

Commission of the European Communities

ESPRIT III

PROJECT NB 6756

CAMAS

COMPUTER AIDED MIGRATION OF APPLICATIONS SYSTEM

CAMAS-TR-2.2.4.6

F2SAD - prediction capabilities

Date: March 1995 — Review 5.0

ACE - Univ. of Amsterdam - ESI SA - ESI GmbH - FECS - PARSYTEC -
Univ. of Southampton

Authors: Berry A.W. van Halderen
Jan de Ronde

March, 1995
University of Amsterdam,
Faculty of Mathematics and Computer Science
Parallel Scientific Computing and Simulation group
Netherlands

Chapter 1

Introduction

This report describes an evaluation of the F2SAD tool with several well known basic algorithms. In this report we consider some sorting algorithms. With this report we respond to the specific request of the review commission to make a more detailed validation of the prediction capabilities of the UvA workbench tools.

The intention of this document is to provide confidence in the ability of the tools (that implement the models developed in SAD and PARASOL) to estimate the execution time of some well known algorithms. The algorithms described here have a performance behaviour that is common knowledge. Despite this fact, we are still able to come up with some points that are interesting and not plain textbook knowledge.

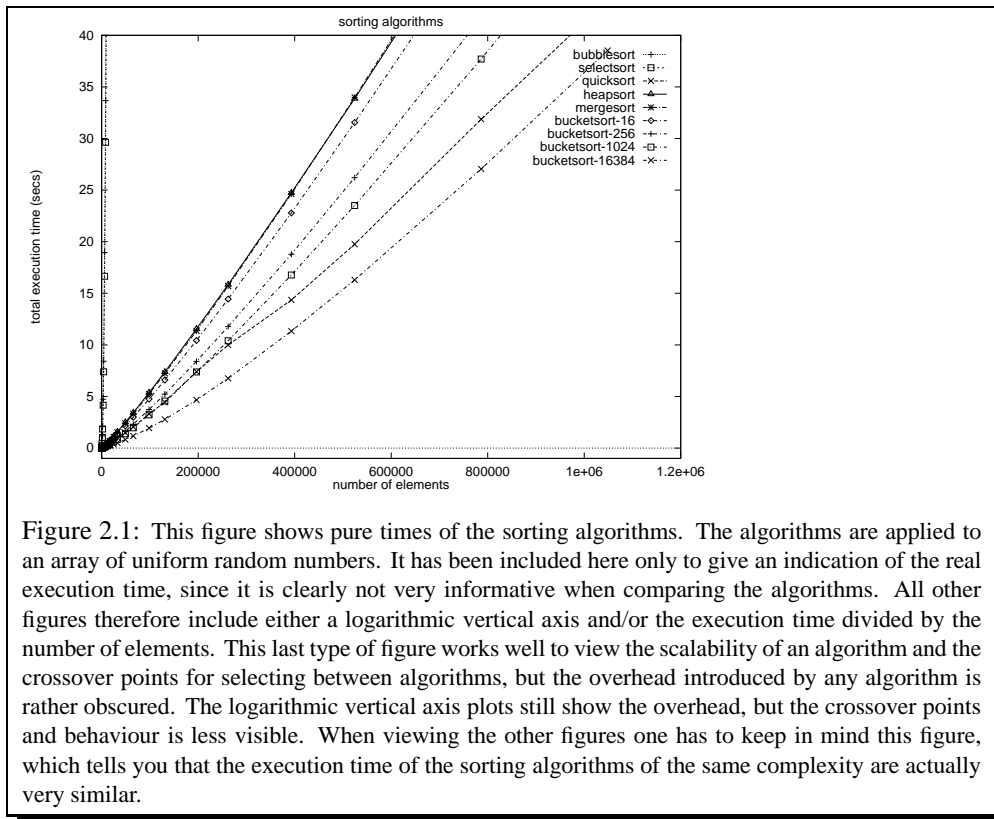
In this report, various figures depicting measured and predicted execution times of Fortran programs will pass. The main part of the annotated algorithms has been included in appendix A.

