

Assignment 2

Sorting

Prof. Darrell Long
Department of Computer Engineering
Spring 2017

Any inaccuracies in this index may be explained by the fact that it has been sorted with the help of a computer.

—Donald Knuth, Vol. III, *Sorting and Searching*

1 Introduction

Putting items into a sorted order is one of the most common tasks in Computer Science, and so as a result there are myriad library routines that will do this task for you—but that does not absolve you of the obligation of understanding how it is done, and in fact it behooves you to understand the various algorithms in order to make wise choices.

The best that can be accomplished (the *lower bound*) for sorting using *comparisons* is $\Omega(n \log n)$ where n is the number of elements to be sorted. If the universe of elements to be sorted is limited (small), then we can do better using a *Count Sort*, where is $O(n)$, where we count the number of occurrences of each element in an array, or a *Radix Sort*, which is also $O(n)$ with a constant proportional to the maximum number of digits in the numbers being sorted.

What is this O and Ω stuff? It's how we talk about the execution time (or space) of programs. We will discuss it in class, and you will see it again in CMPS 101.

1.1 min-Sort

Perhaps the simplest sorting method is to look for the smallest element (the minimum) and move it to the top. We could do this by making a copy of that element, and then shifting everything down by one, and then placing the copy in the top slot. But that seems silly, why move all of those elements? Let's just *exchange* (swap) the top element with the smallest element.

At this point what do we know? We know that the smallest element is at the top. Since that is true and it will not change (we call this an *invariant*), we can forget about it and move on. Let's consider the second element: Why not just do what we did the first time? If we do that, then what do we know? We now know that the top (first) element is the smallest, and the second element is the second smallest (or the same, if there are duplicates). We can repeat this for each element in the array, in succession, up to but not including the last element. By the method of *induction*, we can show that the array is sorted.

Why do we not need to concern ourselves with the last element? The answer is that if it were not the smallest element when we were at the last step, then it was exchanged with the penultimate element, and thus must necessarily be the largest (and consequently, last) element in the array.

To get you started, here is the code for `minSort`. Notice that it is composed of two functions: `minIndex` which finds the *location* of the smallest element, and `minSort` which actually performs the sorting.

```

1 // minIndex: find the index of the least element.
2
3 uint32_t minIndex(uint32_t a[], uint32_t first, uint32_t last)
4 {
5     uint32_t smallest = first; // Assume the first is the least
6     for (uint32_t i = first; i < last; i += 1)
7     {
8         smallest = a[i] < a[smallest] ? i : smallest;
9     }
10    return smallest;
11 }
12
13 // minSort: sort by repeatedly finding the least element.
14
15 void minSort(uint32_t a[], uint32_t length)
16 {
17     for (uint32_t i = 0; i < length - 1; i += 1)
18     {
19         uint32_t smallest = minIndex(a, i, length);
20         if (smallest != i) // It's silly to swap with yourself!
21         {
22             SWAP(a[smallest], a[i]);
23         }
24     }
25    return;
26 }

```

minSort (in C)

What is the time complexity of minSort? The key is to look at the for loop, and we see that the loop is executed `length` times. This would seem to indicate that the sort is $O(n)$, but we need to look further. In the for loop we see that a call is made to `minIndex`, and when we look there we find yet another for loop. This loop is executed `(last - first)` times, which is based on `length`. So we call `minIndex` approximately `length` times, each time requiring approximately `length` operations, thus the time complexity of minSort is $O(n^2)$.

It is often useful to define a macro when a task is done repetitively, and when we also may want to hide the details. One could write a function (an in-line function would be preferred), but you will recall that functions in C always pass their arguments *by value* which is inconvenient for the *swap* operation. A clever person might take advantage of the macro to do some instrumentation.

```

1 # ifdef _INSTRUMENTED
2 # define SWAP(x,y) { uint32_t t = x; x = y; y = x; ; moveCount += 3; }
3 # else
4 # define SWAP(x,y) { uint32_t t = x; x = y; y = x; ; }
5 # endif

```

SWAP macro

1.2 Bubble Sort

Bubble Sort works by examining adjacent pairs of items. If the second item is smaller than the first, exchange them. If you can pass over the entire array and no pairs are out of order, then the array is sorted. You will have noticed

that the largest element falls in a single pass to the bottom of the array. Since it is in fact the largest, we do not need to consider it again, and so the next pass must only consider $n - 1$ pairs of items.

