

HPC With R: The Basics

Drew Schmidt

November 12, 2016



Tutorial Goals

We hope to introduce you to:

- 1 Basic debugging.
- 2 Evaluating the performance of R code.
- 3 Some R best practices to help with performance.
- 4 Basics of parallelism in R.



Exercises

Each section has a complement of exercises to give hands-on reinforcement of ideas introduced in the lecture.

- 1 Later exercises are more difficult than earlier ones.
- 2 Some exercises require use of things not explicitly shown in lecture; look through the documentation mentioned in the slides to find the information you need.



Part I

Basics



- 1 Introduction
- 2 Debugging
- 3 Profiling
- 4 Benchmarking



Resources for Learning R

- *The Art of R Programming* by Norm Matloff: <http://nostarch.com/artofr.htm>
- *An Introduction to R* by Venables, Smith, and the R Core Team: <http://cran.r-project.org/doc/manuals/R-intro.pdf>
- *The R Inferno* by Patrick Burns: http://www.burns-stat.com/pages/Tutor/R_inferno.pdf
- Mathesaurus: <http://mathesaurus.sourceforge.net/>
- R programming for those coming from other languages: http://www.johndcook.com/R_language_for_programmers.html
- *aRrgh: a newcomer's (angry) guide to R*, by Tim Smith and Kevin Ushey: <http://tim-smith.us/arrgh/>



Other Invaluable Resources

- *R Installation and Administration*: <http://cran.r-project.org/doc/manuals/R-admin.html>
- *Task Views*: <http://cran.at.r-project.org/web/views>
- *Writing R Extensions*: <http://cran.r-project.org/doc/manuals/R-exts.html>
- Mailing list archives: <http://tolstoy.newcastle.edu.au/R/>
- The [R] stackoverflow tag.
- The #rstats hastag on Twitter.



- 1 Introduction
- 2 Debugging
 - Debugging R Code
 - The R Debugger
 - Debugging Compiled Code Called by R Code
- 3 Profiling
- 4 Benchmarking



2 Debugging

- Debugging R Code
- The R Debugger
- Debugging Compiled Code Called by R Code



Debugging R Code

- Very broad topic ...
- We'll hit the highlights.
- For more examples, see:
cran.r-project.org/doc/manuals/R-exts.html#Debugging

Object Inspection Tools

- `print()`
- `str()`
- `unclass()`

Object Inspection Tools: print()

Basic printing:

```
1 > x <- matrix(1:10, nrow=2)
2 > print(x)
3      [,1] [,2] [,3] [,4] [,5]
4 [1,]    1    3    5    7    9
5 [2,]    2    4    6    8   10
6 > x
7      [,1] [,2] [,3] [,4] [,5]
8 [1,]    1    3    5    7    9
9 [2,]    2    4    6    8   10
```

Object Inspection Tools: `str()`

Examining the structure of an R object:

```
1 > x <- matrix(1:10, nrow=2)
2 > str(x)
3 int [1:2, 1:5] 1 2 3 4 5 6 7 8 9 10
```

Object Inspection Tools: `unclass()`

Exposing all data with `unclass()`:

```
1 df <- data.frame(x=rnorm(10), y=rnorm(10))
2 mdl <- lm(y~x, data=df) ### That's a "tilde" character
3
4 mdl
5 print(mdl)
6
7 str(mdl)
8
9 unclass(mdl)
```

Try it!

2 Debugging

- Debugging R Code
- The R Debugger
- Debugging Compiled Code Called by R Code



The R Debugger

- `debug()`
- `debugonce()`
- `undebug()`

Using The R Debugger

- 1 Declare function to be debugged: `debug(foo)`
- 2 Call function: `foo(arg1, arg2, ...)`
 - `next`: Enter or n followed by Enter.
 - `break`: Halt execution and exit debugging: Q.
 - `exit`: Continue execution and exit debugging: c.
- 3 Call `undebug()` to stop debugging

Using the Debugger

Example Debugger Interaction

```
1 > f <- function(x){y <- z+1;z <- y*2;z}
2 > f(1)
3 Error in f(1) : object 'z' not found
4 > debug(f)
5 > f(1)
6 debugging in: f(1)
7 debug at #1: {
8     y <- z + 1
9     z <- y * 2
10    z
11 }
12 Browse[2]>
13 debug at #1: y <- z + 1
14 Browse[2]>
15 Error in f(1) : object 'z' not found
16 >
```

2 Debugging

- Debugging R Code
- The R Debugger
- Debugging Compiled Code Called by R Code

Debugging Compiled Code

- Reasonably easy to use gdb and Valgrind.
- See “Writing R Extensions” manual.



- 1 Introduction
- 2 Debugging
- 3 Profiling
 - Why Profile?
 - Profiling R Code
 - Advanced R Profiling
- 4 Benchmarking



3 Profiling

- Why Profile?
- Profiling R Code
- Advanced R Profiling



Performance and Accuracy



Sometimes $\pi = 3.14$ is (a) infinitely faster than the “correct” answer and (b) the difference between the “correct” and the “wrong” answer is meaningless. . . . The thing is, some specious value of “correctness” is often irrelevant because it doesn’t matter. While performance almost always matters. And I absolutely detest the fact that people so often dismiss performance concerns so readily.

— Linus Torvalds, August 8, 2008

Compilers often correct bad behavior...

