

# 15–210: Parallel and Sequential Data Structures and Algorithms

## FINAL (SOLUTIONS)

### Practice Final Exam (Fall 2012)

- There are 10 pages in this examination, comprising 7 questions worth a total of 111 points (including 11 bonus points).
- You have 180 minutes to complete this examination.
- Please answer all questions in the space provided with the question. Clearly indicate your answers.
- You may refer to your one double-sided  $8\frac{1}{2} \times 11$ in sheet of paper with notes, but to no other person or source, during the examination.
- Your answers for this exam must be written in blue or black ink.

Full Name: \_\_\_\_\_ Edsger W. Dijkstra

Andrew ID: \_\_\_\_\_ Section: \_\_\_\_\_

Question	Points	Score
Binary Answers	20	
Costs	12	
Short Answers	18	
Slightly Longer Answers	21	
Longest Contiguous Increasing Subsequence	16	
Median ADT	12	
Geometric Coverage	12	
Total:	111	

**Question 1: Binary Answers** (20 points)

Clearly mark **T** or **F** to the left of each question.

**ANS: F, T, F, F, T, F, F, T, F, T**

- (a) (2 points) The expressions `(Seq.reduce f I A)` and `(Seq.iter f I A)` always return the same result as long as `f` is commutative.
- (b) (2 points) The expressions `(Seq.reduce f I A)` and `(Seq.reduce f I (Seq.reverse A))` always return the same result if `f` is associative and commutative.
- (c) (2 points) Any parallel algorithm for a problem is always faster than a sequential algorithm for the same problem.
- (d) (2 points) Solving recurrences with induction can be used to show both upper and lower bounds?
- (e) (2 points) Let  $p$  be an odd prime. In open address hashing with a table of size  $p$  and given a hash function  $h(k)$ , quadratic probing uses  $h(k, i) = (h(k) + i^2) \bmod p$  as the  $i$ th probe position for key  $k$ . If there is an empty spot in the table quadratic hashing will always find it.
- (f) (2 points) `scan f b L` and `reduce f b L` always have the same asymptotic cost.
- (g) (2 points) If a randomized algorithm has expected  $O(n)$  work, then there exists some constant  $c$  such that the work performed is guaranteed to be at most  $cn$ .
- (h) (2 points) The height of any binary search tree (BST) is  $O(\log n)$ .
- (i) (2 points) Dijkstra's algorithm always terminates even if the input graph contains negative edge weights.
- (j) (2 points) A  $\Theta(n^2)$ -work algorithm always takes longer to run than a  $\Theta(n \log n)$ -work algorithm.

**Question 2: Costs** (12 points)

- (a) (6 points) Solve the following recurrence using the tree method. Show your work.

$$W(n) = 2W(n/2) + n \log n$$

**Solution:** At  $i^{th}$  level there are  $2^i$  subproblems each of which cost is  $\frac{n}{2^i} \log \frac{n}{2^i}$  for total cost of  $n(\log n - i)$ .

$$\begin{aligned} W(n) &= \sum_{i=0}^{\log n - 1} n(\log n - i) \\ &= n \sum_{j=1}^{\log n} j \\ &= n \log n (\log n + 1) / 2 \\ W(n) &\in O(n \log^2 n) \end{aligned}$$

- (b) (6 points) Check the appropriate column for each row in the following table:

	root dominated	leaf dominated	balanced
$W(n) = 2W(n/2) + n^{1.5}$			
$W(n) = \sqrt{n}W(\sqrt{n}) + \sqrt{n}$			
$W(n) = 8W(n/2) + n^2$			

**Solution:**

	root dominated	leaf dominated	balanced
$W(n) = 2W(n/2) + n^{1.5}$	X		
$W(n) = \sqrt{n}W(\sqrt{n}) + \sqrt{n}$		X	
$W(n) = 8W(n/2) + n^2$		X	

### Question 3: Short Answers (18 points)

Answer each of the following questions in the spaces provided.

- (a) (3 points) What simple formula defines the parallelism of an algorithm (in terms of work and span)?