At the time of the CAMAS review 5 we will present additional results on numerical relaxation algorithms.

Chapter 2

Sorting algorithms

Sorting algorithms are symbolic algorithms. Rather than performing heavy computations they compare and manipulate (swap or reorder) data. Sorting algorithms are quite well understood in their complexity behaviour. Despite this fact, few textbooks do actually compare the execution time of the algorithms. Two algorithms can be of the same order of complexity and still differ in their performance because a different number of instructions is executed.



If we classify the sorting algorithms by their average-case time complexity, we can distinguish three classes. The exponential order $O(n^2)$, the logarithmic based sorting algorithms $n \log(n)$ and the linear time sorting algorithms $O(n)$. The exponential order algorithms are sometimes still (unjustifiably) used if a programmer is lazy and works on small arrays or under the disguise of a parallel computer. Shell-sort is such a variant used on parallel computers because of its parallel nature.

The most common sorting algorithms are the $O(n \log n)$ order sorting algorithms like quick-, heap- and mergesort. Although all three have the same average-case order of

<i>algorithm</i>	<i>average</i>	<i>worst-case</i>	<i>implementation</i>
bubblesort	$O(n^2)$	$\Theta(n^2)$	array (list also possible)
selectsort	$O(n^2)$	$\Theta(n^2)$	array (list also possible)
quicksort	$O(n \log n)$	$\Theta(n^2)$	array
heapsort	$O(n \log n)$	$\Theta(n \log n)$	array
mergesort	$O(n \log n)$	$\Theta(n \log n)$	list
bucketsort	$O(n)$	$\Theta(n^2)$	list

Table 2.1: Sorting algorithms; their time complexity as average case time complexity $O(\dots)$ and worst-case time complexity $\Theta(\dots)$. The implementation can either be in linear array or a linked list form.

performance, they are inherently different. Quicksort has very bad worst-case performance, namely $\Theta(n^2)$. It is almost impossible to implement mergesort on arrays and the natural way to implement heapsort is on arrays. It is controversial whether quicksort is the fastest (on average), mergesort or heapsort. Comparing these algorithms is an interesting test case for our tool, F2SAD.

The theoretical complexity limit to sorting is $O(n \log n)$, but also in this case the practical knowledge of constraints to the input of the sorting algorithms can be exploited. To this albeit linear order sorting algorithms have been developed. These include radix, counting and bucketsort, in which the latter is sometimes used as a classifier or find-algorithm. These sorting algorithms do not base their algorithm on comparisons, like the sorting algorithms earlier mentioned, but rather on classification or decision trees.

The predicted times in these sections were produced using F2SAD which will now also incorporate the Parasol II tool in the same program. The tool produces a time-complexity formula in which the machine constants and the algorithmic parameters are still abstract. When not specified, we have used a Sun Classic LX model as test machine.

```
Memory          : 32 MB
Model           : SPARCstation LX
At frequency    : 50 MHz
CPU             : Texas Instruments TMS390S10 (MicroSparc)
Data cache      : 2 Kb blocksize=16 1-way associative
Instruction cache : 4 Kb blocksize=32 1-way associative
```

2.1 Randomly distributed input data

Figure 2.2 shows the measured execution timings of the exponential and $O(n \log n)$ performance as well as the predicted execution timing. The input to the sorting algorithms is uniformly randomly distributed. Also the parameters in the time complexity formula have been set in such a way that they reflect this condition.

The reason for taking random input data is obviously that sorting algorithms will then expose an almost average-case behaviour. Some algorithms have a worst-case behaviour which is of a different complexity order or, less dramatic, have different minor terms in the complexity formula.

How the parameters are set we will come to later, but first we have a look at the measured and estimated execution time and compare them. Figure 2.3 shows the expected and estimated execution time for the bucketsort algorithm, compared to a mergesort implementation, while figure 2.2 gives the same data for the other sorting algorithms.

