

CP 331 Assignment 2: Parallel Merge  
Group 2

Hilts, Vaughan  
hilt2740@mylaurier.ca, 120892740

Morouney, Robert  
moro1422@mylaurier.ca, 069001422

Rusu, David  
rusu0260@mylaurier.ca, 131920260

February 2<sup>nd</sup> 2016

# Contents

0.1	Naive Implementation . . . . .	2
0.1.1	Results . . . . .	3
0.2	Virtual Memory Implementation . . . . .	3
0.2.1	Results . . . . .	4
0.3	Binary Search . . . . .	4
0.3.1	Distances are fairly equal between gaps . . . . .	4
0.3.2	When gaps are not equal . . . . .	5
0.4	Regarding heuristics for algorithms . . . . .	5
0.5	Pooled Memory Implementation . . . . .	5
0.5.1	The Code . . . . .	6

## 0.1 Naive Implementation

This was our first approach at solving the problem. The goal was to get something working fast.

generate 2 sorted lists of size  $N_a$ , and  $N_b$  in two blocks of contiguous memory.

Each processor  $p$  calculates it's block of  $A$  to merge with the equations:

$$i_{a,start,p} = p \left\lfloor \frac{N_a}{P} \right\rfloor + \min(p, N_a \bmod P) \quad (1)$$

$$n_{a,p} = \left\lfloor \frac{N_a}{P} \right\rfloor + \begin{cases} 1 & p < N_a \bmod P \\ 0 & \text{otherwise} \end{cases} \quad (2)$$

Where  $i_{a,start,p}$  is the start index of it's block into list  $A$ , and  $n_{a,p}$  is the length of the block.

Next process  $p$  calculates the block of  $B$  to merge with the block in  $A$ .

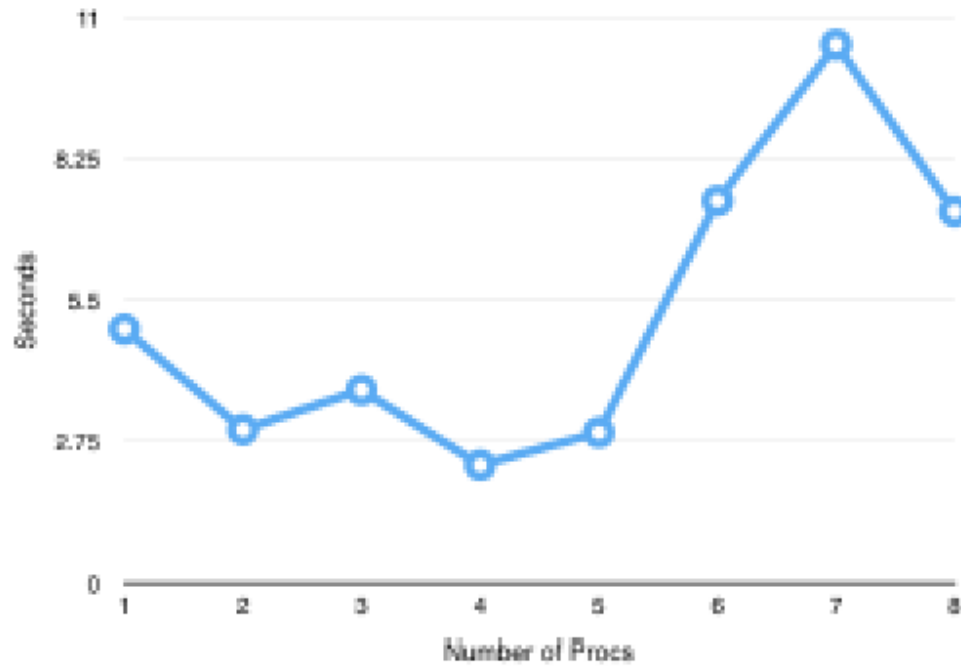
$$i_{b,start,p} = \max\{x \in \mathbb{N} : B_x < A_{i_{a,start,p}-1}\} + 1 \quad (3)$$

$$n_{b,p} = \max\{x \in \mathbb{N} : B_x < A_{i_{a,start,p}+n_{a,p}-1}\} + 1 - i_{b,start,p} \quad (4)$$

We then merge these two lists with a sequential merge

This works well for small lists of sizes smaller than  $2^{32}$  elements but becomes unwieldily with larger lists, since allocating this much contiguous memory will quickly become not practical.

