# Profiling & Parallelization

**Lecture 21**

Dr. Colin Rundel

# Profiling & Benchmarking

# profvis demo

```
1  n = 1e6
2  d = tibble(
3    x1 = rt(n, df = 3),
4    x2 = rt(n, df = 3),
5    x3 = rt(n, df = 3),
6    x4 = rt(n, df = 3),
7    x5 = rt(n, df = 3),
8  ) |>
9    mutate(y = -2*x1 - 1*x2 + 0*x3 + 1*x4 + 2*x5 + rnorm(n))
```

```
1  profvis::profvis({
2    lm(y~., data=d)
3  })
```

# profvis demo 2

```r
profvis::profvis({
  data = data.frame(value = runif(5e4))

  data$sum[1] = data$value[1]
  for (i in seq(2, nrow(data))) {
    data$sum[i] = data$sum[i-1] + data$value[i]
  }
})
```

```r
profvis::profvis({
  x = runif(5e4)
  sum = x[1]
  for (i in seq(2, length(x))) {
    sum[i] = sum[i-1] + x[i]
  }
})
```

# Benchmarking - bench

```
 1  d = tibble(
 2    x = runif(10000),
 3    y = runif(10000)
 4  )
 5
 6  (b = bench::mark(
 7    d[d$x > 0.5, ],
 8    d[which(d$x > 0.5), ],
 9    subset(d, x > 0.5),
10    filter(d, x > 0.5)
11  ))
```

```
# A tibble: 4 × 6
  expression                 min   median `itr/sec` mem_alloc `gc/sec`
  <bch:expr>            <bch:tm> <bch:tm>     <dbl> <bch:byt>    <dbl>
1 d[d$x > 0.5, ]          39.3µs   43.9µs    21703.  235.44KB     60.4
2 d[which(d$x > 0.5), ]   55.1µs   58.9µs    16690.  269.94KB     86.5
3 subset(d, x > 0.5)      63.5µs   67.9µs    14549.  288.01KB     73.8
4 filter(d, x > 0.5)     303.9µs  327.3µs     2960.    1.47MB     16.7
```

# Larger n

```r
1  d = tibble(
2    x = runif(1e6),
3    y = runif(1e6)
4  )
5
6  (b = bench::mark(
7    d[d$x > 0.5, ],
8    d[which(d$x > 0.5), ],
9    subset(d, x > 0.5),
10   filter(d, x > 0.5)
11 ))
```

```
# A tibble: 4 × 6
  expression                  min   median `itr/sec` mem_alloc `gc/sec`
  <bch:expr>             <bch:tm> <bch:tm>     <dbl> <bch:byt>    <dbl>
1 d[d$x > 0.5, ]           2.81ms    2.9ms      315.    13.4MB     242.
2 d[which(d$x > 0.5), ]    6.03ms   6.18ms      147.    24.8MB     264.
3 subset(d, x > 0.5)       6.75ms   6.91ms      129.    24.8MB     249.
4 filter(d, x > 0.5)       4.13ms   4.38ms      186.    24.8MB     299.
```

# bench - relative results

```
1  summary(b, relative=TRUE)
```

```
# A tibble: 4 × 6
  expression                  min median `itr/sec` mem_alloc `gc/sec`
  <bch:expr>                <dbl>  <dbl>     <dbl>     <dbl>    <dbl>
1 d[d$x > 0.5, ]                1      1      2.45         1        1
2 d[which(d$x > 0.5), ]      2.15   2.13      1.14      1.86     1.09
3 subset(d, x > 0.5)         2.41   2.38      1         1.86     1.03
4 filter(d, x > 0.5)         1.47   1.51      1.45      1.86     1.24
```

# t.test

> Imagine we have run 1000 experiments (rows), each of which collects data on 50 individuals (columns). The first 25 individuals in each experiment are assigned to group 1 and the rest to group 2.

