# Scalability

## Preparing for Production

Dr Zak Varty

# Scalability and Production

When put into production code gets used more and on more data.

We will likely have to consider scalability of our methods in

- Computation time

- Memory requirements

When doing so we have to balance a trade-off between development costs and usage costs.

# Example: Bayesian Inference

- MCMC originally takes ~24 hours

- Identifying and amending bottlenecks in code reduced this to ~24 minutes.
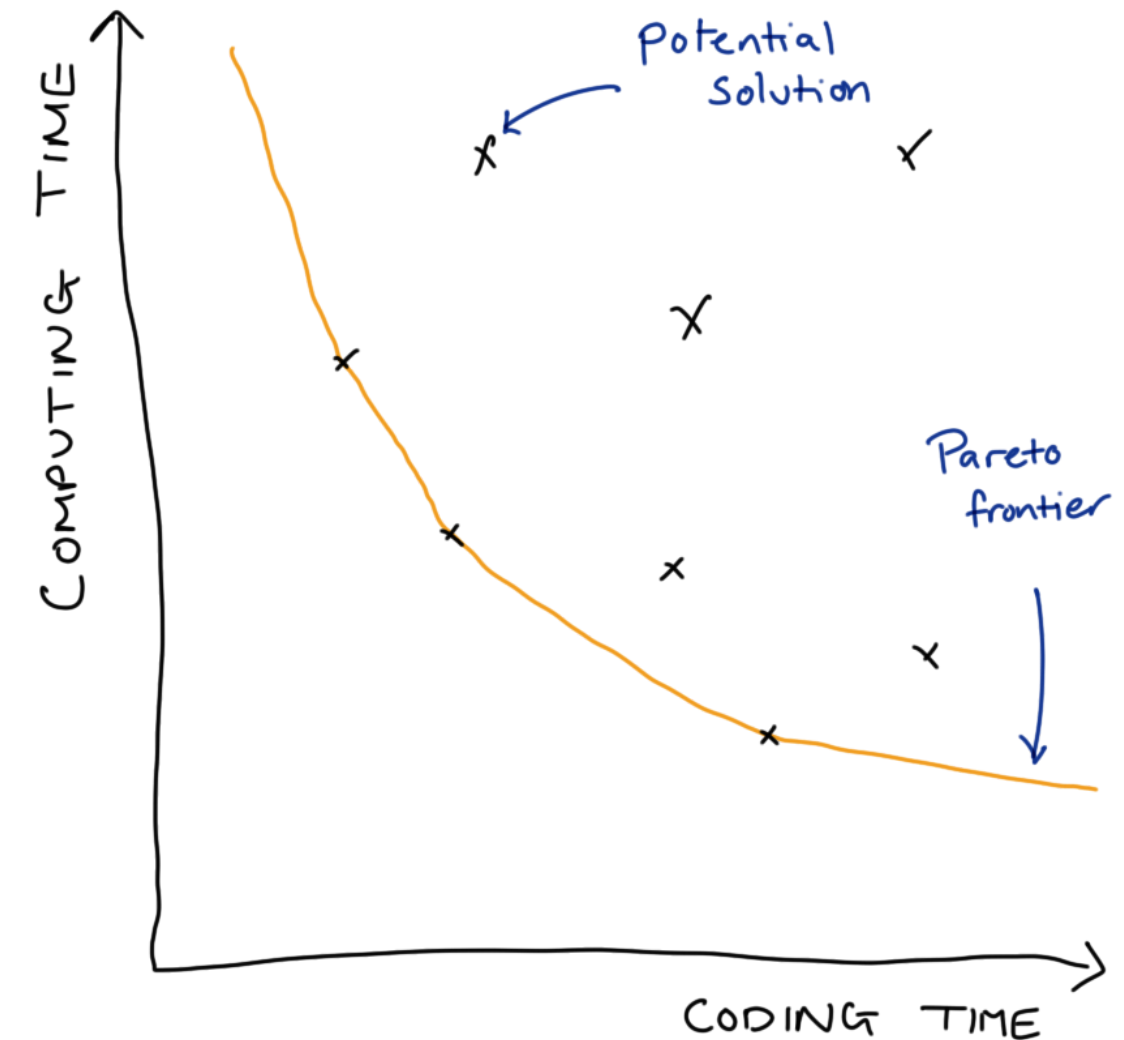
**Is this actually better?**

