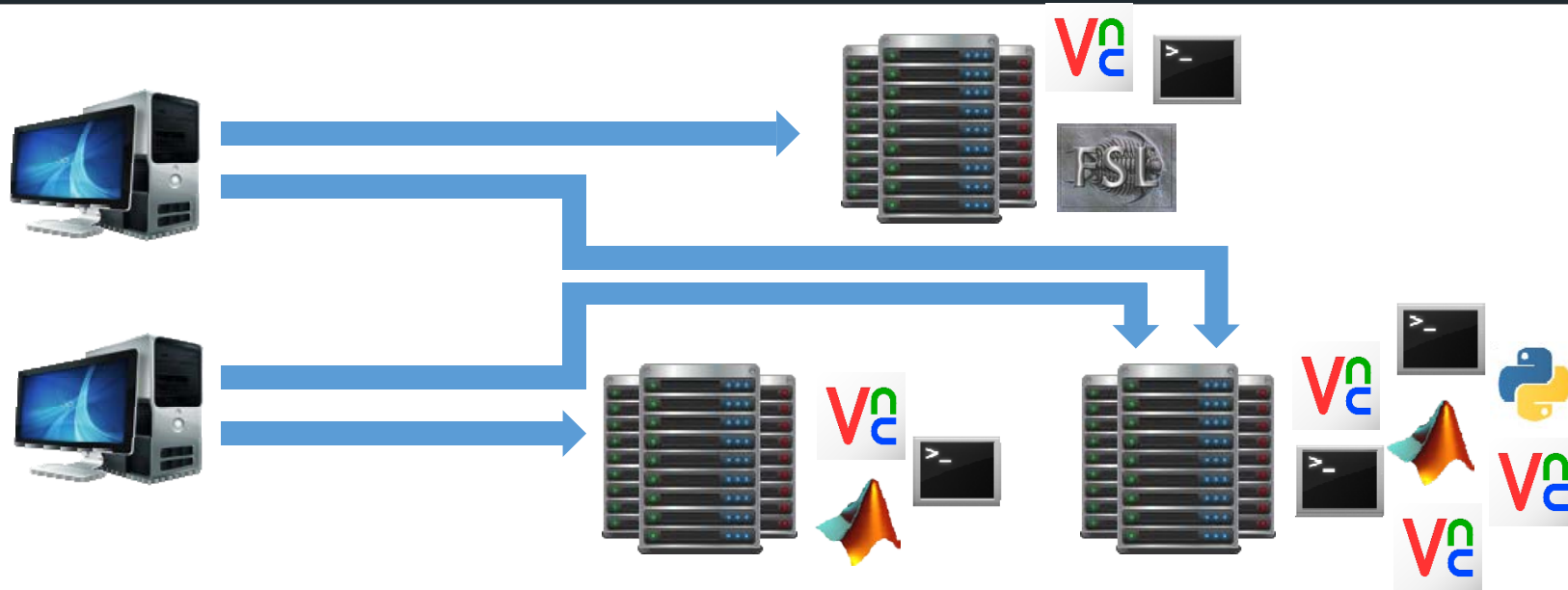
- **Why use a scheduling system?**
- **Submitting jobs**
- **Monitoring and managing jobs**
- **MATLAB**
- **Best Practices**

Overview

- **Why use a scheduling system?**
- Submitting jobs
- Monitoring and managing jobs
- MATLAB
- Best Practices

Why use a scheduling system?

Non-scheduled system (old cluster)