### 0.1.1 Results



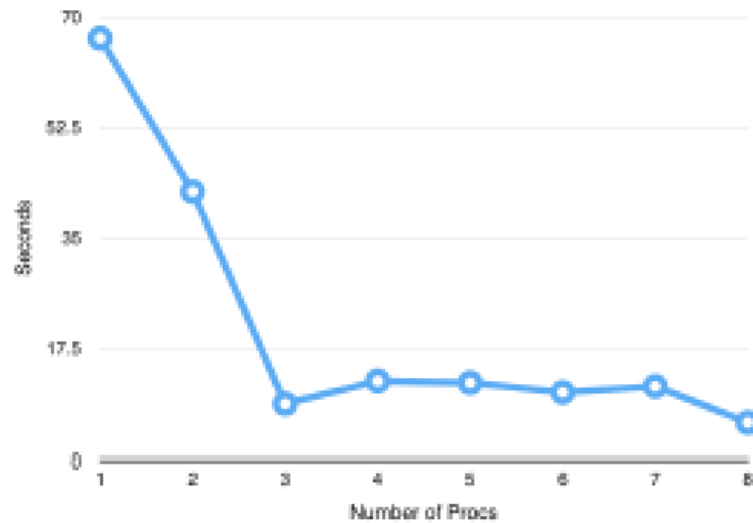
**Figure 1.** Performance on  $10^8$  size lists

## 0.2 Virtual Memory Implementation

To solve the problem of allocating contiguous memory, we instead allocate a few pages of memory and build some abstractions on top of them to give the impression of contiguous memory.

This still allows us to have fast random access into memory but without the requirement of large contiguous blocks of memory. The rest of the algorithm is exactly the same as in our naive approach.

### 0.2.1 Results



**Figure 2.** Performance on  $10^8$  size lists

## 0.3 Binary Search

The linear search in the naive approach of the algorithm is suitable for most array sizes when the set is quite small. However, there are some cases where replacing it with a binary search can be more effective. Specifically, a linear search is used to find the last index of an element that is smaller (when creating partitions of the lists). The results of this are summarized below:

### 0.3.1 Distances are fairly equal between gaps

When the number of partitions and the gaps between the partitions are fairly equal, the binary search performs at an acceptable level. There is a small, but not significant speed up (since only about half the time is spent in this search ? and this is serially, so the less the better). However, this is one of the few cases where the speed up is able to make a difference. When we use Linear Search, the search will remain fairly constant and the CPU cache is able to accurately predict what to store. When we do the same thing with binary search with consistent chunks, we see results that are similar. This would seem reasonable, considering linear search is on the order of  $n$  and binary search is logarithmic in comparison. However, in reality there are some other factors that make this not always true?

### 0.3.2 When gaps are not equal

If your data source gaps are not neatly compacted or the data source has suffered from abnormal randomness, you may have runs where the binary search is significantly slower. Binary search incurs some overhead from its use which is often neglected when studying from a purely Big-O standpoint. There are two things to consider when comparing it with linear search

- In Linear Search, the CPU cache is on your side. The processor has already decided to keep the memory chunks ahead of you in memory and that means access to them will be quick. In contrast, Binary Search does not have this advantage as it may be jumping around quite a few points before settling down on a range.
- There is purely computational overhead in Binary Search. Branching is often neglected and other minor things in Big-O notation which become relevant here.

So, these two things help us draw a conclusion after many benchmarks.

1. If the number of iterations required for linear search is relatively low (10,000 comparisons is a low-figure that is fair from testing), then it can actually be faster to use it instead of Binary Search, despite Binary Search being able to do this theoretically quicker. We have the CPU Cache to blame for this.
2. The other interesting thing to note however, is the implication of a parallelization this would have if it was extended to be parallelized somehow. As the processor count goes up, the effectiveness of a binary search decreases, as the range clearly shrinks. To combat this problem, further heuristics which could select between the two algorithms would be ideal.

