

UNIVERSITY OF CALIFORNIA, DAVIS

ECS 158

PROGRAMMING ON PARALLEL ARCHITECTURES

Parallelization of R's `combn()` Function

Authors:

Trisha FUNTANILLA
Syeda INAMDAR
Eva LI
Jennifer WONG

Professor:

Norm MATLOFF

March 19, 2015

Contents

1	Objective	1
2	Function	1
2.1	Description	1
2.2	Why Speedups Are Possible	2
3	Parallelization	2
3.1	OpenMP	2
3.1.1	Load Balancing	2
3.1.2	Chunk Size	3
3.1.3	Scheduling Policies	3
3.1.4	Cache and Other Performance Tuning Considerations	3
3.1.5	Comparative Analysis	4
4	Timing Comparisons	5
5	Conclusions	5
A	Appendix	6
	References	12

1 Objective

Choose some function `g()` in R, either built-in to base R or in a CRAN package, to parallelize. Do so in Snow, OpenMP and CUDA (or substitute Thrust for either CUDA or OpenMP but not both), and run timing tests.

2 Function

We have chose to parallelize the function `combn()` from the “combinat” package in CRAN.

2.1 Description

Taken from the official documentation:[2]

“Generate all combinations of the elements of `x` taken `m` at a time. If `x` is a positive integer, returns all combinations of the elements of `seq(x)` taken `m` at a time. If argument “`fun`” is not null, applies a function given by the argument to each point. If `simplify` is `FALSE`, returns a list; else returns a vector or an array. “...” are passed unchanged to function given by argument `fun`, if any.”

2.2 Why Speedups Are Possible

We have tried `combn()`, which has the inputs `combn(x, m, fun=NULL, simplify=TRUE, ...)`, with different inputs for the vector `x` and value `m`. We noticed that `combn()` became quite slow once `x` had a size numbering around 100 or higher and `m` around 5 or higher, taking several minutes to compute (but not print) its results. We think that we can achieve speedups by dividing the input `x` equally among the threads/nodes, and parallelizing the work of creating the output matrix or list. The goal is to speed up the function, which, based on its source code, currently computes combinations serially.

3 Parallelization

3.1 OpenMP

The C++ OpenMP program is called through R via `.Call()` using the Rcpp interface. It accepts the same input arguments as the original function, can handle the same types of valid inputs, and through new optional arguments, also allows the user to decide on a scheduling policy and chunk size.

Code highlights:

- Load balancing algorithm [1] (next section)
- Avoidance of critical sections and barriers

Other notes:

- Similar to the original function, the OpenMP implementation has a limit on the input/output size. The program terminates if R decides that it cannot allocate enough space for the output.

3.1.1 Load Balancing

Good load balancing is very critical to our program because for each element i in the input vector x , the number of combinations that can be generated with x_i is larger than that with x_{i+1} . A naive task distribution without proper load balancing could heavily skew the distribution of the workload for each thread. For example, directly assigning an x_i for each thread will result in the threads with lower id's generating far more combinations (thus, more work) than threads with higher id's).

In order to compensate for the load imbalance, the assignment of tasks to the threads was done in a wrap around manner. Using this method, for the first distribution, each thread will work on the element indexed in the same value as their id (e.g. thread 0 on $x[0]$, thread 1 on $x[1]$, and so on). Then, the following distribution depends on the total number of threads, the id of each thread, and the index of its previous assignment. With nth corresponding to the total number of threads and me corresponding to the thread id, the assignment following the initial distribution is determined using the equation:

$$next_task = prev_task + 2 * (nth - me - 1) + 1 \quad (1)$$

The distribution that will then follow the above depends on the id of each thread the index of the previous task. This next distribution is determined by the equation:

$$next_task = prev_task + 2 * me + 1 \quad (2)$$

The succeeding distributions alternates in using the two equations until the last distribution, which is for $n - m + 1$ or the last x_i that can form a combination of size m with the elements succeeding it.

In this algorithm, each thread knows which indexes in x to generate combinations for. The threads do not need to communicate with each other, avoiding huge overhead especially for large values of n .

3.1.2 Chunk Size

The equations from the section above implicitly defines the chunk size so that each thread gets assigned roughly the same number x_i 's to work on. It is essentially the number of times task distributions occur in the load balance algorithm, which is roughly $(n - m + 1)/nth$.