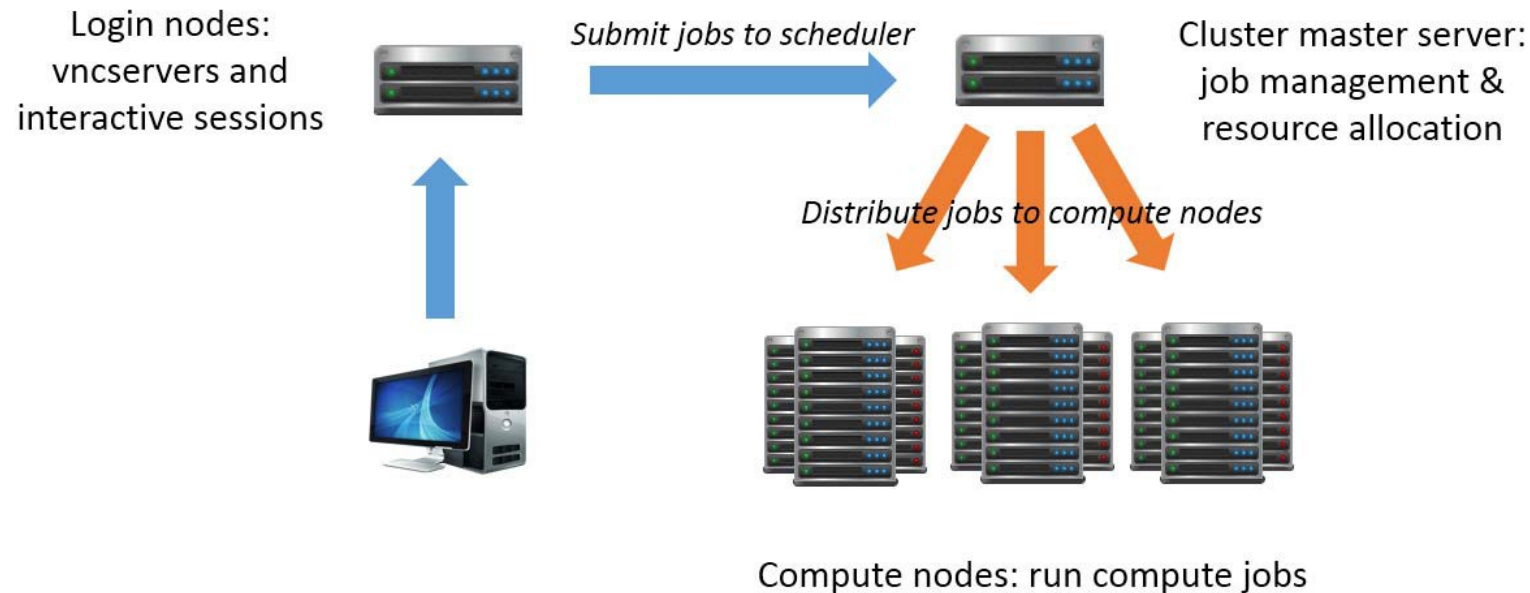


- **Issues**

- No distinction between login and compute nodes
 - → run both interactive sessions and large compute jobs on the same machines
- No management of jobs → e.g. machines can run out of memory

Why use a scheduling system?

Scheduled system (new cluster)



- **Solutions**

- Distinction between login and compute nodes
 - Login and run interactive sessions on a login node
 - Run large compute jobs on compute nodes
- Scheduling system manages allocation of resources to compute jobs

Why use a scheduling system?

Scheduled system (new cluster)

- **Resource management**
 - Dedicated resources for jobs → running jobs do not compete for the same resources
 - Resources are fully utilized, but not overloaded
- **Parallel processing (saving time)**
 - „Obviously” parallel jobs: e.g. processing data from different subjects
 - „Truly” parallel jobs: e.g. via Message Passing Interface (MPI)

Overview

- **Why use a scheduling system?**
- **Submitting jobs**
- Monitoring and managing jobs
- MATLAB
- Best Practices

Submitting jobs

Accessing



ssh →

```
login as: vwxz123
vwxz123@psyclogin.rhul.ac.uk's password:
[vwxz123@psyclogin ~]$
```

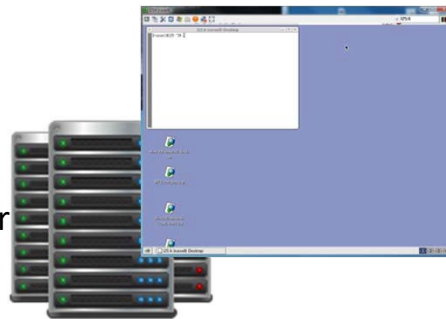


```
[vwxz123@psyclogin ~]$ vncserver :7

Desktop 'TurboVNC: psyclogin:7 (vwxz123)'
started on display psyclogin:7
```



vnc viewer →



1. Log in using SSH¹ (text only)

- Linux: *ssh*
- Win: PuTTY

2. Graphical sessions via VNC²

- *vncserver* → desktop number (ask IT)

