

High Performance Programming - Individual Assignment

Optimizing the quicksort algorithm

Johan Rideg

2020-03-24

1 Preface

For a long time I have wanting to dive deeper into sorting algorithms, especially quicksort and now I have a great excuse to indulge myself. Getting an intuition for why quicksort is so fast and why it even works was not immediately apparent but is something that grew over time as I got deeper into the development.

2 Introduction

Quicksort is a recursive *divide and conquer* algorithm i.e. it works by recursively divide the array into smaller subsets and sort these subsets separately. We will see in detail how the divisions are done which are the very essence of quicksort and why it performs so efficiently. Quicksort is today recognized as one of the fastest sorting algorithms, outperforming other $\mathcal{O}(n \log n)$ algorithms like mergesort if implemented well. It was first developed by Tony Hoare in 1959, and is now used in many standard libraries across multiple programming languages.

Quicksort works by selecting a pivot. The pivot can be any element of the array and there exists many strategies to choose the pivot "in a good way". When a pivot is determined, move all elements less than the pivot to the left of the pivot and elements greater than the pivot to the right of the pivot. This is called *partitioning* and there exists multiple partitioning schemes. When an array is partitioned, it can be split into two arrays where partitioning is performed until we either end up with an array with one element or choose to do something smarter...

All test were performed on a UNIX system using Intel Xeon E5520 @2.27 GHz with the compiler gcc (Ubuntu 7.4.0-1ubuntu1~18.04.1) 7.4.0. The code was compiled using `gcc -g -Wall -O3 -c Quicksort.c` unless stated otherwise.

3 Methology

3.1 Preparatory Work

We must first have some methods for creating, copying and printing arrays which we wish to not be sorted by default. The `generateArray` function allocates memory for an array of size n dynamically and assigns each element an integer, specifically the elements index. This results in an array on the form:

$$\{0, 1, 2, \dots, n - 1\}$$

The array clearly need to be shuffled to be interesting for our purposes. This is done by calling `srand` to set a seed for `rand` function. We then loop over each element, starting from the top and switch its place with a randomly selected element from the left. Decrement until the second element is reached.

The backbone of a sorting algorithm is the `swap` function which takes two numbers and swaps them. We will see this function in many places of the code. In total, the preparatory work resulted in eight simple functions which made the development of quicksort very convenient.

```
1 // Swaps two integers.
2 void swap(int* a, int* b)
3
4 // Shuffles an array using rand().
5 // arr: array.
6 // n: array size.
7 void shuffle(int* arr, const int n)
8
9 // Prints an array to teminal in a column.
10 void printArrVert(int* arr, const int n)
11
12 // Prints an array to terminal horizontally.
13 void printArr(int* arr, const int n)
14
15 // Prints a subset of an arr from lo to hi inclusive.
16 void printSubArr(int* arr, const int lo, const int hi)
17
18 // Generates an ordered array [0, 1, 2, 3, ..., n-1].
19 int* generateArr(const int n)
20
21 // Makes a copy of an array.
22 int* copyArr(int* oldArr, const int n)
23
24 // Verifies that an array is sorted.
25 void verifySorted(int* arr, const int n)
```

3.2 Developing Quicksort

Quicksort is comprised of two main parts: The recursive function and the partition function, which the latter is of particular interest. Being a quicksort novice, I found Lomuto's partitioning scheme easy to both comprehend and to implement. Lomuto's partitioning uses a pivot which is chosen to be the rightmost element of the array. A counter i is used to keep track of the leftmost element greater or equal to the pivot. We then loop over the array from left to right checking if current element is less than the pivot. If it is less than the pivot, swap it with element i and increment i . When the whole array has been iterated over swap the pivot with the leftmost element greater or equal to the pivot. This results in an array where all elements less than the pivot is placed in the left end of the array, followed by the pivot and then all elements greater or equal to the pivot is placed to the right of the pivot. Lastly the index of the pivot is returned and a complete partition has been performed. The result of a partitioning is the following array:

$$\{e_1 < p, \quad e_2 < p, \quad \dots \quad e_{i-1} < p, \quad \textbf{pivot}, \quad e_{i+1} \geq p, \quad e_{i+2} \geq p, \quad \dots \quad e_{n-1} \geq p\}$$

An interesting characteristic of Lomuto's partitioning is that when a partition has been performed, the pivot is at its final index and does not need to be passed on into the recursion.

```
1 // Lomuto partition scheme: Choose last element as pivot.
2 int partitionLomuto(int* arr, const int lo, const int hi)
3 {
4     const int pivot = arr[hi];
5     int i = lo;
6
7     for (int j=lo; j<=hi; j++)
8     {
9         if (arr[j] < pivot)
10        {
11            swap(&arr[i], &arr[j]);
12            i += 1;
13        }
14    }
15
16    swap(&arr[i], &arr[hi]);
17    return i;
18 }
```

Having completed the implementation of the `partitionLomuto` function, the next step is to implement the recursive function `quicksortLomuto`.

```
1 // Quicksort recursion using Lomuto's partitioning scheme.
2 void quicksortLomuto(int* arr, const int lo, const int hi)
3 {
4     if (lo < hi)
5     {
6         int p = partitionLomuto(arr, lo, hi);
7         quicksortLomuto(arr, lo, p-1);
8         quicksortLomuto(arr, p+1, hi);
9     }
10 }
```

We now have a very basic quicksort algorithm on which to apply optimization strategies.