Figure 2.4 shows the error of the predicted versus the expected value.

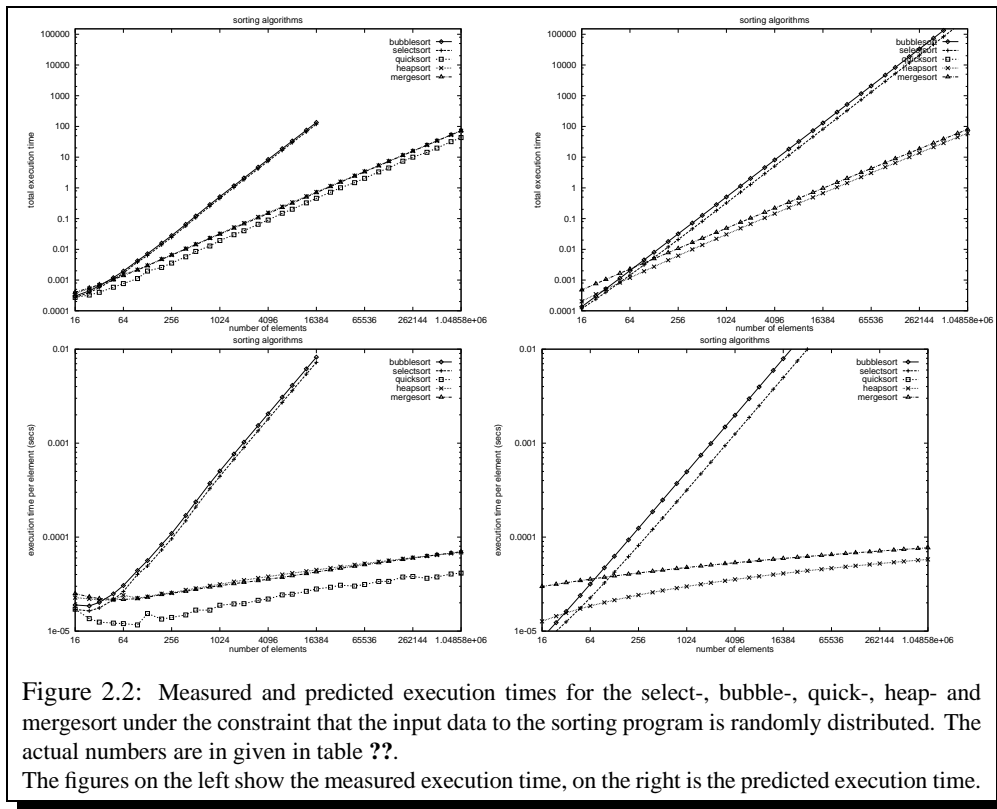


Figure 2.2: Measured and predicted execution times for the select-, bubble-, quick-, heap- and mergesort under the constraint that the input data to the sorting program is randomly distributed. The actual numbers are in given in table ??.

The figures on the left show the measured execution time, on the right is the predicted execution time.

