

# Analysing data in a HPC environment using R

---

**Date:** 14-15 december 2022

**Time:** 9:00 - 12:00

**Zoom:** <https://umu.zoom.us/j/61971917562?pwd=U3NrWkc1MXp6QnYvcGM0Z3QyOTdOZz09>

(<https://umu.zoom.us/j/61971917562?pwd=U3NrWkc1MXp6QnYvcGM0Z3QyOTdOZz09>)

**Instructor:** Pedro Ojeda, Henric Zazzi, Birgitte Brydsö

## Schedule

---

### Wednesday 14 December

Time	Lesson	Tutor
9:00 - 9:10	Welcome	Pedro Ojeda
9:10 - 9:40	What is parallelization	Pedro Ojeda
09:40 - 10:20	Introduction to HPC2N	Birgitte Brydsö
10:20 - 10:35	Coffee break	
10:35 - 10:50	How to use RStudio	Pedro Ojeda
10:50 - 11:15	Lab: Login using RStudio/Thinlinc	Birgitte Brydsö/Pedro Ojeda
11:15 - 11:30	Serial R	Henric Zazzi
11:30 - 12:00	Lab: serial computing in R	Henric Zazzi

### Thursday 15 December

Time	Lesson	Tutor
9:00 - 9:30	Parallel computing in R	Henric Zazzi
9:30 - 10:15	Lab: shared memory computing	Henric Zazzi
10:15 - 10:30	Coffee break	
10:30 - 11:00	Best Practices for R in HPC	Pedro Ojeda
11:00 - 11:30	Usable parallelized R functions	Pedro Ojeda
11:30 - 12:00	Questions and comments	All

# Important links

---

How R packages are installed on HPC2N

[https://www.hpc2n.umu.se/resources/software/user\\_installed/r](https://www.hpc2n.umu.se/resources/software/user_installed/r)

([https://www.hpc2n.umu.se/resources/software/user\\_installed/r](https://www.hpc2n.umu.se/resources/software/user_installed/r)).

Installation information for the course R in HPC

<https://www.hpc2n.umu.se/events/courses/2022/R-in-HPC/setup>

(<https://www.hpc2n.umu.se/events/courses/2022/R-in-HPC/setup>).

Installation information (Ubuntu) for the course “R in an HPC environment”:

<https://umeauniversity.sharepoint.com/:w:/s/HPC2N630/EeLAM89iSrhHuMbduJkzjDwBoqxTHNo5TF6b0TplvaoLPw>

(<https://umeauniversity.sharepoint.com/:w:/s/HPC2N630/EeLAM89iSrhHuMbduJkzjDwBoqxTHNo5TF6b0TplvaoLPw>).

Repository for presentations and labs:

[https://github.com/menzzana/parallel\\_R\\_course/](https://github.com/menzzana/parallel_R_course/) ([https://github.com/menzzana/parallel\\_R\\_course/](https://github.com/menzzana/parallel_R_course/)).

Course evaluation: <https://events.prace-ri.eu/event/1458/surveys/1023> (<https://events.prace-ri.eu/event/1458/surveys/1023>).

Recordings: <https://www.youtube.com/watch?v=5HxbfiqN5kY&list=PL6jMHLEmPVLy63E9RXwHivhaN0epxaEVi>

(<https://www.youtube.com/watch?v=5HxbfiqN5kY&list=PL6jMHLEmPVLy63E9RXwHivhaN0epxaEVi>)

# Questions

---

## What is parallelization

---

### Question

Regarding Amdahls law. Would this be calculated with time complexity or running time?

### Answer

In general we use the running time or the time for the simulation to finish. By timing specific parts of the code, either manually or by using debugger, we can analyze the running time on specific parts of the code.

# Introduction to HPC2N

---

## Question

Do we need the package “clusternor” for the course or always when running R on HPC2N?

## Answer

No, you do not need this package in general. It is just a package that is an example for a function called k-means, which a parallelized version for R already exists. An example of this will be presented tomorrow.

## Question

If I want to use ssh X11 for Rstudio, how do I load Rstudio? (Only from thinlink was shown in the presentation)

## Answer

I would not recommend running RStudio via X11 as then you will be running on the login node and this has the potential to overload the node and hinder other users sharing it. Thinlinc has dedicated nodes on the cluster for running RStudio etc...  
RStudio is not installed on login node either on Kebnekaise or Dardel.

## Question

How can I submit a job with a package I installed myself?

## Answer

You install the package yourself in your \$HOME folder for that specific version of R. When running your job R will see that it is a local package installation and run with it. If R does not detect that set *export R\_LIBS\_USER=<LOCAL INSTALLATION FOLDER>* in your batch script to point to your locally installed packages.

## How to use RStudio

---

## Question

When I click to open the RStudio in the menu will it always open the version for R 4 right

## Answer

Yes. R 4.0.4. If you want a different one you need to start from the command line (and then you also need to load the R module)

## Question

You used both R CMD BATCH and Rscript in your submitscript. is that how we should use Rscript?