3.1.3 Scheduling Policies

The load balancing algorithm described in Section 3.1.1 implements a static scheduling policy since the threads are assigned specific tasks and only work on the tasks assigned to them. For the purposes of timing comparisons, additional optional arguments `sched` and `chunksize` were added to the function in order to test the speed of the program for the

other scheduling policies, namely dynamic and guided. If both arguments are not provided (or set as NULL), the program defaults to the static scheduling policy using the load balancing algorithm. If either dynamic or guided is set, the program uses OpenMP's built-in `schedule` clause to distribute the tasks to the threads. If the chunksize is not provided, the program sets it to the default chunk size of 1.

We performed various tests using the three scheduling policies and experimented with different chunk sizes (for dynamic and guided). It was clear that the load balancing algorithm using static scheduling was the fastest. It eliminates the communication overhead of the dynamic and guided scheduling policies since the threads don't have to repeatedly access the task farm for work, and idleness is still reduced because of the way the tasks are distributed among the threads.

Hence, the static method is the one compared to the original function in the timing comparisons.

3.1.4 Cache and Other Performance Tuning Considerations

The following performance tuning applies only to the load balancing algorithm described in 3.1.1.

- Avoiding instances of false sharing: In order to avoid instances of false sharing (and the resulting cache coherence overheads), we avoided having shared data that could be modified by multiple threads. The data frequently accessed and shared by the threads in the program, such as the input and position vectors are read-only data. Hence, cache lines are not invalidated whenever data from these vectors are accessed. The main shared data structure that is modified by all threads is the output matrix `retmat`, but the threads never have to perform any read on this data.
- Synchronization overhead: The threads are made to work as independent of each other as possible. The only synchronization occurs implicitly at the end of the `parallel` pragma.
- Avoidance of critical sections and barriers

3.1.5 Comparative Analysis

The following test inputs were used for comparing the speeds of the source code and our OpenMP implementation using 8 threads.

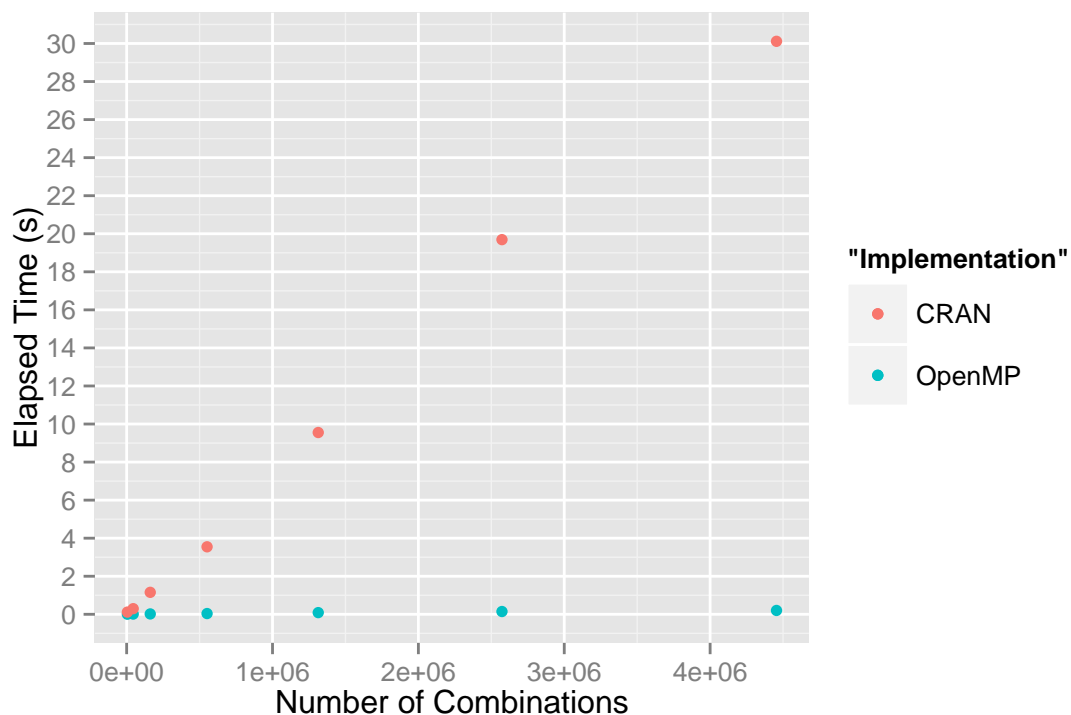
```
> system.time(combn(c(1:100), 2))
> system.time(combn(c(1:300), 2))
> system.time(combn(c(1:100), 3))
> system.time(combn(c(1:150), 3))
> system.time(combn(c(1:200), 3))
> system.time(combn(c(1:300), 3))
> system.time(combn(c(1:250), 3))
```