3.3 Optimizing Quicksort

3.3.1 Choosing an Effective Partitioning Scheme

When tasked with optimizing a well known algorithm such as quicksort, it becomes clear that there exists a smorgasbord of optimization choices available. Another partitioning scheme, namely Hoare's, seemed like a good place to start, promising three times less `swap` calls than Lomuto's on average.

We begin by choosing a pivot, for example the middle element. Hoare uses two indices initialized at the ends of the array and moves them towards each other. This is done by letting the leftmost index, i , move right until it finds an element greater than or equal to the pivot. It then waits for the rightmost index, j , to move left until it encounters an element less than or equal to the pivot. When both indices have found the *inversion*, swap the elements and proceed until the two indices meet. At that point the algorithm returns the meeting point where the pivot is now located.

```
1 // Basic hoare partitioning. No median of three strategy is applied.
2 int partitionBasicHoare(int* arr, const int lo, const int hi)
3 {
4     // Choose middle element as pivot.
5     const int pivot = arr[(hi+lo)/2];
6
7     // Initialize indices at the ends of the array.
8     int i = lo - 1;
9     int j = hi + 1;
10
11    // Repeat until indices meet.
12    while (1)
13    {
14        // Increment i until arr[i] exceeds pivot.
15        do
16            i++;
17        while (arr[i] < pivot);
18
19        // Decrement j until pivot exceeds arr[j].
20        do
21            j--;
22        while (arr[j] > pivot);
23
24        // If indices meet, return j.
25        if (i >= j)
26            return j;
27
28        // Inversion has been found. Swap current elements.
29        swap(&arr[i], &arr[j]);
30    }
31 }
```

A new quicksort version is written. Note that when we call `quicksortHoare` on the partition with values less than the pivot, we also include the pivot in that partition. This is because unlike Lomuto's partitioning scheme, the pivot may not be in its final position.

```
1 // Quicksort recursion using Hoares's partitioning scheme.
2 void quicksortBasicHoare(int* arr, const int lo, const int hi)
3 {
4     if (lo < hi)
5     {
6         int p = partitionBasicHoare(arr, lo, hi);
7         quicksortBasicHoare(arr, lo, p);
8         quicksortBasicHoare(arr, p+1, hi);
9     }
10 }
```

Compiling with only `funroll-loops` as optimization flag and comparing the performance against Lomuto's yield a speedup of about 1.3 times which is lower than expected. Even stranger, when compiling with the `-O3`

optimization flag, both methods perform roughly the same. The test was performed five times on identical arrays whose size ranged from 2,000 to 10,000,000.

3.3.2 Quicksort-Insertionsort Hybrid

One drawback quicksort is that it behaves poorly on "small" (<50) arrays compared to simple $\mathcal{O}(N^2)$ algorithms such as the insertionsort. One approach is that when operating on a large array and the recursion reaches a partition of a small size, let insertionsort perform the sorting instead and back out of the current recursion branch. This is what was hinted at the end of the introduction. Let us create an insertionsort function.

```
1 // Fast insertion sorting algorithm. Used for small partitions.
2 // lo: lowest index of the array;
3 // hi: highest index of the array;
4 void insertionSort(int* arr, const int lo, const int hi)
5 {
6     int tmp;
7     int j;
8     for (int i=lo+1; i<hi+1; i++)
9     {
10         tmp = arr[i];
11         j = i-1;
12
13         while(j>=0 && arr[j]>tmp)
14         {
15             arr[j+1] = arr[j];
16             j = j-1;
17         }
18         arr[j+1] = tmp;
19     }
20 }
```

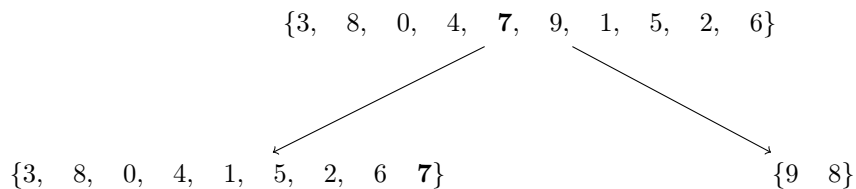
It is important that the `&&`-operator must use short-circuit evaluation which it does in C. If not when we try to evaluate `arr[j]>tmp` with a `j` smaller than 0, we are likely to get an `arrayOutOfBounds` error. With a working insertionsort we can modify our quicksort algorithm.

```
1 // Quicksort recursion using Hoares's partitioning scheme
2 // utilizing insertion sort for small partitions.
3 void quicksortHoare(int* arr, const int lo, const int hi)
4 {
5
6     if (lo < hi)
7     {
8         // If array is small, use insertion sort.
9         if (hi-lo < 50)
10         {
11             insertionSort(arr, lo, hi);
12             return;
13         }
14         int p = partitionHoare(arr, lo, hi);
15         quicksortHoare(arr, lo, p);
16         quicksortHoare(arr, p+1, hi);
17     }
18 }
```

The optimum size for which to switch to insertion sort depends on a number of factors, some of which are the machine, how the code is implemented and how the code is compiled. Finding the optimum size is done empirically by simply iterating through some trial values and check where the execution time dips. A file is written each run with timings. For this particular environment, this indicated that choosing 50 as a cutoff provides the best performance increase, providing a speedup of 1.2 times.