## 0.4 Regarding heuristics for algorithms

A good technique for doing this often deployed in High Performance Computing which allows code to be portable across many platforms is pre-benchmarking. In a pre-benchmarking situation, the software will run a set of tests on various input sizes on the test data and measure the time required for each. Each test is an algorithm or combination of algorithms ? which together make up a strategy. Then, the real execution uses the best strategy to execute on the full set of data. Since certain strategies work on different sets of data and hardware differently, having a few different ones and automatic selection is a big boon. This applies with the linear vs. binary situation as well, in these cases they are both strategies.

To maximize gains in these situations, it would be possible to setup a pre benchmark for these two different strategies and at the start of the launch, pick the best one for the current situation

## 0.5 Pooled Memory Implementation

Merging 2 ordered lists is a fairly simple task with contiguous lists as you can see above. The difficulty comes when trying to merge 2 lists so large that the resulting list is larger than the heap size of the system they are accessed on. This problem has implications in almost every field that uses “Big Data” and finding a solution to this problem was a vital step in completing this assignment.

To solve this we implemented a linked list of buckets where each bucket held a slice of lists A and B which will be merged. With this we can control the size of each allocated block of contiguous memory. To take this a step further we aligned memory on a 32byte boundary each time it was allocated. Because we targeted the Intel Xeon processors on the ORCA cluster of SHARCNET we knew that each processor has a 64 byte cache line size. By investigating further we found that 32 bytes of this are reserved for instructions. By aligning allocated memory on a 32 byte boundary and controlling the size of each bucket to make sure it is divisible by 32 we limited potential L1 cache misses.

This allowed us to achieve the maximum speed when generating the lists. Since the list generation is non-trivial and only needed to be done by one processor we decided for the work division to use  $2^n + 1$  cores. This allowed the master processor to spend time generating the lists then sending them out then receiving them and doing the final compilation. Using this scheme we were able to maintain a processor utilization of above 98% at all times during the run.

Since we used a linked list of buckets, dividing the tasks to the remaining processors was a trivial task. Each processor receives buckets in order and splits their lists recursively into individual elements then merges up to a final list. This final list is then sent back to the master processor which outputs the final result to a file. The result is not as fast as the naive solution however it is far more scaleable when using large sets of data.

### 0.5.1 The Code

---

```

1 #include <stdio.h>
2 #include <stdlib.h>
3 #include <stdint.h>
4 #include <mpi.h>
5 #include <time.h>
6
7 //DEFINITIONS
8 #define NUM_ELEMENTS (1<<32) // 2^32
9 #define MASTER 0
10
11 /// If the size of a pointer is 4 bytes (32bit arch)
12 /// then to align the sizeof(bucket) on a 2^n boundry
13 /// we must reserve enough space for the pointer to the
14 /// next bucket. If it is 8 we have to do a bit more work
15
16 #if __SIZEOF_POINTER__ == 4
17 #define BUCKETSIZ (1 << 4)
18 #elif __SIZEOF_POINTER__ == 8
19 #define BUCKETSIZ (1 << 5)
20 #else
21 #error " [ - ] sizeof(void*) is FUCKED"
22 #endif
23
24 //TYPE DEFINITIONS
25 typedef uint32_t u32;
26 typedef uint64_t u64;
27
28 //
29 //EXPLANATION:
30 // u32 = 4 bytes
31 // struct bucket = 4 or 8 bytes