A Really Dumb Loop

```

1 int main(){
2     int x, i;
3     for (i=0; i<10; i++)
4         x = 1;
5     return 0;
6 }

```

clang -O3 -S example.c

```

main:
    .cfi_startproc
# BB#0:
    xorl    %eax, %eax
    ret

```

clang -S example.c

```

main:
    .cfi_startproc
# BB#0:
    movl    $0, -4(%rsp)
    movl    $0, -12(%rsp)
.LBB0_1:
    cmpl    $10, -12(%rsp)
    jge     .LBB0_4
# BB#2:
    movl    $1, -8(%rsp)
# BB#3:
    movl    -12(%rsp), %eax
    addl    $1, %eax
    movl    %eax, -12(%rsp)
    jmp     .LBB0_1
.LBB0_4:
    movl    $0, %eax
    ret

```



R will not!

Dumb Loop

```
1 for (i in 1:n){
2   tA <- t(A)
3   Y <- tA %*% Q
4   Q <- qr.Q(qr(Y))
5   Y <- A %*% Q
6   Q <- qr.Q(qr(Y))
7 }
8
9 Q
```

Better Loop

```
1 tA <- t(A)
2
3 for (i in 1:n){
4   Y <- tA %*% Q
5   Q <- qr.Q(qr(Y))
6   Y <- A %*% Q
7   Q <- qr.Q(qr(Y))
8 }
9
10 Q
```

Example from a Real R Package

Exerpt from Original function

```
1 while(i<=N){
2   for(j in 1:i){
3     d.k <- as.matrix(x)[l==j,l==j]
4     ...
```

Exerpt from Modified function

```
1 x.mat <- as.matrix(x)
2
3 while(i<=N){
4   for(j in 1:i){
5     d.k <- x.mat[l==j,l==j]
6     ...
```

By changing just 1 line of code, performance of the main method improved by **over 350%**!



Some Thoughts

- R is slow.
- Bad programmers are slower.
- R can't fix bad programming.



3 Profiling

- Why Profile?
- Profiling R Code
- Advanced R Profiling



Timings

Getting simple timings as a basic measure of performance is easy, and valuable.

- `system.time()` — timing blocks of code.
- `Rprof()` — timing execution of R functions.
- `Rprofmem()` — reporting memory allocation in R .
- `tracemem()` — detect when a copy of an R object is created.



Performance Profiling Tools: `system.time()`

`system.time()` is a basic R utility for timing expressions

```
1 x <- matrix(rnorm(20000*750), nrow=20000, ncol=750)
2
3 system.time(t(x) %*% x)
4 #   user  system elapsed
5 #  2.187   0.032   2.324
6
7 system.time(crossprod(x))
8 #   user  system elapsed
9 #  1.009   0.003   1.019
10
11 system.time(cov(x))
12 #   user  system elapsed
13 #  6.264   0.026   6.338
```

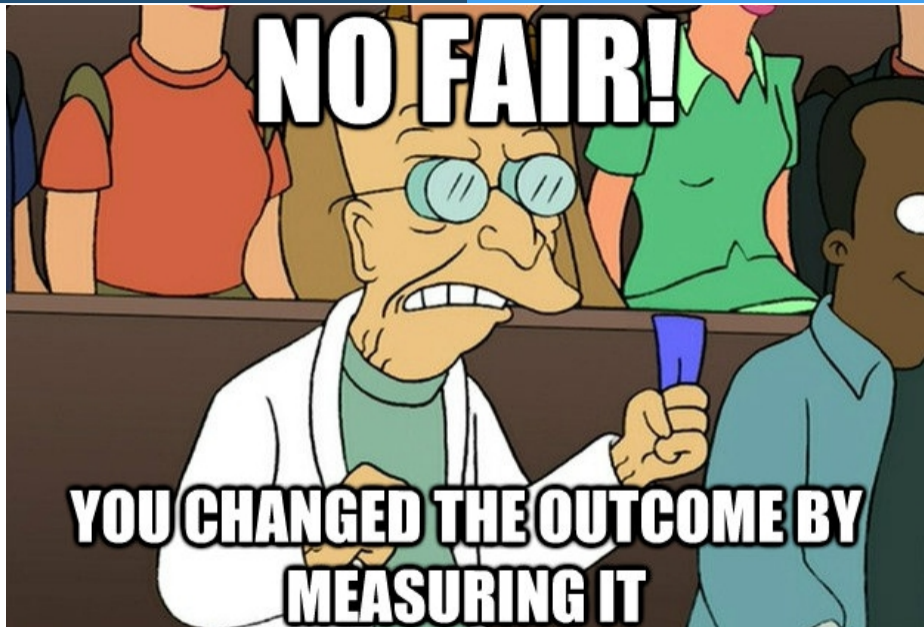
Performance Profiling Tools: `system.time()`

Put more complicated expressions inside of brackets:

```
1 x <- matrix(rnorm(20000*750), nrow=20000, ncol=750)
2
3 system.time({
4   y <- x+1
5   z <- y*2
6 })
7 #   user  system elapsed
8 # 0.057   0.032   0.089
```

Performance Profiling Tools: Rprof()

```
1 Rprof(filename="Rprof.out", append=FALSE, interval=0.02,  
2   memory.profiling=FALSE, gc.profiling=FALSE,  
3   line.profiling=FALSE, numfiles=100L, bufsize=10000L)
```

Performance Profiling Tools: Rprof()

```
1 x <- matrix(rnorm(10000*250), nrow=10000, ncol=250)
2
3 Rprof()
4 invisible(prcomp(x))
5 Rprof(NULL)
6
7 summaryRprof()
8
9 Rprof(interval=.99)
10 invisible(prcomp(x))
11 Rprof(NULL)
12
13 summaryRprof()
```