**Solution:**  $P(n) = \frac{W(n)}{S(n)}$

- (b) (3 points) Name two algorithms we covered in this course that use the greedy method.

**Solution:** Dijkstra's, Prim's, Kruskal's ...

- (c) (3 points) Given a sequence of key-value pairs  $A$ , what does the following code do?

`Table.map length (Table.collect A)`

**Solution:** Makes a table that maps each key to the number values in sequence  $A$  that equal the key.

- (d) (3 points) What is the cut property of graphs that enables MST algorithms such as Kruskal's, Prim's and Borůvka's to work correctly?

**Solution:** For any nonempty connected subset  $U$  of  $V$ , the minimum weight edge between  $U$  and  $V \setminus U$  is in a MST of  $G$ .

- (e) (3 points) What asymptotically efficient parallel algorithm/technique can one use to count the number of trees in a forest (tree and forest have their graph-theoretical meaning)? (*Hint: the ancient saying of "can't see forest from the trees" may or may not be of help.*) Give the work and span for your proposed algorithm.

**Solution:** Run tree contraction over the entire forrest to contract each tree into a single vertex. (You can use either star contract or rake and compress.) Count the number of vertices at the end.

$$W(n) = O(n) \quad S(n) = O(\log^2 n)$$

- (f) (3 points) What are the two ordering invariants of a Treap? (Describe them briefly.)

**Solution: Heap property:** Each node has a higher priority than all of its descendants.

**BST property:** Each node's key is  $\leq$  to the keys in its left subtree and  $\geq$  to the keys in its right subtree.

#### Question 4: Slightly Longer Answers (21 points)

- (a) (7 points) Certain locations on a straight pathway recently built for robotics research have to be covered with a special surface, so CMU hires a contractor who can build arbitrary length segments to cover these locations (a location is covered if there is a segment covering it). The segment between  $a$  and  $b$  (inclusive) costs  $(b - a)^2 + k$ , where  $k$  is a non-negative constant. Let  $k \geq 0$  and  $X = \langle x_1, \dots, x_n \rangle$ ,  $x_i \in \mathbb{R}_+$ , be a sequence of locations that have to be covered. Give an  $O(n^2)$ -work dynamic programming solution to find the cheapest cost of covering these points (all given locations must be covered). Be sure to state the subproblems and give a recurrence, including the base case(s).

**Solution:** Let  $f(i)$  be the cheapest way to cover locations  $i, i + 1, \dots, n$ . Then, we have the following recurrence:

$$f(i) = \begin{cases} 0 & i > n \\ \min_{i \leq j \leq n} \left( f(j + 1) + k + (x_i - x_j)^2 \right) & 1 \leq i \leq n \end{cases}$$

- (b) (7 points) Consider the following variant of the optimal binary search tree (OBST) algorithm given in class:

```

1  fun OBST(A) = let
2    fun OBST'(S, d) =
3      if |S| = 0 then 0
4      else min_{i ∈ {1...|S|}} (OBST'(S_{1,i-1}, d + 1) + d · p(S_i) + OBST'(S_{i+1,|S|}, d + 1))
5    in
6      OBST'(A, 1)
7  end

```

Recall that  $S_{i,j}$  is the subsequence  $\langle S_i, S_{i+1}, \dots, S_j \rangle$  of  $S$ . For  $|A| = n$ , place an asymptotic upper bound on the number of distinct arguments  $OBST'$  will have (a tighter bound will get more credit).

**Solution:** There are  $\binom{n+1}{2} = n(n+1)/2$  possible subsequences of  $S$  and  $d$  is between 1 and  $n$ , so the number of distinct arguments is upper-bounded by  $O(n^3)$ .

- (c) (7 points) Given  $n$  line segments in 2 dimensions, the 3-intersection problem is to determine if any three of them intersect at the same point. Explain how to do this in  $O(n^2)$  work and  $O(\log n)$  span. You can assume the lines are given with integer endpoints (i.e. you can do exact arithmetic and not worry about roundoff errors).

