

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(lm(y~., data=d))
```

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,]	128µs	137µs	7183.	238.81KB	19.3
2 d[which(d\$x > 0.5),]	134µs	148µs	6723.	269.64KB	36.1
3 subset(d, x > 0.5)	166µs	181µs	5441.	287.82KB	26.2
4 filter(d, x > 0.5)	376µs	397µs	2490.	1.47MB	20.6

Larger n

```
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,]	11.7ms	12ms	83.3	13.3MB	58.3
2 d[which(d\$x > 0.5),]	13.2ms	13.3ms	75.2	24.8MB	90.2
3 subset(d, x > 0.5)	17.3ms	17.6ms	56.7	24.8MB	70.8
4 filter(d, x > 0.5)	13.3ms	13.7ms	73.4	24.8MB	122.

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	1.47	1	1
2	d[which(d\$x > 0.5),]	1.13	1.10	1.33	1.86	1.55
3	subset(d, x > 0.5)	1.48	1.46	1	1.86	1.21
4	filter(d, x > 0.5)	1.13	1.14	1.30	1.86	2.10

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.

```
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)
15 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 16.7 12.0  9.77  7.76 10.2 11.5
2     2     1  2.56  8.37  8.65 12.3   6.25  9.51
3     3     1 13.6 12.3 14.3 14.9 13.6  6.11
4     4     1 12.5  8.40 12.5  9.97  9.38  7.14
5     5     1  2.67  9.36  9.86  8.10  8.59 11.1
6     6     1 12.7  8.24 17.6  9.74 10.3  4.64
7     7     1 12.3 11.2 10.8  8.12  7.76 10.3
8     8     1 12.1  8.70 10.8  7.42  6.88  8.58
9     9     1  7.64  4.53 10.9 10.4  8.65 12.1
10    10     1  8.12 11.8  9.79  9.12 12.0  8.57
# i 40 more rows
# i 994 more variables: exp7 <dbl>, exp8 <dbl>,
#   exp9 <dbl>, exp10 <dbl>, exp11 <dbl>,
#   exp12 <dbl>, exp13 <dbl>, exp14 <dbl>.
```

Implementations

```
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  
0.226  0.001  0.228
```

```
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  
0.062  0.000  0.063
```

```
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  0.001  0.055
```


Implementations (cont.)

```
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 }
```

user	system	elapsed
0.014	0.000	0.015

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 )
```

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)	201.74ms	203.25ms	4.90	8.24MB	29.4
2	ttest_for(X, m)	67.09ms	67.59ms	14.6	1.91MB	31.1
3	ttest_apply(X)	60.55ms	60.92ms	16.4	3.49MB	32.8
4	ttest_hand_calc(X)	8.76ms	9.56ms	87.0	3.45MB	29.7

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`
 - `mcpipeline` & `mccollect`

detectCores

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

```
1 detectCores()
```

```
[1] 10
```

pvec

Parallelization of a vectorized function call

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

	user	system	elapsed
	0.095	0.011	0.106

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

	user	system	elapsed
	0.167	0.126	0.242

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

	user	system	elapsed
	0.089	0.170	0.155

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

	user	system	elapsed
	0.018	0.016	0.034

pvec - bench::system_time

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

process	real
60.5ms	59.7ms

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

process	real
164ms	203ms

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

process	real
173ms	196ms

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

process	real
59 μ s	497ms

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

user	system	elapsed
0.000	0.000	0.502

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)
19
20 res
```

	10 ⁶	10 ⁷	10 ⁸
1 cores	0.004	0.024	0.320
4 cores	0.039	0.146	1.879
6 cores	0.032	0.127	1.323
8 cores	0.036	0.138	1.519
10 cores	0.036	0.164	1.434

mclapply

Parallelized version of lapply

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

user	system	elapsed
0.265	0.003	0.270

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

user	system	elapsed
0.312	0.088	0.251

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

user	system	elapsed
0.327	0.088	0.164

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

user	system	elapsed
0.331	0.138	0.162

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

user	system	elapsed
0.369	0.159	0.183

mcparallel

Asynchronously evaluation of an R expression in a separate process

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

```
1 str(m)
```

List of 2

```
$ pid: int 53890  
$ fd : int [1:2] 4 7  
- attr(*, "class")= chr [1:3] "parallelJob" "childProcess" "process"
```

```
1 str(n)
```