The goal is to calculate the t-statistic for each experiment comparing group 1 to group 2.

```r
1   m = 1000
2   n = 50
3   X = matrix(
4     rnorm(m * n, mean = 10, sd = 3),
5     ncol = m
6   ) |>
7     as.data.frame() |>
8     set_names(paste0("exp", seq_len(m))) |>
9     mutate(
10      ind = seq_len(n),
11      group = rep(1:2, each = n/2)
12    ) |>
13    as_tibble() |>
14    relocate(ind, group)
```

```r
1   X
```

```
# A tibble: 50 × 1,002
      ind group  exp1  exp2  exp3  exp4  exp5  exp6
    <int> <int> <dbl> <dbl> <dbl> <dbl> <dbl> <dbl>
 1      1     1  8.66 10.5   5.83  5.35 10.9  12.8
 2      2     1 17.0  13.8  14.2   5.30  7.47 10.5
 3      3     1  6.54  9.09 14.3  10.7  11.4   8.26
 4      4     1  7.04 11.7  14.2   5.76 10.1   1.50
 5      5     1  8.41  5.29  8.96 13.1   8.49  7.66
 6      6     1 11.0   6.49  6.44  9.23  9.80  9.62
 7      7     1  8.15 10.2  11.6  12.0   6.85  8.03
 8      8     1 11.4  12.3   9.82  8.48 11.2   8.08
 9      9     1 12.8   9.49 12.6  17.3   8.99  9.11
10     10     1  9.60  6.23 11.7  11.8   8.27 11.0
# ℹ 40 more rows
# ℹ 994 more variables: exp7 <dbl>, exp8 <dbl>,
#   exp9 <dbl>, exp10 <dbl>, exp11 <dbl>,
#   exp12 <dbl>, exp13 <dbl>, exp14 <dbl>,
#   exp15 <dbl>, exp16 <dbl>, exp17 <dbl>,
#   exp18 <dbl>, exp19 <dbl>, exp20 <dbl>,
```

# Implementations

```r
1  ttest_formula = function(X, m) {
2    for(i in 1:m) t.test(X[[2+i]] ~ X$group)$stat
3  }
4  system.time(ttest_formula(X,m))
```

```
 user   system elapsed
0.197   0.000   0.198
```

```r
1  ttest_for = function(X, m) {
2    for(i in 1:m) t.test(X[[2+i]][X$group == 1], X[[2+i]][X$group == 2])$stat
3  }
4  system.time(ttest_for(X,m))
```

```
 user   system elapsed
0.063   0.000   0.063
```

```r
1  ttest_apply = function(X) {
2    f = function(x, g) {
3      t.test(x[g==1], x[g==2])$stat
4    }
5    apply(X[,-(1:2)], 2, f, X$group)
6  }
7  system.time(ttest_apply(X))
```

```
 user   system elapsed
0.054   0.000   0.055
```

# Implementations (cont.)

```r
1   ttest_hand_calc = function(X) {
2     f = function(x, grp) {
3       t_stat = function(x) {
4         m = mean(x)
5         n = length(x)
6         var = sum((x - m) ^ 2) / (n - 1)
7
8         list(m = m, n = n, var = var)
9       }
10
11      g1 = t_stat(x[grp == 1])
12      g2 = t_stat(x[grp == 2])
13
14      se_total = sqrt(g1$var / g1$n + g2$var / g2$n)
15      (g1$m - g2$m) / se_total
16    }
17
18    apply(X[,-(1:2)], 2, f, X$group)
19  }
20  system.time(ttest_hand_calc(X))
```

```
  user   system elapsed
 0.016    0.000    0.016
```

# Comparison

```
1  bench::mark(
2    ttest_formula(X, m),
3    ttest_for(X, m),
4    ttest_apply(X),
5    ttest_hand_calc(X),
6    check=FALSE
7  )
```

```
Warning: Some expressions had a GC in every iteration; so filtering
is disabled.
```

```
# A tibble: 4 × 6
  expression                  min    median `itr/sec` mem_alloc `gc/sec`
  <bch:expr>             <bch:tm> <bch:tm>      <dbl> <bch:byt>    <dbl>
1 ttest_formula(X, m) 218.73ms    224.1ms       4.47    8.24MB     23.8
2 ttest_for(X, m)       69.59ms    73.2ms       13.7    1.91MB     25.4
3 ttest_apply(X)        63.28ms    68.4ms       14.7    3.48MB     25.8
4 ttest_hand_calc(X)     9.36ms    10.1ms       72.4    3.44MB     21.5
```

# Parallelization

# parallel

Part of the base packages in R

- tools for the forking of R processes (some functions do not work on Windows)
- Core functions:
  - `detectCores`
  - `pvec`
  - `mclapply`
  - `mcparallel` & `mccollect`