Here are the total number of combinations generated for each test above:

```
> nCm(100, 2, 0.100000000000000002)
[1] 4950
> nCm(300, 2, 0.100000000000000002)
[1] 44850
> nCm(100, 3, 0.100000000000000002)
[1] 161700
> nCm(150, 3, 0.100000000000000002)
[1] 551300
> nCm(200, 3, 0.100000000000000002)
[1] 1313400
> nCm(300, 3, 0.100000000000000002)
[1] 4455100
> nCm(250, 3, 0.100000000000000002)
[1] 2573000
```

The elapsed time recorded for each test is the average of three trials. The OpenMP parallelization massively improved the performance of the function. The improvement becomes increasingly invaluable as the input size increases (e.g. `nCm(100, 5)`, the OpenMP implementation still manages to produce an output under 5 seconds, while the original function takes several minutes).

The following plot illustrates the differences in the speeds of the two implementations.



4 Timing Comparisons

5 Conclusions

A Appendix

```
1 #####
2 # R call function for the OpenMP parallelization of combn() from the CRAN
3 # combinat package: http://cran.r-project.org/web/packages/combinat/index.html
4
5 # NOTE: Output is out of order.
6
7 # Function Arguments:
8 # x <- input vector of integers and/or characters\
9 # m <- number of elements in a combination
10 # fun <- function to apply to the resulting output
11 # simplify <- if TRUE print output as a matrix with m rows and nCm columns
12 # where nCm is the total number of combinations generated
13 # ... <- parameters for fun
14
15 # Helper functions for handling characters in input vector x:
16 # is.letter <- # function to check if there's a char in x
17 # asc <- convert char to ASCII decimal value
18 # chr <- convert decimal value to ASCII character
19 # nCm <- calculating the total number of combinations
20 # (taken directly from R combinat package)
21 # inserted into this file since combinat could not be installed in CSIF
22 #####
23
24 combn <- function(x, m, fun = NULL, simplify = TRUE,
25   sched = NULL, chunksize = NULL, ...)
26 {
27   require(Rcpp)
28   dyn.load("combn-final.so")
29
30   # Input checks taken directly from combn source code
31   if(length(m) > 1) {
32     warning(paste("Argument m has", length(m),
33       "elements: only the first used"))
34     m <- m[1]
35   }
36   if(m < 0)
37     stop("m < 0")
38   if(m == 0)
39     return(if(simplify) vector(mode(x), 0) else list())
40   if(is.numeric(x) && length(x) == 1 && x > 0 && trunc(x) == x)
41     x <- seq(x)
42   n <- length(x)
43   if(n < m)
44     stop("n < m")
45
46   nofun <- is.null(fun)
47   count <- nCm(n, m, 0.100000000000000002)
48
49   # Error checks for the scheduling variables: sched and chunksize
50   # R handles the error when 'sched' is not a string/character vector
51
52   # If sched is provided, then sched must be static, dynamic, guided, or NULL
53   if (!grepl('static', sched) && !grepl('dynamic', sched) && !grepl('guided', sched) &&
54     !is.null(sched)) {
55     stop("Scheduling policy must be static, dynamic, or guided.")
56   }
57   # Set to default values depending on what is/are provided
58   if (is.null(sched) && is.null(chunksize)) {
59     sched <- 'static'
60     chunksize <- 1
61   }
62   else if (!is.null(sched) && is.null(chunksize)) {
63     chunksize <- 1
64   }
65   else if (is.null(sched) && !is.null(chunksize)) { # if sched is provided, but chunk
66     size is not
67     sched <- 'static'
```



```

66     warning("'sched' is replaced with default 'static' and 'chunksize' is overridden with
        default value.")
67 }
68
69 # Checks if input vector x has characters
70 # If so, then convert chars to their ASCII decimal values
71 # Operate on the ASCII decimal values for the chars
72 ischarx <- match('TRUE', is.letter(x))
73 if (!is.na(ischarx)) {
74     ischarx_arr <- is.letter(x)
75     for (i in 1:length(charx)) {
76         if (ischarx_arr[i]) {
77             if (length(asc(x[i])) == 1) {
78                 x[i] <- asc(x[i])
79             }
80             else {
81                 x[i] <- as.character(x[i])
82             }
83         }
84     }
85 }
86 x <- strtoi(x, base=10)
87 }
88
89
90 #Calculate positions for output
91 pos <- vector()
92 temp_n <- n
93 for (i in 1:(n-m+1)) {
94     pos <- c(pos, nCm(temp_n-i, m-1))
95 }
96 temp <- pos[1]
97 pos[1] <- 0
98 for (i in 2:length(pos)) {
99     temp2 <- pos[i]
100    pos[i] <- temp
101    temp <- pos[i] + temp2
102 }
103
104 # Initialize output matrix
105 retmat <- matrix(0, m, count)
106 # Call the function through Rcpp
107 retmat <- .Call("combn", x, m, n, count, sched, chunksize, pos)
108
109 # Convert from ASCII decimal values back to chars if necessary
110 if (!is.na(ischarx)) {
111     for (i in 1:length(retmat)) {
112         if ((as.integer(retmat[i]) >= 97 && as.integer(retmat[i]) <= 122)
113             || (as.integer(retmat[i]) >= 65 && as.integer(retmat[i]) <= 90)) {
114             retmat[i] <- chr(retmat[i]);
115         }
116     }
117 }
118
119 # Apply provided function to the output
120 if (!is.null(fun)) {
121     apply(retmat, 2, fun(...))
122 }
123
124 # Format results
125 if (simplify) {
126     out <- retmat
127 }
128 else {
129     out <- list()
130     for (i in 1:ncol(retmat)) {
131         out <- c(out, list(c(retmat[, i])))
132     }
133 }
134 return(out)