3.3.3 Pivot Choice

Naively choosing a pivot, be it the middle element or at random creates inconsistencies in the size of the left and right side of a partitioned array. For example if we have a shuffled array of ten elements as shown below and the pivot is chosen to be the middle element, **7**, we would end up with one significantly larger partition and one smaller as shown below.



Ideally we would like to choose the median element. Calculating the median is however an expensive process and is not something we can realistically do and expect better performance. By creating two equal sized partitions, the recursion should reach the same depth for all partitions, and the time taking partitioning an array on depth d would be cn/d resulting in the best possible performance.

We can employ the *median of three* method of choosing a pivot. A non-deterministic method would pick three elements at random and choose their median value as a pivot. Choosing elements at random also decreases the risks of having bad performance on some specific arrays. This way of choosing a pivot was implemented but chosen not to be used due to `srand` and `rand` functions being too slow. One approach that was not attempted was to build a pseudorandom number generator specific for this problem. Eventually, a deterministic median of three pivot choice was implemented. This worked considerably faster than the non-deterministic one and chooses the pivot by finding the median of the leftmost, middle and rightmost element of the array.

```
1 // Deterministic median of three.
2 // Finds the median of an array by sampling at three points.
3 int medianOfThree(int* arr, const int lo, const int hi)
4 {
5     int left = arr[lo];
6     int mid = arr[(hi+lo)/2];
7     int right = arr[hi];
8
9     if (right < left)
10         swap(&right, &left);
11     if (mid < left)
12         swap(&mid, &left);
13     if (right < mid)
14         swap(&right, &mid);
15
16     return mid;
17 }
```

3.3.4 Parallelizing Quicksort with Pthreads.

At the same recursion level, each process of sorting a partition is independent of each other, which means that they may be performed in parallel. One approach is to partition up until a maximum recursion level and store each partition in a global array of `threadData_t` structures referencing a partition, a thread, and the thread's ID.

```
1 // Preprocess by taking more samples to lower the risks of unbalanced load between threads.
2 // the global args array holds a threadData_t struct for each thread.
3 const int maxLevel = 4;
4 int threadCounter = 0;
5 threadData_t args[16];
6 void preprocess(int* arr, const int lo, const int hi, const int level)
7 {
8     if (level < maxLevel)
9     {
10         int p = partitionHoareBalanced(arr, lo, hi);
11         preprocess(arr, lo, p, level+1);
12         preprocess(arr, p+1, hi, level+1);
13     }
14     else
15     {
16         threadData_t data;
17         data.arr = arr;
18         data.lo = lo;
19         data.hi = hi;
20         data.threadID = threadCounter;
21
22         args[threadCounter] = data;
23         threadCounter++;
24     }
25 }
```

This solution is simple and has its risks. The largest of the partitions at maximum recursion level will restrict minimum execution time. Choosing a maximum recursion level of 4 will result in sixteen partitions to be sorted by each of the 16 threads. Therefore, it is especially important that the pivots chosen are close to the true median to ensure a good load balance. Without a good approximate median, quicksort will produce inconsistent results using this parallelization scheme. At one run it may provide a near optimal speedup while at another run provide close to no speedup or even degrade performance slightly due to the overhead when creating and joining threads.

This is solved by the `partitionHoareBalanced` function. It works the same way as the regular Hoare partitioning with a slight difference. It samples the array at larger number of equidistant points and chooses the median of those values as a pivot. This reduces the risks that the partition sizes deviate significantly from each other. The extra computation required is well worth it since this will only be done $2^{maxLevel}$ times. This proved to drastically help performance consistency without a noticeably impact on the best execution time. The parallelized quicksort first resets the `threadData_t` array and proceeds to create the partitions at maximum recursion level. When all partitions have been created, the function spawns a thread to sort each partition sequentially. Finally, all threads are joined and the global variable `threadCounter` is reset.

```

1 // Creates threads and assigns their workload.
2 void quicksortParall(int* arr, const int lo, const int hi)
3 {
4     // Reset argument array.
5     for (int i=0; i<16; i++)
6     {
7         args[i].thread = 0;
8         args[i].threadID = 0;
9         args[i].arr = NULL;
10        args[i].lo = 0;
11        args[i].hi = 0;
12    }
13
14    // Partition array with low risk of unbalanced load.
15    preprocess(arr, lo, hi, 0);
16
17    // For all leaf partitions let a thread perform sequential quicksort.
18    for(int i=0; i<threadCounter; i++)
19        pthread_create(&args[i].thread, NULL, parallHelper, (void*)&args[i]);
20
21    // Join threads and reset thread counter.
22    void* status;
23    for(int i=0; i<threadCounter; i++)
24        pthread_join(args[i].thread, &status);
25
26    threadCounter = 0;
27
28 }