```

```

32 // data[BUCKETSIZE].a = 4 * BUCKETSIZE BYTES
33 // " " .b = " "
34 //
35 // to align on 2^n boundry to maximize storage
36 // u32_ must be added when there is an 8 byte
37 // pointer (64 bit arch) this adds an extra 4
38 // bytes of padding to keep the sturct aligned in RAM
39 //
40
41 typedef struct bucket {
42     struct bucket *next;
43     u32 count;
44 #if __SIZEOF_POINTER__ == 8
45     u32 _;
46 #endif
47     struct {
48         u32 a; // list a
49         u32 b; // list b
50     } data[BUCKETSIZE];
51 } bucket;
52
53 #define newbucket() do { uintptr_t *start; bucket *b; \
54     b = get_mem((u32)sizeof(bucket)); \
55     b->next = master_list; b->count = 0; \
56     master_list = b; \
57     total_buckets++; } while(0)
58
59 void *get_mem(u32 amt);
60 u32 *merge(u32 *sub_list, u32 block_size, u32 *sub_list2, u32 block_size2);
61 void split_and_merge(u32 *sub_list, u32 min, u32 max);
62
63 int main(int argc, char **argv)
64 {
65     if (MPI_Init(&argc,&argv) != MPI_SUCCESS) {
66         printf(" [ - ] Filled to initialize MPI... Exiting");
67         exit(-1);
68     };
69
70     int id, procs;
71     int TAG_A = 1;
72     int TAG_B = 2;
73     int TAG_C = 3;
74     MPI_Status status;
75
76     MPI_Comm_rank(MPI_COMM_WORLD,&id);
77     MPI_Comm_size(MPI_COMM_WORLD,&procs);
78
79     double wall_time1 = MPI_Wtime();
80
81     // count the total numbers generated and placd in buckets as well as the
82     u64 count = 0;
83
84     // total number of buckets which have been created.
85     u32 total_buckets = 0;
86
87     u32 buckets_per_proc = NUM_ELEMENTS/BUCKETSIZE/procs + 1;

```



```

90     if (id == 0) // MASTER
91     {
92         bucket                *master_list;
93         uintptr_t             *start_addr;
94
95         master_list = NULL;
96         newbucket();
97
98         printf(" [ + ] NOTE - sizeof(void*) = %d\n", __SIZEOF_POINTER__);
99         printf(" [ + ] BUCKETSIZE = %'d\n", BUCKETSIZE);
100        printf(" [ ? ] NUM_ELEMENTS = %u\n", NUM_ELEMENTS);
101        srandom(time(NULL));
102
103
104        double start = clock();
105        while(count < NUM_ELEMENTS)
106        {
107            if(master_list->count >= BUCKETSIZE)
108                newbucket();
109
110            u32 i = master_list->count++;
111
112            /// ALL THE REQUIREMENTS WERE FOR RANDOM!!! AHAHAHA
113            u32 a = NUM_ELEMENTS - (count + (random() % 2)),
114                b = NUM_ELEMENTS - count;
115
116            master_list->data[i].a = a;
117            master_list->data[i].b = b;
118
119            count++;
120        }
121
122        double serial_time = clock() - start;
123
124        printf(" [ + ]   Generated %lu sorted random numbers in 2 lists\n", count*2);
125
126        printf(" [ + ]   This took a total of %.5fseconds\n", serial_time/
CLOCKS_PER_SEC);
127
128        printf(" [ + ]   Total number of buckets generated = %u\n", total_buckets );
129        buckets_per_proc = total_buckets/(procs - 1);
130
131        printf(" [ + ]   Total number of buckets per processor = %u\n",
buckets_per_proc);
132
133        /// SEND BUCKETS TO MPI COMM WORLD
134        MPI_Bcast(&buckets_per_proc, 1, MPI_INT, MASTER, MPI_COMM_WORLD);
135
136        u32 p;
137        u32 *temp_mem_a, *temp_mem_b;
138
139        temp_mem_a = malloc(BUCKETSIZE * sizeof(u32));
140        temp_mem_b = malloc(BUCKETSIZE * sizeof(u32));
141        for( p = 1; p<procs; p++)
142        {
143
144            if (temp_mem_a == NULL || temp_mem_b == NULL){

