

# 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	140µs	7002.	239.84KB	19.4
2 d[which(d\$x > 0.5), ]	137µs	152µs	6388.	271.49KB	33.5
3 subset(d, x > 0.5)	167µs	195µs	4940.	288.85KB	26.3
4 filter(d, x > 0.5)	379µs	432µs	2167.	1.48MB	35.4

# 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, ]	12ms	12.4ms	80.6	13.4MB	60.4
2 d[which(d\$x > 0.5), ]	13.3ms	13.6ms	73.6	24.8MB	210.
3 subset(d, x > 0.5)	17.7ms	17.9ms	55.7	24.8MB	65.0
4 filter(d, x > 0.5)	13.5ms	13.9ms	71.9	24.8MB	56.8

# 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	1.45	1	1.06
2	d[which(d\$x > 0.5), ]	1.11	1.09	1.32	1.86	3.70
3	subset(d, x > 0.5)	1.47	1.45	1	1.86	1.14
4	filter(d, x > 0.5)	1.12	1.12	1.29	1.86	1

# 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  12.1  15.0  12.8   5.54  12.0  11.5
2     2     1   5.61   6.99  14.3   8.47   6.49   5.17
3     3     1  12.5  11.6  15.7  10.3  14.0   7.52
4     4     1   8.33   8.39  10.5  15.6   9.94   8.21
5     5     1  11.1   5.80   8.22   9.48   5.08   2.59
6     6     1   9.48   6.19   6.05   9.16   8.69  17.0
7     7     1  11.2  11.8   5.98   6.05   9.70  12.9
8     8     1   9.56   8.14  14.6  11.4   7.39   9.63
9     9     1  10.1   8.72  10.3   7.71  11.5  12.1
10    10     1  17.1  11.6   8.38  12.3   7.88   8.59
# 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  
5 system.time(ttest_formula(X,m))
```

	user	system	elapsed
	0.183	0.001	0.186

```
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  
5 system.time(ttest_for(X,m))
```

	user	system	elapsed
	0.064	0.001	0.066

```
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  
8 system.time(ttest_apply(X))
```



user	system	elapsed
0.053	0.000	0.054

# 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)	197.05ms	199.89ms	5.01	8.24MB	26.7
2	ttest_for(X, m)	68.75ms	69.1ms	14.3	1.91MB	26.8
3	ttest_apply(X)	58.24ms	62.82ms	16.1	3.49MB	26.8
4	ttest_hand_calc(X)	8.61ms	9.21ms	89.6	3.45MB	25.9

# 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
0.016	0.012	0.028

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

user	system	elapsed
0.168	0.151	0.249

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

user	system	elapsed
0.092	0.164	0.152

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

user	system	elapsed
0.055	0.018	0.072

# pvec - bench::system\_time

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

process	real
25ms	24.8ms

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

process	real
150ms	179ms

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

process	real
191ms	223ms



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

process	real
64 $\mu$ s	497ms

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

user	system	elapsed
0.000	0.000	0.505

# 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	0.032	0.371
4 cores	0.032	0.145	2.006
6 cores	0.024	0.138	1.367
8 cores	0.033	0.127	1.265
10 cores	0.033	0.147	1.548

# mclapply

## Parallelized version of lapply

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

user	system	elapsed
0.262	0.004	0.266

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

user	system	elapsed
0.309	0.095	0.268

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

user	system	elapsed
0.322	0.089	0.161

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

user	system	elapsed
0.336	0.145	0.172

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

user	system	elapsed
0.360	0.150	0.179

# 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 14040  
$ fd : int [1:2] 4 7  
- attr(*, "class")= chr [1:3] "parallelJob" "childProcess" "process"
```

```
1 str(n)
```

List of 2

```
$ pid: int 14041  
$ 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

```
$ 14040: num [1:1000000] -0.7172 -1.8334 -0.0983 0.953 0.1629 ...  
$ 14041: num [1:1000000] 0.235 0.554 0.301 0.977 0.644 ...  
$ 14042: num [1:1000000] 0.448 0.192 0.342 1.386 0.397 ...
```

# mccollect - waiting

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

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

```
$`14043`
```

```
[1] -0.001305267
```

```
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 }  
4  
5 foreach(i = 1:5) %do% {  
6   sqrt(i)  
7 }
```

