

# Boosting R Code Performance

Kun Ren

# About Me

- Hedge fund quant researcher in Shanghai, China
- Low frequency and high frequency trading of equity and futures
- Financial data analysis, predictive statistical modeling, etc.
- Creator of several R packages: formattable, rlist, and pipeR
- Author of *Learning R Programming*

# Performance of R

- R is designed for easiness of use rather than performance
- R is interpreted, dynamic scripting language
- R can be much slower compared with C++ but fast enough for most of its purposes

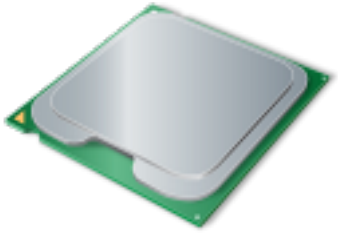
# When performance is important?

- Scenario 1: When data is of billions of rows
- Scenario 2: When an operation needs to be repeated many times
- Scenario 3: When the computing time is limited

# Understanding performance in general

```
1  [|||||||||100.0%]  11 [|||||||||100.0%]  21 [|||||||||100.0%]  31 [|||||||||100.0%]
2  [|||||||||100.0%]  12 [|||||||||100.0%]  22 [|||||||||100.0%]  32 [|||||||||100.0%]
3  [|||||||||100.0%]  13 [|||||||||100.0%]  23 [|||||||||100.0%]  33 [|||||||||100.0%]
4  [|||||||||100.0%]  14 [|||||||||100.0%]  24 [|||||||||100.0%]  34 [|||||||||100.0%]
5  [|||||||||100.0%]  15 [|||||||||100.0%]  25 [|||||||||100.0%]  35 [|||||||||100.0%]
6  [|||||||||100.0%]  16 [|||||||||100.0%]  26 [|||||||||100.0%]  36 [|||||||||100.0%]
7  [|||||||||100.0%]  17 [|||||||||100.0%]  27 [|||||||||100.0%]  37 [|||||||||100.0%]
8  [|||||||||100.0%]  18 [|||||||||100.0%]  28 [|||||||||100.0%]  38 [|||||||||100.0%]
9  [|||||||||100.0%]  19 [|||||||||100.0%]  29 [|||||||||100.0%]  39 [|||||||||100.0%]
10 [|||||||||100.0%]  20 [|||||||||100.0%]  30 [|||||||||100.0%]  40 [|||||||||100.0%]
Mem[|||||100.0%] 35.8G/504G Tasks: 122, 70 thr; 41 running
Swp[|] 111M/357G Load average: 29.05 14.30 7.76
Uptime: 21 days, 17:27:19
```