3. Win: TurboVNC Viewer

- psyclogin.rhul.ac.uk:7

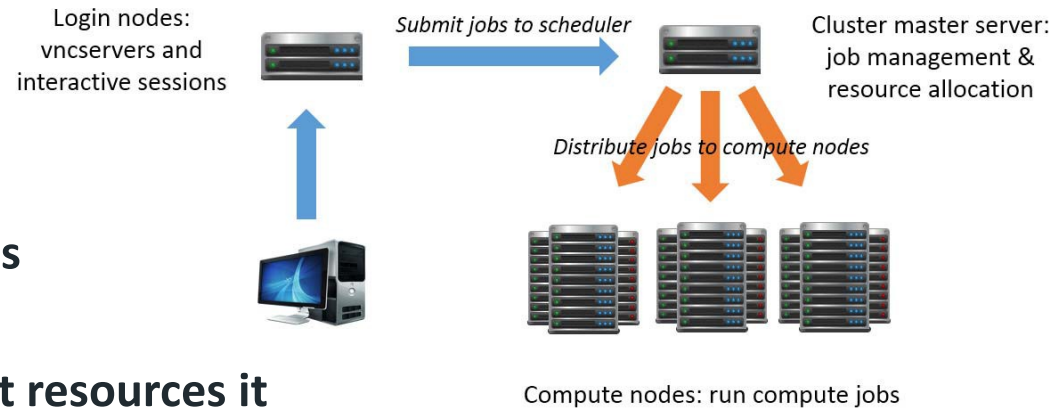
4. Close when finished!³

1. Close TurboVNC Viewer
2. *vncserver -kill :7*
3. Close SSH

Submitting jobs

- **Workflow**

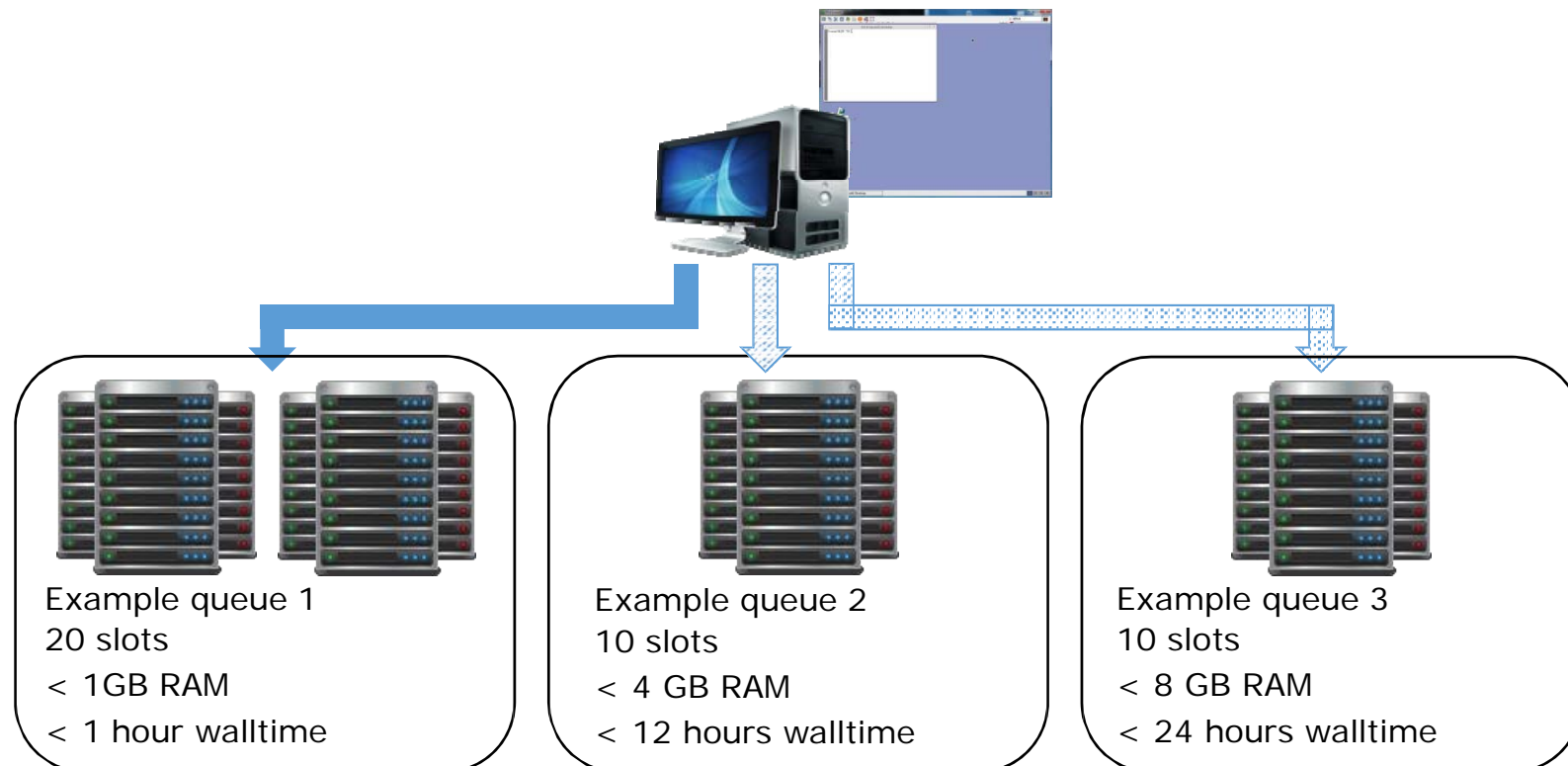
1. Log in to login node via a VNC server
2. Create a batch script to run your analyses
3. Test the batch script and determine what resources it needs (esp. RAM and CPU time)
4. Submit the script to the scheduling system and request specific cluster resources
 - Number of cores (typically 1 per job), RAM
 - Processing time („walltime”)
 - After the walltime elapses, the job is killed.
 - If the job finishes earlier, resources will be released.



Submitting jobs

- **Queues**

- A job is held in a queue and is executed as soon as the requested resources are available.
- Queues are associated with specific sets of resources
 - e.g. maximum RAM, maximum walltime, etc.
- **Choose a queue based on resource requirements!**



Submitting jobs

- **Queues**

Name	Priority	Slots	CPU Time	Memory
fastq	1	24	1h	1GB
normal	5	24	12h	2GB
longq	10	24	24h	4GB
