

**13** (100 PTS.) Not a sorting network.

You are given an array  $A[1 \dots n]$ , but you can not access it directly (or even read the values in it, or compare them directly). Instead, you are give some procedures that can access the array and do certain operations. Your task is to sort the array.

- 13.A.** (30 PTS.) You are given a procedure **sortBlock**( $i$ ), which sorts (in increasing order and in place), all the elements in  $A[i \dots i + u]$ , where  $u \geq 1$  is a prespecified fixed parameter (i.e.,  $u$  is a known fixed number between 1 and  $n$ , but it is not under your control). For example, for the following array:

2.71828	18284	59045	23536	02874	7.1352	66249
1	2	3	4	5	6	7

With  $u = 3$ , the call **sortBlock**(2) would result in:

2.71828	02874	18284	23536	59045	7.1352	66249
1	2	3	4	5	6	7

Describe an algorithm, that uses  $O((n/u)^2)$  calls to **sortBlock**, and sorts the array. What is the running time of your algorithm if calling **sortBlock** takes  $O(1)$  time?

**Solution:**

The idea is to simulate insertion sort. Let  $\Delta = \lfloor (u + 1)/2 \rfloor$ .

```

for  $i = \lfloor n/\Delta \rfloor, \dots, 0$  do
  for  $j = 0, \dots, j$  do
     $s \leftarrow \min(1 + j\Delta, n - u)$ 
    sortBlock( $A[s, \dots, s + u]$ )

```

The bound on the number of calls to **sortBlock** is easy, since it is  $O((n/\Delta)^2) = O(n^2/u^2)$ .

As for correctness, it is easy to argue (by easy but maybe boring induction) that in the end of the  $t$ th outer iteration, the last  $t\Delta$  elements of the array contains the  $t\Delta$  largest elements of the array. Since these elements are in the right order, it follows that in the end of the execution of the algorithm it would be sorting all the elements correctly.

- 13.B.** (50 PTS.) Congratulations! You just got a better sorting primitives **bMerge**, **copy**, and a work array  $W[1, \dots, n]$ .

- (i) **copy** can copy any block of length at most  $u + 1$  between the two arrays (or inside them).
- (ii) **bMerge** is weirder. It takes two blocks  $L$  and  $U$  (both with at most  $u + 1$  elements), treat them as a single block, sort the unified block, and writes the smaller  $|L|$  elements (in sorted order) to  $L$ , and the other (larger)  $|U|$  elements in sorted order to  $U$  (the two blocks do not have to be of the same length).

For example, for

2.71828	18284	59045	23536	02874	7.1352	66249
1	2	3	4	5	6	7

With  $u = 1$ , the call **bMerge**( $A[2 \dots 3]$ ,  $A[6 \dots 7]$ ) would result in:

2.71828	7.1352	18284	23536	02874	59045	66249
1	2	3	4	5	6	7

**bMerge** also returns the number of elements in original block  $L$  that are still in  $L$  after the operation is completed. In this example, since 18284 was in  $L$  before the operation, and it is in  $L$  after the operation is completed, the returned value would be 1.

Note that the blocks given to **bMerge** can not overlap.

Assume that  $A[1 \dots n/2]$  and  $A[n/2 + 1, \dots n]$  are already sorted ( $n$  is even). Describe an algorithm that performs a minimal total number of calls to **sortBlock**, **copy** and **bMerge**, and sorts the array  $A$ . What is the running time of your algorithm if calling **sortBlock**, **copy** and **bMerge** takes  $O(1)$  time? (Prove your bound.)

## Solution:

The idea is of course to implement merge sort. To this end, we copy the prefix, of size  $u$ , of the two parts being merged into the working array, we sort them using a call to **bMerge**. We then consider the first  $u$  elements in the sorted block to be the smallest  $u$  elements in the parts of the arrays yet to be sorted. We add these  $u$  elements to the output, and step forward the two input arrays, according to how many elements were used in the output block. Here is the pseudo-code – we use  $B$  here to make the code clearer (but in reality it would just be a different part of  $A$ ).