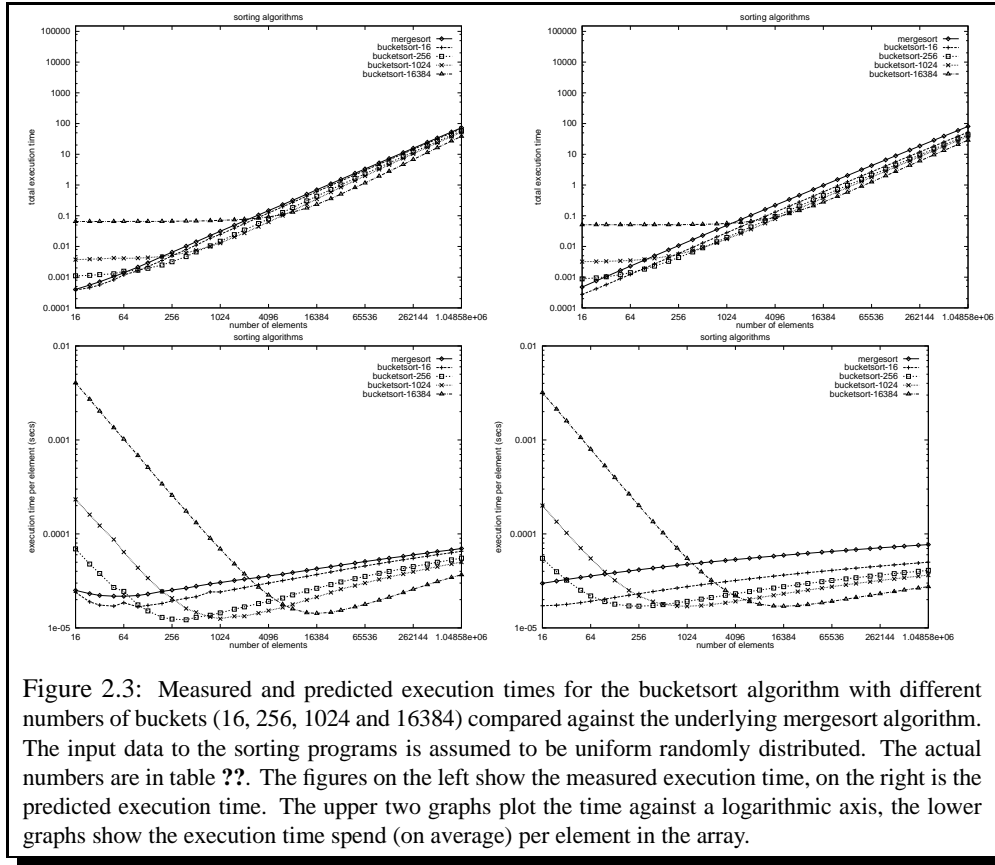


Figure 2.3: Measured and predicted execution times for the bucketsort algorithm with different numbers of buckets (16, 256, 1024 and 16384) compared against the underlying mergesort algorithm. The input data to the sorting programs is assumed to be uniform randomly distributed. The actual numbers are in table ?? . The figures on the left show the measured execution time, on the right is the predicted execution time. The upper two graphs plot the time against a logarithmic axis, the lower graphs show the execution time spend (on average) per element in the array.

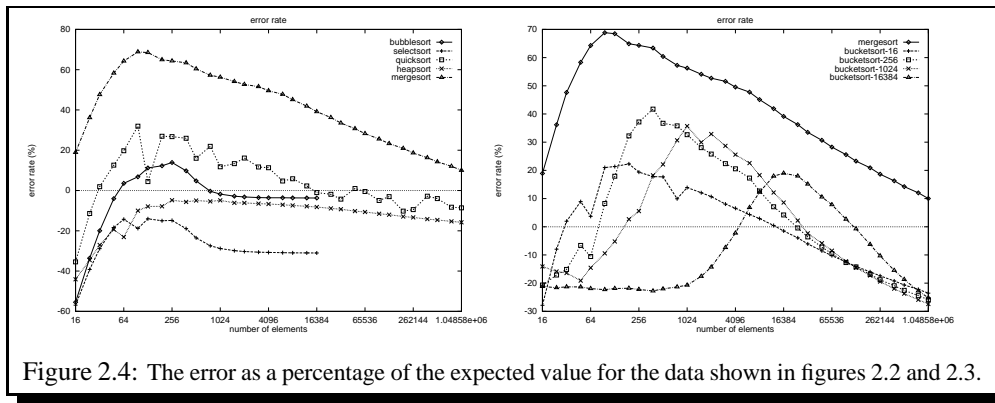


Figure 2.4: The error as a percentage of the expected value for the data shown in figures 2.2 and 2.3.