What then is the expected time complexity of this algorithm? The first pass requires n pairs to be examined; the second pass, $n - 1$ pairs; the third pass $n - 2$ pairs, and so forth. When Carl Friedrich Gauss was a child of 7 (in 1784), he is reported to have amazed his elementary school teacher being a pest he was given the task of summing the integers from 1 to 100. The precocious little Gauss produced the correct answer immediately, having quickly observed that the sum was actually 50 pairs of numbers, with each pair summing to 101, total 5,050. We can then easily see that:

$$n + (n - 1) + (n - 2) + \dots + 1 = \frac{n(n + 1)}{2},$$

and so the *worst case* time complexity is $O(n^2)$ (it may be much better, if for example, the list is already sorted).

```

1 procedure bubbleSort( A : list of sortable items )
2   n = length(A)
3   repeat
4     swapped = false
5     for i = 1 to n-1 inclusive do
6       if A[i-1] > A[i] then
7         swap(A[i-1], A[i])
8         swapped = true
9       end if
10    end for
11    n = n - 1
12  until not swapped
13 end procedure

```

bubbleSort (pseudocode)

1.3 Insertion Sort

Insertion Sort is another in-place sort, it functions by taking an item and then inserting it into its correct position in the array. It consumes one input element each repetition, and growing the sorted portion of the list. Each iteration, insertion sort removes one element from the input unsorted portion and finds its location in the sorted portion of the list.

```

1 procedure insertionSort( A : list of sortable items )
2   for i = 1 to length(A)
3     tmp = A[i]
4     j = i - 1
5     while j >= 0 and A[j] > tmp
6       A[j + 1] = A[j]
7       j = j - 1
8     end while
9     A[j + 1] = tmp
10  end for
11 end procedure

```

insertionSort (pseudocode)

What is the expected time complexity of Insertion sort? We look and observe that there two two nested loops, each operating on approximately the length of the array (it grows shorter at each iteration, but as we learned from young Gauss, $\sum_{i=1}^n = \frac{n(n+1)}{2}$). Nested loops mean that we *multiply* the execution times, and so it is also $O(n^2)$.

1.4 QuickSort (recursive)

One of the most important skills that you will need to develop is the ability to read and understand, though perhaps not program in, languages with which you are unfamiliar. Below, you will find the code for QuickSort written in Python. It would be a questionable choice for you to simply try to translate this into C, but there are some interesting elements to it.

First, you see that the code partitions the array into three parts: (i) a part that is *lesser* in value than the pivot element; (ii) *equal* in value to the pivot point, and (iii) *greater* in value than the pivot point. The choice of pivot element is arbitrary.

Second, you will see that we first recursively call quickSort first on the left partition, and then on the right partition. We then join these together to form a sorted array.

Third, you may notice that aside from the arbitrarily chosen pivot element, there is no array indexing. One could, in principle, use this algorithm to sort a *linked list*.

You will want to write a helper function called *partition* that will divide an array into three parts, as described earlier. You should do this *in place*, in other words you do not need to create three arrays. Instead, you will *swap* elements so that those less than or equal to the pivot element are on the left, while those greater than it are on the right. Why was this acceptable in Python? Python is a very different language that has lists (which is can treat as arrays) as a native data type. Different language enable different choices, and writing this in Python leads to the (surprising to some) insight that QuickSort can be implemented on linked lists.

```
1 def quickSort(a):
2     if len(a) > 0:
3         pivot = a[0]
4         left = []
5         mid = []
6         right = []
7         for x in a:
8             if x == pivot:
9                 mid.append(x)
10            elif x < pivot:
11                left.append(x)
12            else:
13                right.append(x)
14        return quickSort(left) + mid + quickSort(right)
15    else:
16        return []
```

quickSort.py

1.5 Merge Sort

Merge Sort is a different type of sort in that it works through the array sequentially. It is well-suited for the case when you *do not have* random access, such as when working with magnetic tapes, or linked lists.

There are two main varieties of Merge Sort: *binary merge sort* (which may be recursive) and *natural merge sort*. The recursive binary merge sort was originally discovered by John von Neumann, one of the great polymaths of the previous century.

A single item (or, if you prefer to start your induction with zero, and empty list) is by definition sorted. Now, consider the case where we have two unsorted lists: either (i) the item at the front of list 1 is the smallest, or (ii) the item at the front of list 2 is the smallest, but in either case *greater than or equal to* the item at the end of list 3. Move that item to the end of the the (initially empty) list 3. Repeat this process until both of list 1 or list 2 have elements at their head less than the element at the end of list 3 or are empty. What do we know at this point? We know that

(i) the third list contains a *sorted subsequence*, and (ii) that sequence is equal in length to the sum of the sorted subsequences that we pulled from lists 1 and 2. This sorted subsequence is called a *run*. We repeat this process until we have exhausted either list 1 or list 2, at which point we append the remaining list to list 3.