The code handles carefully the case where the number of remaining elements is smaller than  $u$ , thus the code looks a bit messier than desired.

```

mSort(  $A[1, \dots, N], B[1 \dots, M]$  ):
   $a \leftarrow 1$ 
   $b \leftarrow 1$ 
   $o \leftarrow 1$ 
  while ( $a \leq N$  and  $b \leq M$ ) do
     $a' = \min(i + u, N)$ 
     $\Delta_A = a' - a + 1$ 
     $b' = \min(i + u, M)$ 
     $\Delta_B = b' - b + 1$ 
    if  $\Delta_A \geq \Delta_B$  then
       $o' = o + \Delta_A$ 
      copy(  $A[a, \dots, a'] \Rightarrow W[o, \dots, o + \Delta_A - 1]$ )
      copy(  $B[b, \dots, b'] \Rightarrow W[o', \dots, o' + \Delta_B - 1]$ )
       $\lambda \leftarrow \mathbf{bMerge}(W[o, \dots, o + \Delta_A - 1], W[o', \dots, o' + \Delta_B - 1])$ 
       $\lambda' \leftarrow \Delta_A - \lambda$ 
       $a \leftarrow a + \lambda$ 
       $b \leftarrow b + \lambda'$ 
       $o \leftarrow o + \Delta_A$ 
    else
      OtherCase
      if  $\Delta_A < u$  and  $\Delta_B < u$  then break loop
  Using  $\lceil (N + M)/u \rceil$  calls to copy, copy the sorted array back into  $A$  and  $B$ .

```

**OtherCase:**

$$o' = o + \Delta_B$$

**copy**(  $B[b, \dots, b'] \Rightarrow W[o, \dots, o + \Delta_B - 1]$ )

**copy**(  $A[a, \dots, a'] \Rightarrow W[o', \dots, o' + \Delta_A - 1]$ )

$\lambda \leftarrow \mathbf{bMerge}(W[o, \dots, o + \Delta_B - 1], W[o', \dots, o' + \Delta_A - 1])$

$$\lambda' \leftarrow \Delta_B - \lambda$$

$$b \leftarrow b + \lambda$$

$$a \leftarrow a + \lambda'$$

$$o \leftarrow o + \Delta_B$$

Clearly, in each outer iteration, the algorithm add at least  $u$  elements to the output array. In the last iteration, the two remaining parts are small enough, that they are now sorted correctly, an the algorithm can bail out.

As such, the running time of the algorithm

$$\mathbf{mSort}(A[1 \dots n/2], A[n/2 + 1, \dots n])$$

is  $O(n/u + 1)$ , and this also bounds the number of calls to the primitives provided.

- 13.C.** (20 PTS.) Building on (B) and expanding on it, describe a sorting algorithm using these primitives that sort the given array  $A$  (that is initially not sorted). What is the running time of your algorithm if calling **sortBlock**, **copy** and **bMerge** takes  $O(1)$  time? Naturally, the faster the better (Prove your bound).

### Solution:

Just implement merge sort. If the subarray to be sorted is of size smaller than  $u$ , then sort directly using **sortBlock**, otherwise, break it into two subarrays, sort them recursively, and merge them using the procedure from (B). The running time is

$$T(n) = 2T(n/2) + O(n/u + 1),$$

where  $T(n) = O(1)$  if  $n \leq u$ . As such, the solution to this recurrence is  $O((n/u + 1)(1 + \log(n/u)))$  (which is essentially  $O((n/u) \log(n/u))$ ), and this bounds the running time and number of calls to the primitives.

## **14** (100 PTS.) Not a sorting question.

Consider an array  $A[0 \dots n - 1]$  with  $n$  distinct elements. Each element is an  $\ell$  bit string representing a natural number between 0 and  $2^\ell - 1$  for some  $\ell > 0$ . The only way to access any element of  $A$  is to use the function **FetchBit**( $i, j$ ) that returns the  $j$ th bit of  $A[i]$  in  $O(1)$  time.

- 14.A.** (20 PTS.) Suppose  $n = 2^\ell - 1$ , i.e. exactly one of the  $\ell$ -bit strings does not appear in  $A$ . Describe an algorithm to find the missing bit string in  $A$  using  $\Theta(n \log n)$  calls to **FetchBit** without converting any of the strings to natural numbers.