# detectCores

Surprisingly, detects the number of cores of the current system.

```
1  detectCores()
```

```
[1] 32
```

# pvec

Parallelization of a vectorized function call

```
1 system.time(pvec(1:1e7, sqrt, mc.cores = 1))
```

```
  user  system elapsed
 0.028   0.038   0.068
```

```
1 system.time(pvec(1:1e7, sqrt, mc.cores = 4))
```

```
  user  system elapsed
 0.211   0.253   0.324
```

```
1 system.time(pvec(1:1e7, sqrt, mc.cores = 8))
```

```
  user  system elapsed
 0.101   0.288   0.176
```

```
1 system.time(sqrt(1:1e7))
```

```
  user  system elapsed
 0.024   0.042   0.066
```

# pvec - bench::system_time

```
1  bench::system_time(pvec(1:1e7, sqrt, mc.cores = 1))
```

```
process     real
 125ms     126ms
```

```
1  bench::system_time(pvec(1:1e7, sqrt, mc.cores = 4))
```

```
process     real
 220ms     250ms
```

```
1  bench::system_time(pvec(1:1e7, sqrt, mc.cores = 8))
```

```
process     real
 272ms     286ms
```

```
1 bench::system_time(Sys.sleep(.5))
```

```
process      real
 54.4µs 500.8ms
```

```
1 system.time(Sys.sleep(.5))
```

```
  user    system elapsed
 0.000    0.000   0.501
```

# Cores by size

```
1  cores = c(1,4,6,8,10)
2  order = 6:8
3  f = function(x,y) {
4    system.time(
5      pvec(1:(10^y), sqrt, mc.cores = x)
6    )[3]
7  }
8
9  res = map(
10   cores,
11   function(x) {
12     map_dbl(order, f, x = x)
13   }
14 ) |>
15   do.call(rbind, args = _)
16
17 rownames(res) = paste0(cores," cores")
18 colnames(res) = paste0("10^",order)
```

```
1  res
```

```
         10^6  10^7  10^8
1 cores  0.004 0.126 0.745
4 cores  0.033 0.171 2.190
6 cores  0.038 0.173 1.847
8 cores  0.045 0.183 1.827
10 cores 0.050 0.187 1.852
```

# mclapply

implements a parallelized version of `lapply`

```
1  system.time(rnorm(1e7))
```

```
 user   system elapsed
0.197   0.020   0.219
```

```
1  system.time(unlist(mclapply(1:10, function(x) rnorm(1e6), mc.cores = 2)))
```

```
 user   system elapsed
0.247   0.200   0.260
```

```
1  system.time(unlist(mclapply(1:10, function(x) rnorm(1e6), mc.cores = 4)))
```

```
 user   system elapsed
0.231   0.216   0.172
```

```
1  system.time(unlist(mclapply(1:10, function(x) rnorm(1e6), mc.cores = 8)))
```

```
 user   system elapsed
0.238   0.271   0.149
```

```
1  system.time(unlist(mclapply(1:10, function(x) rnorm(1e6), mc.cores = 10)))
```

```
 user   system elapsed
0.246   0.332   0.144
```

# mcparallel

Asynchronously evaluation of an R expression in a separate process

```r
1  m = mcparallel(rnorm(1e6))
2  n = mcparallel(rbeta(1e6,1,1))
3  o = mcparallel(rgamma(1e6,1,1))
```

```r
1  str(m)
```

```
List of 2
 $ pid: int 787800
 $ fd : int [1:2] 6 9
 - attr(*, "class")= chr [1:3] "parallelJob" "childProcess" "process"
```

```r
1  str(n)
```

```
List of 2
 $ pid: int 787801
 $ fd : int [1:2] 7 11
 - attr(*, "class")= chr [1:3] "parallelJob" "childProcess" "process"
```

# mccollect

Checks `mcparallel` objects for completion

```
1  str(mccollect(list(m,n,o)))
```

```
List of 3
 $ 787800: num [1:1000000] 1.113 -1.375 0.254 -1.055 0.641 ...
 $ 787801: num [1:1000000] 0.5526 0.0744 0.9768 0.9385 0.1238 ...
 $ 787802: num [1:1000000] 2.00498 4.60403 0.00115 0.58452 0.48502 ...
```