2.2 Setting the parameters

For the algorithmic parameters —the number of times loops and conditions are taken— the F2SAD/ Parasol II toolset has basically two ways actualizing. One is by using a profile files to determine the parameters and the other is by defining them by hand. The programs which have been analyzed in this report have been deliberately analyzed by hand. Below we give an example of theoretical complexity parameters for the selection sorting algorithm.

2.2.1 Selection Sort

In appendix A, the algorithm considered here, can be found. The number of times the outer most loop at line 9 is executed is clearly $n - 1$, the size of the array to be sorted. The inner loop at line 12 has different properties. In the first iteration of the outermost loop it is executed $n - 2$ times, the second time $n - 3$ times continuing until it is executed only once. This leads to the summation:

$$\sum_{i=1}^{n-2} i = \frac{1}{2}(n-2)((n-2)+1) = \frac{1}{2}(n-1)(n-2)$$

Since the outer loop iterates $n - 1$ times, the inner loop will iterate $\frac{1}{2}(n - 2)$ times on average.

The conditional on line 13 is true whenever an element a_i in the list $a_0, a_1 \dots a_n$ is smaller than all its predecessors ($a_{i-1}, a_{i-2} \dots a_0$) For $i = 0$ the conditional is always true, for $i = 1$ with a uniform random list this will be $\frac{1}{2}$, for $i = 2$ it will be $\frac{1}{4}$. We will not go into any detail, but the the idea between this logic is that each predecessor has a probability of $\frac{1}{2}$ to be smaller and in this way each predecessor will half the chance that the element a_i is smaller than all its predecessors.

Each element a_i will be subject to the conditional i times, which leads to the following formula for the chance that the condition evaluates to true:

$$\frac{\frac{1}{i} + \frac{1}{i-1} + \dots + \text{frac}11}{n - i}$$

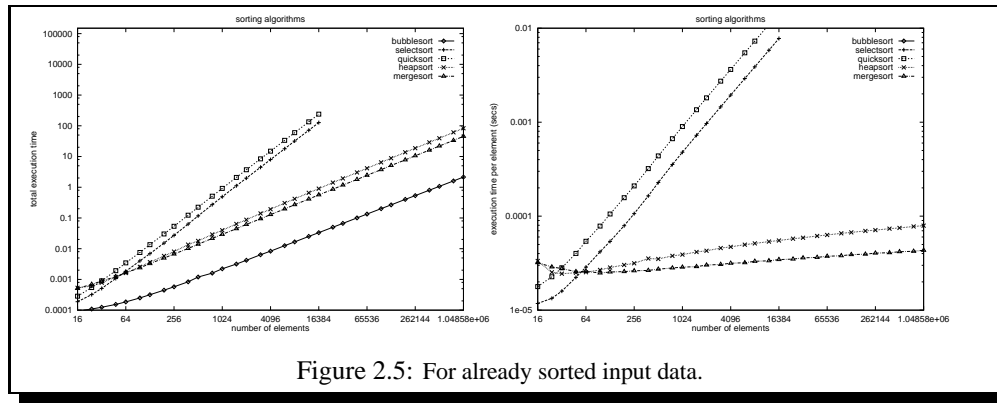
The nominator approaches 2 thus we get:

$$\sum_{i=0}^{n-1} \frac{2}{n - i} = 2 \ln(n) + C$$

In which we can ignore C , and since there are n elements we have to divide this by n

To recapulate we are left with the following parameters:

N.1	$n-1$	(the outer loop)
N.2	$0.5*(n-2)$	(the inner loop)
P.1	$2*\ln(n)/n$	(the constant)



We have used here the notation $N . x$ and $P . x$, for the control flow parameters, which is also used by F2SAD. The numbers after the $N .$ and $P .$ have no real meaning, but are distributed according to the flow of the program. They are the same each time the program is run through F2SAD.

For sorted data the conditional $P . 2$ will be nearly 0 (actually it will be $\frac{n}{\frac{1}{2}(n-1)(n-2)}$, all other parameters remain the same.

Figure 2.5 gives results in the hypothetical case that all input data is sorted. In that case obviously for example the bubblesort algorithm displays a very friendly execution time behaviour.

