

Exercises

a) Log in to Taito either with your training or CSC user account, either from a terminal (with X11 forwarding) or using NX client

b) Go to working directory and untar the exercises file (which you have copied/retrieved earlier)

```
$ cd $WRKDIR
$ tar xf <exercise-file-name.tar>
```

Note. The not bolded dollar sign "\$" at the beginning of some lines is not to be typed, but marks the command prompt. What follows is to be typed on the command prompt. Text in *italics* is not a command and you can choose what to put there (but you need to be consistent later).

1. Simple batch job script

a) Create a batch job script that prints the compute node on which it is running.

Using nano editor (use whichever editor you like):

```
$ nano test_hostname.sh
```

Put this into the file.

```
#!/bin/bash -l
#SBATCH -J print_hostname
#SBATCH -o output.txt
#SBATCH -e errors.txt
#SBATCH -t 00:01:00
#SBATCH -p test
#SBATCH -n 1
#SBATCH --mem=500
echo "This job runs on the host: "; hostname
```

To exit from the nano editor:

```
CTRL+O; CTRL+X (confirm save)
```

Submit the batch script to Taito

```
$ sbatch test_hostname.sh
```

b) Check the job status.

In the following command replace *<your username>* with the training account or your CSC username, which ever you used to log in to Taito. If you are not sure which it is, you can check it with `whoami` or with this command `echo $USERNAME`).

```
$ squeue -u <your username>
```

c) What and where did the job print out?

```
$ less output.txt (type q to quit)
$ less errors.txt (type q to quit)
```

2. Simple R job

Run a simple R job from a script. The script will fit a straight line through a file containing some x,y value pairs.

a) Create a script containing the R commands.

Go to the directory "r-job", where you have the data (a file called `data.csv`). In that directory, create the following R script file (R commands to be executed). Name the file `fit.R`

```
mydata <- read.csv("data.csv")
fit <- lm(y~x,mydata)
fit$coefficients
```

b) Set up the R environment with the module command

```
$ module load r-env
```

c) Run the script with

```
$ R --no-save --no-restore -f fit.R
```

d) Did the job succeed? What are the fit coefficients?

3. Simple R job as a batch job

Run the previous R script as a batch job.

a) Create a batch job script, which will submit the job to the queue.

Start with the batch script file from the first exercise. Copy it to the current folder and edit it as follows. Add to the batch script file the `module load ...` command and then the command to run the R-script (the commands you gave in the command line in the previous exercise, you can remove the hostname command if you like). The other SLURM requirements can be as in the previous example.

b) Submit the batch script with

```
$ sbatch your_script_name
```

c) Check the job results.

Did the job succeed? Where are the fit constants?

4. Run tens of R batch jobs as an array job

In this example, we will repeat the previous fitting job for 20 datasets using the array job functionality of SLURM.

a) Prepare a list of files to process.

Go to the folder named **r-array**. Create there a file called **datanames.txt**. This file will contain the names of all those files that will be used as input to the fitting. Run the following commands to create it.

```
$ cd data_dir
$ ls * > ../datanames.txt
```

b) Write the R script, that will do the fitting.

Go back to the r-array folder, create a script named **modelscript.R** and put the following commands to it (you can copy the previous script and edit that, or start from scratch).

```
dataname <- commandArgs(trailingOnly = TRUE)
mydata <- read.csv(paste0("data_dir/",dataname))
fit <- lm(y~x,mydata)
write(fit$coefficients,file=paste0("result_dir/",dataname,"_result.txt"))
```

The first line will extract from the batch command the name of the dataset to be fitted. The next line reads that data into the variable mydata. Then we fit, like in the previous example, and finally write the coefficients into a file.

c) Create a batch script to submit the job.

Name it **R_array.sh**. Copy the contents from the previous example. Add the following line among the other lines starting #SBATCH:

```
#SBATCH -a 1-20
```

It will ask SLURM to run an array of 20 jobs. Edit the output and error files to go to their own directories and files by editing/adding:

```
#SBATCH -o out/output%a.txt
#SBATCH -e err/errors%a.txt
```

After the line with module **load r-env/default**, add the following line

```
dataname=$(sed -n "$SLURM_ARRAY_TASK_ID"p datanames.txt)
```

and replace the line to run the R command into:

```
R --no-save --no-restore -f modelscript.R --args $dataname
```

You should now have:

- 1) `datanames.txt`, which has the names of your datafiles
- 2) `modelscript.R`, which contains the R code to do the fitting
- 3) `R_array.sh`, which is the batch script to submit the job
- 4) (and the folders `out`, `err`, `data_dir`, `result_dir` which were there already)

d) run the batch script with

```
$ sbatch R_array.sh
```

You should get the fit coefficients in separate files in the `result_dir`. Let's now use interactive R to look at the results.

e) Initialise R and Rstudio and start RStudio

```
$ module load r-env # (if you did this already, no need to repeat)
$ module load rstudio
$ rstudio
```

f) Collect the results and plot them.