```

The `parallHelper` function takes a `void*` arguments, casts to a `threadData_t` structure and calls the sequential algorithm.

```

1 // Intermediate step to call sequential quicksort with a pthread.
2 void* parallHelper(void* args)
3 {
4     threadData_t* data = (threadData_t*)args;
5     int* arr = (*data).arr;
6     int lo = (*data).lo;
7     int hi = (*data).hi;
8     quicksortHoare(arr, lo, hi);
9     pthread_exit(NULL);
10 }

```

4 Performance analysis

4.1 General Optimization

The sorting works in-place meaning that the array which is passed will be sorted instead of making a copy and sorting that instead. While it may or may not have been easier to implement an algorithm which creates a copy it would either way have been inefficient in regards to memory. The performance is expected to be worse aswell because the pointer needs to move inbetween a memory span twice the size of the in-line version.

The choice to restrict to integer elements was made for two main reasons. First, integer comparisons and arithmetic are relatively quick compared to other datatypes and still provide a large span of possible numbers (-32,768 to 32,767). Secondly, a `int` occupies less memory than its long or double counterparts, meaning a shorter span of memory is allocated for the array. This leads to shorter pointer moves and thus, better performance.

Declaring constants as `const` may help the compiler to make some optimizations at assembly level. The `const` declaration by itself should not yield any speedup, however, it is good practice to help prevent the developer introducing errors and fascilitates readability.

Setting the attribute `inline` to functions, where applicable, is a great method for a noticeable speedup for this specific problem since most of the code is segmented into small, easaly-read, subroutines, this is essential to achieve a good performance. The performance increase comes from the removal of function overhead by essentially replacing a call to an `inline` function by its code. Care must be taken as this may vary between different codes. `inlining` when the executable is large may worsen the performance because the system spends more time fetching the next code chunk. This is done automatically if the `-O3` optimization flag is passed.

No manual loop unrolling is performed but is done by the compiler by passing the `funroll-loops` flag. This increases the size of the executable. The performance increase comes from the elimination of the loop control instruction at assembly level. This is done automatically if the `-O3` optimization flag is passed.

4.2 Parallel performance

Comparing the parallel performance increase is done by generating an array with a size n and making a copy. The optimized sequential quicksort is timed sorting the original and the optimized parallel quicksort is timed sorting the copy using 16 threads. This is repeated as n is incremented from 2,000 to 100,000,000 with steps of 2000. When all timings have been performed, a result file is written. This is repeated five times and the lowest sorting times for both the sequential and parallel quicksort are chosen at each n to compare against both methods. Looking at figure 1, it becomes apparent that at small n , the sequential performs better than the parallelized version due to decreased overhead when spawning threads. At higher n , the parallel speedup is consistently around 4 times.

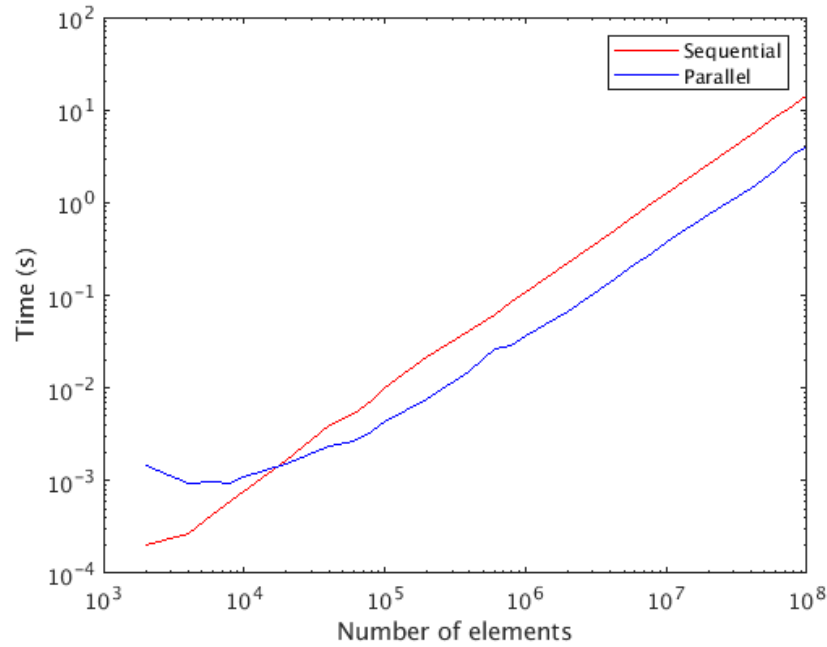


Figure 1: *LogLog plot of the best recorded execution times for the optimized sequential (red) quicksort and the parallelized (blue) version. Best recordings was sampled from a set of five executions.*

4.3 Final performance

Comparing the first version, presented in 3.2, to the final parallelized and optimized quicksort is interesting to see how much of a performance increase was gained in total. The test was performed similarly to when comparing the optimized sequential and parallel algorithms. When compiling with no optimization flags, the speedup is slightly below 5 times for large arrays. With `-O3` optimization flag, the speedup is closer to 4 times. This is reasonable as the optimized quicksort algorithm is not expected to do as many function calls, such as `swap`, and the benefit of removed function overhead is not as impactful for the optimized quicksort than it is for the basic Lomuto algorithm.