# Limited resources



CPU-bound



Memory-bound



I/O-bound

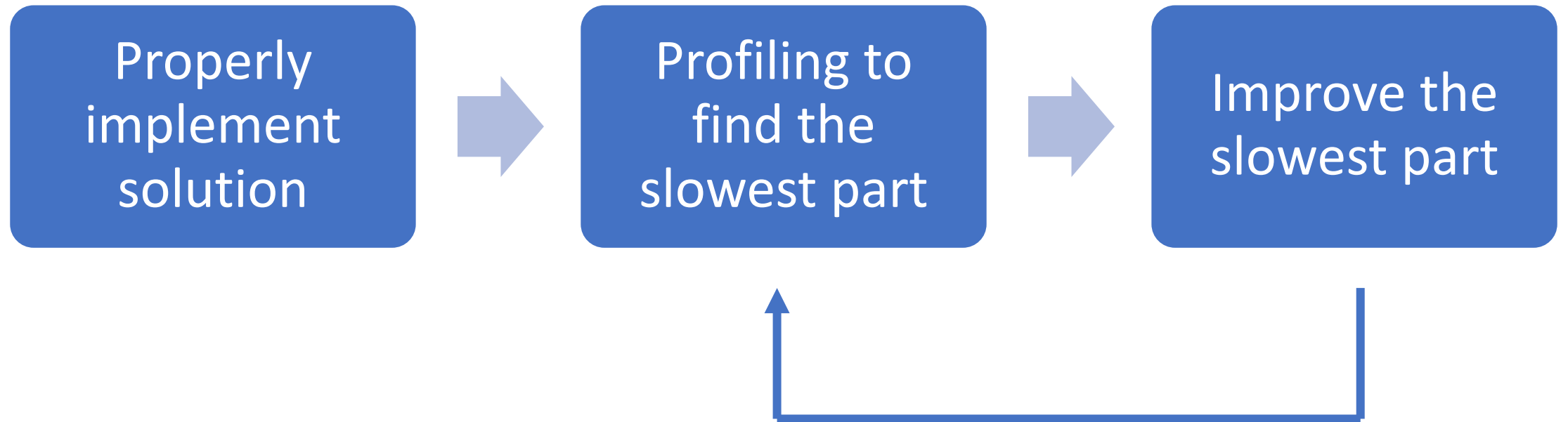


Network-bound

# Why can R be slow?

- R as an interpreted, dynamic scripting language
- Overhead of function calling
- Inefficient use of data structures
- Frequent copy of data: growing vectors, converting between data frame and matrix, etc

# A road map to speed up R code





# A rolling regression problem

date	ret	mret
1	-0.560475647	-0.4941738794
2	-0.230177489	1.1275934662
3	1.558708314	-1.1469495486
4	0.070508391	1.4810185971
5	0.129287735	0.9161912125
6	1.715064987	0.3351310171
7	0.460916206	0.5746753285
8	-1.265061235	0.2036196619
9	-0.686852852	-0.4470411920
10	-0.445661970	-0.3435259276
11	1.224081797	-0.6038357965
12	0.359813827	1.2386725185
13	0.400771451	0.5994920533
14	0.110682716	-0.1087337576
15	-0.555841135	-1.1388569648

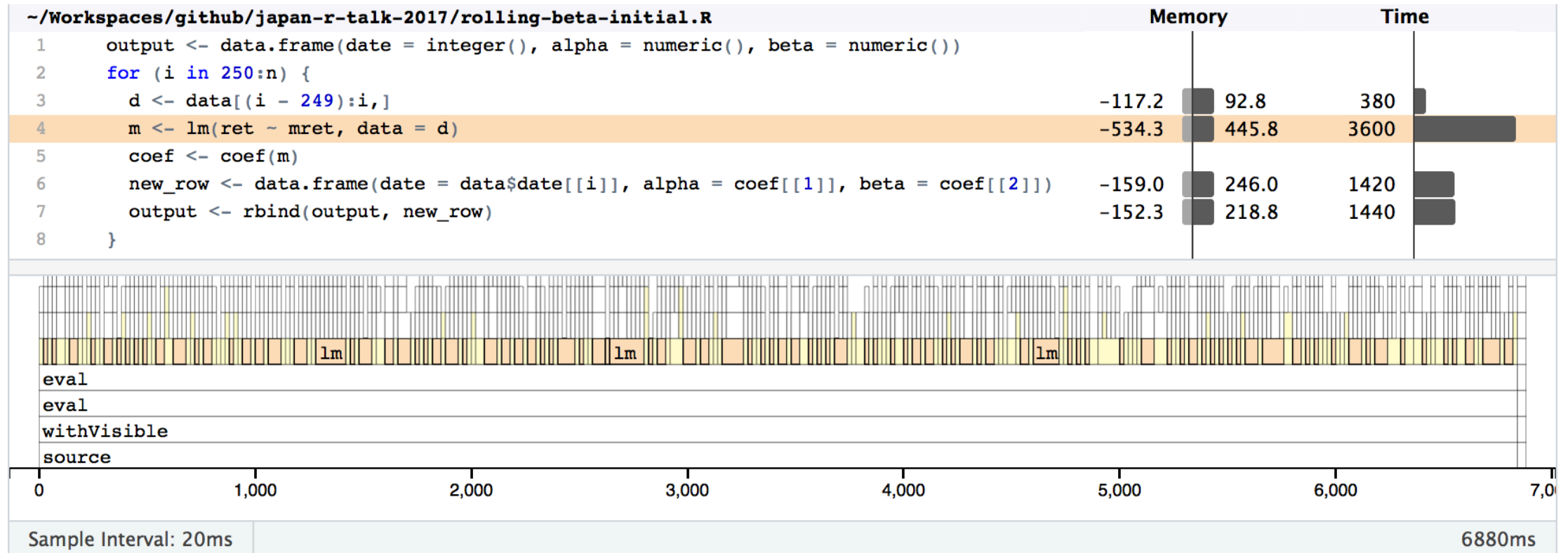
```
1 set.seed(123)
2 n <- 250 * 20
3 data <- data.frame(date = 1:n, ret = rnorm(n), mret = rnorm(n))
```

