# Table of Contents

# Parallelization in R

Here comes a brief introduction on how to run an R code in parallel on multiple cores. Depends, of course, whether your computer's processor (CPU) has more than one core ([here](#) you can find how to figure out the number of cores on Windows).

Parallel calculation in R works in this way (hugely simplified description):

1. Initiate the `library (parallel)` (or other libraries supporting parallel backend).
2. Created several "workers" (usually one worker = one core).
3. The workers are a bit like a newly launched instance of R - they don't know anything, their global environment has no variables and no uploaded libraries, so whatever you want them to do, you need to supply them (objects, libraries etc.)
4. Use some of the functions which are able to run the loop in parallel. E.g. most of the functions like `apply`, `lapply`, `sapply` have their parallel alternatives: `parApply`, `parLapply`, `parSapply`. These functions will distribute the job evenly among workers, collect results from them and output them out.
5. When you finish using parallel backend and you don't need workers any more, stop them, otherwise, they will keep hanging in the memory.

If we stick to using `library (parallel)`, then the simplified description of the script looks like this:

1. ```r
   library (parallel)
   ```

2. ```r
   cl <- makeCluster (4)
   ```

   - this creates four workers (clusters), and assigns their description to variable `cl`, which will be important further - all functions able to run parallel need this description so as they know where the workers actually are and how many of them;

3. ```r
   clusterExport (cl, varlist = c("object1", "object2"))
   ```

   - this function exports the variables (here named `object1` and `object2`, but it could be whatever) from current R instance into the global environment of newly created workers - so as they could use them. Note that the first argument is the description of workers `cl`)

4. ```r
   result <- parSapply (cl, some.vector, FUN = function (i)
   {some.function1; some.function2})
   ```

   - this is one of the functions which is running in parallel, analogy of `sapply`. Note that the first argument is again the description of workers. It takes elements in `some.vector` and evenly ditributes them among workers; then it process the `some.function1` and `some.function2` on these workers, and collects results - here into variable `results`.

5. ```r
   stopCluster (cl)
   ```

- this stops the workers and closes the parallel backend.

Let's demonstrate the parallelization on the example of a function which generates a large number ($10^6$) of random values from the normal distribution and then calculates their mean; this is replicated many times (100 times here).