A Appendix - Complete code

```
1 #include <stdio.h>
2 #include <stdlib.h>
3 #include <string.h>
4 #include <math.h>
5 #include <time.h>
6 #include <pthread.h>
7
8 // For timings.
9 double get_wall_seconds()
10 {
11     struct timespec tspec;
12     clock_gettime(CLOCK_REALTIME, &tspec);
13
14     double seconds = tspec.tv_sec + (double)tspec.tv_nsec/(double)1000000000;
15     return seconds;
16 }
17
18 // Swaps two integers.
19 void swap(int* a, int* b)
20 {
21     int temp = *a;
22     *a = *b;
23     *b = temp;
24 }
25
26 // Shuffles an array.
27 // arr: array.
28 // n: array size.
29 void shuffle(int* arr, const int n)
30 {
31     // Set seed for rand().
32     srand(time(NULL));
33
34     // Shuffle array.
35     for (int i = n - 1; i > 0; i--)
36     {
37         int j = rand() % (i + 1);
38         swap(&arr[i], &arr[j]);
39     }
40 }
41
42 // Prints an array to terminal in a column.
43 void printArrVert(int* arr, const int n)
44 {
45     for (int i = 0; i < n; i++)
46     {
47         printf("idx %i: %i\n", i, arr[i]);
48     }
49 }
50
51 // Prints an array to terminal horizontally.
52 void printArr(int* arr, const int n)
53 {
54     printf("[");
55     for (int i = 0; i < n-1; i++)
56     {
57         printf("%i, ", arr[i]);
58     }
59     printf("%i\n", arr[n-1]);
60 }
61
62 // Prints a subset of an arr from lo to hi inclusive.
```

```

63 void printSubArr(int* arr, const int lo, const int hi)
64 {
65     printf("[");
66     for (int i = lo; i < hi; i++)
67     {
68         printf("%i, ", arr[i]);
69     }
70     printf("%i\n", arr[hi]);
71 }
72
73 // Generates an ordered array [0, 1, 2, 3, ..., n-1].
74 int* generateArr(const int n)
75 {
76     int* arr = (int*)malloc(sizeof(int) * n);
77     for (int i = 0; i < n; i++)
78     {
79         arr[i] = i;
80     }
81
82     return arr;
83 }
84
85 // Makes a copy of an array.
86 int* copyArr(int* oldArr, const int n)
87 {
88     int* newArr = (int*)malloc(sizeof(int)*n);
89     for (int i=0; i<n; i++)
90     {
91         newArr[i] = oldArr[i];
92     }
93     return newArr;
94 }
95
96 // Verifies that an array is sorted.
97 void verifySorted(int* arr, const int n)
98 {
99     for (int i=1; i<n; i++)
100     {
101         if (arr[i]<arr[i-1])
102         {
103             printf("*** Array is NOT sorted. ***\n");
104             return;
105         }
106     }
107     printf("*** Array is sorted. ***\n");
108 }
109
110
111
112 // === QUICKSORT FUNCTIONS ===
113
114 // Fast insertion sorting algorithm. Used for small partitions.
115 // lo: lowest index of the array;
116 // hi: highest index of the array;
117 void insertionSort(int* arr, const int lo, const int hi)
118 {
119     int tmp;
120     int j;
121     for (int i=lo+1; i<hi+1; i++)
122     {
123         tmp = arr[i];
124         j = i-1;
125
126         while(j>=0 && arr[j]>tmp)
127         {

```



```

128         arr[j+1] = arr[j];
129         j = j-1;
130     }
131     arr[j+1] = tmp;
132 }
133 }
134
135 // Lomuto partition scheme: Choose last element as pivot.
136 int partitionLomuto(int* arr, const int lo, const int hi)
137 {
138     const int pivot = arr[hi];
139     int i = lo;
140
141     for (int j=lo; j<=hi; j++)
142     {
143         if (arr[j] < pivot)
144         {
145             swap(&arr[i], &arr[j]);
146             i += 1;
147         }
148     }
149
150     swap(&arr[i], &arr[hi]);
151     return i;
152 }
153
154 // Deterministic median of n elements. array is a small array.
155 // Used when partitioning in the first levels of recursion
156 // when running multiple threads. Lowers risk of bad load balance.
157 int medianOfN(int* arr, const int lo, const int hi)
158 {
159     insertionSort(arr, lo, hi);
160     return arr[(hi+lo)/2];
161 }
162
163 // Deterministic median of three.
164 // Finds the median of an array by sampling at three points.
165 int medianOfThree(int* arr, const int lo, const int hi)
166 {
167     int left = arr[lo];
168     int mid = arr[(hi+lo)/2];
169     int right = arr[hi];
170
171     if (right<left)
172         swap(&right, &left);
173     if (mid<left)
174         swap(&mid, &left);
175     if (right<mid)
176         swap(&right, &mid);
177
178     return mid;
179 }
180
181
182 // Non-deterministic median of three. Returns median of a subset of an array.
183 // Note: RAND_MAX = 2147483647 on ITC1515 system so we should be safe.
184 // arr: array for which the median value is to be found in selected interval.
185 // lo: lower interval bound.
186 // hi: upper interval bound.
187 int medianOfThreeRand(int* arr, const int lo, const int hi)
188 {
189     srand(time(NULL));
190     int vals[3];
191     int index;
192