# mccollect - waiting

```r
p = mcparallel(mean(rnorm(1e5)))
```

```r
mccollect(p, wait = FALSE, 10)
```

```
$`787803`
[1] -0.0008566918
```

```r
mccollect(p, wait = FALSE)
```

```
Warning in selectChildren(jobs, timeout): cannot wait for child
787803 as it does not exist

NULL
```

```r
mccollect(p, wait = FALSE)
```

```
Warning in selectChildren(jobs, timeout): cannot wait for child
787803 as it does not exist

NULL
```

# doMC & foreach

# doMC & foreach

Packages by Revolution Analytics that provides the `foreach` function which is a parallelizable `for` loop (and then some).

- Core functions:

    - `registerDoMC`

    - `foreach`, `%dopar%`, `%do%`

# registerDoMC

Primarily used to set the number of cores used by `foreach`, by default uses `options("cores")` or half the number of cores found by `detectCores` from the parallel package.

```
1  options("cores")
```

```
$cores
NULL
```

```
1  detectCores()
```

```
[1] 32
```

```
1  getDoParWorkers()
```

```
[1] 1
```

```
1  registerDoMC(4)
2  getDoParWorkers()
```

```
[1] 4
```

# foreach

A slightly more powerful version of base `for` loops (think `for` with an `lapply` flavor). Combined with `%do%` or `%dopar%` for single or multicore execution.

```r
for(i in 1:10) {
  sqrt(i)
}
```

```r
foreach(i = 1:5) %do% {
  sqrt(i)
}
```

```
[[1]]
[1] 1

[[2]]
[1] 1.414214

[[3]]
[1] 1.732051

[[4]]
[1] 2

[[5]]
[1] 2.236068
```

# **foreach** - iterators

foreach can iterate across more than one value, but it doesn't do length coercion

```
1  foreach(i = 1:5, j = 1:5) %do% {
2    sqrt(i^2+j^2)
3  }
```

```
1  foreach(i = 1:5, j = 1:2) %do% {
2    sqrt(i^2+j^2)
3  }
```

[[1]]
[1] 1.414214

[[2]]
[1] 2.828427

[[3]]
[1] 4.242641

[[4]]
[1] 5.656854

[[5]]
[1] 7.071068

[[1]]
[1] 1.414214

[[2]]
[1] 2.828427

# **foreach** - combining results

```r
foreach(i = 1:5, .combine='c') %do% {
  sqrt(i)
}
```

```
[1] 1.000000 1.414214 1.732051 2.000000 2.236068
```

```r
foreach(i = 1:5, .combine='cbind') %do% {
  sqrt(i)
}
```

```
     result.1 result.2 result.3 result.4 result.5
[1,]        1 1.414214 1.732051        2 2.236068
```

```r
foreach(i = 1:5, .combine='+') %do% {
  sqrt(i)
}
```

```
[1] 8.382332
```

# foreach - parallelization

Swapping out `%do%` for `%dopar%` will use the parallel backend.

```
1  registerDoMC(4)
2  system.time(foreach(i = 1:10) %dopar% mean(rnorm(1e6)))
```

```
  user  system elapsed
 0.164   0.051   0.090
```

```
1  registerDoMC(8)
2  system.time(foreach(i = 1:10) %dopar% mean(rnorm(1e6)))
```

```
  user  system elapsed
 0.176   0.098   0.063
```

```
1  registerDoMC(10)
2  system.time(foreach(i = 1:10) %dopar% mean(rnorm(1e6)))
```

```
  user  system elapsed
 0.202   0.139   0.070
```

# furrr / future

```
1  system.time( purrr::map(c(1,1,1), Sys.sleep) )
```

```
  user   system elapsed
 0.000    0.000   3.003
```

```
1  system.time( furrr::future_map(c(1,1,1), Sys.sleep) )
```

```
  user   system elapsed
 0.040    0.003   3.064
```

```
1  future::plan(future::multisession) # See also future::multicore
2  system.time( furrr::future_map(c(1,1,1), Sys.sleep) )
```

```
  user   system elapsed
 0.392    0.009   1.700
```