## Solution:

The idea is to simply count the number of '1's and '0's in the  $j$ th bit of all strings. If all strings were there, we would have  $2^{\ell-1}$  ones and  $2^{\ell-1}$  zeros. Whichever is less, is the missing bit. We will store the bits of the missing string in  $S[0, \dots, \ell - 1]$ .

```
for  $j = 0, \dots, \ell - 1$  do
     $c = 0$ 
    for  $i = 0, \dots, n - 1$  do
         $c = c + \text{FetchBit}(i, j)$ 
    if  $c < 2^{\ell-1}$ 
         $S[j] = 1$ 
    else
         $S[j] = 0$ 
```

The algorithm calls **FetchBit**  $\ell \times n$  times and  $\ell = \log n$ . Hence,  $O(n \log n)$  calls.

- 14.B. (40 PTS.) Suppose  $n = 2^\ell - 1$ . Describe an algorithm to find the missing bit string in  $A$  using only  $O(n)$  calls to **FetchBit**.

## Solution:

We will use the same the idea as the first part but we will keep track of the previous bit and recursively search half of the array for the next missing bit.

```
findMissing( $A[0, \dots, n - 1], B[0, \dots, 2^{m+1} - 2], S[0, \dots, \ell - 1], m$ )
    if  $m \geq 0$ 
         $z = 0$ 
         $o = 0$ 
        for  $i = 0, \dots, 2^{m+1} - 2$  do
            if FetchBit( $B[i], m$ ) = 1
                 $B_1[o] = B[i]$ 
                 $o = o + 1$ 
            else
                 $B_0[z] = B[i]$ 
                 $z = z + 1$ 
        if  $o < 2^m$ 
             $S[m] = 1$ 
            findMissing( $A, B_1, S, m - 1$ )
        else
             $S[m] = 0$ 
            findMissing( $A, B_0, S, m - 1$ )
```

To find the missing bit string, we call the above algorithm with the following input **findMissing**( $A, [0, \dots, n - 1], S, \ell - 1$ ). The below figure shows a run of the algorithm for  $\ell = 3$  and  $A$  missing the string 011.

0 1 0	0 1 0	0 1 0
1 1 1	1 1 1	1 1 1
1 1 0	1 1 0	1 1 0
0 0 1	0 0 1	0 0 1
1 0 0	1 0 0	1 0 0
0 0 0	0 0 0	0 0 0
1 0 1	1 0 1	1 0 1
↑	↑	↑
$m = 2$	$m = 1$	$m = 0$
$B = [0, \dots, 6]$	$B = [0, 3, 5]$	$B = [0]$
$z = 3, o = 4$	$z = 2, o = 1$	$z = 1, o = 0$
$S[0] = 0$	$S[1] = 1$	$S[2] = 1$

The number of calls to **FetchBit** is  $2^{m+1} - 1$  and it decreases by half with every recursive call starting from  $n$  in the first recursive call. Hence, the number of calls to **FetchBit** follows the following recurrence:

$$\begin{aligned}
 T(n) &= T(n/2) + n \\
 &= T(n/4) + n/2 + n \\
 &= 1 + 2 + 4 + \dots + n/2 + n \\
 &= O(n)
 \end{aligned}$$

- 14.C. (40 PTS.) Suppose  $n = 2^\ell - k$ , i.e. exactly  $k$  of the  $\ell$ -bit strings do not appear in  $A$ . Describe an algorithm to find the  $k$  missing bit strings in  $A$  using only  $O(n \log k)$  calls to **FetchBit**.

### Solution:

We will again use the same idea. However, instead of recursing on one branch (either zeros or ones), we will recurse on both as long as there is a missing bit. We will use a  $k \times \ell$  array  $S$  to store the missing strings and an array  $C$  to keep track of which missing strings each recursive call is looking searching for. To find the  $k$  missing strings, we call the below algorithm with the following input **findMissingk**( $A, [0, \dots, n-1], S, [0, \dots, k-1], \ell-1$ ). The figure below also shows a run of the algorithm for  $\ell = 3$ ,  $k = 3$ , and  $A$  missing the strings 110, 011, 001.