- human hours invested

- frequency of use

- safe / stable / general / readable

- trade for scalability

# Knowing when to worry

Sub-optimal optimisation can be worse than doing nothing

> … programmers have spent far too much time worrying about efficiency in **the wrong places** and at **the wrong times**; premature optimisation is the root of all evil (or at least most of it) in programming. - Donald Knuth

# This Lecture

- Basic profiling to find bottlenecks.

- Some simple solutions

- Strategies for scalable (R) code

- Signpost advanced methods & further reading

# Profiling your code: basics

## R as a stopwatch

```r
1  t_start <- Sys.time()
2  Sys.sleep(0.5) # YOUR CODE
3  t_end <- Sys.time()
4
5  t_end - t_start
```

```
Time difference of 0.5098979 secs
```

```r
1  library(tictoc)
2
3  tic()
4  Sys.sleep(0.5) # YOUR CODE
5  toc()
```

```
0.505 sec elapsed
```

## With {tictoc} we can get fancy

```r
 1  tic("total")
 2  tic("first, easy part")
 3  Sys.sleep(0.5)
 4  toc(log = TRUE)
 5  ## first, easy part: 0.505 sec elapsed
 6  tic("second, hard part")
 7  Sys.sleep(3)
 8  toc(log = TRUE)
 9  ## second, hard part: 3.008 sec elapsed
10  toc()
11  ## total: 3.519 sec elapsed
```