```

```

193     index = (rand()%(hi-lo+1))+lo;
194     vals[0] = arr[index];
195     index = (rand()%(hi-lo+1))+lo;
196     vals[1] = arr[index];
197     index = (rand()%(hi-lo+1))+lo;
198     vals[2] = arr[index];
199
200     if (vals[2]<vals[0])
201         swap(&vals[2], &vals[0]);
202     if (vals[1]<vals[0])
203         swap(&vals[1], &vals[0]);
204     if (vals[2]<vals[1])
205         swap(&vals[2], &vals[1]);
206
207     return vals[1];
208 }
209
210 // Basic hoare partitioning. No median of three strategy is applied.
211 int partitionBasicHoare(int* arr, const int lo, const int hi)
212 {
213     // Choose middle element as pivot.
214     const int pivot = arr[(hi+lo)/2];
215
216     // Initialize indices at the ends of the array.
217     int i = lo - 1;
218     int j = hi + 1;
219
220     // Repeat until indices meet.
221     while (1)
222     {
223         // Increment i until arr[i] exceeds pivot.
224         do
225             i++;
226         while (arr[i] < pivot);
227
228         // Decrement j until pivot exceeds arr[j].
229         do
230             j--;
231         while (arr[j] > pivot);
232
233         // If indices meet, return j.
234         if (i >= j)
235             return j;
236
237         // Inversion has been found. Swap current elements.
238         swap(&arr[i], &arr[j]);
239     }
240 }
241
242 // Hoare partition scheme: Choose middle element as pivot.
243 int partitionHoare(int* arr, const int lo, const int hi)
244 {
245     // Choose pivot as the median of three elements.
246     const int pivot = medianOfThree(arr, lo, hi);
247
248     // Initialize indices at the ends of the array.
249     int i = lo - 1;
250     int j = hi + 1;
251
252     while (1)
253     {
254         // Increment i until arr[i] exceeds pivot.
255         do
256             i++;
257         while (arr[i] < pivot);

```

```

258         // Decrement j until pivot exceeds arr[j].
259         do
260             j--;
261         while (arr[j] > pivot);
262
263         // If indices meet, return j.
264         if (i >= j)
265             return j;
266
267         // Inversion has been found. Swap current elements.
268         swap(&arr[i], &arr[j]);
269     }
270 }
271
272 }
273
274 // Hoare's partitioning scheme: sample the array at regular intervals
275 // and choose their median as pivot
276 int partitionHoareBalanced(int* arr, const int lo, const int hi)
277 {
278     // Choose median of any nSamples number of elements.
279     const int nSamples = 50;
280     const int stepsize = (hi-lo)/(nSamples-1);
281     int* pivotArr = (int*)malloc(sizeof(int)*nSamples);
282     pivotArr[0] = arr[lo];
283     for (int i=1; i<nSamples-1; i++)
284     {
285         // Sample array at regular intervals.
286         pivotArr[i] = arr[lo+i*stepsize];
287     }
288     pivotArr[nSamples-1] = arr[hi];
289     const int pivot = medianOfN(pivotArr, 0, nSamples-1);
290     free(pivotArr);
291
292     // Initialize indices at the ends of the array.
293     int i = lo - 1;
294     int j = hi + 1;
295
296     while (1)
297     {
298         // Increment i until arr[i] exceeds pivot.
299         do
300             i++;
301         while (arr[i] < pivot);
302
303         // Decrement j until pivot exceeds arr[j].
304         do
305             j--;
306         while (arr[j] > pivot);
307
308         // If indices meet, return j.
309         if (i >= j)
310             return j;
311
312         // Inversion has been found. Swap current elements.
313         swap(&arr[i], &arr[j]);
314     }
315 }
316
317 // Quicksort recursion using Lomuto's partitioning scheme.
318 void quicksortLomuto(int* arr, const int lo, const int hi)
319 {
320     if (lo < hi)
321     {