# Example - Bootstraping

Bootstrapping is a resampling scheme where the original data is repeatedly reconstructed by taking a samples of size *n* (with replacement) from the original data, and using that to repeat an analysis procedure of interest. Below is an example of fitting a local regression (`loess`) to some synthetic data, we will construct a bootstrap prediction interval for this model.
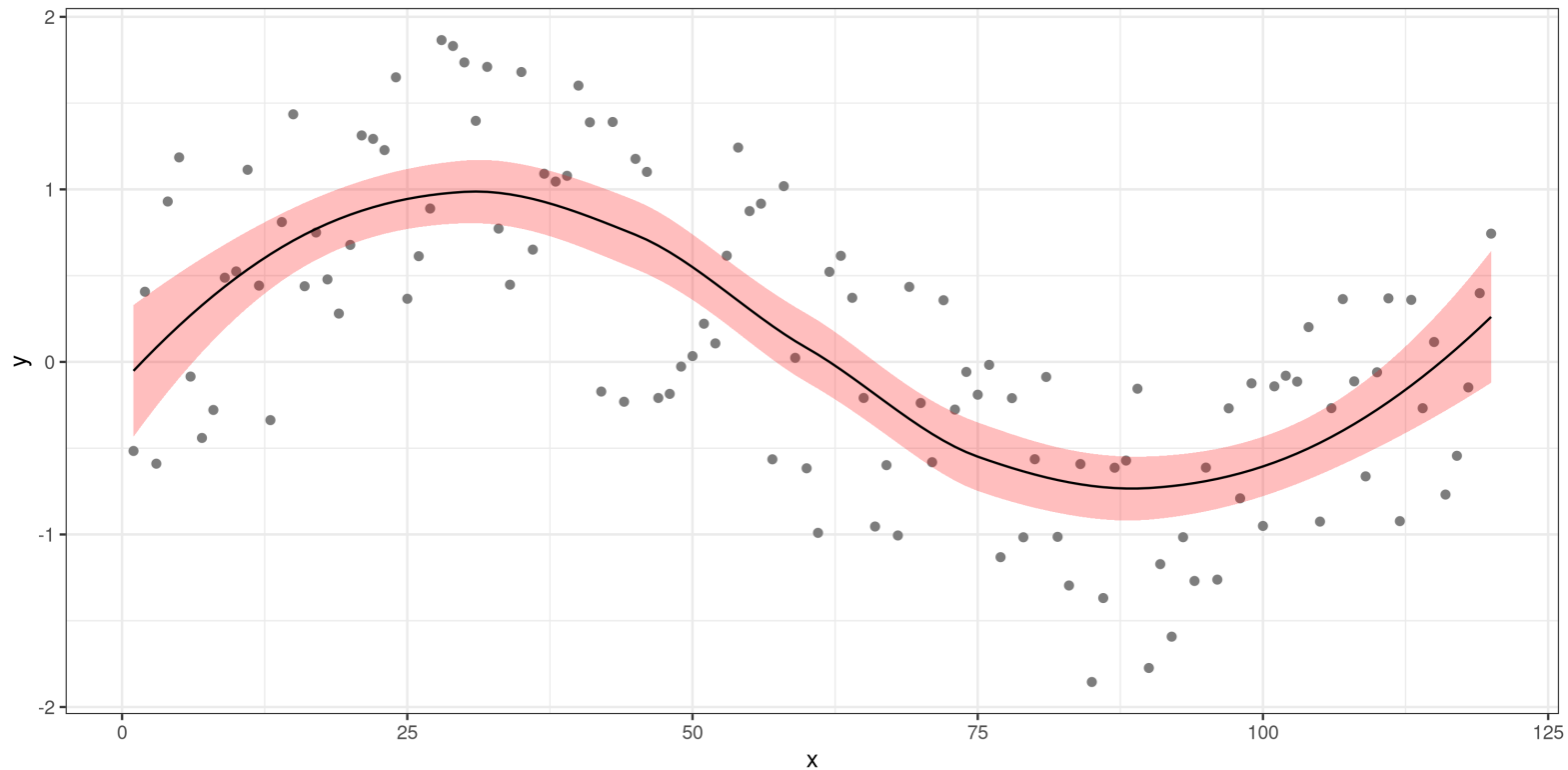
```r
1  set.seed(3212016)
2  d = data.frame(x = 1:120) |>
3      mutate(y = sin(2*pi*x/120) + runif(length(x),-1,1))
4
5  l = loess(y ~ x, data=d)
6  p = predict(l, se=TRUE)
7
8  d = d |> mutate(
9    pred_y = p$fit,
10   pred_y_se = p$se.fit
11 )
```

```r
ggplot(d, aes(x,y)) +
  geom_point(color="gray50") +
  geom_ribbon(
    aes(ymin = pred_y - 1.96 * pred_y_se,
        ymax = pred_y + 1.96 * pred_y_se),
    fill="red", alpha=0.25
  ) +
  geom_line(aes(y=pred_y)) +
  theme_bw()
```