verylongq	15	4	INFINITE	4GB

- One can run up to 32 jobs only (also subject to availability).
- Can submit more jobs, but they remain at the bottom of the queue.

Submitting jobs

qsub

- **Submit jobs in linux terminal:**

```
qsub <arguments> <command to run>
```

- Arguments:

- -n job name
- -q which queue to use
- -o path to file where standard output should be redirected
- -e path to file where standard error output should be redirect
- -l request specific resources (e.g. RAM, CPU time, number of cores)
- -v variable to pass to job batch script
- -V pass environment (e.g. FSL configuration) to worker

- E.g.:

```
qsub -q fastq -V -v \  
    RAWDIR=/.../201706131237_19810218EIJO/Series_002_MPRAGE,\  
    OUTDIR=/MRIWork/MRIWork09/cluster/demo/bet,\  
    SUBJID=1 \  
    /MRIWork/MRIWork09/cluster/demo/bet_script.sh
```

Submitting jobs

qsub

- **Submit jobs in linux terminal:**

```
qsub -q fastq -V -v \  
    RAWDIR=/.../201706131237_19810218EIJO/Series_002_MPRAGE,\  
    OUTDIR=/MRIWork/MRIWork09/cluster/demo/bet,\  
    SUBJID=1 \  
    /MRIWork/MRIWork09/cluster/demo/bet_script.sh
```