```

```

135 }
136
137 # function to check if there's a char in x
138 is.letter <- function(x) grepl("[[:alpha:]]", x)
139 # convert char to ascii decimal value
140 asc <- function(x) { strtoi(charToRaw(x),16L) }
141 # convert decimal value to ascii character
142 chr <- function(n) { rawToChar(asc.raw(n)) }
143
144 # n choose m - calculates the total number of combinations for a given input
145 "nCm"<-
146 function(n, m, tol = 9.9999999999999984e-009)
147 {
148 #   DATE WRITTEN:  7 June 1995                LAST REVISED:  10 July 1995
149 #   AUTHOR:  Scott Chasalow
150 #
151 #   DESCRIPTION:
152 #       Compute the binomial coefficient ("n choose m"), where n is any
153 #       real number and m is any integer. Arguments n and m may be vectors;
154 #       they will be replicated as necessary to have the same length.
155 #
156 #       Argument tol controls rounding of results to integers. If the
157 #       difference between a value and its nearest integer is less than tol,
158 #       the value returned will be rounded to its nearest integer. To turn
159 #       off rounding, use tol = 0. Values of tol greater than the default
160 #       should be used only with great caution, unless you are certain only
161 #       integer values should be returned.
162 #
163 #   REFERENCE:
164 #       Feller (1968) An Introduction to Probability Theory and Its
165 #       Applications, Volume I, 3rd Edition, pp 50, 63.
166 #
167   len <- max(length(n), length(m))
168   out <- numeric(len)
169   n <- rep(n, length = len)
170   m <- rep(m, length = len)
171   mint <- (trunc(m) == m)
172   out[!mint] <- NA
173   out[m == 0] <- 1 # out[mint & (m < 0 | (m > 0 & n == 0))] <- 0
174   whichm <- (mint & m > 0)
175   whichn <- (n < 0)
176   which <- (whichm & whichn)
177   if(any(which)) {
178     nnow <- n[which]
179     mnow <- m[which]
180     out[which] <- ((-1)^mnow) * Recall(mnow - nnow - 1, mnow)
181   }
182   whichn <- (n > 0)
183   nint <- (trunc(n) == n)
184   which <- (whichm & whichn & !nint & n < m)
185   if(any(which)) {
186     nnow <- n[which]
187     mnow <- m[which]
188     foo <- function(j, nn, mm)
189     {
190       n <- nn[j]
191       m <- mm[j]
192       iseq <- seq(n - m + 1, n)
193       negs <- sum(iseq < 0)
194       ((-1)^negs) * exp(sum(log(abs(iseq))) - lgamma(m + 1))
195     }
196     out[which] <- unlist(lapply(seq(along = nnow), foo, nn = nnow,
197       mm = mnow))
198   }
199   which <- (whichm & whichn & n >= m)
200   nnow <- n[which]
201   mnow <- m[which]
202   out[which] <- exp(lgamma(nnow + 1) - lgamma(mnow + 1) - lgamma(nnow -
203     mnow + 1))
204   nna <- !is.na(out)

```

```

205 outnow <- out[nna]
206 rout <- round(outnow)
207 smalldif <- abs(rout - outnow) < tol
208 outnow[smalldif] <- rout[smalldif]
209 out[nna] <- outnow
210 out
211 }