## Answer

No. It is just 2 examples on how to submit your scripts

## Question

When I type “rstudio”, it says that the command is not known...  
Through the terminal!

## Answer

You should run via Thinlinc as RStudio is not available via terminal.

Okay thanks!!

## Question

In what situation would R CMD BATCH be the preferred option over Rscript? It seems Rscript allows more flexibility?

## Answer

You can always use Rscript if you want as it allows more flexibility.

## Question

Now in Thinlinc, when I use the rstudio command, a message pops up, saying “Path to R not specified, and no module binary specified; Invalid R module ()”...

## Answer

Did you load the R modules?

## Question

Out of curiosity. Is there a way to connect to jupyterlabs using an R notebook via thinlinc or the command line in Kebnekaise?

## Answer

Send a mail to HPC2N support and they will help you set it up.

([support@hpc2n.umu.se](mailto:support@hpc2n.umu.se) (<mailto:support@hpc2n.umu.se>))

## Question

I can't log in on think link

I use server: [kebnekaise-tl.hpc2n.umu.se](http://kebnekaise-tl.hpc2n.umu.se) (<http://kebnekaise-tl.hpc2n.umu.se>).

Username: jnord

Password same as i used to log in with ssh

only one time with wrong password

Banned forever or temporarily?

Same error msg again "Connection to kbenekaise"... , timed out

## Answer

If you tried several times with wrong passwords, it could be that your account was banned.

Temporarily if this happens. I am checking with my colleagues on the status.

## Question

I downloaded a Thinlinc installer for Mac and double clicked it but it does not do anything... Why?

The usual installer window does not pop up. Nothing happens after I double click the downloaded installer. I don't see breakout room #2? OK, see you tomorrow.

## Answer

I am in breakout room 1.

Could you take a look at the launchpad and see if the icon of ThinLinc is there?

## Serial R

---

## Question

I ran R CMD BATCH --no-save --no-restore lab\_serial.R but the results did not show up in the terminal like when I used Rscript. Where is the result stored?

## Answer

I know the answer now. It is lab\_serial.Rout  
Good!

## Parallel computing in R

---

### Question

Two things; is the recordings up from yesterday? Also, and this is just a suggestion - I think it would be better to put the Q/A stuff in a separate document. Scrolling down past all the other info is troublesome, especially since the page jumps all the time. You could just put the link to that page at the top of the Q/A page

### Answer

Thanks for the suggestion we will take it into account.

The videos are being uploaded here: <https://www.youtube.com/@HPC2N/playlists>  
(<https://www.youtube.com/@HPC2N/playlists>)  
under R in HPC playlist

### Question

```
input <- list(1:100000)  
res <- mclapply(input, calcpi, mc.cores = n)
```

I am wondering why this “res” only have one number, I thought it will have a list of n

### Answer

Yes, that is an error as it creates a list of 1. I will correct it.  
In reality you should do

```
input <- 1:1000000
```

so remove the list command. It has now been corrected in the answers as well. Thanks for noticing.

You could however also create a big list using

```
input <- vector(mode='list', length=100000)
```

As mclapply can have as input either List, vector or data frame

## Question

is this correct: `/pfs/stor10/users/home/j/jnord/documents/r/clang/dag2lab1/lab_parallel.Rout.1` ?

Also:

when module-loading openBlas i got no extra speedup. Why?

## Answer

I think that the installation on Kebnekaise for R already has OpenBlas as it was installed with it.

Yes, OpenBlas is the one installed on Kebnekaise. You can check it with the **`sessionInfo()`** function

## Question

How does tidyverse `map()` and `future()` functions compare to `apply()` and `parallel()` etc? Similar speeds?

## Answer

`map()` and `apply()` are very similar in the functioning. There are some slight differences with more capabilities of `map()` but here an additional package tidyverse needs to be installed. `future()` is used for asynchronous computing where different (maybe plenty) are generated but not executed. Execution only occurs when tasks are explicitly called for finalizing. `parallel()` is a “backend” for many parallelized packages.

## Question

I am trying Lab 3 but the time elapsed is much much higher when doing this (7-8 s, not related to the number of cores). I might be specifying the `registerDoParallel(n)` in the wrong place when I do it inside the loop?

```
...  
no_cores <- detectCores() - 1  
...  
registerDoParallel(n)  
res <- foreach(i=1:100000) %do% calcpi(i)  
...
```

`%dopar%` improved things but still quite slow, improved with the number of cores though. Will paste code

```
library(parallel)  
library(foreach)  
library(doParallel)
```

```

a<-NULL
calcp_i <- function(no) {
y <- runif(100)
x <- runif(100)
z <- sqrt(x2+y2)
length(which(z<=1))*4/length(z)
}
no_cores <- detectCores() - 1
for (n in 1:no_cores) {
print(n)
start_time <- Sys.time()
registerDoParallel(n)
res <- foreach(i=1:100000) %dopar% calcp_i(i)
vres <- unlist(res)
print(mean(vres))
print(Sys.time() - start_time)
a[n]<-(Sys.time() - start_time)
}
plot(a)

```

should there be makeCluster when using foreach too? I thought registerDoParallel would cover that? Could you provide a tweak or hint?