- For each date, estimate the linear coefficients for the recent 250 days
- $r = \alpha + \beta r_m + \epsilon$
- `lm(ret ~ mret)`

# An initial solution

```
1 output <- data.frame(date = integer(), alpha = numeric(), beta = numeric())
2 for (i in 250:n) {
3   d <- data[(i - 249):i,]
4   m <- lm(ret ~ mret, data = d)
5   coef <- coef(m)
6   new_row <- data.frame(date = data$date[[i]], alpha = coef[[1]], beta = coef[[2]])
7   output <- rbind(output, new_row)
8 }
```

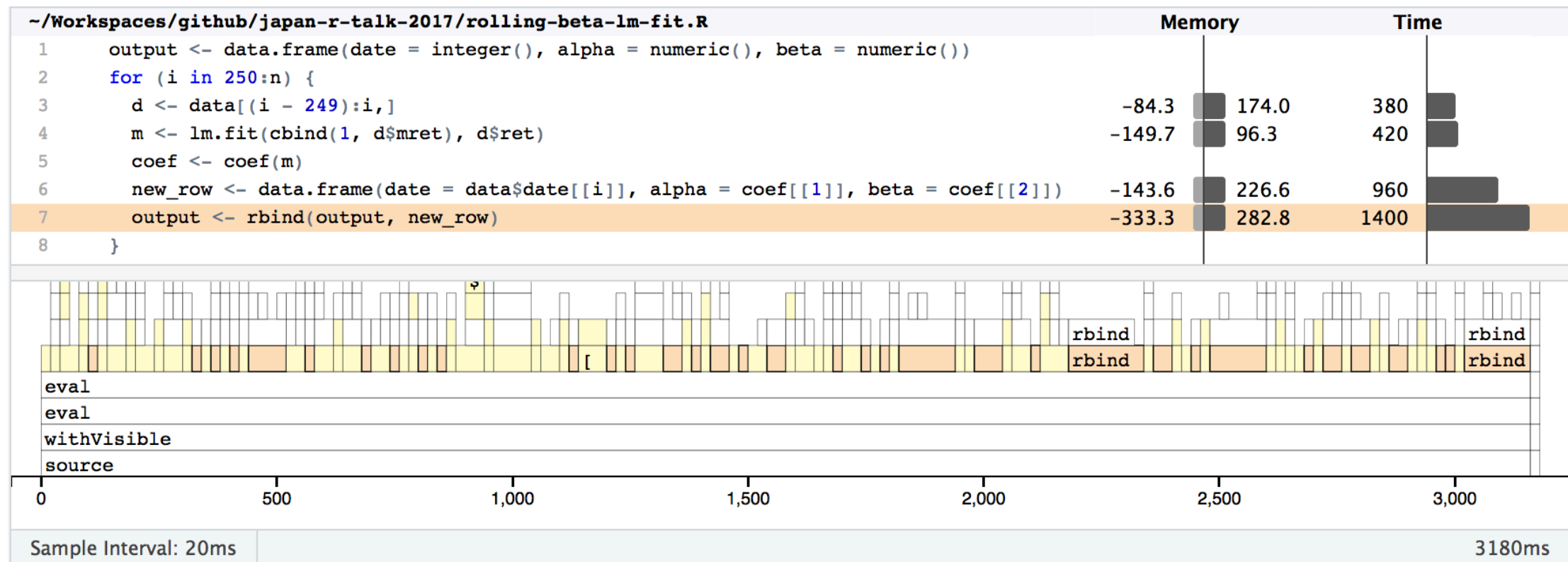
# Profiling initial solution (6880ms)



# Replace `lm` with `lm.fit`

```
1 output <- data.frame(date = integer(), alpha = numeric(), beta = numeric())
2 for (i in 250:n) {
3   d <- data[(i - 249):i,]
4   m <- lm.fit(cbind(1, d$mret), d$ret)
5   coef <- coef(m)
6   new_row <- data.frame(date = data$date[[i]], alpha = coef[[1]], beta = coef[[2]])
7   output <- rbind(output, new_row)
8 }
```