Performance Profiling Tools: Rprof()

```

1 $by.self
2           self.time self.pct total.time total.pct
3 "La.svd"      0.68   69.39      0.72   73.47
4 "%*%"        0.12   12.24      0.12   12.24
5 "aperm.default" 0.04    4.08      0.04    4.08
6 "array"       0.04    4.08      0.04    4.08
7 "matrix"      0.04    4.08      0.04    4.08
8 "sweep"       0.02    2.04      0.10   10.20
9 ### output truncated by presenter
10
11 $by.total
12           total.time total.pct self.time self.pct
13 "prcomp"      0.98   100.00      0.00    0.00
14 "prcomp.default" 0.98   100.00      0.00    0.00
15 "svd"         0.76   77.55      0.00    0.00
16 "La.svd"      0.72   73.47      0.68   69.39
17 ### output truncated by presenter
18
19 $sample.interval
20 [1] 0.02
21
22 $sampling.time
23 [1] 0.98

```

Performance Profiling Tools: Rprof()

```
1 $by.self
2 [1] self.time self.pct total.time total.pct
3 <0 rows> (or 0-length row.names)
4
5 $by.total
6 [1] total.time total.pct self.time self.pct
7 <0 rows> (or 0-length row.names)
8
9 $sample.interval
10 [1] 0.99
11
12 $sampling.time
13 [1] 0
```

3 Profiling

- Why Profile?
- Profiling R Code
- Advanced R Profiling



Other Profiling Tools

- perf, PAPI
- fpmi, mpiP, TAU
- pbdPROF
- pbdPAPI

See forthcoming paper *Analyzing Analytics: Advanced Performance Analysis Tools for R* for more details.



- 1 Introduction
- 2 Debugging
- 3 Profiling
- 4 Benchmarking



Benchmarking

- There's *a lot* that goes on when executing an R function.
- Symbol lookup, creating the abstract syntax tree, creating promises for arguments, argument checking, creating environments, ...
- Executing a second time can have dramatically different performance over the first execution.
- Benchmarking several methods fairly requires some care.

Benchmarking tools: rbenchmark

rbenchmark is a simple package that easily benchmarks different functions:

```
1 x <- matrix(rnorm(10000*500), nrow=10000, ncol=500)
2
3 f <- function(x) t(x) %*% x
4 g <- function(x) crossprod(x)
5
6 library(rbenchmark)
7 benchmark(f(x), g(x), columns=c("test", "replications", "elapsed", "relative"))
8
9 #   test replications elapsed relative
10 # 1 f(x)           100   13.679    3.588
11 # 2 g(x)           100    3.812    1.000
```

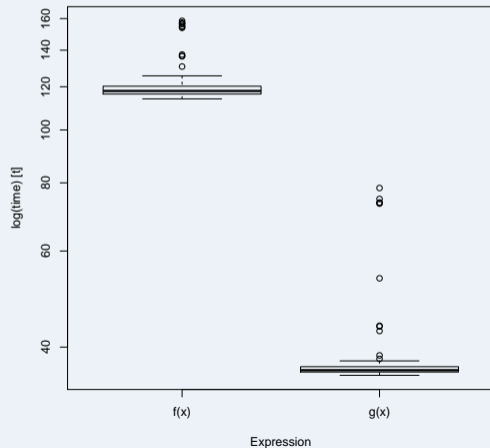
Benchmarking tools: microbenchmark

microbenchmark is a separate package with a slightly different philosophy:

```
1 x <- matrix(rnorm(10000*500), nrow=10000, ncol=500)
2
3 f <- function(x) t(x) %*% x
4 g <- function(x) crossprod(x)
5
6 library(microbenchmark)
7 microbenchmark(f(x), g(x), unit="s")
8
9 # Unit: seconds
10 #   expr      min       lq      mean     median      uq      max  neval
11 #   f(x) 0.11418617 0.11647517 0.12258556 0.11754302 0.12058145 0.17292507   100
12 #   g(x) 0.03542552 0.03613772 0.03884497 0.03668231 0.03740173 0.07478309   100
```

Benchmarking tools: microbenchmark

```
1 bench <- microbenchmark(f(x), g(x), unit="s")  
2 boxplot(bench)
```



Part II

Improving R Performance



5 Free Improvements

- Packages
- The Bytecode Compiler
- Choice of BLAS Library

6 Writing Better R Code



5 Free Improvements

- Packages
- The Bytecode Compiler
- Choice of BLAS Library

Packages

- Many high-quality “application” packages exist.
- Data manipulation: dplyr, data.table
- Modeling/math: Many! Try the CRAN taskviews.
- Parallelism: Discussed separately.

5 Free Improvements

- Packages
- The Bytecode Compiler
- Choice of BLAS Library



The Compiler Package

- Released in 2011 (Tierney)
- Bytecode: sort of like machine code for interpreters. . .
- Improves R code speed by 2-5% generally.
- Does best on loops.



Bytecode Compilation

- Non-core packages not (bytecode) compiled by default.
- “Base” and “recommended” (core) packages are.
- Downsides:
 - (slightly) larger install size
 - (much!) longer install process
 - doesn't fix bad code
- Upsides: slightly faster.