- **Content of cluster_bet.sh**

```
FNAME=S$(printf "%02d" $SUBJID)  
RAWFILE1=${RAWDIR}/${echo $(ls ${RAWDIR}) | awk '{print $1}'}  
  
mri_convert ${RAWFILE1} ${OUTDIR}/${FNAME}.nii.gz > ${OUTDIR}/log_convert.txt  
fslswapdim ${OUTDIR}/${FNAME} z -x -y ${OUTDIR}/${FNAME} >> ${OUTDIR}/log_convert.txt  
bet ${OUTDIR}/${FNAME} ${OUTDIR}/${FNAME}_brain -R -v;
```

Submitting jobs

Multiple jobs (Job arrays)

- **Submit jobs in linux terminal:**

```
qsub -q fastq -V -t 1-8 -v \  
    RAWDIR=/.../201706131237_19810218EIJO/Series_002_MPRAGE,\  
    OUTDIR=/MRIWork/MRIWork09/cluster/demo/bet,\  
    /MRIWork/MRIWork09/cluster/demo/bet_script.sh
```

- **Content of cluster_bet.sh**

```
FNAME=S$(printf "%02d" $SGE_TASK_ID)  
RAWFILE1=${RAWDIR}/${echo $(ls ${RAWDIR}) | awk '{print $1}'}  
  
mri_convert ${RAWFILE1} ${OUTDIR}/${FNAME}.nii.gz > ${OUTDIR}/log_convert.txt  
fslswapdim ${OUTDIR}/${FNAME} z -x -y ${OUTDIR}/${FNAME} >> ${OUTDIR}/log_convert.txt  
bet ${OUTDIR}/${FNAME} ${OUTDIR}/${FNAME}_brain -R -v;
```

Overview

- Why use a scheduling system?
- Submitting jobs
- Monitoring and managing jobs
- MATLAB
- Best Practices

Monitoring jobs

qstat

- For information about the cluster, queues, jobs:

```
qstat <options>
```

- Arguments:
 - -g c cluster queue summary
 - -q <queue name> selected queue
 - -f full output
 - -F [resource] full output and show (selected) resources

 - -u <user ID> jobs of selected users
 - -j <job ID> selected job

Monitoring jobs

qstat

- E.g.: Queues
 - List queues

```
[udjt005@psyclogin cluster]$ qstat -g c
```

CLUSTER QUEUE	CQLOAD	USED	RES	AVAIL	TOTAL	aoACDS	cdsuE
fastq	0.01	0	0	16	16	0	0
longq	0.01	0	0	16	16	0	0
normalq	0.01	0	0	16	16	0	0
verylongq	0.01	0	0	4	4	0	0

- Instances of „fastq” queue

```
[udjt005@psyclogin cluster]$ qstat -q fastq -f
```

queuename	qtype	resv/used/tot.	load_avg	arch	states
fastq@psycmri01.rhul.ac.uk	BP	0/0/8	0.09	lx-amd64	
fastq@psycmri02.rhul.ac.uk	BP	0/0/8	0.05	lx-amd64	

Monitoring jobs

qstat

- E.g.: Queues
 - Requestable resources of „fastq” queue

```
[udjt005@psyclogin cluster]$ qstat -q fastq -F
queuename                qtype resv/used/tot. load_avg arch          states
-----
fastq@psycmri01.rhul.ac.uk  BP    0/0/8          0.07    lx-amd64
qf:h_cpu=01:00:00
qf:h_rss=1.000G

hl:mem_total=31.357G
hl:swap_total=6.000G
hl:virtual_total=37.357G

hl:mem_free=30.631G
hl:swap_free=5.906G
hl:virtual_free=36.537G

hl:mem_used=743.750M
hl:swap_used=95.770M
hl:virtual_used=839.520M
...
```