**Solution:** First, we compute all possible intersection points between pairs of line segments. There can be at most  $O(n^2)$  points. Then, insert these points into a hash table, checking if any collision seen is a result of 3 lines intersecting at the same point. This meets the time bound since hashing  $O(n^2)$  points takes  $O(n^2)$  work and  $O(\log n^2) \subseteq O(\log n)$  span.

### Question 5: Longest Contiguous Increasing Subsequence (16 points)

Given a sequence of numbers, the *longest contiguous increasing subsequence* problem is to find the largest number of contiguous increases in a sequence of numbers. For example,

LCIS(<7, 2, 3, 4, 1, 8>)

will return 2 since there are 2 increases in a row in the contiguous subsequence <2, 3, 4>. Note that this is different from the longest increasing subsequence problem discussed in recitation.

- (a) (4 points) The LCIS problem can be solved in linear work by strengthening the problem (inductive hypothesis) and solved using divide and conquer by splitting the sequence in half and solving each half. Describe what values you would return from the recursive calls to efficiently construct the solution.

**Solution:** longest increasing prefix, longest contiguous increasing subsequence, longest increasing suffix

- (b) (4 points) Fill in the following SML code for your recursive divide-and-conquer algorithm:

```
fun LCIS (S : int seq) : int =
  let
    fun LCIS'(S : int seq) : int * int * int =
      case (showt S) of
        EMPTY => (0, 0, 0)
      | ELT(x) => (0, 0, 0)
      | NODE(L,R) => (* fill in below *)
        let val (L_pre, L_lcis, L_suf) = LCIS' L
            val (R_pre, R_lcis, R_suf) = LCIS' R
            val bridge = Seq.nth L (Seq.length L -1) < Seq.nth R 0
            val pre = if (L_pre = (Seq.length L -1) andalso bridge)
                      then L_pre + 1 + R_pre
                      else L_pre
            val suf = if (R_suf = (Seq.length R -1) andalso bridge)
                      then L_suf + 1 + R_suf
                      else R_suf
            val lcis = max (max(pre, suf), max(L_lcis, R_lcis))
        in
          (pre, lcis, suf)
        end
      in (* fill in below *)
        case (LCIS' S) of (_,X,_) => X
      end
  end
```

- (c) (4 points) Assuming a tree-based implementation of sequences in which `showt`, and `nth` take  $O(\log n)$  work, write recurrences for the work and span of `LCIS'` and state the solutions of the recurrences.

**Solution:**

$$W(n) = 2W\left(\frac{n}{2}\right) + O(\log n) = O(n)$$

$$S(n) = S\left(\frac{n}{2}\right) + O(\log n) = O(\log^2 n)$$

- (d) (4 points) The problem can also be solved with a scan. Here is the code.

```
datatype Dir = UP of int | MIX of int

fun LCIS(S : int seq) =
  let
    fun up i = if (nth S (i+1)) > (nth S i) then UP(1) else MIX(0)
    val Sup = tabulate up ((length S)-1)
    val (R,MIX(v)) = scan binop (MIX(0)) Sup
    val R' = map (fn MIX(x) => x) R
  in
    Int.max(v, reduce Int.max 0 R')
  end
```

Fill in the following code for `binop`.

```
fun binop(_ , MIX b) =      MIX b

  | binop(MIX a , UP b) =   MIX (a + b)

  | binop(UP a , UP b) =    UP (a + b)
```

**Question 6: Median ADT** (12 points)

The *median* of a set  $C$ , denoted by  $\text{median}(C)$ , is the value of the  $\lceil n/2 \rceil$ -th smallest element (counting from 1). For example,

$$\text{median}(\{1, 3, 5, 7\}) = 3$$

$$\text{median}(\{4, 2, 9\}) = 4$$

In this problem, you will implement an abstract data type **medianT** that maintains a collection of integers (possibly with duplicates) and supports the following operations:

<b>insert</b> ( $C, v$ )	: $\text{medianT} \times \text{int} \rightarrow \text{medianT}$	add the integer $v$ to $C$ .
<b>median</b> ( $C$ )	: $\text{medianT} \rightarrow \text{int}$	return the median value of $C$ .
<b>fromSeq</b> ( $S$ )	: $\text{int Seq.seq} \rightarrow \text{medianT}$	create a <b>medianT</b> from $S$ .