2.3 A note on linear sorting

As was mentioned above the linear order sorting algorithms use some other sorting algorithm to sort the classes they have build. As we have seen, the overhead and the usage of an other algorithm do not make it attractive for sorting purposes, since it is only very slightly better than the underlying sort. But, the linear order sorting algorithms have also a very different purpose. If it is necessary to classify the input in ranges, resulting in a list of only roughly sorted lists, there is no need for the underlying comparison sort mechanism. And therefor these algorithms have their separate usefulness, especially in parallel computers in which data has be redistributed. The bucket method can be used to classify the data, and to distribute each class to a processor.

Appendix A

Source code

This appendix includes all the source code of the algorithms studied in this report. The main program is not included since it is generated in order to provide multiple input data sets.

A.1 Bubblesort

```
1      SUBROUTINE bubblesort(usize, a)
2      IMPLICIT NONE
3      INTEGER usize
4      DOUBLE PRECISION a
5      DIMENSION a(*)
6      DOUBLE PRECISION swap
7      INTEGER i, size
8      LOGICAL flag
9
10     size = usize
11 10   flag = .FALSE.
12     DO 20, i=1, size-1
13         IF(a(i) .GT. a(i+1)) THEN
14             PRINT *, i, a(i), a(i+1)
15             swap = a(i)
16             a(i) = a(i+1)
17             a(i+1) = swap
18             flag = .TRUE.
19         END IF
20 20   CONTINUE
21     size = size - 1
22     IF(flag) GOTO 10
23     END
```

A.2 Selectsort

```
1      SUBROUTINE selectsort(size, a)
2      IMPLICIT NONE
3      INTEGER size
4      DOUBLE PRECISION a
5      DIMENSION a(*)
6      DOUBLE PRECISION smallest
7      INTEGER i, j, index
8
9      DO 20, i=1, size-1
10     index = i
11     smallest = a(index)
12     DO 10, j=i+1, size
13         if(smallest .GT. a(j)) THEN
14             index = j
15             smallest = a(index)
16         END IF
17 10     CONTINUE
18     a(index) = a(i)
19     a(i) = smallest
20 20     CONTINUE
21     END
```

A.3 Heapsort

The heapify routine is the key to the heapsort algorithm. The parameters to the heapify routine are an array A and an index i into that array. The precondition for the heapify routine is that the left binary subtree and the right binary subtree are both heaps. $A(i)$ however may be larger than the elements in both subtrees, thus violating the heap property. The heapify routine will “sift down” this element $A(i)$ and by this way both subtrees and $A(i)$ will become one larger heap.

```

 $l \leftarrow \text{Left}(i)$ 
 $r \leftarrow \text{Right}(i)$ 
if  $l \leq \text{HeapSize}[A]$  and  $A[l] > A[i]$ 
    then  $\text{largest} \leftarrow l$ 
    else  $\text{largest} \leftarrow i$ 
if  $r \leq \text{HeapSize}[A]$  and  $A[r] > A[\text{largest}]$ 
    then  $\text{largest} \leftarrow r$ 
if  $\text{largest} \neq i$ 
    then exchange  $A[i] \leftrightarrow A[\text{largest}]$ 
        Heapify( $A, \text{largest}$ )

```

```

1
2     SUBROUTINE heapify(size, a, parent)
3     IMPLICIT NONE
4     c     left(index) = index*2
5     c     right(index) = index*2 + 1
6     DOUBLE PRECISION a
7     INTEGER size, parent
8     DIMENSION a(*)
9     INTEGER i, l, r, largest
10    DOUBLE PRECISION swap
11
12    i = parent
13    10   l = i*2                                left(i)
14        r = i*2+1                              right(i)
15    IF ((l .LE. size) .AND. (a(l) .GT. a(i))) THEN
16        largest = l
17    ELSE
18        largest = i
19    END IF
20    IF ((r .LE. size) .AND. (a(r) .GT. a(largest))) THEN
21        largest = r
22    END IF
23    IF (largest .NE. i) THEN
24        a(i) = a(largest)
25        a(largest) = swap
26        i = largest
27        GOTO 10
28    END IF
29
30    END
31