```
[[1]]
```

```
[1] 1
```

```
[[2]]
```

```
[1] 1.414214
```

```
[[3]]
```

```
[1] 1.732051
```

```
[[ 4]]  
[1] 2
```

# 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) {  
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 foreach(i = 1:5, j = 1:2) {  
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	0.028	0.109

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

user	system	elapsed
0.302	0.032	0.076

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

user	system	elapsed
0.325	0.042	0.069



# 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
0.053	0.008	3.097

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

user	system	elapsed
0.206	0.007	1.451

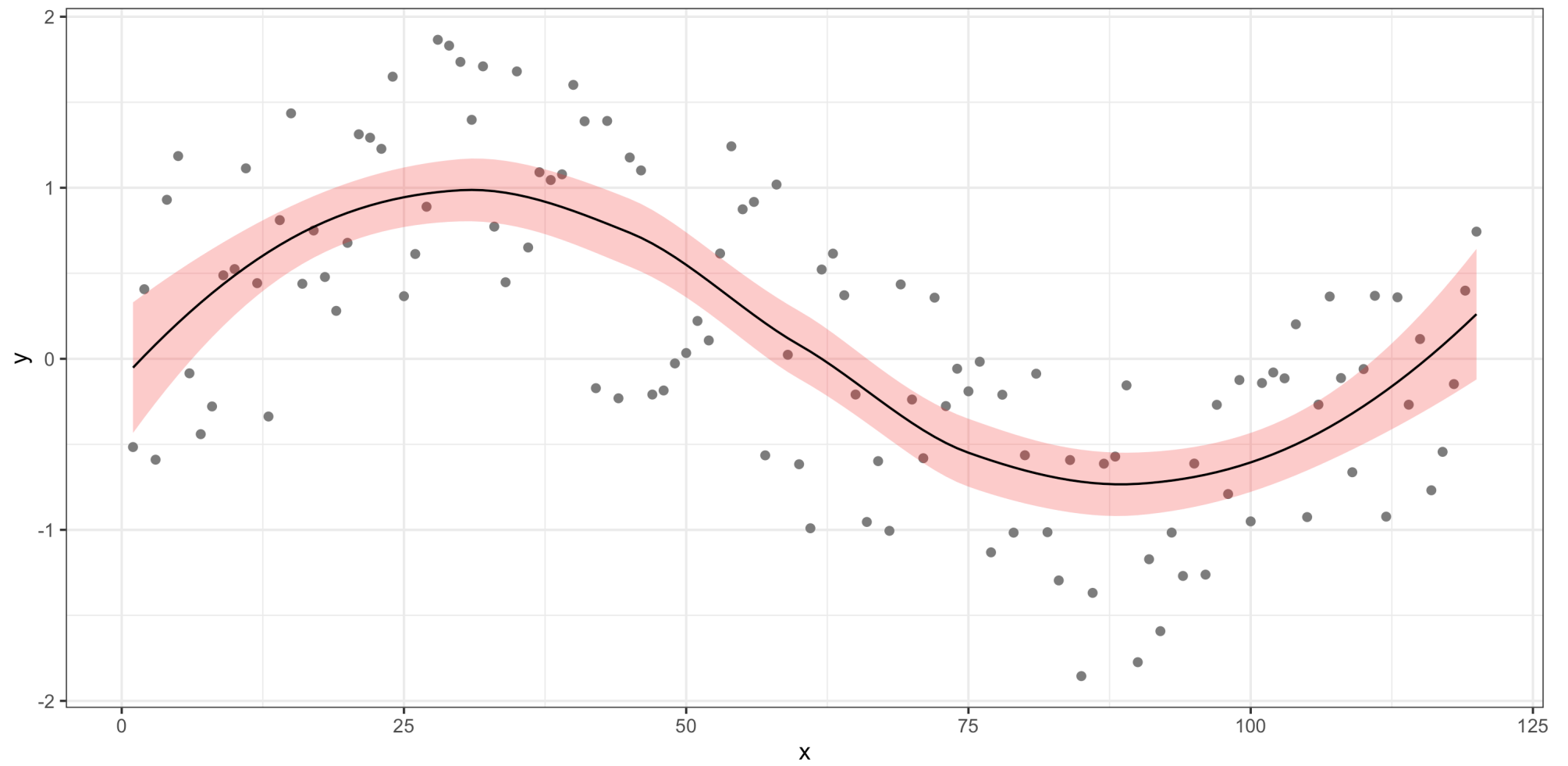


# Example - Bootstrapping

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:

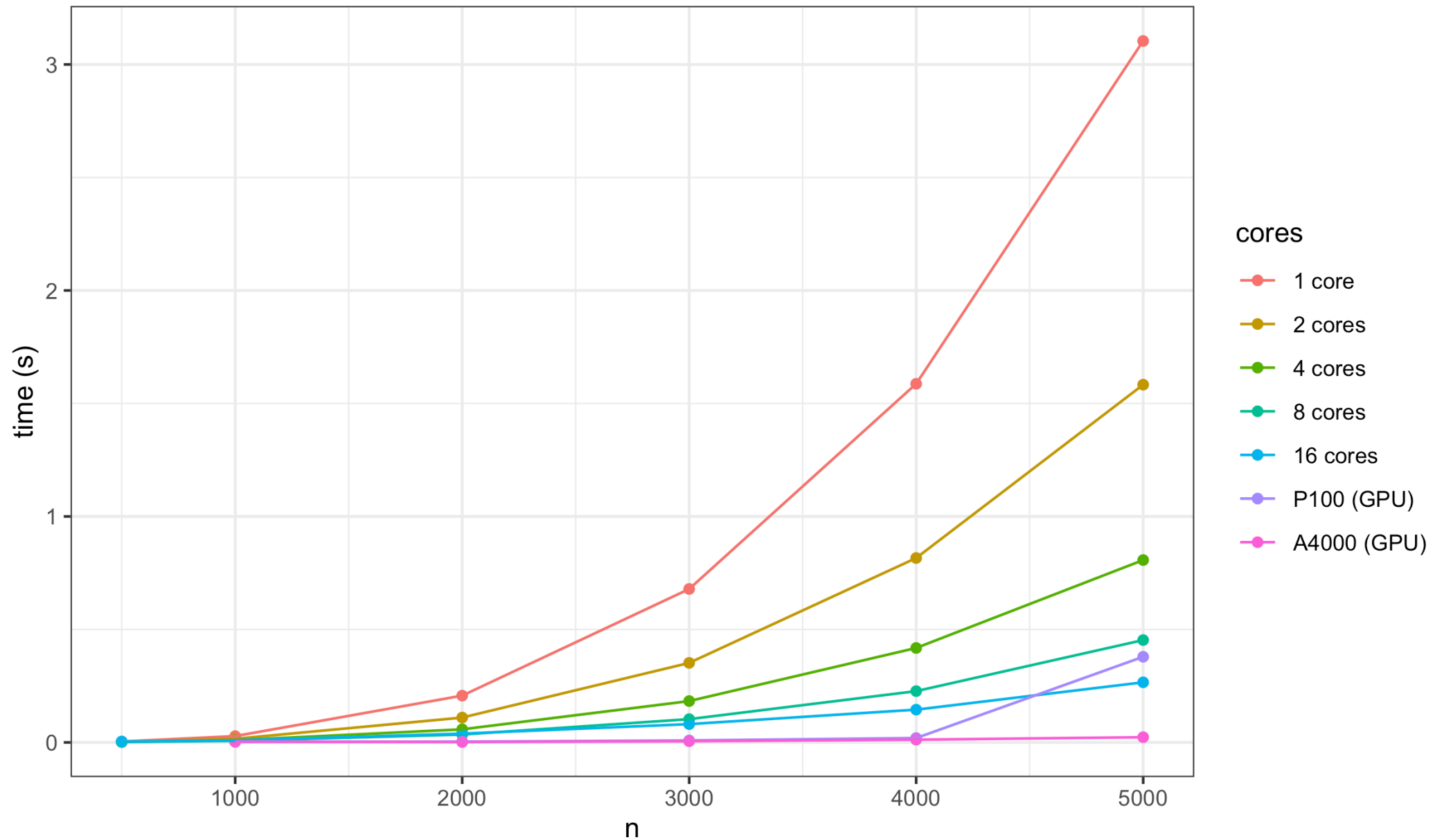
- ATLAS - Automatically Tuned Linear Algebra Software
- OpenBLAS - fork of GotoBLAS from TACC at UTexas
- Intel MKL - Math Kernel Library, part of Intel's commercial compiler tools
- cuBLAS / Magma - GPU libraries from Nvidia and UTK respectively
- Accelerate / vecLib - Apple's framework for GPU and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