It might be tempting to assume that **findMissingk** recursively calls itself twice on half the input size leading to the following recurrence on the number of calls to **FetchBits**:

$$T(n) \leq 2T(n/2) + O(n)$$

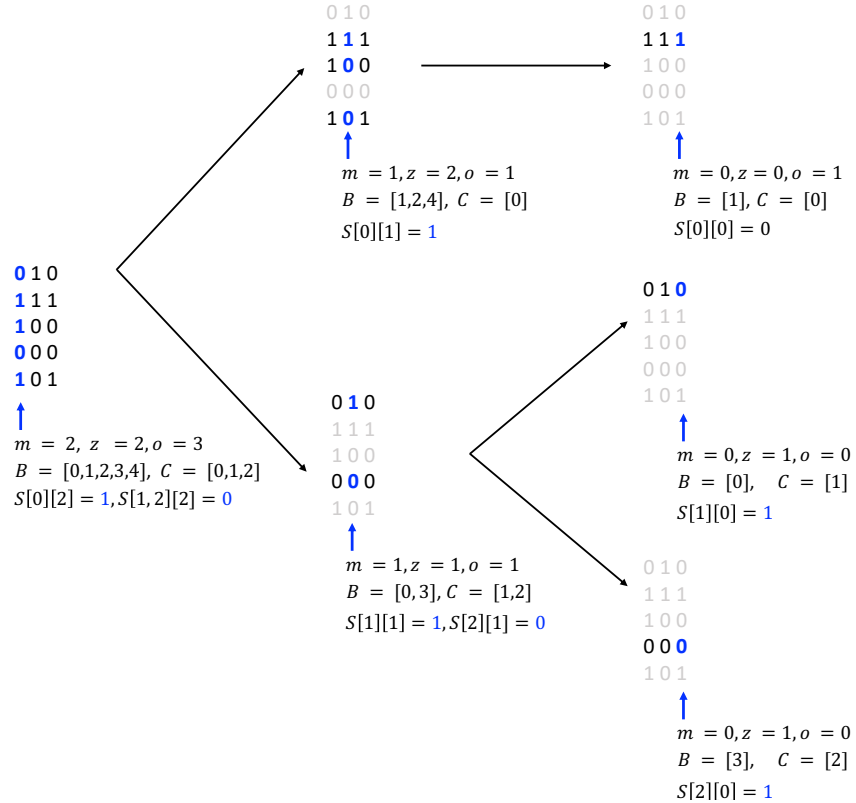
which leads on an  $O(n \log n)$  bound. However, this bound is not tight. Notice that the recursion tree for **findMissingk** will have at most  $k$  nodes at any level of the recursion. Hence, at level  $j$ , the total work is at most:  $\min\{2^j, k\} \times n/(2^j)$ . Thus,

$$\begin{aligned}
 T(n) &\leq \sum_{j=0}^{\ell-1} \min\{2^j, k\} \times n/(2^j) = \sum_{j=0}^{\log k} n + \sum_{j=\log k+1}^{\ell-1} kn/2^j = n \log k + (n - k) \\
 &\leq O(n \log k)
 \end{aligned}$$

```

findMissing( $A[0, \dots, n-1], B[0, \dots, 2^{m+1}-1-k], S, C[0, \dots, k-1], m$ )
  if  $m \geq 0$ 
     $z = 0$ 
     $o = 0$ 
    for  $i = 0, \dots, 2^{m+1}-1-k$  do
      if FetchBit( $B[i], m$ ) = 1
         $B_1[o] = B[i]$ 
         $o = o + 1$ 
      else
         $B_0[z] = B[i]$ 
         $z = z + 1$ 
     $k_1 = 2^m - o$ 
     $k_0 = 2^m - z$ 
    if  $k_1 > 0$ 
      for  $i = 0, \dots, (k_1 - 1)$  do
         $S[C[i]][m] = 1$ 
         $C_1[i] = C[i]$ 
      findMissing( $A, B_1, S, C_1, m-1$ )
    if  $k_0 > 0$ 
      for  $i = 0, \dots, (k_0 - 1)$  do
         $S[C[i+k_1]][m] = 0$ 
         $C_0[i] = C[i+k_1]$ 
      findMissing( $A, B_0, S, C_0, m-1$ )

```



**Alternative Proof:**

We can write the two parameter recurrence for the number of calls to **FetchBits** as:

$$T(n, k) = T(n/2, k_1) + T(n/2, k - k_1) + n$$

where  $k_1$  is the number of strings in which the current missing bit is 1 and  $0 \leq k_1 \leq k \leq n$ . Observe that  $T(n, 0) = 0$ ,  $T(n, 1) = O(n)$  (from part B), and  $T(1, 1) = O(1)$ .