# Bootstraping Demo

# What to use when?

Optimal use of parallelization / multiple cores is hard, there isn't one best solution

- Don't underestimate the overhead cost

- Experimentation is key

- Measure it or it didn't happen

- Be aware of the trade off between developer time and run time

# BLAS and LAPACK

# Statistics and Linear Algebra

An awful lot of statistics is at its core linear algebra.

For example:

- Linear regession models, find

$$\hat{\beta} = (\mathrm{X}^\mathrm{T}\mathrm{X})^{-1}\,\mathrm{X}^\mathrm{T}\mathrm{y}$$

- Principle component analysis

  - Find $\mathrm{T} = \mathrm{X}\mathrm{W}$ where $\mathrm{W}$ is a matrix whose columns are the eigenvectors of $\mathrm{X}^\mathrm{T}\mathrm{X}$.

  - Often solved via SVD - Let $\mathrm{X} = \mathrm{U}\Sigma\mathrm{W}^\mathrm{T}$ then $\mathrm{T} = \mathrm{U}\Sigma$.

# Numerical Linear Algebra

Not unique to Statistics, these are the type of problems that come up across all areas of numerical computing.

- Numerical linear algebra ≠ mathematical linear algebra

- Efficiency and stability of numerical algorithms matter

  - Designing and implementing these algorithms is hard

- Don't reinvent the wheel - common core linear algebra tools (well defined API)

# BLAS and LAPACK

Low level algorithms for common linear algebra operations

## BLAS

- **B**asic **L**inear **A**lgebra **S**ubprograms

- Copying, scaling, multiplying vectors and matrices

- Origins go back to 1979, written in Fortran

## LAPACK

- **L**inear **A**lgebra **Pack**age

- Higher level functionality building on BLAS.

- Linear solvers, eigenvalues, and matrix decompositions

- Origins go back to 1992, mostly Fortran (expanded on LINPACK, EISPACK)

# Modern variants?

Most default BLAS and LAPACK implementations (like R's defaults) are somewhat dated