# Replace `lm` with `lm.fit` (3180ms)



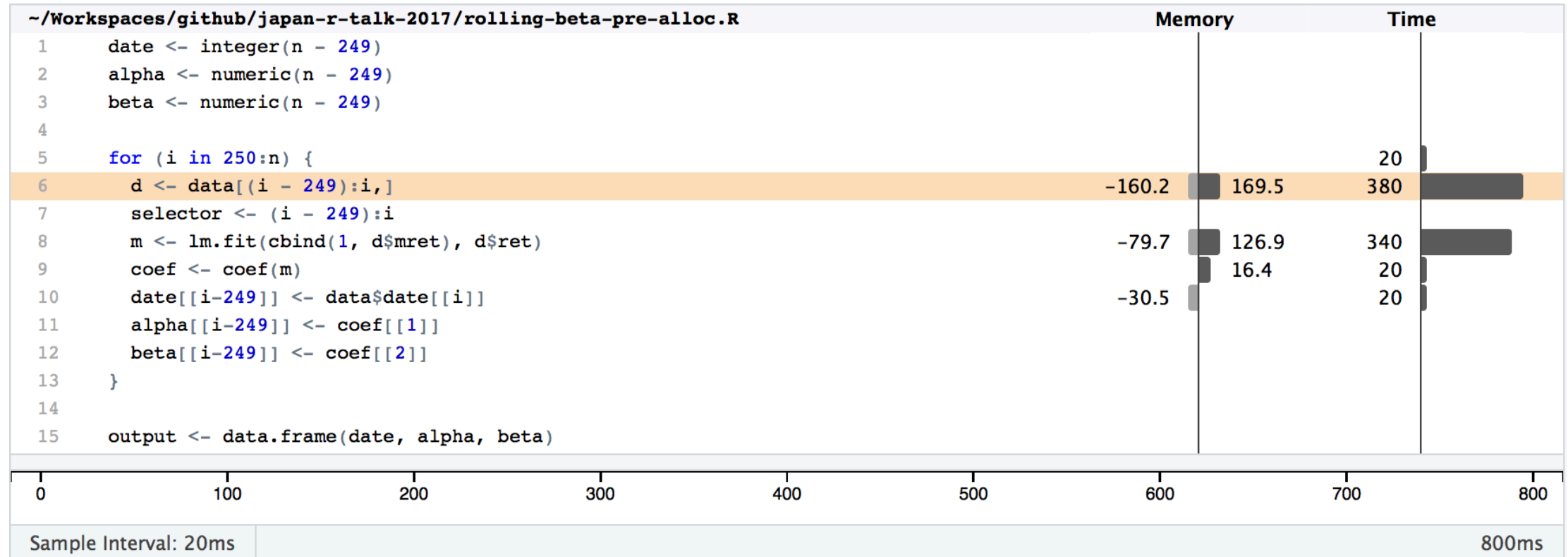
# Pre-allocate vectors

```
1 date <- integer(n - 249)
2 alpha <- numeric(n - 249)
3 beta <- numeric(n - 249)
4
5 for (i in 250:n) {
6   d <- data[(i - 249):i,]
7   selector <- (i - 249):i
8   m <- lm.fit(cbind(1, d$mret), d$ret)
9   coef <- coef(m)
10  date[[i-249]] <- data$date[[i]]
11  alpha[[i-249]] <- coef[[1]]
12  beta[[i-249]] <- coef[[2]]
13 }
14
15 output <- data.frame(date, alpha, beta)
```

Pre-allocation of vectors

Avoid copying vectors

# Pre-allocate vectors (800ms)



# Avoid row selection from `data.frame` (440ms)

~/Workspaces/github/japan-r-talk-2017/rolling-beta-vector-selector.R		Memory		Time
1	date <- integer(n - 249)			
2	alpha <- numeric(n - 249)			
3	beta <- numeric(n - 249)			
4				
5	for (i in 250:n) {			20
6	selector <- (i - 249):i			
7	m <- lm.fit(cbind(1, data\$mret[selector]), data\$ret[selector])	-90.5	97.1	360
8	coef <- coef(m)		4.2	20
9	date[[i-249]] <- data\$date[[i]]	-33.1	12.4	40
10	alpha[[i-249]] <- coef[[1]]			
11	beta[[i-249]] <- coef[[2]]			
12	}			
13				
14	output <- data.frame(date, alpha, beta)			

