

# Recitation 1

## Priority Queues and Hashing

### 1.1 Announcements

- PASLLab is due this **Friday, December 9**.
- The final exam is **Tuesday, December 13**.
- A review session for the final is on **Sunday, December 11, 4-6pm, in Rashid**.
- A practice final and its solutions are up on the website and on Piazza.

## 1.2 Leftist Heaps

**Task 1.1.** *Identify the defining properties of a leftist heap.*

A leftist heap is a binary tree given by

**datatype** tree = Leaf | Node **of** key × tree × tree

which satisfies

- (a) the *heap property*, requiring that the key stored at each node is smaller<sup>1</sup> than any descendent key, and
- (b) the *leftist property*, requiring that for every  $\text{Node}(\_, L, R)$ , we have  $\text{rank}(L) \geq \text{rank}(R)$ . We define the *rank* of a heap to be the number of nodes in its right spine, i.e.,

$$\text{rank}(\text{Leaf}) = 0$$

$$\text{rank}(\text{Node}(\_, L, R)) = 1 + \text{rank}(R)$$

**Task 1.2.** *What is an upper bound on the rank of the root of a leftist heap?*

For a leftist heap containing  $n$  entries, the rank of the root is at most  $\log_2(n + 1)$ .

---

<sup>1</sup>We assume a min-heap. In a max-heap, each key is larger than its descendents.

## 1.2.1 Building A Leftist Heap

Consider the following pseudo-SML code implementing leftist heaps.

### Data Structure 1.3. *Leftist Heap*

```

1  datatype PQ = Leaf | Node of int × key × PQ × PQ
2
3  fun rank Q =
4    case Q of
5      Leaf ⇒ 0
6      | Node (r,_,_,_) ⇒ r
7
8  fun makeLeftistNode (k,A,B) =
9    if rank A < rank B
10   then Node (1 + rank A, k, B, A)
11   else Node (1 + rank B, k, A, B)
12
13 fun meld (A,B) =
14   case (A,B) of
15     (_, Leaf) ⇒ A
16     | (Leaf, _) ⇒ B
17     | (Node (_,ka,La,Ra), Node (_,kb,Lb,Rb)) ⇒
18       if ka < kb
19       then makeLeftistNode (ka, La, meld (Ra,B))
20       else makeLeftistNode (kb, Lb, meld (A,Rb))
21
22 fun singleton k = Node (1,k,Leaf,Leaf)
23
24 fun insert (Q,k) = meld (Q, singleton k)
25
26 fun fromSeq S = Seq.reduce meld Leaf (Seq.map singleton S)
27
28 fun deleteMin Q =
29   case Q of
30     Leaf ⇒ (NONE, Q)
31     | Node (_,k,L,R) ⇒ (SOME k, meld (L,R))

```

**Task 1.4.** Diagram the process of executing the code

`fromSeq`  $\langle 3, 5, 2, 1, 4, 6, 7, 8 \rangle$

3 5 2 1 4 6 7 8

3 1 4 7  
/  
5 2 6 8

1 4  
/ \ / \  
2 3 6 7  
/ /  
5 8

1  
/ \  
3 2  
/ \  
4 5  
/ \  
6 7  
/  
8

**Task 1.5.** What are the work and span of  $(\text{fromSeq } S)$  in terms of  $|S| = n$ ?

Notice that `meld` only traverses the right spines of its arguments, each of which are logarithmic in length, and therefore `meld`( $A, B$ ) requires  $O(\log |A| + \log |B|)$  work and span and returns a heap of size  $|A| + |B|$ . This suggests the recurrences

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

$$S(n) = S(n/2) + O(\log n)$$

both of which we have seen before; they solve to  $O(n)$  work and  $O(\log^2 n)$  span, respectively.

## 1.3 Removing Duplicates

Removing duplicates is a crucial substep of many interesting algorithms. For example, in BFS, consider the step where we construct a new frontier. One viable method would be to generate the sequence of all out-neighbors, and then remove duplicates:

$$F' = \text{removeDuplicates} \langle v : u \in F, v \in N_G^+(u) \rangle$$

So, how fast is it to remove duplicates? Can we do it in parallel?