Throughout this problem, let  $n$  denote the size of the collection at the time, i.e.,  $n = |C|$ .

- (a) (5 points) Describe how you would implement the **medianT** ADT using (balanced) binary search trees so that **insert** and **median** take  $O(\log n)$  work and span.

**Solution:** As described in the augmented tree lecture, we keep a balanced BST where each node is augmented with the size of the subtree, so that the  $(n/2)$ -th element can be found in  $O(\log n)$ ; inserting an element also takes  $O(\log n)$  because we simply need to “update” the size information on a relevant path, which has length  $O(\log n)$ .

- (b) (7 points) Using some other data structure, describe how to improve the work to  $O(\log n)$ ,  $O(1)$  and  $O(|S|)$  for the three operations respectively. The **fromSeq S** function needs to run in  $O(\log^2 |S|)$  span. (*Hint: think about maintaining the median, the elements less than the median, and the elements greater than the median separately.*)

**Solution:** Keep a max heap of values smaller than the median, the current median, and a min heap of values bigger than the median. When inserting, put the new value in the correct heap, rebalancing as necessary. The function **fromSeq** can be easily supported since using quick select to the initial median only requires  $O(n)$  expected work and  $O(\log^2 n)$  span.



### Question 7: Geometric Coverage (12 points)

For points  $p_1, p_2 \in \mathbb{R}^2$ , we say that  $p_1 = (x_1, y_1)$  *covers*  $p_2 = (x_2, y_2)$  if  $x_1 \geq x_2$  and  $y_1 \geq y_2$ . Given a set  $S \subseteq \mathbb{R}^2$ , the *geometric cover number* of a point  $q \in \mathbb{R}^2$  is the number of points in  $S$  that  $q$  covers. Notice that by definition, every point covers itself, so its cover number must be at least 1.

In this problem, we'll compute the geometric cover number for every point in a given sequence. More precisely:

**Input:** a sequence  $S = \langle s_1, \dots, s_n \rangle$ , where each  $s_i \in \mathbb{R}^2$  is a 2-d point.

**Output:** a sequence of pairs each consisting of a point and its cover number. Each point must appear exactly once, but the points can be in any order.

Assume that we use the `ArraySequence` implementation for sequences.

- (a) (4 points) Develop a brute-force solution `gcnBasic` (in pseudocode or Standard ML). Despite being a brute-force solution, your solution should not do more work than  $O(n^2)$ .

**Solution:**

```
fun gcnBasic s =  
  let  
    fun covers ((x1,y1),(x2,y2)) = (x1 >= x2) andalso (y1 >= y2)  
  in  
    map (fn p => length (filter (fn p' => covers(p, p')) s)) s  
  end
```

- (b) (4 points) In words, outline an algorithm `gcnImproved` that has  $O(n \log n)$  work. You may assume an implementation of `OrderedTable` in which `split`, `join`, and `insert` have  $O(\log n)$  cost (i.e., work and span), and `size` and `empty` have  $O(1)$  cost.

**Solution:** We'll keep an ordered table  $T$  of points ordered by their  $x$  values. Initially,  $T$  is empty. To compute the cover number for every point, we'll first sort these points by their  $y$  values. Then, for each of these points, we insert them one by one into  $T$ —and the cover number of this point can be found by splitting  $T$  using its  $x$  value and taking the size of the left side.

- (c) (4 points) Show that the work bound cannot be further improved by giving a lower bound for the problem.

**Solution:** We'll reduce comparison-based sorting to GCN, which means that GCN cannot be solved in less than  $\Omega(n \log n)$  work. The reduction is as follows: for a given input sequence  $s = \langle s_1, \dots, s_n \rangle$ , we create a sequence of points

$$P = \langle (s_1, s_1), (s_2, s_2), \dots, (s_n, s_n) \rangle$$

(using map in  $O(n)$  work and  $O(1)$  span). Running GCN on this  $P$  gives the “rank” of each element, which we can then use as indices to inject and get a sorted sequence.