## Follow-up question

so this is the tweak that I should make? It unfortunately slowed down to about 12s

```

...
nproc <- makeCluster(n)
registerDoParallel(nproc)
res <- foreach(i=1:100000) %dopar% calcp_i(i)
stopCluster(nproc)
...

```

## Follow-up 2

I do get Warning messages that I am unsure why they are caused (9 for this run):

1: In list(...) : closing unused connection 26 (<-localhost:11370)

## Follow-up answer 2

Good question, I also do get this from time to time. Usually it is because you have not stopped the cluster, but I see that you have done that, which I also have, so probably there are some orphaned processes that are closed badly when the script terminates. In my case it is intermittent meaning that it just happens from time to time, but it does not affect the actual script



## Final follow-up

The trick was to move the vector out of the foreach call (like you did in the example). When I tweaked the code like this (below) then I obtained the same speed as the example

...

```
input <- 1:10000
```

```
res <- foreach(i=input) %dopar% calcpi(i)
```

...

## Answer

how the timings look like? Could you paste the code?

## Answer 2

You are missing the command `makeCluster` in this example.

Yes, `registerDoParallel` covers that if you are running on all cores, If you are running on a defined set of cores, then you need to *makeCluster* before and *stopCluster* after. See slide 22.

## Answer 3

Yes, that looks correct. Also I have not tried this example in full and this could be a good example that the amount of overhead of setting this up (parallel distributed processes etc...) takes longer time than the actual computation as you need to set it up for each process. This example might be too simple to show that and if you had heavy calculations for each foreach steps then parallelisation with foreach would make it much faster.

In general you can see that *mclapply* is quite faster than running *foreach*.

## Question

when I run lab3 script on cluster, then then my front en session is busy while the job is running. Why ist this? Isn't the job running on the node in the background, and I should get control back to my command line?

OK

How do I use `mc.cores=...` when I use foreach. That was not spelled out in the slider about foreach (i think?). Same as with `mcapply`?

```
makeCluster
```

Aha. OK slide 22 . I see. Great.

## Answer

If you run on the cluster with `detectCores()-1` then it will run on all cores on that node, and you are sharing the node with everyone else. Please restrict it more like running on 1-10 cores for example. The question that was posted as answer for the labs is better suited for your laptop that

you are not sharing.

It was specified in the next slide in the presentation after the foreach example, meaning that you need to use makeCluster and stopCluster commands. See slide 22

## Question

I tested the lab\_parallel\_mclapply.R script on the cluster but got longer and longer time lapses the more core it used. The same with the parapply script.

Do I need to change the script to make it optimized?

## Answer

But it is faster on your laptop, or have you only tried it on Kebne?

## Answer

I only tried on Kebne cluster

## Answer

---

It should be faster so probably you are not requesting the number of cores in your SLURM script correctly. Also I have tested the code so there should be a decrease in time.

## Followup

Ok, how should you request it correctly in the SLURM script?

## Question

Naive question. If my function needed libraries, should I called them inside the function or outside? What is the best practice?

## Answer

Call them outside of the function. I always put all needed libraries at the beginning of every script so they are loaded on initiation of the script.

Thanks! :)

# Best Practices for R in HPC

---

## Question

Should we run in our home directories or on scratch from batch?

(how do I stop hackmd from jumping and deleting what I write? I had to rewrite three times!)

## Answer

You should run in either your home directory or in a project storage directory if you need more space. You can use /scratch on the compute nodes, but the space there is node-specific and you need to copy it elsewhere before the job ends (copy from inside the job) as the /scratch is cleared between jobs.

As for HackMD I have the same problem intermittently. It seems to be worse if someone else is writing further up in the document, so if you see that happening try and wait a little while before trying again.

## Question

Can we do checkpointing in R?

## Answer

There are several meanings to checkpoint but in this question I guess by checkpointing you mean setting breakpoints and checking the values of variables. In this case you can do that well in RStudio.

Breakpoints can be added to the script (Click on the left hand side) and then pull up the environment tab.

## Follow up

I mean in batch scripts. Sometimes if I run a very long job there could be some issue and the job crashes (node goes down) so the work is gone. Can I checkpoint to save some of my work and continue the calculation. I can't use Rstudio for batch right?

OK thank you

# Usable parallelized R functions

---

## Question

Do all parallel R functions work with vectorized input?

## **Answer**

That depends on the function. Most of them work with vectorized input, but there are some that use single variables, or lists, ...

If you talk about the parallel functions we talked about today, (mcXapply, parXlapply, foreach) then the answer is yes.

Thanks!

## **Question**

Is there a similar utility of microbenchmark that can be used in the cluster directly? maybe in the script submitted to slurm instead of the R code?

## **Answer**

I am aware of the packages: Rslurm and slurmR for interacting with slurm. I have tried them but in my opinion using batch scripts \*.sh give you more flexibility.

## **Question**

## **Answer**

## **Question**

## **Answer**

## **Question**

## **Answer**

## **Question**

## **Answer**

**Question**

**Answer**