Monitoring jobs

qstat

- E.g.: Queues
 - Requestable walltime of all queue instances

```
[udjt005@psyclogin cluster]$ qstat -F h_cpu
queuename                qtype resv/used/tot. load_avg arch          states
-----
fastq@psycmri01.rhul.ac.uk  BP    0/0/8          0.19    lx-amd64
    qf:h_cpu=01:00:00
-----
fastq@psycmri02.rhul.ac.uk  BP    0/0/8          0.06    lx-amd64
    qf:h_cpu=01:00:00
-----
normalq@psycmri01.rhul.ac.uk BP    0/0/8          0.19    lx-amd64
    qf:h_cpu=12:00:00
-----
normalq@psycmri02.rhul.ac.uk BP    0/0/8          0.06    lx-amd64
    qf:h_cpu=12:00:00
-----
longq@psycmri01.rhul.ac.uk  BP    0/0/8          0.19    lx-amd64
    qf:h_cpu=24:00:00
-----
longq@psycmri02.rhul.ac.uk  BP    0/0/8          0.06    lx-amd64
    qf:h_cpu=24:00:00
...

```

Monitoring jobs

qstat

- E.g.: Jobs
 - All jobs of user abcd123

```
[udjt005@psyclogin cluster]$ qstat -u abcd123
```

job-ID	prior	name	user	state	submit/start	at	queue	slots	ja-task-ID
4928	0.05500	Job1	abcd123	Eqw	04/03/2017	10:09:52		12	
4929	0.05500	Job4	abcd123	Eqw	04/03/2017	10:12:23		12	
4931	0.05500	Job1	abcd123	Eqw	04/03/2017	11:02:53		12	
4933	0.05500	Job1	abcd123	Eqw	04/03/2017	15:58:48		12	
4937	0.05500	Job1	abcd123	Eqw	04/04/2017	11:58:44		12	
4938	0.05500	Job1	abcd123	Eqw	04/04/2017	12:11:47		12	
4939	0.05500	Job1	abcd123	Eqw	04/04/2017	12:28:04		12	
4943	0.05500	Job1	abcd123	Eqw	04/04/2017	14:59:11		12	
4944	0.05500	Job1	abcd123	Eqw	04/04/2017	15:24:57		12	
4970	0.05500	Job1	abcd123	Eqw	04/06/2017	15:40:59		12	
4971	0.05500	Job1	abcd123	Eqw	04/06/2017	15:50:42		12	
4972	0.05500	Job1	abcd123	Eqw	04/06/2017	16:03:00		12	
4974	0.05500	Job1	abcd123	Eqw	04/06/2017	16:32:37		12	
4975	0.05500	Job1	abcd123	Eqw	04/06/2017	16:54:32		12	
4976	0.05500	Job1	abcd123	Eqw	04/06/2017	17:06:40		12	
4979	0.05500	Job1	abcd123	Eqw	04/06/2017	17:53:28		12	
4980	0.05500	Job1	abcd123	Eqw	04/06/2017	18:05:16		12	
4983	0.05500	Job1	abcd123	Eqw	04/07/2017	08:21:11		12	
4984	0.05500	Job1	abcd123	Eqw	04/07/2017	08:33:15		12	

Monitoring jobs

qstat

- **E.g.: Particular job**

- Details of job 4928¹

```
[udjt005@psyclogin cluster]$ qstat -j 4928
=====
job_number:                4928
submission_time:           Tue Apr  4 12:11:47 2017
sges_o_workdir:            /MRIWork/.../CoSMoMVPA-master/mvpa
hard_resource_list:        h_cpu=7200
stdout_path_list:          NONE:NONE:/MRIWork/.../CoSMoMVPA-master/mvpa/Job1/Job1.log
script_file:               /usr/local/apps/matlab/R2015b/toolbox/local/communicatingJobWrapper.sh
error_reason                1:      04/04/2017 13:21:23 [795555711:29570]: error: can't open
output file "/MRIWork/.../CoSMoMVPA-master/mvpa/Job1/Job1.log": Stale file handle
scheduling info:           Job is in error state
...
```

Monitoring jobs

Output and error messages

- Execution of a program locally (in a terminal)

```
n=1
bet \
/MRIWork/MRIWork09/tibor_auer/temp/T1_$(printf "%02d" $n).nii.gz \
/MRIWork/MRIWork09/tibor_auer/temp/T1_$(printf "%02d" $n)_brain.nii.gz -R -v
```