List of 2

```
$ pid: int 53891  
$ fd : int [1:2] 5 9  
- attr(*, "class")= chr [1:3] "parallelJob" "childProcess" "process"
```

mccollect

Checks `mcpaallel` objects for completion

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

List of 3

```
$ 53890: num [1:1000000] -0.144 0.786 2.704 0.157 -0.558 ...  
$ 53891: num [1:1000000] 0.538 0.691 0.593 0.932 0.171 ...  
$ 53892: num [1:1000000] 1.319 0.17 1.548 0.137 0.151 ...
```

mccollect - waiting

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

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

```
$`53893`
```

```
[1] 0.0009924475
```

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

```
NULL
```

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

```
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] 10
```

```
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.

```
1 for(i in 1:10) {  
2   sqrt(i)  
3 }
```

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

```
[[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]]  
[1] 1.414214
```

```
[[2]]  
[1] 2.828427
```

```
[[3]]  
[1] 4.242641
```

```
[[4]]  
[1] 5.656854
```

```
[[5]]  
[1] 7.071068
```

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

```
[[1]]  
[1] 1.414214
```

```
[[2]]  
[1] 2.828427
```

foreach - combining results

```
1 foreach(i = 1:5, .combine='c') %do% {  
2   sqrt(i)  
3 }
```

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

```
1 foreach(i = 1:5, .combine='cbind') %do% {  
2   sqrt(i)  
3 }
```

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