Compiling a Function

```
1 test <- function(x) x+1
2 test
3 # function(x) x+1
4
5 library(compiler)
6
7 test <- cmpfun(test)
8 test
9 # function(x) x+1
10 # <bytecode: 0x38c86c8>
11
12 disassemble(test)
13 # list(.Code, list(7L, GETFUN.OP, 1L, MAKEPROM.OP, 2L, PUSHCONSTARG.OP,
14 #   3L, CALL.OP, 0L, RETURN.OP), list(x + 1, '+', list(.Code,
15 #   list(7L, GETVAR.OP, 0L, RETURN.OP), list(x)), 1))
```

Compiling Packages

From R

```
1 install.packages("my_package", type="source", INSTALL_opts="--byte-compile")
```

From The Shell

```
1 export R_COMPILE_PKGS=1
2 R CMD INSTALL my_package.tar.gz
```

Or add the line: `ByteCompile: yes` to the package's DESCRIPTION file.



The Compiler: How much does it help *really*?

```
1 f <- function(n) for (i in 1:n) 2*(3+4)
2
3
4 library(compiler)
5 f_comp <- cmpfun(f)
6
7
8 library(rbenchmark)
9
10 n <- 100000
11 benchmark(f(n), f_comp(n), columns=c("test", "replications", "elapsed",
12   "relative"),
13   order="relative")
14 #           test replications elapsed relative
15 # 2 f_comp(n)          100    2.604    1.000
16 # 1      f(n)           100    2.845    1.093
```



The Compiler: How much does it help *really*?

```
1 g <- function(n){
2   x <- matrix(runif(n*n), nrow=n, ncol=n)
3   min(colSums(x))
4 }
5
6 library(compiler)
7 g_comp <- cmpfun(g)
8
9
10 library(rbenchmark)
11
12 n <- 1000
13 benchmark(g(n), g_comp(n), columns=c("test", "replications", "elapsed",
14   "relative"),
15   order="relative")
16 #           test replications elapsed relative
17 # 2 g_comp(n)           100    6.854    1.000
18 # 1 g(n)                 100    6.860    1.001
```

5 Free Improvements

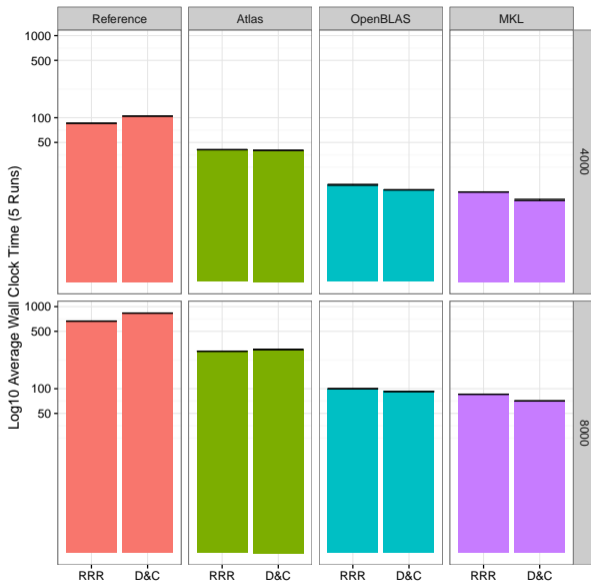
- Packages
- The Bytecode Compiler
- Choice of BLAS Library



The BLAS

- Basic Linear Algebra Subprograms.
- Basic numeric matrix operations.
- Used in linear algebra and many statistical operations.
- Different implementations available.
- Several multithreaded BLAS libraries exist.





Comparing Symmetric Eigenvalue Performance <http://bit.ly/2f49Sop>



5 Free Improvements

6 Writing Better R Code

- Loops
- Ply Functions
- Vectorization
- Loops, Plys, and Vectorization



6 Writing Better R Code

- **Loops**
- Ply Functions
- Vectorization
- Loops, Plys, and Vectorization



Loops

- `for`
- `while`
- No `goto`'s or `do while`'s.
- They're *really* slow.



Loops: Best Practices

- *Profile, profile, profile.*
- Mostly try to avoid.
- Evaluate practicality of rewrite (plys, vectorization, compiled code)
- Always preallocate storage; don't grow it dynamically.



6 Writing Better R Code

- Loops
- **Ply Functions**
- Vectorization
- Loops, Plys, and Vectorization



“Ply” Functions

- R has functions that apply other functions to data.
- In a nutshell: loop sugar.
- Typical *ply's:
 - `apply()`: apply function over matrix “margin(s)” .
 - `lapply()`: apply function over list/vector.
 - `mapply()`: apply function over multiple lists/vectors.
 - `sapply()`: same as `lapply()`, but (possibly) nicer output.
 - Plus some other mostly irrelevant ones.
- Also `Map()` and `Reduce()`.