In the GUI that opens, click "File", "Open File", choose `analyse.R` and accept. The script contents will appear in the top left window. To execute the commands move your mouse to the line, and press "ctrl - Enter" to run them. Run them in order. The original data was created by calculating the y values by $y=2x$ + some random noise. How do the fit coefficients match that?

5. Copying files to `hpc_archive`

You can access `hpc_archive` only from Taito and Sisu. After installing some tools, you can access IDA, which also uses iRods technology, also from your computer.

Log in to Taito. Check the contents of your `hpc_archive`:

```
$ ils
```

Show the directory that you're in in `hpc_archive`

```
$ ipwd
```

Show the directory that you're in in `taito`:

```
$ pwd
```

Create a directory in hpc_archive

```
$ mkdir test
```

Move to test directory in hpc_archive

```
$ cd test
```

Confirm where you are in hpc_archive

```
$ pwd
```

Copy (put) a file to the test directory in hpc_archive:

```
$ iput <filename>
```

Confirm that the file is in hpc_archive

```
$ ls
```

Copy the file back from hpc_archive, but to a local folder called **localtest**

```
$ mkdir localtest
$ cd localtest
$ iget <filename>
```

6. Archive a file

Make a new directory in hpc_archive home directory (not under test) called **mysafe**. Copy there your files, e.g., **test_hostname.sh** and **R_array.sh** from previous exercises. Compress and archive a file, e.g., **fit.R** using a command **tar -zcvf**, and copy it to **mysafe**.

7. Batch job with thread parallelization

Some applications can be run in parallel to speed them up. In this example you run the HMMER software to describe and analyze related or similar sequence areas both in serial and parallel to see if the jobs speed up.

HMMER uses a database that is already installed, but the protein sequences you want to study need to be copied first to be used as input:

```
$ cp /appl/bio/hmmer/example.fasta .
```

Let's first run the job with just one core. Copy one of the old batch scripts to current directory, and change / add the following items in it:

- 1) Output to **out_%j.txt**

2) error to `err_%j.txt`

3) run time 10 minutes

4) load the `hmmmer` -module

5) run command:

```
hmmmscan $HMMERDB/Pfam-A.hmm example.fasta > example_1.result
```

Submit the job with:

```
$ sbatch jobscript.sh
```

Submitting the job echoes the SLURM job id number to the screen, but that is also shown in the output and error filenames (`out_<SLURM_JOBID>.out`). Check if the job is running with

```
$ squeue -u <your username>
```

Or

```
$ squeue -j <SLURM_JOBID>
```

Once the job is finished you can check how much memory and time it used:

```
$ sacct -j <SLURM_JOBID> -o elapsed,reqmem,maxrss
```

Did you reserve a good amount of memory? (not excessively too much, but enough to not be close to memory running out and terminating the job).

Now, let's try with 4 cores. At this point we'll also switch to using the environment variable `$SLURM_CPUS_PER_TASK` to avoid mistakes and the need to change that in many places. Add this line to the batch script:

```
#SBATCH --cpus-per-task=4
```

Change the run command to:

```
hmmmscan --cpu $SLURM_CPUS_PER_TASK $HMMERDB/Pfam-A.hmm \
  example.fasta > example_${SLURM_CPUS_PER_TASK}.result
```

As you've asked for 4 cpus per task, the environment variable `$SLURM_CPUS_PER_TASK` evaluates to 4 when the script is run, and you only need to change the number on the `#SBATCH` line.

Submit the job and check with the `sacct` command how long it took to run the hmmer job and how did the memory usage change and try to answer these questions:

- a) Does it make sense to use 4 cores instead of 1?
- b) Was the memory reservation ok?
- c) Does it make sense to use more than 4 cores?
- d) How to speed up the job?

8. Batch job memory consumption

Create a new R-script (like in exercise 3) named `mem-test.R`. It should have the following contents:

```
dim=10
mat <- matrix(rnorm(dim*dim), dim)
mat3 <- t(mat)
```

Variable `dim` is the dimension of the square matrix, which will be filled with normal distributed random numbers. Make a new batch script to run your R-script `mem-test.R`. In addition to the resource requests your script needs to load the R-environment and then run the R-script (as in exercise 3). Submit the script with the `sbatch` command.

Once the job has been completed (how can you check if it is running or queuing?), check with `sacct` how much memory was used as in the previous exercise. If the job completed successfully, increase the `dim` variable in your script *i.e.* make a bigger matrix and rerun the job. Note also the time it takes to run the job. How does the time and memory needed by the job depend on the number of elements in the array?

dim	# of elements	Time [s] (Elapsed)	Memory used (MaxRSS)

- a) How big matrix is needed to exceed the default allowed memory of the batch job?