```

```

145         fprintf(stderr, " [ - ] Error getting memory for temp_lists\n");
146         exit(-1);
147     }
148     u32 j;
149
150     for ( j = 0; j < buckets_per_proc; j++ ) {
151         u32 i = 0;
152
153         for ( i = 0; i < BUCKETSIZ; i++){
154             temp_mem_a[i] = master_list->data[i].a;
155             temp_mem_b[i] = master_list->data[i].b;
156
157         }
158
159         // SEND THE LISTS TO PROC P
160         MPI_Send(temp_mem_a,BUCKETSIZ,MPI_INT,p,TAG_A,MPI_COMM_WORLD);
161         MPI_Send(temp_mem_b,BUCKETSIZ,MPI_INT,p,TAG_B,MPI_COMM_WORLD);
162
163         if(master_list->next != NULL)
164             master_list = master_list->next;
165
166     }
167
168     printf(" [ + ] Buckets sent to Proc %d\n", p);
169
170 }
171
172 free(temp_mem_a);
173 free(temp_mem_b);
174 free(master_list);
175
176 } else // SLAVE
177
178 {
179     u32      *list_a , *list_b , *list_c , *list_t ;
180
181     // THIS IS CUTE... IT PAUSES THE SLAVES UNTIL THEY HEAR FROM THEIR MASTER
182     MPI_Bcast(&buckets_per_proc , 1 , MPI_INT, MASTER, MPI_COMM_WORLD);
183
184     u32 min = NUM_ELEMENTS, max = 0;
185
186     list_a = malloc(BUCKETSIZ*buckets_per_proc * sizeof(u32));
187     list_b = malloc(BUCKETSIZ*buckets_per_proc * sizeof(u32));
188     list_c = malloc(2 * buckets_per_proc * BUCKETSIZ * sizeof(u32));
189     u32 i;
190     u32 a=0, b=0;
191     for(i = 0; i < buckets_per_proc; i++)
192     {
193         list_t = malloc(BUCKETSIZ * sizeof(u32));
194         u32 m;
195         MPI_Recv(list_t ,BUCKETSIZ,MPI_INT,MASTER,TAG_A,MPI_COMM_WORLD,&status);
196         split_and_merge(list_t ,0,BUCKETSIZ-1);
197         for(m=0;m<BUCKETSIZ;m++){
198             list_a[a++] = list_t[m];
199         }
200
201         MPI_Recv(list_t ,BUCKETSIZ,MPI_INT,MASTER,TAG_B,MPI_COMM_WORLD,&status);

```

```

202     split_and_merge(list_t, 0, BUCKETSIZ-1);
203     for (m=0; m<BUCKETSIZ; m++){
204         list_b[b++] = list_t[m];
205     }
206
207     free(list_t);
208 }
209
210 ////////////////////////////////////////////////////
211 //
212     list_c = merge(list_a, a, list_b, b);
213 //
214 ////////////////////////////////////////////////////
215 // IN THE MEAN TIME HERE'S A PICTURE OF SEXY A DUCK
216 //
217 //      ,--',-'.
218 //      ,--',-'.
219 //      ;(((---)))
220 //      ;((#) (#)
221 //      | \_/_/_/_/_/|
222 //      " ,--',-'.
223 //      ( ( ,--',-'. )--',-'.
224 //      ,--',-'. \_/_/_/_/_/ ,--',-'.
225 //      ,--',-'. )--',-'.
226 //      ,--',-'. /_/_/_/_/_/ ,--',-'.
227 //      ,--',-'.
228 //      \_/_/_/_/_/
229 //      ,--',-'.
230 //      ,--',-'.
231 //      ,--',-'.
232 //      ,--',-'.
233 //      ,--',-'.
234 //      ,--',-'.
235 //      ,--',-'.
236 //      QUACK!
237 ////////////////////////////////////////////////////
238
239
240 MPI_Send(list_c, a+b, MPI_INT, MASTER, TAG_C, MPI_COMM_WORLD);
241 printf(" [ + ] Processor %d is sending its merged list to MASTER\n", id);
242 free(list_a);
243 free(list_b);
244 free(list_c);
245
246 }
247
248 if (id == 0) // MASTER
249 {
250     FILE* f;
251     f = fopen("c.list", "w");
252
253     u32 p;
254     for (p=1; p<procs; p++)
255     {
256         u32 *list_c;
257         u32 c_size = 2 * buckets_per_proc * BUCKETSIZ;
258         list_c = malloc(c_size * sizeof(u32));