source

050100150200250300350400450

Sample Interval: 20ms

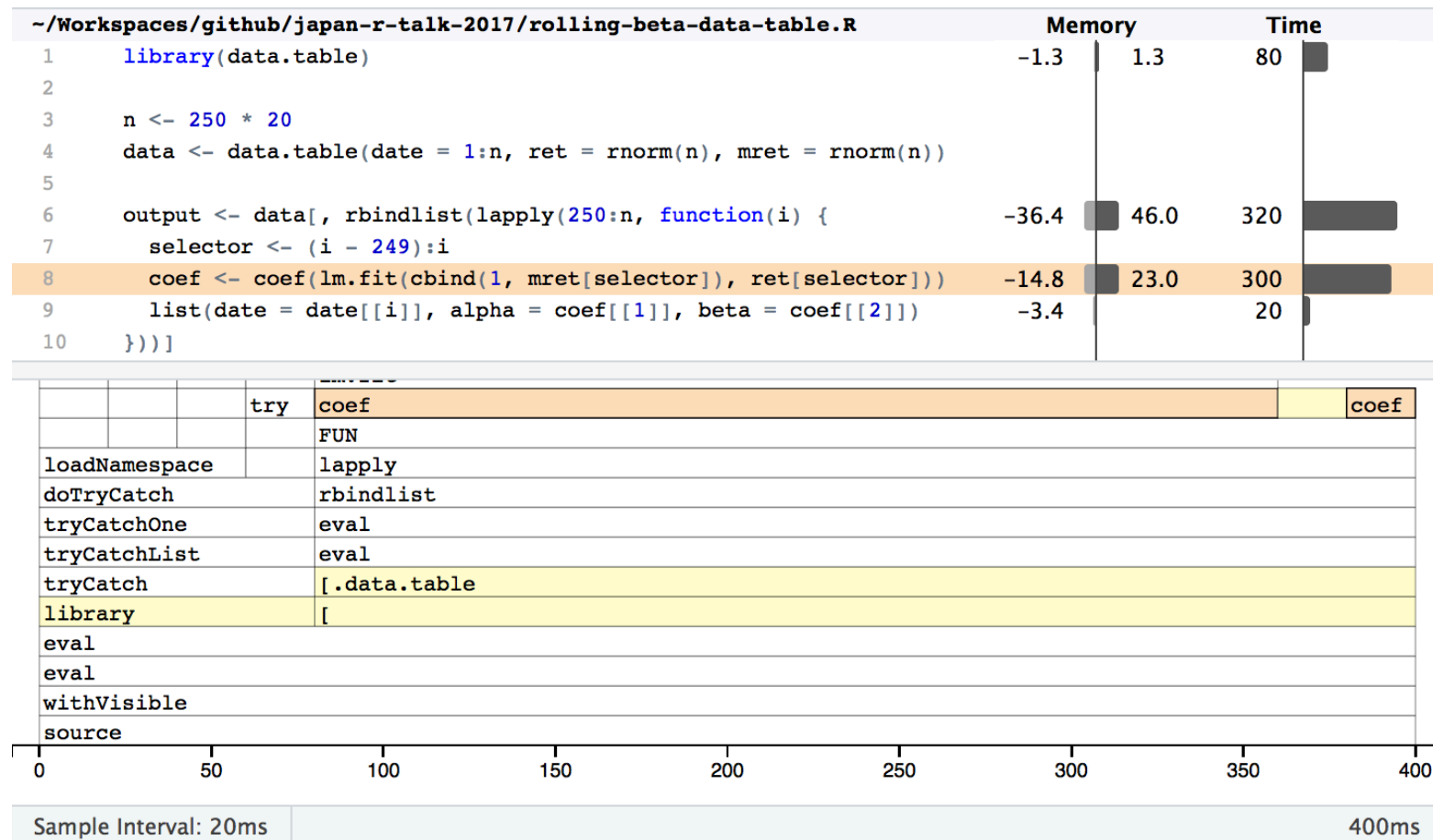
460ms



# Rewrite with `data.table`

```
1 library(data.table)
2
3 n <- 250 * 20
4 data <- data.table(date = 1:n, ret = rnorm(n), mret = rnorm(n))
5
6 output <- data[, rbindlist(lapply(250:n, function(i) {
7   selector <- (i - 249):i
8   coef <- coef(lm.fit(cbind(1, mret[selector]), ret[selector]))
9   list(date = date[[i]], alpha = coef[[1]], beta = coef[[2]])
10 })))
```