```

```

323     int p = partitionLomuto(arr, lo, hi);
324     quicksortLomuto(arr, lo, p-1);
325     quicksortLomuto(arr, p+1, hi);
326 }
327 }
328
329 // Quicksort recursion using Hoares's partitioning scheme.
330 void quicksortBasicHoare(int* arr, const int lo, const int hi)
331 {
332
333     if (lo < hi)
334     {
335         int p = partitionBasicHoare(arr, lo, hi);
336         quicksortBasicHoare(arr, lo, p);
337         quicksortBasicHoare(arr, p+1, hi);
338     }
339 }
340
341 // Quicksort recursion using Hoares's partitioning scheme with median of three
342 // strategy applied and utilized insertion sort for small partitions.
343 void quicksortHoare(int* arr, const int lo, const int hi)
344 {
345
346     if (lo < hi)
347     {
348         // If array is small, use insertion sort.
349         if (hi-lo < 50)
350         {
351             insertionSort(arr, lo, hi);
352             return;
353         }
354         int p = partitionHoare(arr, lo, hi);
355         quicksortHoare(arr, lo, p);
356         quicksortHoare(arr, p+1, hi);
357     }
358 }
359
360
361
362
363
364 // === PARALLELIZED QUICKSORT ===
365
366 // Struct used to call parallHelper.
367 typedef struct threadData
368 {
369     pthread_t thread;
370     int threadID;
371     int* arr;
372     int lo;
373     int hi;
374 } threadData_t;
375
376 // Intermediate step to call sequential quicksort with a pthread.
377 void* parallHelper(void* args)
378 {
379     threadData_t* data = (threadData_t*)args;
380     int* arr = (*data).arr;
381     int lo = (*data).lo;
382     int hi = (*data).hi;
383     quicksortHoare(arr, lo, hi);
384     pthread_exit(NULL);
385 }
386
387 // Preprocess by taking more samples to lower the risks of unbalanced load between threads.

```

```

388 // the global args array holds a threadData_t struct for each thread to be used later by
    quicksortParall.
389 const int maxLevel = 4;
390 int threadCounter = 0;
391 threadData_t args[16];
392 void preprocess(int* arr, const int lo, const int hi, const int level)
393 {
394     if (level < maxLevel)
395     {
396         int p = partitionHoareBalanced(arr, lo, hi);
397         preprocess(arr, lo, p, level+1);
398         preprocess(arr, p+1, hi, level+1);
399     }
400     else
401     {
402         threadData_t data;
403         data.arr = arr;
404         data.lo = lo;
405         data.hi = hi;
406         data.threadID = threadCounter;
407
408         args[threadCounter] = data;
409         threadCounter++;
410     }
411 }
412 }
413
414 // Creates threads and assigns their workload.
415 void quicksortParall(int* arr, const int lo, const int hi)
416 {
417     // Reset argument array.
418     for (int i=0; i<16; i++)
419     {
420         args[i].thread = 0;
421         args[i].threadID = 0;
422         args[i].arr = NULL;
423         args[i].lo = 0;
424         args[i].hi = 0;
425     }
426
427     // Partition array with low risk of unbalanced load.
428     preprocess(arr, lo, hi, 0);
429
430     // For all leaf partitions let a thread perform sequential quicksort.
431     for(int i=0; i<threadCounter; i++)
432         pthread_create(&args[i].thread, NULL, parallHelper, (void*)&args[i]);
433
434     // Join threads and reset thread counter.
435     void* status;
436     for(int i=0; i<threadCounter; i++)
437         pthread_join(args[i].thread, &status);
438
439     threadCounter = 0;
440 }
441 }
442
443 // Used for finding the optimal array size to switch to insertion sort.
444 int cutoff = 0;
445 void quicksortInsertTest(int* arr, const int lo, const int hi)
446 {
447     if (lo < hi)
448     {
449         // If array is small, use insertion sort.
450         if (hi-lo < cutoff)
451         {

```

```

452         insertionSort(arr, lo, hi);
453         return;
454     }
455     int p = partitionHoare(arr, lo, hi);
456     quicksortInsertTest(arr, lo, p);
457     quicksortInsertTest(arr, p+1, hi);
458 }
459 }
460
461 // Used to approximate what size of an array it is more beneficial
462 // to use insertion sort than to keep recursing.
463 void optimalInsertCutoff()
464 {
465     const int N = 1000000;
466     int* arr;
467     int* arrCopy;
468     double startTime;
469     double insertTime;
470     double noInsertTime;
471
472     // Generate two identical arrays.
473     arr = generateArr(N);
474     shuffle(arr, N);
475     arrCopy = copyArr(arr, N);
476
477     // Do a measurement where no insertion sort is used
478     cutoff = 0;
479     startTime = get_wall_seconds();
480     quicksortInsertTest(arr, 0, N-1);
481     noInsertTime = get_wall_seconds() - startTime;
482     arr = arrCopy;
483     arrCopy = copyArr(arr, N);
484
485     // Compare to when insertion sort is used for different cutoff values and write results
486     // to a file.
487     FILE *outputFile = fopen("VaryingInsertionCutoffTimings.txt", "w");
488     fprintf(outputFile, "nElements,cutoff,noInsertTime,InsertTime,speedup\n");
489     for (int i=1; i<3; i++)
490     {
491         for (int j=1; j<6; j++)
492         {
493             cutoff = 2*j*10*pow(10, i-1);
494             printf("Cutoff: %i\n", cutoff);
495
496             startTime = get_wall_seconds();
497             if (cutoff == 1000000)
498                 printf("Hello!");
499             quicksortInsertTest(arr, 0, N-1);
500             insertTime = get_wall_seconds() - startTime;
501             verifySorted(arr, N);
502             fprintf(outputFile, "%i,%i,%3.6f,%3.6f,%3.6f\n", N, cutoff, noInsertTime,
503                     insertTime, noInsertTime/insertTime);
504
505             free(arr);
506             arr = arrCopy;
507             arrCopy = copyArr(arr, N);
508         }
509     }
510
511 // Used to perform timings between different versions of quicksort.
512 // Writes a result file.
513 void doTimings()
514 {