Ply Examples: `apply()`

```
1 x <- matrix(1:10, 2)
2
3 x
4 #      [,1] [,2] [,3] [,4] [,5]
5 # [1,]    1    3    5    7    9
6 # [2,]    2    4    6    8   10
7
8 apply(X=x, MARGIN=1, FUN=sum)
9 # [1] 25 30
10
11 apply(X=x, MARGIN=2, FUN=sum)
12 # [1]  3  7 11 15 19
13
14 apply(X=x, MARGIN=1:2, FUN=sum)
15 #      [,1] [,2] [,3] [,4] [,5]
16 # [1,]    1    3    5    7    9
17 # [2,]    2    4    6    8   10
```



Ply Examples: `lapply()` and `sapply()`

```
1 lapply(1:4, sqrt)
2 # [[1]]
3 # [1] 1
4 #
5 # [[2]]
6 # [1] 1.414214
7 #
8 # [[3]]
9 # [1] 1.732051
10 #
11 # [[4]]
12 # [1] 2
13
14 sapply(1:4, sqrt)
15 # [1] 1.000000 1.414214 1.732051 2.000000
```

Transforming Loops Into Ply's

```
1 vec <- numeric(n)
2 for (i in 1:n){
3   vec[i] <- my_function(i)
4 }
```

Becomes:

```
1 sapply(1:n, my_function)
```

Ply's: Best Practices

- Most ply's are just shorthand/higher expressions of loops.
- Generally not much faster (if at all), especially with the compiler.
- Thinking in terms of `lapply()` can be useful however...



Ply's: Best Practices

- With ply's and lambdas, can do some fiendishly crafty things.
- But don't go crazy...

```
1 cat(sapply(letters, function(a) sapply(letters, function(b) sapply(letters,  
  function(c) sapply(letters, function(d) paste(a, b, c, d, letters, "\n",  
    sep=" "))))))
```



6 Writing Better R Code

- Loops
- Ply Functions
- **Vectorization**
- Loops, Plys, and Vectorization



Vectorization

- `x+y`
- `x[, 1] <- 0`
- `rnorm(1000)`



Vectorization

- Same in R as in other high-level languages (Matlab, Python, ...).
- Idea: use pre-existing compiled kernels to avoid interpreter overhead.
- Much faster than loops and plys.

```
1 ply <- function(x) lapply(rep(1, 1000), rnorm)
2 vec <- function(x) rnorm(1000)
3
4 library(rbenchmark)
5 benchmark(ply(x), vec(x))
6 #      test replications elapsed relative
7 # 1 ply(x)           100   0.348   38.667
8 # 2 vec(x)           100   0.009    1.000
```

6 Writing Better R Code

- Loops
- Ply Functions
- Vectorization
- Loops, Plys, and Vectorization



Putting It All Together

- Loops are slow.
- `apply()`, `Reduce()` are just for loops.
- `Map()`, `lapply()`, `sapply()`, `mapply()` (and most other core ones) are *not* for loops.
- *Ply functions are not vectorized.*
- Vectorization is fastest, but often needs lots of memory.



Squares

Let's compute the square of the numbers 1–100000, using

- for loop without preallocation
- for loop with preallocation
- `sapply()`
- vectorization



Squares

```
1 square_sapply <- function(n) sapply(1:n, function(i) i^2)
2
3 square_vec <- function(n) (1:n)*(1:n)
```

```
1 library(rbenchmark)
2 n <- 100000
3
4 benchmark(square_loop_noinit(n), square_loop_withininit(n), square_sapply(n),
5           square_vec(n))
6 #           test replications elapsed relative
7 # 1 square_loop_noinit(n)           100 17.296 2470.857
8 # 2 square_loop_withininit(n)         100  0.933  133.286
9 # 3 square_sapply(n)                 100  1.218  174.000
10 # 4 square_vec(n)                   100  0.007   1.000
```



Part III

Parallelism



- 7 An Overview of Parallelism
- 8 Shared Memory Parallelism in R
- 9 Distributed Memory Parallelism with R
- 10 Distributed Matrices



Parallel Programming Packages for R

Shared Memory

Examples: **parallel**, **snow**, **foreach**,
gputools, **HiPLARM**

Distributed

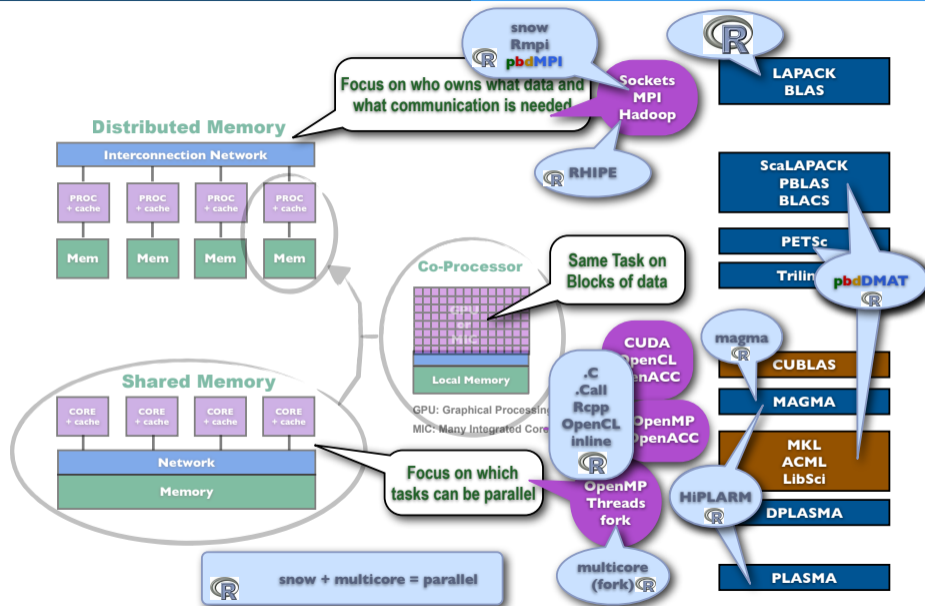
Examples: **pbdR**, **Rmpi**, **RHadoop**, **RHIPE**

CRAN HPC Task View

For more examples, see:

<http://cran.r-project.org/web/views/HighPerformanceComputing.html>





Portability

Many parallel R packages break on Windows

Windows

A fatal exception 0E has occurred at 0028:C0011E36 in UXD UMM(01) + 00010E36. The current application will be terminated.