# Rewrite with `data.table` (320ms)



# Multi-symbol rolling regression

symbol	date	ret	mret
1	1	-0.560475647	-0.079625050
1	2	-0.230177489	0.045941483
1	3	1.558708314	0.126138891
1	4	0.070508391	-0.001914877
1	5	0.129287735	-0.058701079
1	6	1.715064987	-0.014730471
1	7	0.460916206	0.014168433
1	8	-1.265061235	0.028467197
1	9	-0.686852852	0.083210147
1	10	-0.445661970	-0.190119626
1	11	1.224081797	0.033166720
1	12	0.359813827	0.043090255
1	13	0.400771451	0.092460001
1	14	0.110682716	-0.112767268
1	15	-0.555841135	-0.217610923
1	16	1.786913137	0.096349661
1	17	0.497850478	-0.052769189

```
1 library(data.table)
2 set.seed(123)
3 dt <- expand.grid(symbol = 1:100, date = 1:(250 * 20))
4 setDT(dt, key = c("symbol", "date"))
5 dt[, ret := rnorm(.N)]
6 dt[, mret := mean(ret), by = date]
```

- 100 symbols, 20 years of history
- setDT to transform dt as data.table
- (symbol, date) as key to ensure order
- For each symbol and each date, estimate the linear coefficients for the recent 250 days of that symbol

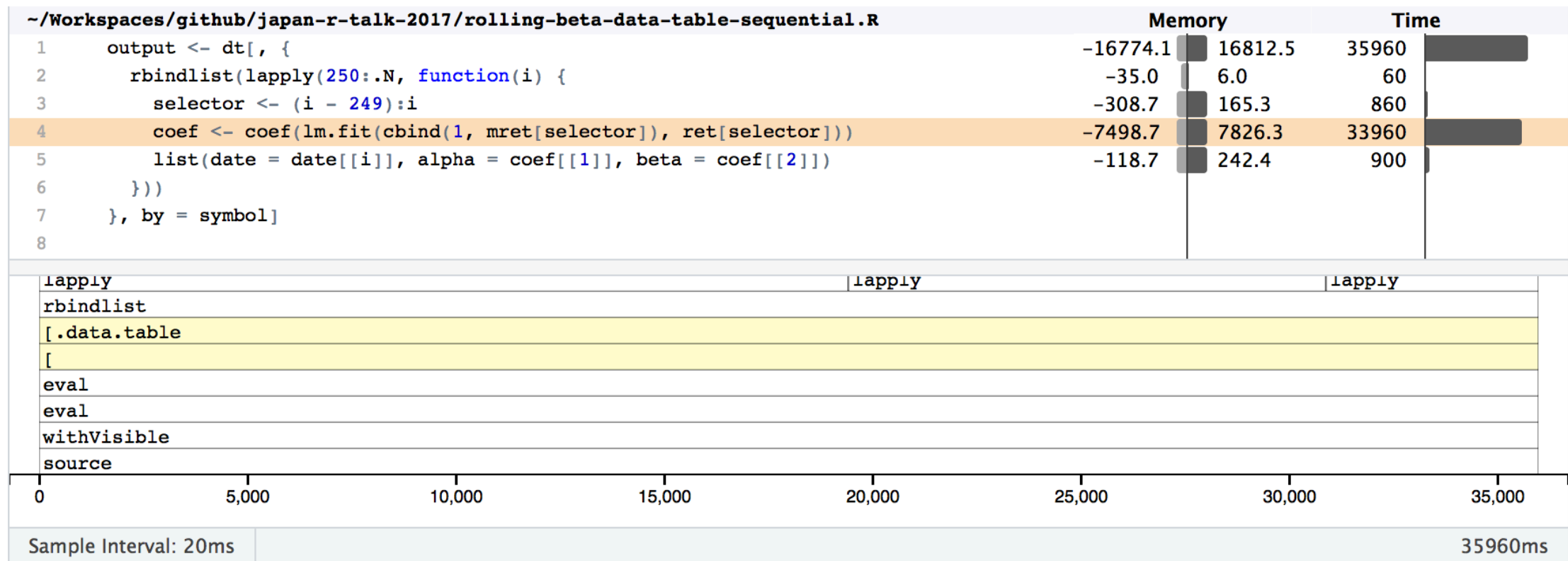
# The data.table solution

```
1 output <- dt[, {  
2   rbindlist(lapply(250:.N, function(i) {  
3     selector <- (i - 249):i  
4     coef <- coef(lm.fit(cbind(1, mret[selector]), ret[selector]))  
5     list(date = date[[i]], alpha = coef[[1]], beta = coef[[2]])  
6   })))  
7 }, by = symbol]
```