```

```

1  /*****
2  OpenMP (C++) implementation of R's combn() function from the combinat package
3
4  Called from R using .Call1() through Rcpp
5  *****/
6
7  #include <Rcpp.h>
8  #include <omp.h>
9
10 using namespace std;
11 using namespace Rcpp;
12
13 // Computes the indices of the next combination to generate
14 // The indices then get mapped to the actual values from the input vector
15 int next_comb(int *comb, int m, int n)
16 {
17     int i = m - 1;
18     ++comb[i];
19
20     while ((i >= 0) && (comb[i] >= n - m + 1 + i)) {
21         --i;
22         ++comb[i];
23     }
24
25     if (comb[0] == 1) {
26         return 0;
27     }
28
29     for (i = i + 1; i < m; ++i) {
30         comb[i] = comb[i - 1] + 1;
31     }
32
33     return 1;
34 }
35
36 RcppExport SEXP combn(SEXP x_, SEXP m_, SEXP n_, SEXP nCm_, SEXP sched_, SEXP chunksize_,
37                       , SEXP pos_, SEXP out)
38 {
39     // Convert SEXP variables to appropriate C++ types
40     NumericVector x(x_); // input vector
41     NumericVector pos(pos_); // position vector for the combinations so that the output is
42     // sorted
43     int m = as<int>(m_), n = as<int>(n_), nCm = as<int>(nCm_), chunksize = as<int>(
44     chunksize_);
45     string sched = as<string>(sched_);
46
47     NumericMatrix retmat(m, nCm);
48
49     // OpenMP schedule clauses
50     if (sched == "dynamic") {
51         omp_set_schedule(omp_sched_dynamic, chunksize);
52     }
53     else if (sched == "guided") {
54         omp_set_schedule(omp_sched_guided, chunksize);
55     }
56
57     if (sched == "static") { // use the load balancing algorithm
58         #pragma omp parallel
59         {
60             // this thread id, total number of threads, combination indexes array
61             int me, nth, *comb;

```

```

59     nth = omp_get_num_threads();
60     me = omp_get_thread_num();
61
62     // array that will hold all of the possible combinations
63     // of size m of the indexes
64     comb = new int[m];
65
66     // initialize comb array
67     for (int i = 0; i < m; ++i) {
68         comb[i] = i;
69     }
70
71
72     int chunkNum = 1; // the number of chunk that has been distributed
73     int mypos; // variable for the output position
74
75     // each thread gets assign a chunk to work on
76     // each thread will have about the same number of chunks
77     // to work on throughout the lifetime of the program
78     for(int current_x = me; current_x < n-m+1; current_x+=1) {
79         int temp;
80         mypos = pos[current_x];
81         for (int i = 0; i < m; ++i) {
82             temp = comb[i] + current_x;
83             retmat(i, mypos) = x[temp];
84         }
85         mypos++;
86         while(next_comb(comb, m, n-current_x)) {
87             int temp;
88             for (int i = 0; i < m; ++i) {
89                 temp = comb[i] + current_x;
90                 retmat(i, mypos) = x[temp];
91             }
92             mypos++;
93         }
94
95         // reset comb array for the next chunk this thread will work on
96         for(int i = 0; i < m; i++) {
97             comb[i] = i;
98         }
99
100         chunkNum++; // increment chunkNum for the next chunk distribution
101         // determine which element this thread will work on
102         if (chunkNum % 2 == 0) {
103             current_x = current_x + 2 * (nth - me - 1);
104         }
105         else {
106             current_x = current_x + 2 * me;
107         }
108     }
109 }
110 }
111 else { // dynamic or guided
112     int mypos;
113     #pragma omp parallel
114     {
115         int *comb = new int[m];
116         for (int i = 0; i < m; ++i) {
117             comb[i] = i;
118         }
119
120         #pragma omp for schedule(runtime)
121         for(int current_x = 0; current_x < (n - m + 1); current_x++) {
122             int temp;
123             mypos = pos[current_x];
124             for (int i = 0; i < m; ++i) {
125                 temp = comb[i] + current_x;
126                 retmat(i, mypos) = x[temp];
127             }
128             mypos++;

```

```
129     while(next_comb(comb, m, n-current_x)) {
130         int temp;
131         for (int i = 0; i < m; ++i) {
132             temp = comb[i] + current_x;
133             retmat(i, mypos) = x[temp];
134         }
135         mypos++;
136     }
137
138     for(int i = 0; i < m; i++) {
139         comb[i] = i;
140     }
141 }
142 }
143 }
144 return retmat;
145 }
```

References

- [1] Junior Barrera, Alfredo Goldman, and Martha Torres
A Parallel Algorithm for Enumerating Combinations
<http://www.ime.usp.br/~gold/ipp03v3.pdf>
- [2] Scott Chasalow
Package ‘combinat’
<http://cran.r-project.org/web/packages/combinat/combinat.pdf>