- * Press any key to terminate the current application.
- * Press CTRL+ALT+DEL again to restart your computer. You will lose any unsaved information in all applications.

Press any key to continue _

RNG's in Parallel

- Be careful!
- Aided by **rlecuyer**, **rsprng**, and **doRNG** packages.

Parallel Programming: In Theory



Parallel Programming: In Practice



- 7 An Overview of Parallelism
- 8 Shared Memory Parallelism in R
 - The parallel Package
 - The foreach Package
- 9 Distributed Memory Parallelism with R
- 10 Distributed Matrices



- 8 Shared Memory Parallelism in R
 - The parallel Package
 - The foreach Package

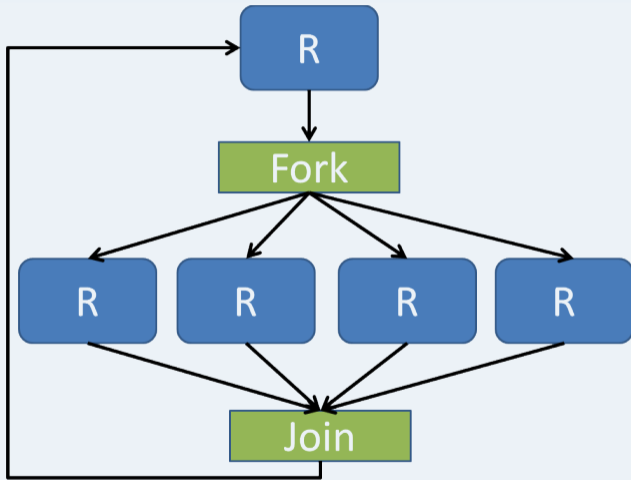
The parallel Package

- Comes with R \geq 2.14.0
- Has 2 disjoint interfaces.

parallel = snow + multicore



The parallel Package: multicore



Operates on fork/join paradigm.

The parallel Package: multicore

- + Data copied to child on write (handled by OS)
- + Very efficient.
 - No Windows support.
 - Not as efficient as threads.



The parallel Package: multicore

```
1 mclapply(X, FUN, ...,
2   mc.preschedule=TRUE, mc.set.seed=TRUE,
3   mc.silent=FALSE, mc.cores=getOption("mc.cores", 2L),
4   mc.cleanup=TRUE, mc.allow.recursive=TRUE)
```

```
1 x <- lapply(1:10, sqrt)
2
3 library(parallel)
4 x.mc <- mclapply(1:10, sqrt)
5
6 all.equal(x.mc, x)
7 # [1] TRUE
```

The parallel Package: multicore

```
1 simplify2array(mclapply(1:10, function(i) Sys.getpid(), mc.cores=4))
2 # [1] 27452 27453 27454 27455 27452 27453 27454 27455 27452 27453
3
4 simplify2array(mclapply(1:2, function(i) Sys.getpid(), mc.cores=4))
5 # [1] 27457 2745
```

The parallel Package: snow

- ? Uses sockets.
- + Works on all platforms.
- More fiddley than `mclapply()`.
- Not as efficient as forks.

The parallel Package: snow

```
1 ### Set up the worker processes
2 cl <- makeCluster(detectCores())
3 cl
4 # socket cluster with 4 nodes on host localhost
5
6 parSapply(cl, 1:5, sqrt)
7
8 stopCluster(cl)
```

The parallel Package: Summary

All

- `detectCores()`
- `splitIndices()`

multicore

- `mclapply()`
- `mcmapply()`
- `mcparallel()`
- `mccollect()`
- and others...

snow

- `makeCluster()`
- `stopCluster()`
- `parLapply()`
- `parSapply()`
- and others...



- 8 Shared Memory Parallelism in R
 - The parallel Package
 - The foreach Package



The foreach Package

- On Cran (Revolution Analytics).
- Main package is **foreach**, which is a single interface for a number of “backend” packages.
- Backends: **doMC**, **doMPI**, **doParallel**, **doRedis**, **doRNG**, **doSNOW**.



The foreach Package: The Idea



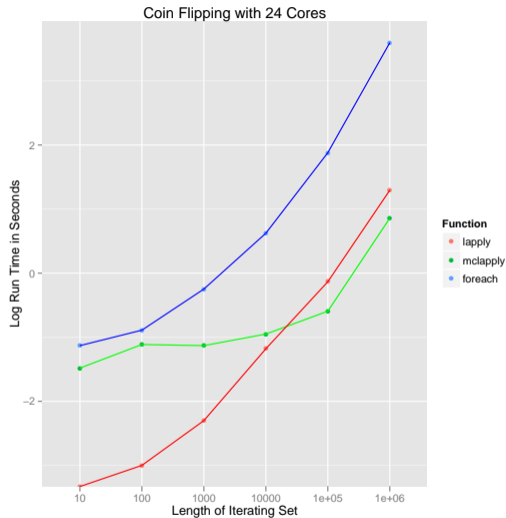
Unify the disparate interfaces.