We now take list 3 and copy elements from it as long the elements continue to increase to list 1, when the first decreasing element occurs we switch to copying the run to list 2. We repeat this, switching between lists 1 and 2, until list 3 is depleted.

We return to merging elements from lists 1 and 2 to increasing length runs on list 3. We do this until there is a single run, at which point the list is sorted.

It is much easier to think about this if you do it *recursively*. Take the list, split it into two lists (left and right). Call `mergeSort(left)` and `mergeSort(right)` and merge the results. For the base case, the recursion will encounter a list of a single element, which by definition is sorted.

What is the execution time of this algorithm? Let's assume that we have runs of length 1 on tapes 1 and 2. We merge these onto list 3 and we now have (at worst) runs of length 2. We processed n elements during this pass. We now need to ask, How many times can we double starting with 1 until the value we get is greater than n ? The answer is $\log n$. The worst case execution time of our algorithm is $O(n \log n)$.

```
1 def mergeSort(items):
2     if len(items) > 1:
3         middle = len(items) / 2 # Split at the mid-point
4         leftList = items[0:middle] # Left half
5         rightList = items[middle:] # Right half
6
7         mergeSort(leftList) # Sort the leftList list
8         mergeSort(rightList) # Sort the RightList list
9
10    # Merge the sorted lists
11    l, r = 0, 0
12    for i in range(len(items)):
13        # Both lists hold elements
14        if l < len(leftList) and r < len(rightList):
15            # The left is smaller
16            if leftList[l] < rightList[r]:
17                items[i] = leftList[l]
18                l += 1
19            # The right is smaller
20            else:
21                items[i] = rightList[r]
22                r += 1
23        # Only the left has any elements
24        elif l < len(leftList):
25            items[i] = leftList[l]
26            l += 1
27        # Only the right has any elements
28        else:
29            items[i] = rightList[r]
30            r += 1
```

mergeSort.py

2 Your Task

You task is to:

- Implement a testing harness for sorting algorithms. You will do this using `getopt`:

```
1  GETOPT(3)                                Linux Programmer's Manual
   GETOPT(3)
2
3
4
5  NAME
6      getopt, getopt_long, getopt_long_only, optarg, optind,
   opterr, optopt -
7      Parse command-line options
8
9  SYNOPSIS
10     #include <unistd.h>
11
12     int getopt(int argc, char * const argv[],
13                const char *optstring);
14
15     extern char *optarg;
16     extern int optind, opterr, optopt;
17
18     #include <getopt.h>
19
20     int getopt_long(int argc, char * const argv[],
21                     const char *optstring,
22                     const struct option *longopts, int *longindex);
   getopt
```

- Implement *five* specified sorting algorithms.
- Gather statistics about their performance.

3 Specifics

You must use `getopt` to parse the command line arguments. To get you started, here is a hint.

```
1      while ((c = getopt(argc, argv, "AmbiqMp:r:n:")) != -1)
```

- -A means employ all sorting algorithms.
- -m means enable minSort.
- -b means enable bubbleSort.
- -i means enable insertionSort.
- -q means enable quickSort.

- -M means enable mergeSort.
- -p n means print the first n elements of the array. The *default* value is 100.
- -r s means set the random seed to s. The *default* is 8222022.
- -n c means set the array size to c. The *default* value is 100.

It is important to read this carefully. None of these options is exclusive of any other (you may specify any number of them, including *zero*).

- Your random numbers should be 24 bits, no larger ($2^{24} - 1 = 16777215$).
- You must use `rand()` and `srand()`, not because they are good (they are not), but because they are what is specified by the C99 standard.
- Your program *must* be able to sort any number of random integers *up to the memory limit of the computer*. That means that you will need to dynamically allocate the array using `calloc()`.
- Your program should have no *memory leaks*.