```

Most paramters of the heapsort are

```

32   for  $i \leftarrow \frac{Size[a]}{2}$  downto 1
      do   Heapify(A,i)

33   SUBROUTINE buildheap(size, a)
34   IMPLICIT NONE
35   DOUBLE PRECISION a
36   INTEGER size
37   DIMENSION a(*)
38   INTEGER i
39
40   DO 10, i=size/2, 1, -1
41       CALL heapify(size, a, i)
42   10 CONTINUE
43
44   END
45
46
47   BuildHeap(A)
48   for  $i \leftarrow length[A]$  downto 2
49       do   exchange A[1]  $\leftrightarrow$  A[i]
50           decrease HeapSize by 1
51           Heapify(A,1)
52
53   SUBROUTINE heapsort(asize, a)
54   IMPLICIT NONE
55   INTEGER asize
56   DOUBLE PRECISION a
57   DIMENSION a(*)
58   DOUBLE PRECISION swap
59   INTEGER i, size
60
61   size = asize
62   CALL buildheap(size, a)
63   DO 10, i=size, 2, -1
64       swap = a(1)
65       a(1) = a(i)
66       a(i) = swap
67       size = size - 1
68       CALL heapify(size, a, 1)
69   10 CONTINUE
70
71   END

```

SX/2

SX-1

A.4 Mergesort

```

1      SUBROUTINE mergelsort(a, lstptr, head, tail)
2      IMPLICIT NONE
3      INTEGER stacksize
4      PARAMETER (stacksize = 256)
5      INTEGER lstptr, head(*), tail(*)
6      DOUBLE PRECISION a(*)
7      INTEGER stackindex, stack(stacksize)
8      INTEGER size, list1, list2, run, hsize, x, y, z
9
10     x = 0
11     y = 0
12     z = 0
13
14     size = 0
15     run = lstptr
16 10   IF(run .GT. 0) THEN
17         size = size + 1
18         run = tail(run)
19         GOTO 10
20     END IF
21     sv = size
22
23     stack(1) = 0
24     stackindex = 2
25
26 1    IF(size .LE. 1) GO TO 4
27
28     list1 = lstptr
29     list2 = lstptr
30     hsize = size/2
31 20   IF(hsize .GT. 0) THEN
32         hsize = hsize - 1
33         list2 = tail(list2)
34         GOTO 20
35     END IF
36
37     stack(stackindex) = size
38     stack(stackindex+1) = list2
39     stack(stackindex+2) = 1
40     stackindex = stackindex + 3
41     size = size/2
42     lstptr = list1
43     GO TO 1
44 2    list1 = lstptr
45     stackindex = stackindex - 3
46     size = stack(stackindex)
47     list2 = stack(stackindex+1)
48
49     stack(stackindex) = size
50     stack(stackindex+1) = list1
51     stack(stackindex+2) = 2
52     stackindex = stackindex + 3
53     size = size - size/2
54     lstptr = list2
55     GO TO 1
56 3    list2 = lstptr
57     stackindex = stackindex - 3
58     size = stack(stackindex)
59     list1 = stack(stackindex+1)
60