The foreach Package

- + Works on all platforms (if backend does).
- + Can even work serial with minor notational change.
- + Write the code once, use whichever backend you prefer.
 - Really bizarre, non-R-ish syntax.
 - Efficiency issues if you aren't careful!





```

1  ### Bad performance
2  foreach(i=1:len) %dopar%
3      tinyfun(i)
4
5  ### Expected performance
6  foreach(i=1:ncores) %dopar% {
7      out <- numeric(len/ncores)
8      for (j in 1:(len/ncores))
9          out[i] <- tinyfun(j)
10     out
11 }

```

The foreach Package: General Procedure

- Load **foreach** and your backend package.
- Register your backend.
- Call `foreach`



Using foreach: serial

```
1 library(foreach)
2
3 ### Example 1
4 foreach(i=1:3) %do% sqrt(i)
5
6 ### Example 2
7 n <- 50
8 reps <- 100
9
10 x <- foreach(i=1:reps) %do% {
11   sum(rnorm(n, mean=i)) / (n*reps)
12 }
```

Using foreach: Parallel

```
1 library(foreach)
2 library(<mybackend>)
3
4 register<MyBackend>()
5
6 ### Example 1
7 foreach(i=1:3) %dopar% sqrt(i)
8
9 ### Example 2
10 n <- 50
11 reps <- 100
12
13 x <- foreach(i=1:reps) %dopar% {
14   sum(rnorm(n, mean=i)) / (n*reps)
15 }
```

foreach backends

multicore

```
1 library(doParallel)
2 registerDoParallel(cores=ncores)
3 foreach(i=1:2) %dopar% Sys.getpid()
```

snow

```
1 library(doParallel)
2 cl <- makeCluster(ncores)
3 registerDoParallel(cl=cl)
4
5 foreach(i=1:2) %dopar% Sys.getpid()
6 stopCluster(cl)
```



- 7 An Overview of Parallelism
- 8 Shared Memory Parallelism in R
- 9 Distributed Memory Parallelism with R
 - Distributed Memory Parallelism
 - Rmpi
 - pbdMPI vs Rmpi
- 10 Distributed Matrices

9 Distributed Memory Parallelism with R

- Distributed Memory Parallelism
 - Rmpi
 - pbdMPI vs Rmpi

Why Distribute?

- Nodes only hold so much ram.
- Commodity hardware: $\approx 32 - 64$ gib.
- With a few exceptions (**ff**, **bigmemory**), R does computations in memory.
- If your problem doesn't fit in the memory of one node. . .

Packages for Distributed Memory Parallelism in R

- **Rmpi**, and **snow** via **Rmpi**.
- **RHIPE** and **RHadoop** ecosystem.
- **pbdR** ecosystem.



9 Distributed Memory Parallelism with R

- Distributed Memory Parallelism
- Rmpi
- pbdMPI vs Rmpi

Rmpi Hello World

```
1 mpi.spawn.Rslaves(nslaves=2)
2 #           2 slaves are spawned successfully. 0 failed.
3 # master (rank 0, comm 1) of size 3 is running on: wootabega
4 # slave1 (rank 1, comm 1) of size 3 is running on: wootabega
5 # slave2 (rank 2, comm 1) of size 3 is running on: wootabega
6
7 mpi.remote.exec(paste("I am",mpi.comm.rank(),"of",mpi.comm.size()))
8 # $slave1
9 # [1] "I am 1 of 3"
10 #
11 # $slave2
12 # [1] "I am 2 of 3"
13
14 mpi.exit()
```

Using Rmpi from snow

```
1 library(snow)
2 library(Rmpi)
3
4 cl <- makeCluster(2, type = "MPI")
5 clusterCall(cl, function() Sys.getpid())
6 clusterCall(cl, runif, 2)
7 stopCluster(cl)
8 mpi.quit()
```

Rmpi Resources

- **Rmpi** tutorial: <http://math.acadiau.ca/ACMMaC/Rmpi/>
- **Rmpi** manual: <http://cran.r-project.org/web/packages/Rmpi/Rmpi.pdf>

9 Distributed Memory Parallelism with R

- Distributed Memory Parallelism
- Rmpi
- pbdMPI vs Rmpi



pbdMPI vs Rmpi

- **Rmpi** is interactive; **pbdMPI** is exclusively batch.
- **pbdMPI** is easier to install.
- **pbdMPI** has a simpler interface.
- **pbdMPI** integrates with other pbdR packages.

Example Syntax

Rmpi

```
1 # int
2 mpi.allreduce(x, type=1)
3 # double
4 mpi.allreduce(x, type=2)
```

pbdMPI

```
1 allreduce(x)
```

Types in R

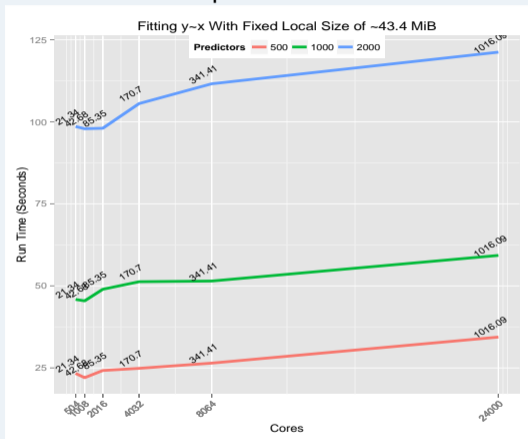
```
1 > typeof(1)
2 [1] "double"
3 > typeof(2)
4 [1] "double"
5 > typeof(1:2)
6 [1] "integer"
```