```

```

259 MPI_Recv(list_c, c_size, MPI_INT, p, TAG_C, MPI_COMM_WORLD, &status);
260 printf(" [ + ] MASTER recieved Buckets sent by Proc %d\n", p);
261 u32 i;
262
263 if(list_c == NULL){
264     fprintf(stderr, " [ - ] Error receiving list!\n");
265     exit(-1);
266 }
267
268 printf(" [ + ] PRINTING TO FILE!\n");
269 for ( i = 0; i < c_size; i++){
270     fprintf(f, "%u\n", list_c[i]);
271 }
272 free(list_c);
273
274 }
275 fclose(f);
276 double wall_time2 = MPI_Wtime();
277
278 printf(" [ + ] Program took %.5f seconds to merge %lu elements\n",
279     wall_time2 - wall_time1, (u64) NUM_ELEMENTS*2);
280 printf(" [ + ] Program finished successfully!\n");
281 }
282
283 MPI_Finalize();
284
285 return 0;
286 }
287
288
289 /// get memory aligned to 32byte boundry. Intel Xeon processors on the ORCA cluster
290 /// have a cache-line size of 64bytes. 32 bytes are for program memory and 32 bytes
291 /// are for instructions. By aligning the memory on a 32byte boundry it makes it
292 /// easier for the processor to split and process the chunks of allocated memory.
293 void *get_mem(u32 amt)
294 {
295     uintptr_t mem_start;
296
297     mem_start = (uintptr_t)malloc((size_t)(amt + 31));
298     if((void *)mem_start == NULL)
299     {
300         fprintf(stderr, " [ - ] Failed to Allocated memory in get_mem function after
301             \
302                 allocating %fMB on the heap\n",
303                 (double)total_buckets*sizeof(bucket)/1024.0/1024.0);
304         exit(-1);
305     }
306     /// accidentally had this aligned on a 33 byte boundry and it had a 2 - 3x slow
307     /// down depending on the node
308     mem_start = (mem_start + 31) & ~(uintptr_t)31;
309     return (void *)mem_start;
310 }
311
312 /// merge the two lists while keeping the integrety of the order
313 u32 *merge(u32 *sub_list, u32 block_size, u32 *sub_list2, u32 block_size2)
314 {
315     u32 i1, i2, merged_index;

```

```

316     u32 *merged_list;
317     u32 merged_size = block_size + block_size2;
318
319     i1 = 0;
320     i2 = 0;
321     merged_index = 0;
322
323     merged_list = (u32 *)malloc(merged_size*sizeof(u32));
324
325     while ((i1 < block_size) && (i2 < block_size2)) {
326         if (sub_list[i1] <= sub_list2[i2]) {
327             merged_list[merged_index] = sub_list[i1];
328             merged_index++; i1++;
329         } else {
330             merged_list[merged_index] = sub_list2[i2];
331             merged_index++; i2++;
332         }
333     }
334
335     u32 i;
336     if (i1 >= block_size)
337         for (i = merged_index; i < merged_size; i++, i2++)
338             merged_list[i] = sub_list2[i2];
339     else if (i2 >= block_size2)
340         for (i = merged_index; i < merged_size; i++, i1++)
341             merged_list[i] = sub_list[i1];
342
343     for (i = 0; i < block_size; i++)
344         sub_list[i] = merged_list[i];
345     for (i = 0; i < block_size2; i++)
346         sub_list2[i] = merged_list[block_size+i];
347
348     return merged_list;
349 }
350
351
352 void split_and_merge(u32 *sub_list, u32 min, u32 max)
353 {
354     u32 *_; // stupid recursion
355
356     u32 middle = (min+max)/2,
357         lowerCount = middle - min + 1,
358         upperCount = max - middle;
359
360     // check if sub_list is already sorted
361     if (max == min) {
362         return;
363     } else {
364         //sort first half of sub_list
365         split_and_merge(sub_list, min, middle);
366         //sort second half of sub_list
367         split_and_merge(sub_list, middle+1, max);
368         /* Now merge the two halves */
369         _ = merge(sub_list + min, lowerCount, sub_list + middle + 1, upperCount);
370     }
371 }

```

Listing 1: Parallel Merge