```

```

515
516     int* arr1;
517     int* arr2;
518     int N;
519
520     double startTime;
521     double seqTime;
522     double parTime;
523
524     FILE *outputFile = fopen("run1.txt", "w");
525     fprintf(outputFile, "nElements,seqTime,parTime,speedup\n");
526
527     // from 1000 to 1000000000
528     for (int i=1; i<6; i++)
529     {
530         for (int j=1; j<6; j++)
531         {
532             // Update N.
533             N = (2*j)*1000*pow(10, i-1);
534
535             // Do first timing.
536             printf("*** Sorting %i elements\n", N);
537             printf(">>> Sorting using sequential quicksort...\n");
538             arr1 = generateArr(N);
539             shuffle(arr1, N);
540             arr2 = copyArr(arr1, N);
541             printf("Original: ");
542             verifySorted(arr1, N);
543             startTime = get_wall_seconds();
544             quicksortHoare(arr1, 0, N-1);
545             seqTime = get_wall_seconds() - startTime;
546             printf("Sorted: ");
547             verifySorted(arr1, N);
548
549             // Do second timing.
550             printf(">>> Sorting using parallelized quicksort...\n");
551             printf("Original: ");
552             verifySorted(arr2, N);
553             startTime = get_wall_seconds();
554             quicksortParall(arr2, 0, N-1);
555             parTime = get_wall_seconds() - startTime;
556             printf("Sorted: ");
557             verifySorted(arr2, N);
558
559             // Print to output file.
560             fprintf(outputFile, "%i,%3.6f,%3.6f,%3.6f\n", N, seqTime, parTime, seqTime/
parTime);
561             printf("Done\n\n");
562             free(arr1);
563             free(arr2);
564             arr1 = NULL;
565             arr2 = NULL;
566         }
567     }
568     fclose(outputFile);
569 }
570
571 // Main driver code.
572 int main(int argc, char *argv[])
573 {
574     int* arr1;
575     int* arr2;
576
577     // Testing parallel sorting.
578     int N = 50000000;

```

```

579 printf("\n\n>>> PARALLEL TIMINGS:\n");
580 printf("    Using array size: %i.\n", N);
581
582 arr1 = generateArr(N);
583 shuffle(arr1, N);
584 arr2 = copyArr(arr1, N);
585
586 printf("arr1: ");
587 verifySorted(arr1, N);
588 printf("arr2: ");
589 verifySorted(arr2, N);
590 printf("\n");
591
592 double startTime;
593 printf("Sorting arr1 using sequential quicksort.\n");
594 printf("Original arr1: ");
595 verifySorted(arr1, N);
596 startTime = get_wall_seconds();
597 quicksortHoare(arr1, 0, N-1);
598 printf("Sorting time: %3.6f seconds.\n", get_wall_seconds() - startTime);
599 printf("Verifying...\n");
600 printf("Sorted arr1:  ");
601 verifySorted(arr1, N);
602 printf("\n\n");
603 free(arr1);
604
605 printf("Sorting arr2 using parallelized quicksort.\n");
606 printf("Original arr2: ");
607 verifySorted(arr2, N);
608 startTime = get_wall_seconds();
609 quicksortParall(arr2, 0, N-1);
610 printf("Sorting time: %3.6f seconds.\n", get_wall_seconds() - startTime);
611 printf("Verifying...\n");
612 printf("Sorted arr2:  ");
613 verifySorted(arr2, N);
614 printf("\n\n");
615 free(arr2);
616
617 pthread_exit(NULL);
618 }

```


B Appendix - Makefile

```
1 CFLAGS=-Wall -O3
2 LDFLAGS=-lm -pthread
3
4 Quicksort: Quicksort.o
5     gcc -g -o Quicksort Quicksort.o $(LDFLAGS)
6
7 Quicksort.o: Quicksort.c
8     gcc -g $(CFLAGS) -c Quicksort.c
9
10 run: Quicksort
11     ./Quicksort
12
13 debug:
14     gdb ./Quicksort
15
16 valgrind:
17     valgrind --leak-check=full ./Quicksort
18
19 clean:
20     rm -f ./Quicksort *.o
```