- 7 An Overview of Parallelism
- 8 Shared Memory Parallelism in R
- 9 Distributed Memory Parallelism with R
- 10 Distributed Matrices



Distributed Matrices and Statistics with pbdDMAT

Least Squares Benchmark



```
x <- ddmatrix("rnorm", nrow=m, ncol=n)
y <- ddmatrix("rnorm", nrow=m, ncol=1)
mdl <- lm.fit(x=x, y=y)
```

pbsR Scripts

- They're just R scripts.
- Can't run interactively (with more than 1 rank).
- We can use **pbsinline** to get “pretend interactivity”.



ddmatrix: 2-dimensional Block-Cyclic with 6 Processors

$$X = \begin{bmatrix} X_{11} & X_{12} & X_{13} & X_{14} & X_{15} & X_{16} & X_{17} & X_{18} & X_{19} \\ X_{21} & X_{22} & X_{23} & X_{24} & X_{25} & X_{26} & X_{27} & X_{28} & X_{29} \\ X_{31} & X_{32} & X_{33} & X_{34} & X_{35} & X_{36} & X_{37} & X_{38} & X_{39} \\ X_{41} & X_{42} & X_{43} & X_{44} & X_{45} & X_{46} & X_{47} & X_{48} & X_{49} \\ X_{51} & X_{52} & X_{53} & X_{54} & X_{55} & X_{56} & X_{57} & X_{58} & X_{59} \\ X_{61} & X_{62} & X_{63} & X_{64} & X_{65} & X_{66} & X_{67} & X_{68} & X_{69} \\ X_{71} & X_{72} & X_{73} & X_{74} & X_{75} & X_{76} & X_{77} & X_{78} & X_{79} \\ X_{81} & X_{82} & X_{83} & X_{84} & X_{85} & X_{86} & X_{87} & X_{88} & X_{89} \\ X_{91} & X_{92} & X_{93} & X_{94} & X_{95} & X_{96} & X_{97} & X_{98} & X_{99} \end{bmatrix}_{9 \times 9}$$

$$\text{Processor grid} = \begin{vmatrix} 0 & 1 & 2 \\ 3 & 4 & 5 \end{vmatrix} = \begin{vmatrix} (0,0) & (0,1) & (0,2) \\ (1,0) & (1,1) & (1,2) \end{vmatrix}$$



Understanding ddmatrix: Local View

$$\begin{array}{c}
 \left[\begin{array}{cc|cc}
 X_{11} & X_{12} & X_{17} & X_{18} \\
 X_{21} & X_{22} & X_{27} & X_{28} \\
 \hline
 X_{51} & X_{52} & X_{57} & X_{58} \\
 X_{61} & X_{62} & X_{67} & X_{68} \\
 \hline
 X_{91} & X_{92} & X_{97} & X_{98}
 \end{array} \right]_{5 \times 4}
 \quad
 \left[\begin{array}{cc|c}
 X_{13} & X_{14} & X_{19} \\
 X_{23} & X_{24} & X_{29} \\
 \hline
 X_{53} & X_{54} & X_{59} \\
 X_{63} & X_{64} & X_{69} \\
 \hline
 X_{93} & X_{94} & X_{99}
 \end{array} \right]_{5 \times 3}
 \quad
 \left[\begin{array}{c|c}
 X_{15} & X_{16} \\
 X_{25} & X_{26} \\
 \hline
 X_{55} & X_{56} \\
 X_{65} & X_{66} \\
 \hline
 X_{95} & X_{96}
 \end{array} \right]_{5 \times 2} \\
 \\
 \left[\begin{array}{cc|cc}
 X_{31} & X_{32} & X_{37} & X_{38} \\
 X_{41} & X_{42} & X_{47} & X_{48} \\
 \hline
 X_{71} & X_{72} & X_{77} & X_{78} \\
 X_{81} & X_{82} & X_{87} & X_{88}
 \end{array} \right]_{4 \times 4}
 \quad
 \left[\begin{array}{cc|c}
 X_{33} & X_{34} & X_{39} \\
 X_{43} & X_{44} & X_{49} \\
 \hline
 X_{73} & X_{74} & X_{79} \\
 X_{83} & X_{84} & X_{89}
 \end{array} \right]_{4 \times 3}
 \quad
 \left[\begin{array}{c|c}
 X_{35} & X_{36} \\
 X_{45} & X_{46} \\
 \hline
 X_{75} & X_{76} \\
 X_{85} & X_{86}
 \end{array} \right]_{4 \times 2}
 \end{array}$$

$$\text{Processor grid} = \left| \begin{array}{ccc|c}
 0 & 1 & 2 & \\
 3 & 4 & 5 &
 \end{array} \right| = \left| \begin{array}{cc|c}
 (0,0) & (0,1) & (0,2) \\
 (1,0) & (1,1) & (1,2)
 \end{array} \right|$$



Methods for class `dmatrix`

pbDMAT has over 100 methods with *identical* syntax to R:

- ``[, rbind(), cbind(), ...`
- `lm.fit(), prcomp(), cov(), ...`
- ``%*%`, solve(), svd(), norm(), ...`
- `median(), mean(), rowSums(), ...`

Serial Code

```
1 cov(x)
```

Parallel Code

```
1 cov(x)
```

Part IV

Wrapup



~Thanks!~

Questions?



Email: wmathematics@gmail.com



GitHub: <https://github.com/wmathematics>



Web: <http://wmathematics.info>



Twitter: [@wmathematics](https://twitter.com/wmathematics)