# Profiling your code in detail

To diagnose scaling issues you have to understand what your code is doing.
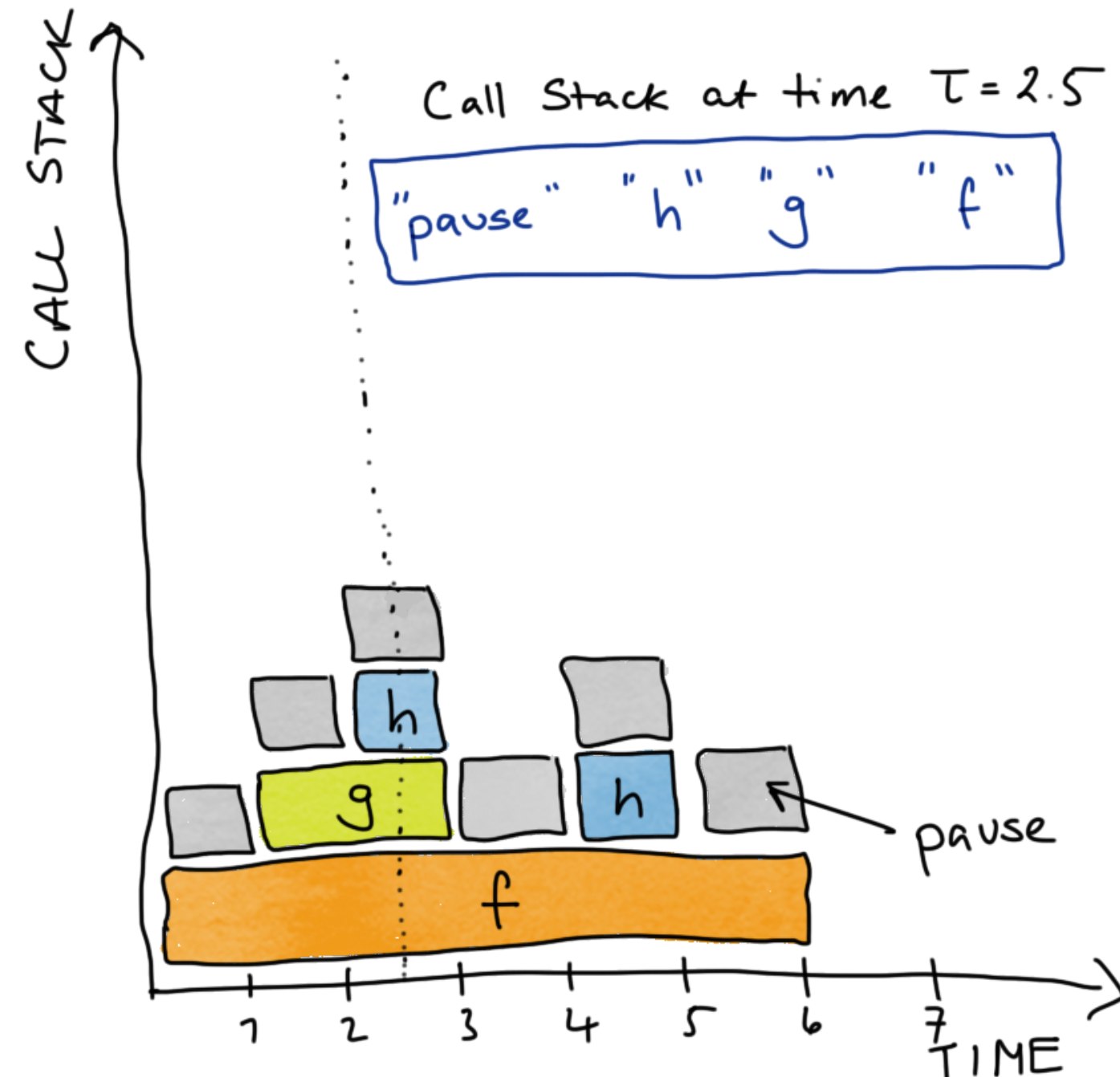
- Stop the code at time $\tau$ and examine the **call-stack**.

    - The current function being evaluated, the function that called that, the function that called that, ..., top level function.

- Do this a lot and you can measure (estimate) the proportion of working memory (RAM) uses over time and the time spent evaluating each function.