- Stderr in case of error (also in the terminal)

```
bash: bet: command not found
```

- Stdout after fsl configured (also in the terminal)

```
IN=/MRIWork/MRIWork09/tibor_auer/temp/T1_01
OUT=/MRIWork/MRIWork09/tibor_auer/temp/T1_01_brain
...
min 0 thresh2 0 thresh 97.8208 thresh98 978.208 max 2779
c-of-g 134.204 88.4442 150.463 mm
radius 90.7149 mm
median within-brain intensity 281
self-intersection total 307.343 (threshold=4000.0)
```

Monitoring jobs

Output and error messages

- Execution of a program as job (on the cluster)

```
qsub -q fastq -V -v outdir=/MRIWork/MRIWork09/tibor_auer/temp,n=1 cluster_bet.sh
```

- Jobs running on the cluster are not interactive, i.e. outputs are not sent to the terminal but redirected to files saved in home by default.
 - Stdout: <home>/<script name>.o<job ID>[.<task ID>]¹
 - Stderr: <home>/<script name>.e<job ID>[.<task ID>]¹

- E.g.:

- Content of cluster_bet.e5038

```
/usr/local/share/sgе/default/spool/psycmri02/job_scripts/5038: line 2: bet:  
command not found
```

- Content of cluster_bet.o5039

```
IN=/MRIWork/MRIWork09/tibor_auer/temp/T1_01  
OUT=/MRIWork/MRIWork09/tibor_auer/temp/T1_01_brain  
...  
min 0 thresh2 0 thresh 97.8208 thresh98 978.208 max 2779  
c-of-g 134.204 88.4442 150.463 mm  
radius 90.7149 mm  
median within-brain intensity 281  
self-intersection total 307.343 (threshold=4000.0)
```



Monitoring jobs

Output and error messages

- **Execution of a program as job (on the cluster)**

```
qsub -q fastq -V -v outdir=/MRIWork/MRIWork09/tibor_auer/temp,n=1 cluster_bet.sh
```

- Jobs running on the cluster are not interactive, i.e. outputs are not sent to the terminal but redirected to files saved in home by default.
 - Stdout: <home>/<script name>.o<job ID>[.<task ID>]¹
 - Stderr: <home>/<script name>.e<job ID>[.<task ID>]¹

- Can specify output locations

```
outdir=<project directory>
```

```
qsub -q fastq -o ${outdir}/job1_out.txt -e ${outdir}/job1_err.txt -V -v n=1 cluster_bet.sh
```


Manging jobs

- **Delete jobs**

```
qdel <job id>
```

- Can delete all jobs at once

```
qselect -u `whoami` | xargs qdel
```

Overview

- Why use a scheduling system?
- Submitting jobs
- Monitoring and managing jobs
- MATLAB
- Best Practices

Submitting MATLAB jobs

- **MATLAB Distributed Computing Server (DCS) and Parallel Computing Toolbox (PCT):**
 - Collection of functions for running matlab jobs in parallel environment
 - DCS supports 3rd party schedulers at various degree (PBS/TORQUE – well, SGE – limited)
 - MATLAB translates jobs into a series of qsub commands with options according to the configuration
- **Workflow**
 1. Create analysis script
 2. Create scheduler object using DCS/PCT functions
 3. Configure scheduler (log file location, walltime, memory)
 4. Create jobs to run
 5. Submit jobs
- **Options**
 - MATLAB PCT functions with RHUL-Psycho wrapper
 - parfor loops with RHUL-Psycho cluster profiler (you can use the wrapper, too)
 - **Automatic Analysis**

Submitting MATLAB jobs

RHUL-Psycho wrapper

- **PsychoSGE constructor**
 - In /usr/local/apps/psycapps/cluster/PsychoSGE.m
 - Provides
 - Sets up PsychoSGE cluster profile¹
 - Configures cluster (logfile location, walltime, memory)
 - Create pool for parfor

```
%% Config
C = PsychoSGE;
% Defaults
% C.Walltime = 2;    % Hours
% C.Memory = 2:     % GB
% C.Logs = pwd;

cluster = C.getCluster;

pool = C.getPool(8); % 8 jobs in parallel
```