```
1 foreach(i = 1:5, .combine='+') %do% {  
2   sqrt(i)  
3 }
```

```
[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
0.298	0.022	0.108

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

user	system	elapsed
0.303	0.038	0.078

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

user	system	elapsed
0.324	0.050	0.066



furrr / future

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

user	system	elapsed
0.000	0.000	3.012

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

user	system	elapsed
0.050	0.006	3.084

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

user	system	elapsed
0.188	0.004	1.445

Example - Bootstrapping

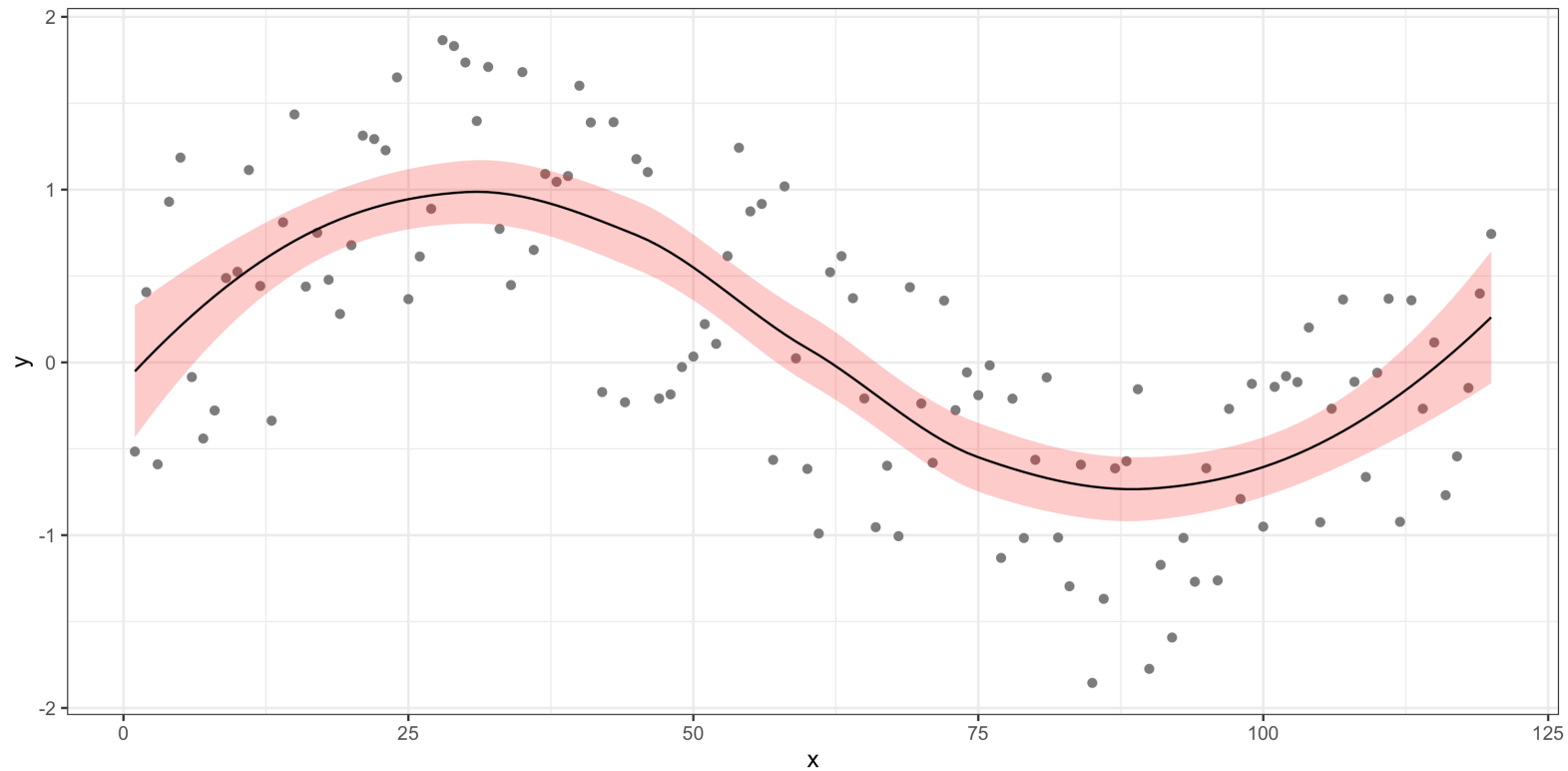
Bootstrapping is a resampling scheme where the original data is repeatedly reconstructed by taking a samples of size n (with replacement) from the original data, and using that to repeat an analysis procedure of interest. Below is an example of fitting a local regression (`loess`) to some synthetic data, we will construct a bootstrap prediction interval for this model.

```
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 )
```

```

1 ggplot(d, aes(x,y)) +
2   geom_point(color="gray50") +
3   geom_ribbon(
4     aes(ymin = pred_y - 1.96 * pred_y_se,
5         ymax = pred_y + 1.96 * pred_y_se),
6     fill="red", alpha=0.25
7   ) +
8   geom_line(aes(y=pred_y)) +
9   theme_bw()

```



Bootstrapping 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 regression models, find

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

- Principle component analysis
 - Find $T = XW$ where W is a matrix whose columns are the eigenvectors of $X^T X$.
 - Often solved via SVD - Let $X = U\Sigma W^T$ then $T = U\Sigma$.

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 \neq 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 **P**ackage
- 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:

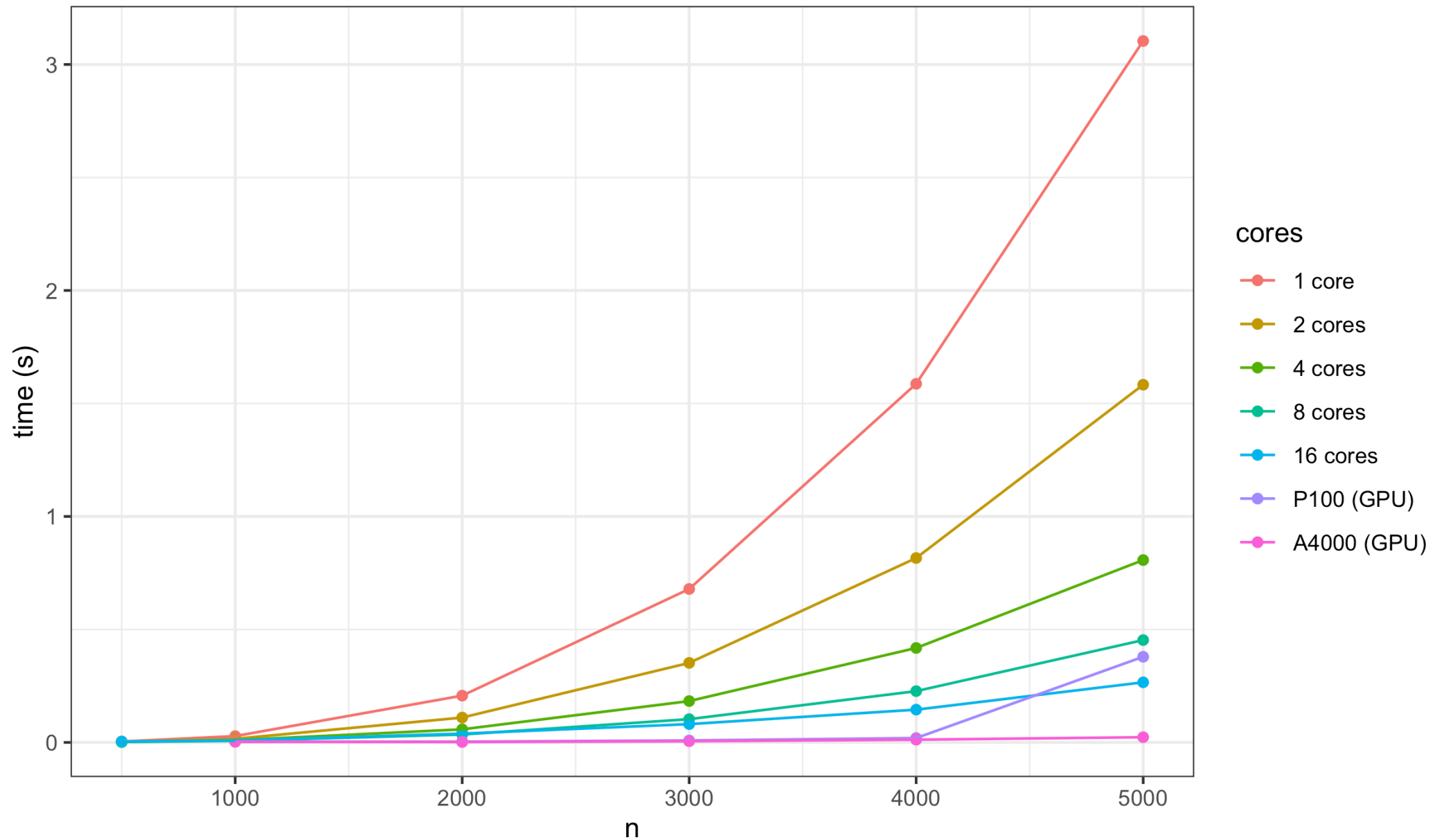
- ATLAS - Automatically Tuned Linear Algebra Software
- OpenBLAS - fork of GotoBLAS from TACC at UTexas
- Intel MKL - Math Kernel Library, part of Intel's commercial compiler tools
- cuBLAS / Magma - GPU libraries from Nvidia and UTK respectively
- Accelerate / vecLib - Apple's framework for GPU and multicore computing

OpenBLAS Matrix Multiply Performance

```
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   {
10     flexiblas::flexiblas_set_num_threads(n_cores)
11     sapply(
12       sizes,
13       function(s)
14       {
15         y = x[1:s,1:s]
16         system.time(y %*% y)[3]
17       }
18     )
19   }
20 )
```


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