A large part of this assignment is understanding and comparing the performance of various sorting algorithms. Consequently, you *must* collect some simple statistics on each algorithm. In particular,

- The *size* of the array,
- The number of *moves* required (each time you transfer an element in the array, that counts), and
- The number of *comparisons* required (comparisons *only* count for elements, not for logic).

```
1 Valhalla:sorting darrell$ make
2 gcc -Wall -Werror -Wextra -pedantic -O3 -DMASK=0x00ffffff -c sorting.c
3 gcc -Wall -Werror -Wextra -pedantic -O3 -c bv.c
4 gcc -Wall -Werror -Wextra -pedantic -O3 -c bubblesort.c
5 gcc -Wall -Werror -Wextra -pedantic -O3 -c minsort.c
6 gcc -Wall -Werror -Wextra -pedantic -O3 -c insertionsort.c
7 gcc -Wall -Werror -Wextra -pedantic -O3 -c quicksort.c
8 gcc -Wall -Werror -Wextra -pedantic -O3 -c mergesort.c
9 gcc -o sorting sorting.o bv.o bubblesort.o minsort.o insertionsort.o
10 quicksort.o mergesort.o
```

```
1 Valhalla:sorting darrell$ ./sorting
2 Valhalla:sorting darrell$
```

```
1 Valhalla:sorting darrell$ ./sorting -n 10 -r 1234 -i
2 Insertion Sort
3 10 elements
4 34 moves
5 25 compares
6      667041      694558      3789884      3962622      6238630      10272683      10916661
7      11017812      11838633      13274474
```

```

1 Valhalla:sorting darrell$ ./sorting -b -q -n 100000 -p 21
2 Bubble Sort
3 100000 elements
4 7489129170 moves
5 2496376390 compares
6      176      685      805      1174      1217      1618      1643
7      1861      1897      2192      2381      2437      2465      3023
8      3058      3225      3285      3348      3425      3909      3928
9
10 Quick Sort
11 100000 elements
12 3539469 moves
13 1046295 compares
14      176      685      805      1174      1217      1618      1643
15      1861      1897      2192      2381      2437      2465      3023
16      3058      3225      3285      3348      3425      3909      3928

```

```

1 Valhalla:sorting darrell$ ./sorting -A
2 Min Sort
3 100 elements
4 288 moves
5 5049 compares
6      47320      272862      813325      931036      1300149      1451478      1482116
7      1599339      1886666      1926530      1999678      2177992      2338129      2381540
8      2504306      2752788      2996335      2998169      3171889      3590861      3829897
9      3967508      4022167      4231676      4422014      4702258      4742819      4841493
10     4915648      5160950      5607401      5792480      5812062      5967954      5968964
11     6053032      6195346      6301066      6375315      6389795      6781764      6797856
12     7260963      7261776      7443963      7614058      7685040      7685284      7739256
13     7989314      8075396      8085360      8249909      8352629      8737396      9017069
14     9214736      9256511      9320602      9414691      9489446      9526978      9657331
15     10177720      10578757      10623131      10856398      10978990      11215204      11268563
16     11339313      11464349      11556747      11680696      11697687      11774380      11989481
17     12291309      12547422      12617960      12814415      12844019      12899314      12912456
18     13059818      13226807      13494413      13782527      14314674      14404909      14506627
19     14956880      15232419      15509585      15563419      15846433      15950040      16608657
20     16703976      16712977
21 Bubble Sort
22 100 elements
23 7074 moves
24 2358 compares
25     47320      272862      813325      931036      1300149      1451478      1482116
26     1599339      1886666      1926530      1999678      2177992      2338129      2381540
27     2504306      2752788      2996335      2998169      3171889      3590861      3829897
28     3967508      4022167      4231676      4422014      4702258      4742819      4841493
29     4915648      5160950      5607401      5792480      5812062      5967954      5968964