Submitting MATLAB jobs

RHUL-Psycho wrapper

- **Example:**
 - Calculate eigenvalues of a square matrix of random numbers
 - Serial execution

```
N = 2000;  
  
%%  
for j = 1:8  
    out(:,j) = eig(rand(N));  
end
```

Submitting MATLAB jobs

RHUL-Psycho wrapper

- **Example:**
 - Calculate eigenvalues of a square matrix of random numbers (with benchmarking)
 - Simple job submission

```
N = 2000;  
  
C = PsychoSGE;  
C.Walltime = 0.5; % Hours  
C.Memory = 1; % GB  
C.Logs = pwd;  
cluster = C.getCluster;  
  
for j = 1:8  
    job(j) = cluster.createJob;  
    job(j).createTask(@eig, 1, {rand(N)});  
    job(j).submit;  
end  
  
while ~all(strcmp({job.State}, 'finished')), pause(1); end  
  
for j = 1:8  
    out(j) = job(j).fetchOutputs;  
end
```

Create and configure the cluster

Create and submit jobs and tasks

Wait for all jobs to finish

Retrieve output(s)

Submitting MATLAB jobs

RHUL-Psycho wrapper

- **Example:**
 - Calculate eigenvalues of a square matrix of random numbers (with benchmarking)
 - Simple job submission
 - Create and submit jobs and tasks

```
for j = 1:8
    job(j) = cluster.createJob;
    job(j).createTask(@eig, 1, {rand(N)});
    job(j).submit;
end
```

- *createTask(F, N, {inputs})*
 - F A handle to the function that is called.
 - N Number of output arguments to be returned.
 - {inputs} A row cell array specifying the input arguments. Each element in the cell array will be passed as a separate input argument.

Submitting MATLAB jobs

RHUL-Psycho wrapper

- **Example:**

- Calculate eigenvalues of a square matrix of random numbers (with benchmarking)
- Simple job submission
 - Wait for all jobs to finish

```
while ~all(strcmp({job.State},'finished')), pause(1); end
```

- *job.State*: status of the job ('pending', 'queued', 'running', 'finished', 'failed')

- Retrieve output(s)

```
for j = 1:8  
    out(j) = job(j).fetchOutputs;  
end
```

- *out = job.fetchOutputs*
 - Returns the result(s) from the tasks of a finished job.
 - If the scalar job has M tasks, each row of the M-by-N cell array data contains the output arguments for the corresponding task in the job. Each row has N elements, where N is the greatest number of output arguments from any one task in the job.

Submitting MATLAB jobs

RHUL-Psycho wrapper

- **Example:**
 - Calculate eigenvalues of a square matrix of random numbers (with benchmarking)
 - Parfor

```
N = 2000;  
  
C = PsychoSGE;  
C.Walltime = 0.5; % Hours  
C.Memory = 1; % GB  
C.Logs = pwd;  
pool = C.getPool(8);  
  
parfor j = 1:8  
    out(:,j) = eig(rand(N));  
end  
  
C.closePool;
```

Create and configure the pool

Run usual functions¹

Close pool

Overview

- **Why use a scheduling system?**
- **Submitting jobs**
- **Monitoring and managing jobs**
- **MATLAB**
- **Best Practices**

Best Practices

- **Write and debug your scripts locally before submitting!**
- **Requesting the appropriate resources allows the scheduling system to operate most efficiently.**
 - Make a note of the required resources (top) – especially memory and CPU time
 - Under-requesting (e.g. 1GB RAM when you need 4GB) can cause the job to crash.
 - Over-requesting (e.g. 4GB RAM when you only need 1GB) means
 - Fewer jobs will run simultaneously.
 - Jobs may wait for longer for resources.

Further Information

- **Cluster wiki on the CUBIC website**

http://www.cubic.rhul.ac.uk/wiki/doku.php?id=cluster:cluster_root