```
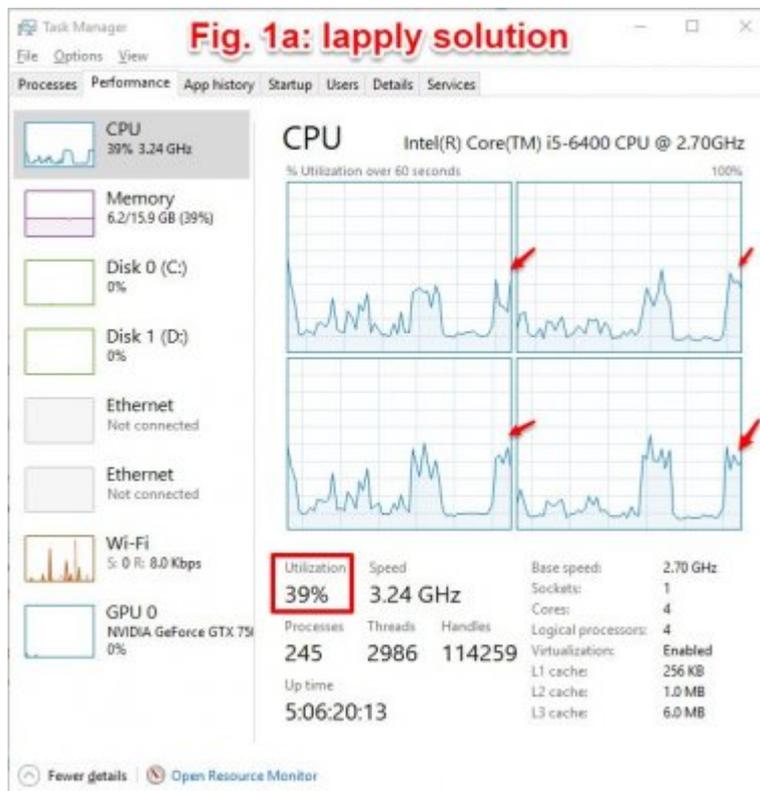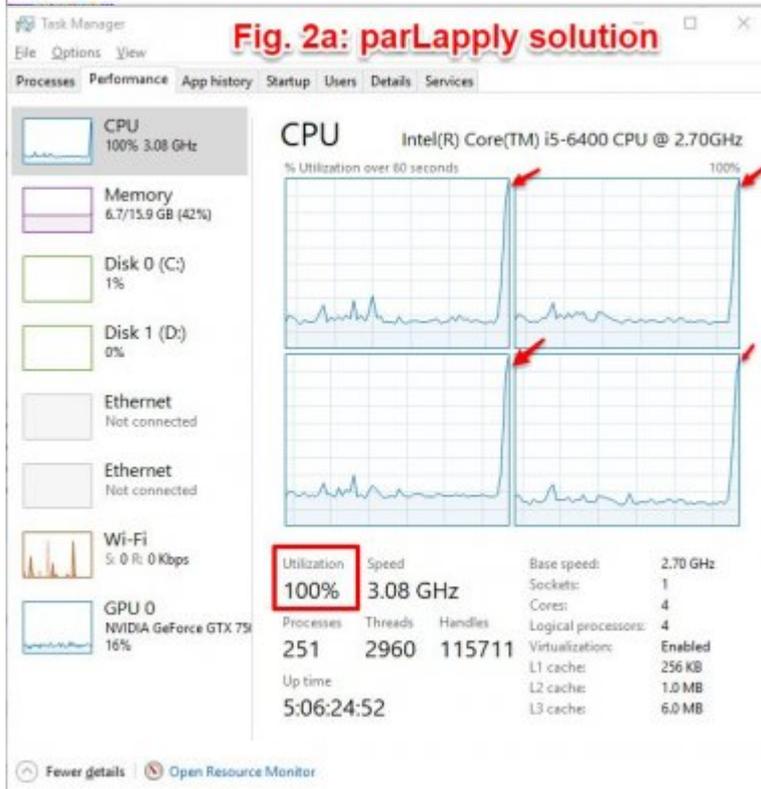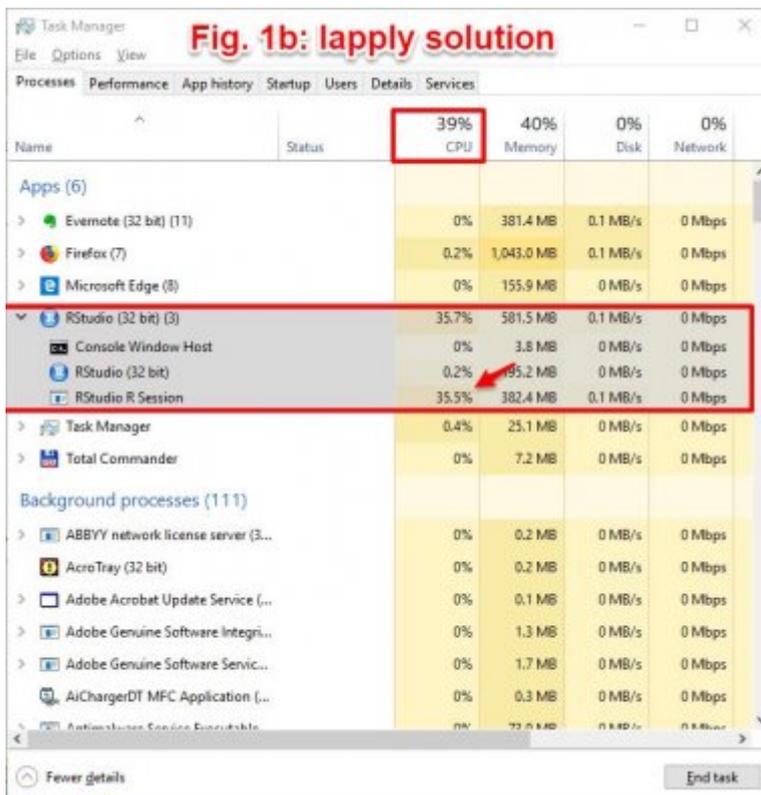lapply (1:100, FUN = function (x) mean (rnorm (1000000)))
```

Parallel version (using library `parallel`) looks like this:

```
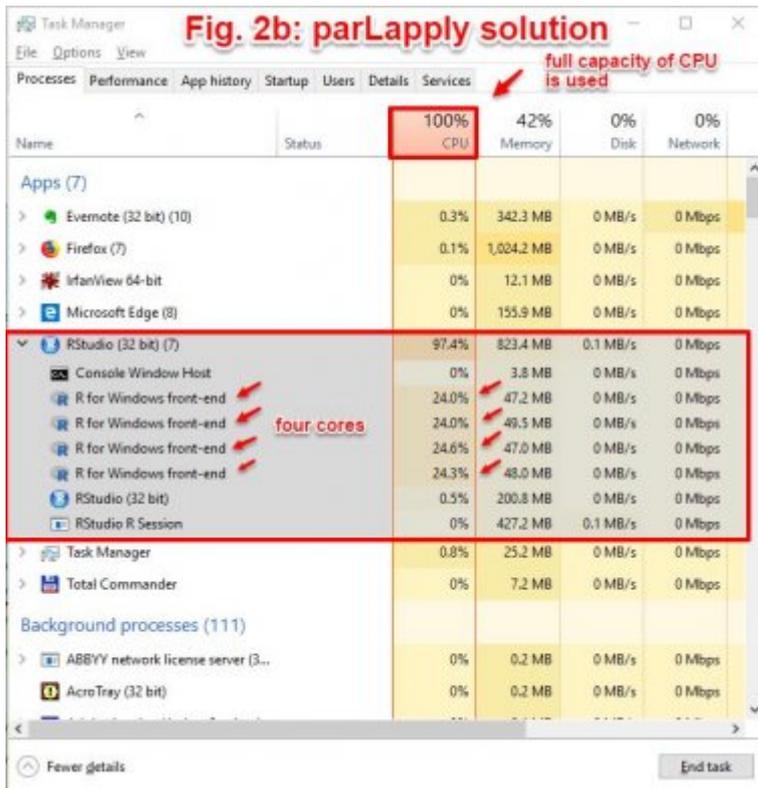library (parallel)
cl <- makeCluster (4)
res <- parLapply (cl, X = 1:100, fun = function (x) mean (rnorm (1000000)))
stopCluster (cl)
```

Both functions return the list where each item is the mean value (all are close to zero since this is mean of the normal distribution with zero mean and unit standard deviation).

What is happening at the computer if you run `lapply` vs `parLapply` solution? For this, I had a looked into the *Windows Task Manager*, the folder *Processes* and *Performance* (you get there by right-clicking on the bottom toolbar of the Windows and choosing the option Task manager, see e.g. here). The `lapply` solution uses only part of the computer capacity - in my example, it is using ~ 39% of CPU, which (given that if I do any calculation the CPU runs on 10-15%) represents approximately one-fourth of the CPU capacity (my CPU has 4 cores, so it is like if it runs only on a single core, Figs 1a and 1b below). In case of `parLapply` the situation is different - first, four instances of R program are launched (which is possible to see in the Process window of the Task manager, Fig. 2a), and CPU is running on 100% of capacity, with each R instance running on roughly 25% (Fig. 2b).

Fig. 1b: lapply solution


Fig. 2a: parLapply solution

This makes the calculation considerably faster; result of benchmarking using function `microbenchmark` from library `microbenchmark`, show that the non-parallel version is in average 2.5 times slower than parallel (column mean, 7.56 sec vs 2.95 sec; note that I modified argument `times` in `microbenchmark` function to only 10 (default is 100), which means that the function repeats each version 10-times and calculates the time statistics):

```
mb <- microbenchmark::microbenchmark (
  {
    lapply (1:100, FUN = function (x) mean (rnorm (1000000)))
  },
  {
    library (parallel)
    cl <- makeCluster (4)
    res <- parLapply (cl, X = 1:100, fun = function (x) mean (rnorm
(1000000)))
    stopCluster (cl)
  },
  times = 10)
mb
```

```
Unit: seconds
...
      min       lq     mean   median       uq      max neval cld
 7.389548 7.522466 7.566548 7.585431 7.605311 7.703006    10   b
 2.853429 2.890022 2.954747 2.943975 2.968527 3.114184    10   a
```

The reason why the parallel version is not truly four times faster than non-parallel version (although it is running on four cores in the case of my computer) is that parallel processes spend some time overhead by managing the parallelization - splitting the data, sending them to individual R workers, collecting them and merging them together. The best use of parallelization is when a single operation

takes considerable time, so the time spent by calculation is much higher than time spent by communication of R with individual cores. Indeed, if I increase the number of random values from which the mean is calculated from $10^6$ to $10^7$, the relative speed of parallel version increases to almost four (83.8 sec non-parallel vs 21.5 sec parallel)(do not try to run the code below unless you have a fairly powerful computer, since the calculation already takes some time; and note that I decreased the `times = 5`, otherwise it takes ages):

```
mb <- microbenchmark::microbenchmark (
  {
    lapply (1:100, FUN = function (x) mean (rnorm (10000000)))
  },
  {
    library (parallel)
    cl <- makeCluster (4)
    res <- parLapply (cl, X = 1:100, fun = function (x) mean (rnorm
(10000000)))
    stopCluster (cl)
  },
  times = 5)
mb
```

```
Unit: seconds
...
      min       lq     mean   median       uq       max neval cld
 83.08273 83.82933 83.95855 83.97395 84.39401 84.51273     5    b
 21.42050 21.43552 21.58001 21.49912 21.58116 21.96373     5   a
```