```

```

61
62     IF(a(head(list1)) .LT. a(head(list2))) THEN
63         lstptr = list1
64     ELSE
65         lstptr = list2
66     END IF
67
68     run = 0
69 30   IF(list1 .GT. 0 .AND. list2 .GT. 0) THEN
70         z = z + 1
71         IF(a(head(list1)) .LT. a(head(list2))) THEN
72             IF(run .GT. 0) THEN
73                 tail(run) = list1
74             ELSE
75                 lstptr = list1
76             END IF
77             run = list1
78             list1 = tail(list1)
79         ELSE
80             IF(run .GT. 0) THEN
81                 tail(run) = list2
82             ELSE
83                 lstptr = list2
84             END IF
85             run = list2
86             list2 = tail(list2)
87         END IF
88         GOTO 30
89     END IF
90     IF(list1 .GT. 0) THEN
91         IF(run .GT. 0) THEN
92             tail(run) = list1
93         ELSE
94             lstptr = list1
95         ENDIF
96     ELSE IF(list2 .GT. 0) THEN
97         IF(run .GT. 0) THEN
98             tail(run) = list2
99         ELSE
100            lstptr = list2
101        ENDIF
102    END IF
103
104
105 4   IF(size .EQ. 1) tail(lstptr) = 0
106
107    F2C isn't able to process vector-if statements, that is why the following
108    IF-statement is commented out and two replacement IF's are dropped
109    in.
110    y = y + 1
111    IF(stack(stackindex-1).EQ.1) GO TO 2
112    x = x + 1
113    IF(stack(stackindex-1).EQ.2) GO TO 3
114
115    END

```

A.5 Quicksort

```

1      SUBROUTINE quicksort(asize, a)
2      INTEGER stacksize
3      PARAMETER (stacksize = 256)
4      INTEGER asize
5      DOUBLE PRECISION a
6      DIMENSION a(*)
7      DOUBLE PRECISION aux
8      INTEGER start, size, front, back, stack, idx
9      DIMENSION stack(stacksize)
10
11     size = asize
12     idx = 0
13     start = 1
14 20    IF(size .GT. 1) THEN
15         front = start+1
16         back = start+size-1
17
18 30    IF(front .LE. back) THEN
19         IF (a(start) .GT. a(front)) THEN
20             front = front + 1
21         ELSE IF(a(start) .LE. a(back)) THEN
22             back = back - 1
23         ELSE
24             aux = a(front)
25             a(front) = a(back)
26             a(back) = aux
27         END IF
28         GOTO 30
29     END IF
30
31     IF(front .NE. start+1) THEN
32         aux = a(start)
33         a(start) = a(front-1)
34         a(front-1) = aux
35         stack(idx+1) = front - start - 1
36         stack(idx+2) = start
37         idx = idx + 2
38     END IF
39     size = size - back + start - 1
40     start = front
41     GOTO 20
42 END IF
43
44 IF(idx .GT. 0) THEN
45     idx = idx - 2
46     size = stack(idx+1)
47     start = stack(idx+2)
48     GOTO 20
49 END IF
50
51 END

```


A.6 Bucketsort

```

1
2     SUBROUTINE bucketlsort(a, lstptr, head, tail)
3     IMPLICIT NONE
4     DOUBLE PRECISION a(*)
5     INTEGER lstptr, head(*), tail(*)
6     INTEGER numbuckets
7     PARAMETER (numbuckets = 16384)
8     INTEGER buckets(numbuckets), bucketnum, i, aux
9
10    DO 10, i=1, numbuckets
11        buckets(i) = 0
12    10 CONTINUE
13    20 IF(lstptr .GT. 0) THEN
14        bucketnum = INT(a(head(lstptr)) * numbuckets) + 1
15        aux          = tail(lstptr)
16        tail(lstptr) = buckets(bucketnum)
17        buckets(bucketnum) = lstptr
18        lstptr          = aux
19        GO TO 20
20    END IF
21
22    DO 30, i=1, numbuckets
23        CALL mergelsort(a, buckets(i), head, tail)
24    30 CONTINUE
25
26    lstptr = 0
27    DO 40, i=1, numbuckets
28        IF(buckets(i) .GT. 0) THEN
29            IF(lstptr .EQ. 0) THEN
30                lstptr = buckets(i)
31            ELSE
32                tail(aux) = buckets(i)
33            END IF
34            aux = buckets(i)
35    50    IF(tail(aux)) THEN
36            aux = tail(aux)
37            GO TO 50
38        END IF
39    END IF
40    40 CONTINUE
41
42    END
43

```