We will prove that  $T(n, k) = O(n \log k)$  for  $k \geq 1$  this by induction on  $k$ .

**Base Case:**  $k = 1$

$$T(n, 1) = O(n) = O(n \log c) \text{ for some constant } c > 1.$$

**Inductive Hypothesis:** Assume for any  $1 \leq k' < k$ ,  $T(n, k') = O(n \log k')$

**Proof:**

$$T(n, k) = T(n/2, k_1) + T(n/2, k - k_1) + n$$

If  $1 \leq k_1 < k$ , then we can apply by inductive hypothesis:

$$\begin{aligned} T(n, k) &\leq O(n \log(k_1)) + O(n \log(k - k_1)) + n \\ &\leq O(n \log(k_1(k - k_1))) \\ &\leq O(n \log(k^2)) \\ &\leq O(n \log k) \end{aligned}$$

If  $k_1 = 0$  or  $k_1 = k$ , then  $T(n, k) = T(n/2, k) + O(n)$ . However, since the missing strings are distinct, there must exist a different bit that will trigger two recursive calls at some level of the recursion. Let  $0 \leq j \leq \ell - \log k - 1$  be the first level that triggers two recursive calls. Then,

$$\begin{aligned} T(n, k) &= n + n/2 + \dots + n/2^j + T(n/2^{j+1}, k'_1) + T(n/2^{j+1}, k - k'_1) \\ &\leq O(n) + O(n/2^{j+1} \log k) \\ &\leq O(n \log k) \end{aligned}$$

## 15 (100 PTS.) Don't want to walk too much.

You are given a set of  $n$  distinct points on a line with x-coordinates  $x_1, x_2, \dots, x_n$ . The points are not sorted and their values are stored in an array  $X$  where  $X[i] = x_i$ . Each point is associated with a positive weight  $w_i$  such that  $\sum w_i = 1$ . The weights are also stored in an array  $W$  where  $W[i] = w_i$ . Our goal is to find  $x_j$  that minimizes the weighted distance given by:  $\sum_i w_i |x_j - x_i|$ .

**15.A.** (10 PTS.) (Easy.) Show that if all the weights are equal,  $x_j$  is the median of  $X$ .

### Solution:

Since all the weights are equal minimizing the weighted distance is the same as minimizing:  $p(x) = \sum_i |x - x_i|$ . Let  $x_M$  be the median. Assume the point that minimizes  $p(x)$  is  $x_j \neq x_M$ . Then,  $p(x_M) - p(x_j) > 0$ . Without loss of generality, assume  $x_j > x_M$ . Let  $\Delta = x_j - x_M$  then with perseverance we can write:

$$\begin{aligned}
 p(x_M) - p(x_j) &= \sum_{x_i \leq x_M} (x_M - x_i) + \sum_{x_i > x_M} (x_i - x_M) - \sum_{x_i \leq x_j} (x_j - x_i) - \sum_{x_i > x_j} (x_i - x_j) \\
 &= - \sum_{x_i \leq x_M} \Delta + \sum_{x_i > x_j} \Delta + \sum_{x_M < x_i \leq x_j} (x_i - x_M) - \sum_{x_M < x_i \leq x_j} (x_j - x_i) \\
 &\leq - \sum_{x_i \leq x_M} \Delta + \sum_{x_i > x_M} \Delta - \sum_{x_M < x_i \leq x_j} (x_j - x_i) \\
 &\leq 0 - \sum_{x_M < x_i \leq x_j} (x_j - x_i) \\
 &< 0
 \end{aligned}$$

which is a contradiction! You can get a similar result by assuming  $x_j < x_M$ .

**15.B.** (20 PTS.) In general, show that  $x_j$  is the point that satisfies the following property:

$$\sum_{x_i < x_j} w_i < \frac{1}{2} \quad \text{and} \quad \sum_{x_i > x_j} w_i \leq \frac{1}{2}.$$

### Solution:

We will first sort the points  $x_i$  and rearrange the weights accordingly. (This step is not necessary. However, it simplifies the notation in our proof.)