```


30	6053032	6195346	6301066	6375315	6389795	6781764	6797856
31	7260963	7261776	7443963	7614058	7685040	7685284	7739256
32	7989314	8075396	8085360	8249909	8352629	8737396	9017069
33	9214736	9256511	9320602	9414691	9489446	9526978	9657331
34	10177720	10578757	10623131	10856398	10978990	11215204	11268563
35	11339313	11464349	11556747	11680696	11697687	11774380	11989481
36	12291309	12547422	12617960	12814415	12844019	12899314	12912456
37	13059818	13226807	13494413	13782527	14314674	14404909	14506627
38	14956880	15232419	15509585	15563419	15846433	15950040	16608657
39	16703976	16712977					
40	Insertion Sort						
41	100 elements						
42	2457 moves						
43	2358 compares						
44	47320	272862	813325	931036	1300149	1451478	1482116
45	1599339	1886666	1926530	1999678	2177992	2338129	2381540
46	2504306	2752788	2996335	2998169	3171889	3590861	3829897
47	3967508	4022167	4231676	4422014	4702258	4742819	4841493
48	4915648	5160950	5607401	5792480	5812062	5967954	5968964
49	6053032	6195346	6301066	6375315	6389795	6781764	6797856
50	7260963	7261776	7443963	7614058	7685040	7685284	7739256
51	7989314	8075396	8085360	8249909	8352629	8737396	9017069
52	9214736	9256511	9320602	9414691	9489446	9526978	9657331
53	10177720	10578757	10623131	10856398	10978990	11215204	11268563
54	11339313	11464349	11556747	11680696	11697687	11774380	11989481
55	12291309	12547422	12617960	12814415	12844019	12899314	12912456
56	13059818	13226807	13494413	13782527	14314674	14404909	14506627
57	14956880	15232419	15509585	15563419	15846433	15950040	16608657
58	16703976	16712977					
59	Quick Sort						
60	100 elements						
61	1542 moves						
62	384 compares						
63	47320	272862	813325	931036	1300149	1451478	1482116
64	1599339	1886666	1926530	1999678	2177992	2338129	2381540
65	2504306	2752788	2996335	2998169	3171889	3590861	3829897
66	3967508	4022167	4231676	4422014	4702258	4742819	4841493
67	4915648	5160950	5607401	5792480	5812062	5967954	5968964
68	6053032	6195346	6301066	6375315	6389795	6781764	6797856
69	7260963	7261776	7443963	7614058	7685040	7685284	7739256
70	7989314	8075396	8085360	8249909	8352629	8737396	9017069
71	9214736	9256511	9320602	9414691	9489446	9526978	9657331
72	10177720	10578757	10623131	10856398	10978990	11215204	11268563
73	11339313	11464349	11556747	11680696	11697687	11774380	11989481
74	12291309	12547422	12617960	12814415	12844019	12899314	12912456
75	13059818	13226807	13494413	13782527	14314674	14404909	14506627
76	14956880	15232419	15509585	15563419	15846433	15950040	16608657
77	16703976	16712977					
78	Merge Sort						
79	100 elements						

```

80 1344 moves
81 316 compares
82      47320      272862      813325      931036      1300149      1451478      1482116
83      1599339      1886666      1926530      1999678      2177992      2338129      2381540
84      2504306      2752788      2996335      2998169      3171889      3590861      3829897
85      3967508      4022167      4231676      4422014      4702258      4742819      4841493
86      4915648      5160950      5607401      5792480      5812062      5967954      5968964
87      6053032      6195346      6301066      6375315      6389795      6781764      6797856
88      7260963      7261776      7443963      7614058      7685040      7685284      7739256
89      7989314      8075396      8085360      8249909      8352629      8737396      9017069
90      9214736      9256511      9320602      9414691      9489446      9526978      9657331
91     10177720     10578757     10623131     10856398     10978990     11215204     11268563
92     11339313     11464349     11556747     11680696     11697687     11774380     11989481
93     12291309     12547422     12617960     12814415     12844019     12899314     12912456
94     13059818     13226807     13494413     13782527     14314674     14404909     14506627
95     14956880     15232419     15509585     15563419     15846433     15950040     16608657
96     16703976     16712977

```

Submission

You *must* turn in your assignment in the following manner:

1. Have file called `Makefile` that when the grader types `make` will compile your program. At this point you will have learned about `make` and can create your own `Makefile`.
 - `CFLAGS=-Wall -Wextra -Werror -pedantic` must be included.
 - `CC=gcc` must be specified.
 - `make clean` must remove all files that are compiler generated.
 - `make` should build your program, as should `make all`.
2. Your program *must* have the source and header files:
 - `minsort.h` specifies the interface to `minSort()`.
 - `minsort.c` implements `minSort()`.
 - Each sorting method will have its own pair of header file and source file.
 - `sorting.c` contains `main()` and *may* contain the other functions necessary to complete the assignment.
3. You may have other source and header files, but *do not try to be overly clever*.
4. A plain text file called `README` that describes how your program works.
5. The executable file produced by the compiler *must be called* `sorting`.
6. These files must be in the directory `assignment2`.
7. You must `commit` and `push` the directory and its contents using `git`.

Points will be assigned according to the difficulty of the sort involved.

- 5% – min-Sort
- 10% – Bubble Sort

- 15% – Insertion Sort
- 20% – Quick Sort (recursive)
- 30% – Merge Sort

A sort is not considered to be implemented if it does not sort *correctly every time*.

Additional criteria are:

- 10% – Code quality and correctness.
- 10% – Completeness: which includes things like the `Makefile`.