# Profiling: Toy Example

```
1  h <- function() {
2    profvis::pause(1)
3  }
4
5  g <- function() {
6    profvis::pause(1)
7    h()
8  }
9
10 f <- function() {
11   profvis::pause(1)
12   g()
13   profvis::pause(1)
14   h()
15 }
```
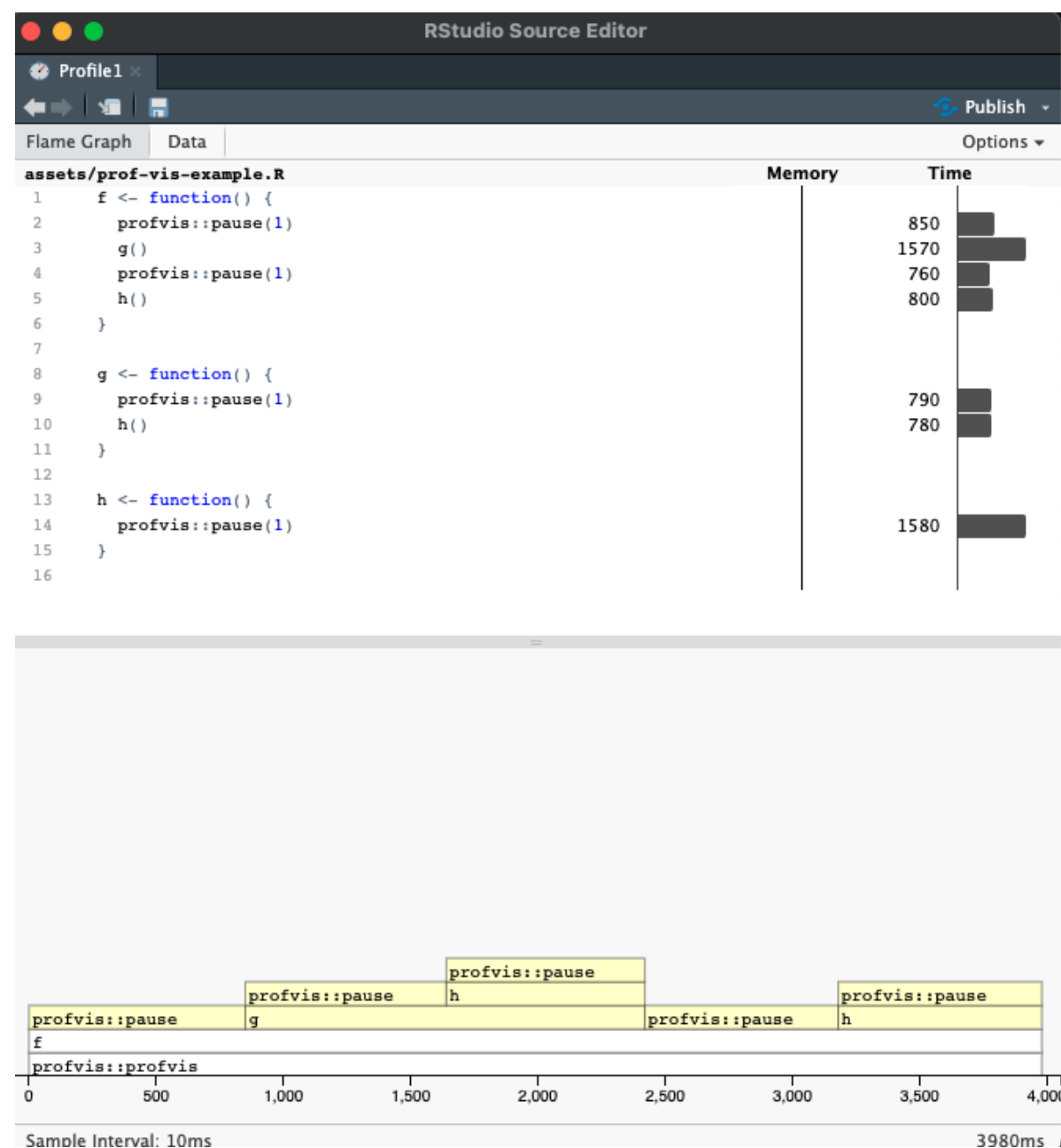
# Profiling: How To

```
1  source("assets/prof-vis-example.R")
2  profvis::profvis(f())
```

# Notes on Time Profiling

- Will get slightly different results each time you run the function

  - Changes to internal state of computer

  - Usually not a big deal, mainly effects fastest parts of code

  - Be careful with stochastic simulations

  - Use `set.seed()` to make a fair comparison over many runs.

# Notes on Profiling

## Function Source

```
1  pad_with_NAs
```
```
function(x, n_left, n_right){
  c(rep(NA, n_left), x, rep(NA, n_right))
}
```