Rolling by each symbol  
Combine with rbindlist

Use by for group operation

# The data.table solution (35960ms)



# mcLapply + data.table solution (8640ms)

```
1 library(parallel)
2 symbols <- dt[, sort(unique(symbol))]
3 output <- mclapply(symbols, function(sym) {
4   dt[symbol == sym, {
5     rbindlist(lapply(250:.N, function(i) {
6       selector <- (i - 249):i
7       coef <- coef(lm.fit(cbind(1, mret[selector]), ret[selector]))
8       list(date = date[[i]], alpha = coef[[1]], beta = coef[[2]])
9     })))
10   }, by = symbol]
11 }, mc.cores = 8)
12 output <- rbindlist(output)
13 setkey(output, symbol, date)
```

Automatically use binary  
search on symbol

user	system	elapsed
52.637	2.097	8.640

# Solutions to R performance issues

- Built-in or lower-level functions (`lm` -> `lm.fit`)
- Vectorization (`cumsum`, `cumprod`, ...)
- Matrix algebra (`%*%`, `solve`, `crossprod`, ...)
- High performance packages (`data.table`, `fst`, `RcppRoll`, ...)
- Rcpp and Rcpp template packages (`RcppArmadillo`, `RcppEigen`, etc.)


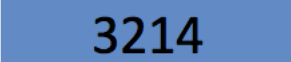


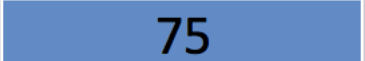

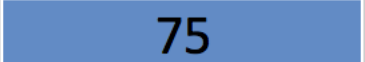

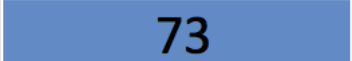

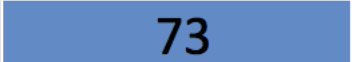

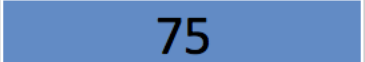

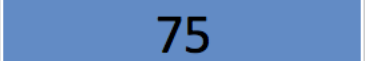





# Package for performance: data.table

- Algorithm: Fast merge based on binary search
- OpenMP-powered multi-threaded subsetting
- Control copy behavior: In-place modification (`:=`)
- High performance helper functions
  - `setnames`, `setcolorder`, `setkey`, `setindex`, ...
  - `rbindlist`, `dcast`, `melt`, `frank`, ...
  - `%between%`, `%chin%`, ...
  - `fread`, `fwrite`
- The data.table cheatsheet by DataCamp



# Package for performance: fst

<https://github.com/fstpackage/fst>: Lightning Fast Serialization of Data Frames for R  
Chinese Stock Market 1-minute data (2014-2017): 542M+ rows, 23 columns

Function	Compression	File size (GB)		Time (s)
saveRDS	Yes		22	
readRDS	Yes		22	
saveRDS	No		75	
readRDS	No		75	
data.table::fwrite	No		73	
data.table::fread	No		73	
fst::write_fst	No		75	
fst::read_fst	No		75	
fst::write_fst	100 (40 threads)		25	
fst::read_fst	100 (40 threads)		25	

# Parallel computing

## Multi-threading



☐ OpenBLAS

☐ Intel MKL

☐ OpenMP

Drop-in solutions

## Multi-processing



☐ Fork-based (Linux, macOS)

☐ Socket-based (Windows)

Platform dependent

## Multi-node



☐ Hadoop

☐ Spark

☐ etc.