- Written in Fortran and designed for a single cpu core

- Certain (potentially non-optimal) hard coded defaults (e.g. block size).
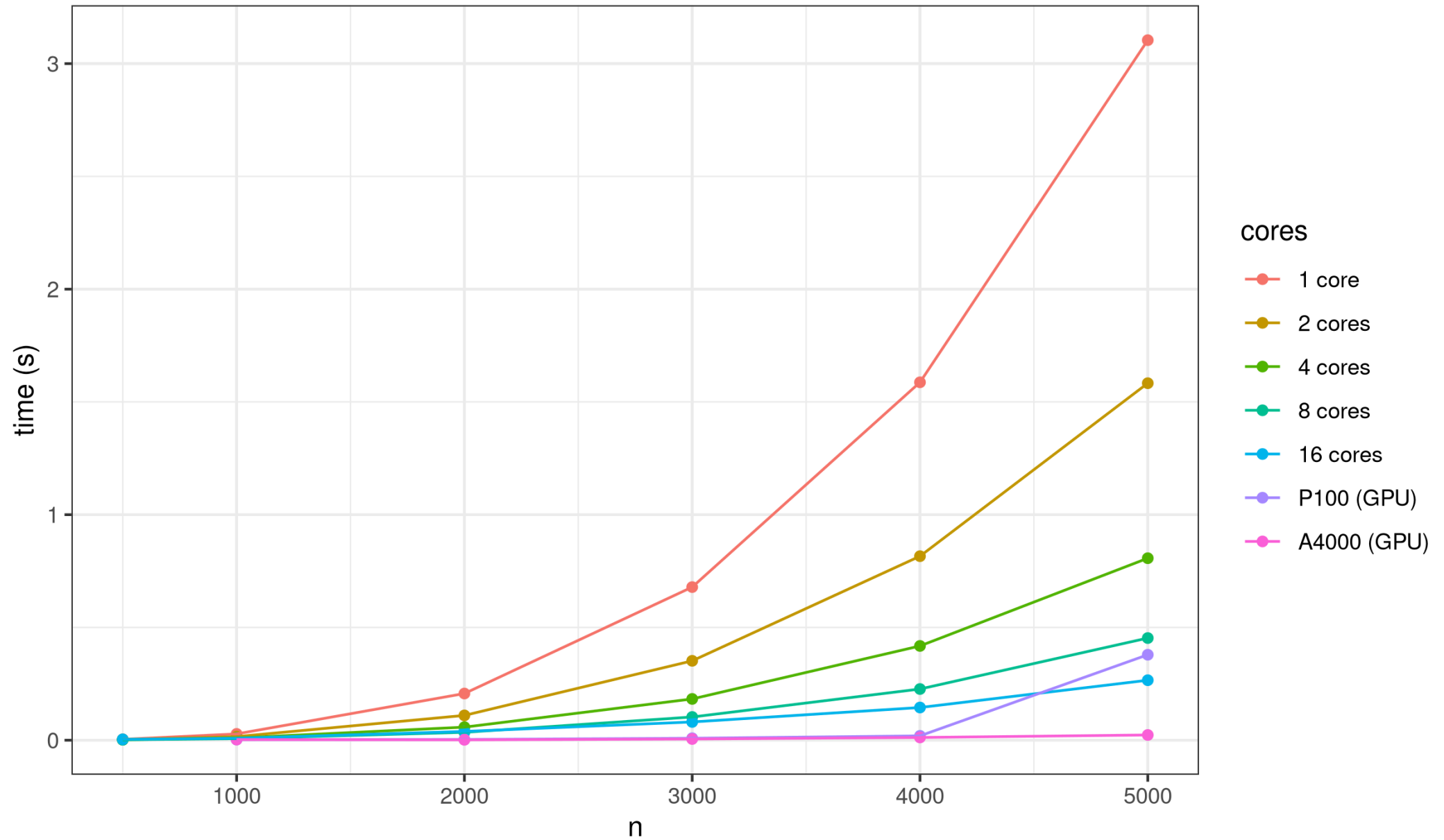
Multithreaded alternatives:

- ATLAS - Automatically Tuned Linear Algebra Software

- OpenBLAS - fork of GotoBLAS from TACC at UTexas

- Intel MKL - Math Kernel Library, part of Intel's commercial compiler tools

- cuBLAS / Magma - GPU libraries from Nvidia and UTK respectively

- Accelerate / vecLib - Apple's framework for GPU and multicore computing

# OpenBLAS Matrix Multiply Performance

```r
1  x=matrix(runif(5000^2),ncol=5000)
2
3  sizes = c(100,500,1000,2000,3000,4000,5000)
4  cores = c(1,2,4,8,16)
5
6  sapply(
7    cores,
8    function(n_cores) {
9      flexiblas::flexiblas_set_num_threads(n_cores)
10     sapply(
11       sizes,
12       function(s) {
13         y = x[1:s,1:s]
14         system.time(y %*% y)[3]
15       }
16     )
17   }
18 )
```

| n | 1 core | 2 cores | 4 cores | 8 cores | 16 cores |
|---|---|---|---|---|---|
| 100 | 0.000 | 0.000 | 0.000 | 0.000 | 0.000 |
| 500 | 0.004 | 0.003 | 0.002 | 0.002 | 0.004 |
| 1000 | 0.028 | 0.016 | 0.010 | 0.007 | 0.009 |
| 2000 | 0.207 | 0.110 | 0.058 | 0.035 | 0.039 |
| 3000 | 0.679 | 0.352 | 0.183 | 0.103 | 0.081 |
| 4000 | 1.587 | 0.816 | 0.418 | 0.227 | 0.145 |
| 5000 | 3.104 | 1.583 | 0.807 | 0.453 | 0.266 |

# Matrix Multiply of (n x n) matrices

Matrix Multiply of (n x n) matrices