## Compiled Function

```
1  mean
```
```
function (x, ...)
UseMethod("mean")
<bytecode: 0x7fabbac6d718>
<environment: namespace:base>
```
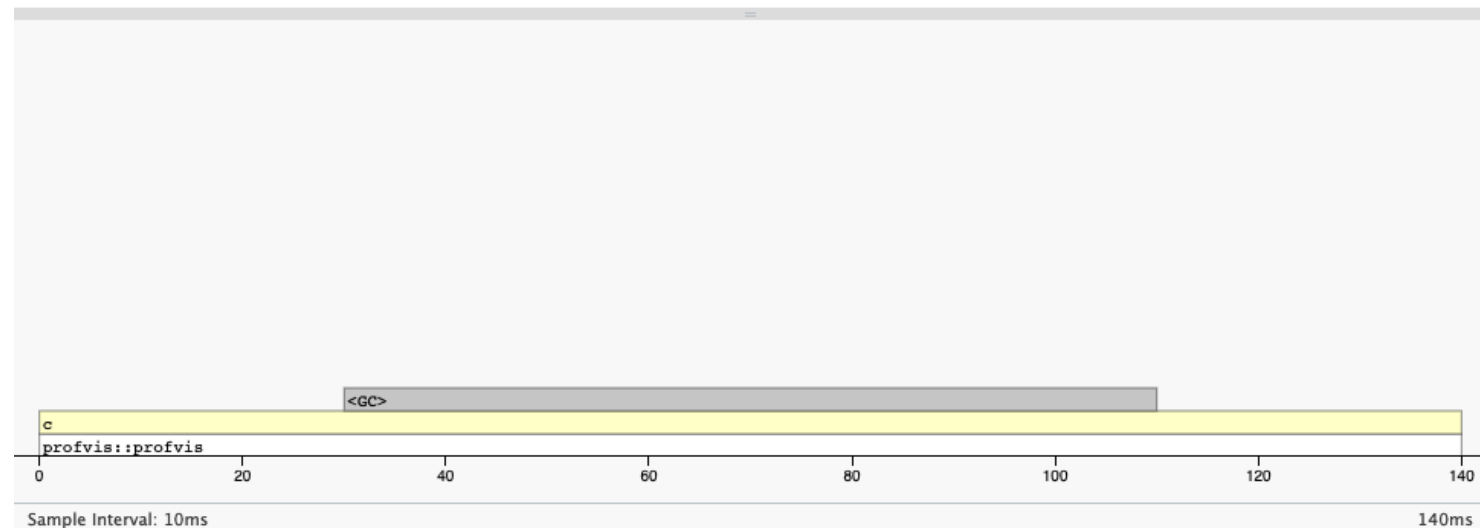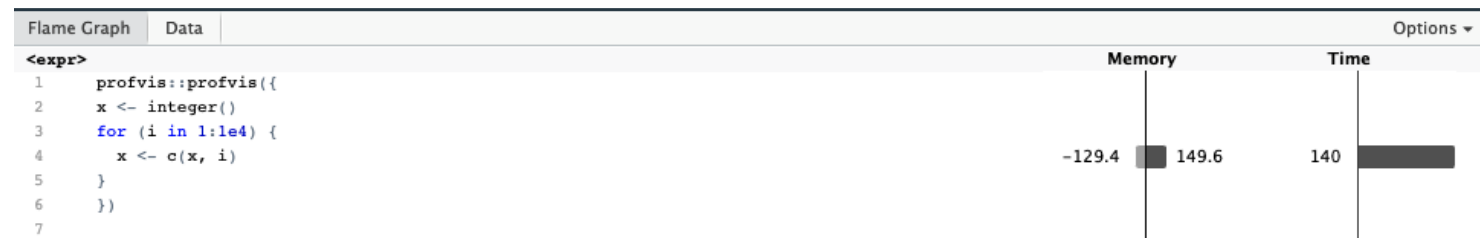
- Compiled functions have no R source code.

- Profiler does not extend into compiled code, see {jointprof} if you really need this.

# Memory Profiling

`profvis()` can similarly measure the memory usage of your code.

```
1  x <- integer()
2  for (i in 1:1e4) {
3    x <- c(x, i)
4  }
```



- Copy-on-modify behaviour makes growing objects slow.

- Pre-allocate storage where possible.

- Strategies and structures, see R inferno and Effecient R.