Let  $x_j$  be the point that minimizes the weighted distance.  $p(x) = \sum_i w_i |x - x_i|$ . Then,  $p(x_j) - p(x_{j+1}) < 0$  and  $x_{j+1} > x_j$ .

$$\begin{aligned}
 p(x_j) - p(x_{j+1}) &= \sum_i w_i |x_j - x_i| - \sum_i w_i |x_{j+1} - x_i| \\
 &= \sum_{i < j} w_i (x_j - x_i) + \sum_{i \geq j} w_i (x_i - x_j) - \sum_{i < j+1} w_i (x_{j+1} - x_i) - \sum_{i \geq j+1} w_i (x_i - x_{j+1}) \\
 &= \left( \sum_{i \geq j} w_i - \sum_{i < j} w_i \right) (x_{j+1} - x_j) < 0
 \end{aligned}$$

$$\text{Thus, } \sum_{i \geq j} w_i - \sum_{i < j} w_i < 0 \quad \text{and} \quad \sum_i w_i = 1 \Rightarrow \sum_{i > j} w_i \leq \frac{1}{2}$$

Similarly,  $p(x_j) - p(x_{j-1}) > 0$  and  $x_{j-1} < x_j$  which will lead to  $\sum_{i < j} w_i < \frac{1}{2}$ .



**Alternative Solution:**

$p(x)$  is continuous. We show that  $\frac{dp(x)}{dx} > 0$  for  $x > x_j$  and  $\frac{dp(x)}{dx} < 0$  for  $x < x_j$  to conclude that  $x_j$  is a global minimum among the  $x_i$  points.

- 15.C. (20 PTS.) Given,  $X$  and  $W$ , describe in few lines an algorithm to find  $x_j$ . What is the running time of your algorithm.

**Solution:**

The idea is to sort  $X$  and rearrange the elements of  $W$  accordingly. Then, sum the weights until we exceed  $1/2$ . The point at which we exceed  $1/2$  is  $x_j$ . The runtime is  $O(n \log n) + O(n) = O(n \log n)$ .

```

Sort( $X, W$ )
 $s \leftarrow 0$ 
 $i \leftarrow 1$ 
while  $s < 1/2$  do
     $s = s + W[i]$ 
     $i = i + 1$ 
 $x_j \leftarrow X[i - 1]$ 

```

- 15.D. (50 PTS.) Given,  $X$  and  $W$ , describe an algorithm to find  $x_j$  in  $O(n)$  time. Prove the bound on the running time of your algorithm.

**Solution:**

The idea is to use binary search. We find the median using the linear time selection algorithm from class. We then compute the weights of all elements smaller than the median. If it is greater than  $1/2$ , we recurse on elements smaller than the median. If it is smaller than  $1/2$ , we recurse on elements larger than the median while keeping track of the sum of weights of points smaller than the median. We run the below algorithm with input **WeightedMedian**( $X, W, 0$ ).

```

WeightedMedian( $X[1, \dots, n], W[1, \dots, n], t$ )
if  $n = 1$ 
    return  $X[1]$ 
else
     $s \leftarrow 0, l \leftarrow 1, g \leftarrow 1$ 
     $y \leftarrow \text{Select}(X, \lfloor (n+1)/2 \rfloor)$ 
    for  $i = 1, \dots, n$  do
        if  $X[i] < y$ 
             $s = s + W[i]$ 
             $X_{<}[l] = X[i], W_{<}[l] = W[i]$ 
             $l = l + 1$ 
        else
             $X_{>}[g] = X[i], W_{>}[g] = W[i]$ 
             $g = g + 1$ 
    if  $t + s > 1/2$ 
        WeightedMedian( $X_{<}, W_{<}, t$ )
    else
        WeightedMedian( $X_{>}, W_{>}, s$ )

```

The algorithm requires  $O(n)$  to find the median and another  $O(n)$  to divide all elements of  $X$  into  $X_{<}$  and  $X_{>}$ . Then the algorithm recursively calls itself on  $X_{<}$  or  $X_{>}$  both of which are of size  $n/2$ . Hence, for some constant  $c > 1$ , the number of operations is:

$$\begin{aligned}
 T(n) &= T(n/2) + cn \\
 &= T(n/4) + cn/2 + cn \\
 &= T(n/2^k) + cn \sum_{i=0}^{k-1} 1/2^i \\
 &= T(1) + cn \sum_{i=0}^{\log n - 1} 1/2^i \\
 &= O(n)
 \end{aligned}$$