# Use parallel: fork vs. cluster

- Platform dependent:
  - Windows: socket-based master-slave cluster
  - Linux and macOS: mcp parallel and mc collect
- Inter-process communication: serialization: slow and high memory peak
- Multi-threading techniques do not work well with fork-based parallelism.
- Use RhpcBLASctl to force number of threads in OpenMP and BLAS to avoid deadlock in fork process

# Rcpp

- R data structures + C++ data structures
- Multi-threading by OpenMP
- Template packages: RcppArmadillo, RcppEigen, etc.
- Rcpp now used by 1000+ packages

# data.table + Rcpp + RcppArmadillo (3700ms)

C++

```
1 // [[Rcpp::depends(RcppArmadillo)]]
2 #include <RcppArmadillo.h>
3 #include <Rcpp.h>
4 using namespace Rcpp;
5
6 // [[Rcpp::export]]
7 List roll_lm(const IntegerVector& date, const arma::mat& x, const arma::colvec& y, int n) {
8   int size = date.size();
9   IntegerVector _date(size - n + 1);
10  NumericVector alpha(size - n + 1), beta(size - n + 1);
11  for (int i = n - 1; i < size; i++) {
12    const arma::mat& xi = x.rows(i - n + 1, i);
13    const arma::colvec& yi = y.subvec(i - n + 1, i);
14    const arma::colvec& coef = solve(xi, yi);
15    _date[i - n + 1] = date[i];
16    alpha[i - n + 1] = coef[0];
17    beta[i - n + 1] = coef[1];
18  }
19  return List::create(_["date"] = _date, _["alpha"] = alpha, _["beta"] = beta);
20 }
```

R

```
10 output <- dt[, roll_lm(date, cbind(1, mret), ret, 250), by = symbol]
```

user	system	elapsed
3.669	0.029	3.700

# data.table + Rcpp + RcppArmadillo + OpenMP (840ms)

C++

```
1 // [[Rcpp::depends(RcppArmadillo)]]
2 // [[Rcpp::plugins(openmp)]]
3 #include <omp.h>
4 #include <RcppArmadillo.h>
5 #include <Rcpp.h>
6 using namespace Rcpp;
7
8 // [[Rcpp::export]]
9 List par_roll_lm(const IntegerVector& date, const arma::mat& x, const arma::colvec& y, int n, int threads) {
10   int size = y.size();
11   IntegerVector _date(size - n + 1);
12   NumericVector alpha(size - n + 1), beta(size - n + 1);
13   #pragma omp parallel for shared(date, x, y), num_threads(threads)
14   for (int i = n - 1; i < size; i++) {
15     const arma::mat& xi = x.rows(i - n + 1, i);
16     const arma::colvec& yi = y.subvec(i - n + 1, i);
17     const arma::colvec& coef = solve(xi, yi);
18     INTEGER(_date)[i - n + 1] = date[i];
19     REAL(alpha)[i - n + 1] = coef[0];
20     REAL(beta)[i - n + 1] = coef[1];
21   }
22   return List::create(_["date"] = _date, _["alpha"] = alpha, _["beta"] = beta);
23 }
```

Use OpenMP header file  
and update compiler flags

OpenMP transforms the  
for loop into parallel

Use pointer to update R  
vectors

R

```
14 output <- dt[, par_roll_lm(date, cbind(1, mret), ret, 250, threads = 8), by = symbol]
```

user	system	elapsed
6.375	0.065	0.840

# Performance summary

Solution	Time	Performance
<b>Initial solution (plain for loop)</b>	6880ms	1x
Replace lm with lm.fit	3180ms	2x
Pre-allocate vectors	800ms	9x
Avoid row selection from data.frame	440ms	16x
Rewrite with data.table	320ms	21x
<b>Increase to 100 symbols</b>	35960ms	21x
mclapply + data.table	8640ms	87x (mc.core = 8)
data.table + Rcpp + RcppArmadillo	3700ms	204x
data.table + Rcpp + RcppArmadillo + OpenMP	840ms	900x (threads = 8)

- Kun Ren <[renkun@outlook.com](mailto:renkun@outlook.com)>
- GitHub: renkun-ken
- Twitter: @renkun\_ken
- Slides and code: <https://github.com/renkun-ken/japan-r-talk-2017>