Effective Data Science: Production - Scalability - Zak Varty

# Tips to work at scale

# Vectorise

```r
1  x <- 1:10
2  y <- 11:20
3  z <- rep(NA, length(x))
4
5  for (i in seq_along(x)) {
6    z[i] <- x[i] * y[i]
7  }
```

```r
1  x <- 1:10
2  y <- 11:20
3  z <- x * y
```

Use and write functions with vectorised inputs.

```r
1  rnorm(n = 100, mean = 1:10, sd = rep(1, 10))
```

Be careful of recycling!

# Special vectors: Linear Algebra

```r
1  X <- diag(x = c(2, 0.5))
2  y <- matrix(data = c(1, 1), ncol = 1)
3
4  X %*% y
```

```
     [,1]
[1,]  2.0
[2,]  0.5
```

More on vectorising: Noam Ross Blog Post

# For loops in disguise: the apply family

Functional programming equivalent of a for loop. [`apply()`, `mapply()`, `lapply()`, ...]

Apply a function to each element of a list-like object.

```
1  A <- matrix(data = 1:12, nrow = 3, ncol = 4)
2  A
```

```
     [,1] [,2] [,3] [,4]
[1,]    1    4    7   10
[2,]    2    5    8   11
[3,]    3    6    9   12
```

```
1  # MARGIN = 1 => rows,  MARGIN = 2 => columns
2  apply(X = A, MARGIN = 1, FUN = sum)
```

```
[1] 22 26 30
```

Generalises functions from `{matrixStats}`

```
1  rowSums(A)
```

```
[1] 22 26 30
```

# For loops in disguise: purrr::map

**Iterate over a single object with `map()`:**

```r
1  mu <- c(-10, 0, 10)
2  purrr::map(.x = mu, .f = rnorm, n = 5)
```
```
[[1]]
[1] -10.630941 -10.536954 -10.393322  -8.395591  -8.770291

[[2]]
[1] -0.2173831  0.4013440 -2.2716535 -1.5199591  0.2817534

[[3]]
[1]  8.217122  7.209396  9.051645 10.628293 11.263494
```

**Iterate over multiple objects `map2()` and `pmap()`:**

```r
1  mu <- c(-10, 0, 10)
2  sigma <- c(0, 0.1, 0)
3  purrr::map2(.x = mu, .y = sigma, .f = rnorm, n = 5)
```
```
[[1]]
[1] -10 -10 -10 -10 -10

[[2]]
[1] -0.097381610 -0.104544473 -0.063586494 -0.002777205
-0.105194011

[[3]]
[1] 10 10 10 10 10
```

For more details and variants see Advanced R chapters 9-11 on functional programming.

# Easy parallelisation with furrr

- {parallel} and {futures} allow parallel coding over multiple cores.

- Powerful, but steep learning curve.

- {furrr} makes this very easy, just add future_ to purrr verbs.

```
1   mu <- c(-10, 0, 10)
2   furrr::future_map(
3     .x = mu,
4     .f = rnorm,
5     .options = furrr::furrr_options(seed = TRUE),
6     n = 5)
```

```
[[1]]
[1] -11.210327  -7.806065 -10.330897  -9.628868 -11.735650

[[2]]
[1]  0.1394846 -2.5554590  0.8892910  1.0211497  1.4515885

[[3]]
[1] 10.22833  9.81417 10.17171 11.94347  9.74026
```

Need to be very careful handling RNG. See R-bloggers for more details.

# Sometimes R doesn't cut it



- An API for running C++ code in R

  - Loops that need to be run in order

  - Lots of function calls (e.g. deep recursion)

  - Fast data structures

- Beyond our scope but good to know exists. Starting point: Advanced R Chapter 25.

# Wrapping up

## Summary

1. Pick you battles wisely

2. Target your energy with profiling

3. Scale loops with vectors

4. Scale loops in parallel processing

5. Scale in another language

## Help!

- Articles and blog links

- The R inferno (Circles 2-4)

- Advanced R (Chapters 23-25),

- Efficient R (Chapter 7).