### 1.3.1 Sequential

Before we think about parallelism, we should acquaint ourselves with a good sequential algorithm solving the same problem. This way, we know what to shoot for in terms of work bounds, since we want our parallel algorithm to be asymptotically work-efficient.

**Task 1.6.** *Describe a sequential algorithm which performs expected  $O(n)$  work to remove duplicates from a sequence of length  $n$ . Also argue that  $\Omega(n)$  work is necessary in order to solve this problem, and conclude that your algorithm is asymptotically optimal.*

*Hint: try hashing elements one at a time.*

We can iterate left-to-right across the sequence, maintaining a set of all elements seen so far. At each element, we check to see if it's present in the set. If it is, we ignore it. If it isn't, we insert it into the set and also write it to the output. Using a hash set, we can check for membership and do insertions both in expected constant time. So, this algorithm has expected  $O(n)$  work.

Removing duplicates requires at least  $\Omega(n)$  work because we have to inspect every element. Here's a sketch of a proof by contradiction: suppose there was an algorithm for removing duplicates which used only  $o(n)$  work. Then there must be at least one element of the input which the algorithm did not inspect. If the algorithm would blindly include this element in the output, then we can adversarially choose that element to be a duplicate. If instead the algorithm blindly excludes the element from the output, then we can adversarially choose the element to be distinct. Thus clearly the given algorithm does not properly remove duplicates.

Since  $\Omega(n)$  work is necessary and our algorithm has  $O(n)$  work, we know it is asymptotically optimal in terms of runtime.

### 1.3.2 Parallel

**Task 1.7.** *Implement a function*

**val** `removeDuplicates` :  $(\alpha \times \text{int} \rightarrow \text{int}) \rightarrow \alpha \text{ Seq.t} \rightarrow \alpha \text{ Seq.t}$

where  $(\text{removeDuplicates } h \ S)$  returns a sequence of all unique elements of  $S$ , given that  $h(e, m)$  hashes the element  $e$  to a uniform random integer in the range  $[0, m)$  (thus the probability of collision for any two distinct elements is  $1/m$ ).

*Hint: as a first attempt, try simultaneously hashing as many elements as possible all at the same time. What do you do when elements collide?*

We can use `inject` to simultaneously insert as many keys as possible into an initially empty hash table  $T$  of some size  $m > n$  (we'll decide on a value for  $m$  later). Specifically, for every  $0 \leq i < |S|$ , we attempt to insert the pair  $(i, S[i])$  at the location  $T[h(S[i], m)]$ .

Every pair  $(i, S[i])$  then compares itself against the value  $(j, S[j])$  stored at  $T[h(S[i], m)]$ . If  $i = j$ , then  $S[i]$  is “accepted,” and will be included in the output. Otherwise, if  $S[i] \neq S[j]$ , then  $S[i]$  is passed on to be retried in the next round. We continue retrying until no elements remain.

Why is this algorithm correct? Consider some key  $k$ ; we never discard  $k$  until  $k$  is accepted, so we only need to argue that we never accept the same key twice. Consider some round and two indices  $i \neq j$  such that  $S[i] = S[j]$ . If  $S[i]$  is accepted then  $S[j]$  won't be (since  $i \neq j$ ), and furthermore it won't be retried on the next round. Therefore  $S[j]$  will never be accepted.

How many elements are retried each round? Consider:

$$\begin{aligned} \Pr[S[i] \text{ is retried}] &= \Pr[\exists j. S[i] \neq S[j] \wedge h(S[i], m) = h(S[j], m)] \\ &\leq \sum_{S[i] \neq S[j]} \Pr[h(S[i], m) = h(S[j], m)] \\ &\leq \frac{|S|}{m} \end{aligned}$$

If we chose  $m = 3|S|/2$ , then  $|S|/m = 2/3$ , and by linearity of expectation, we have that the number of retried elements is at most  $2|S|/3$  in expectation.

Since we need  $O(|S|)$  work and  $O(\log |S|)$  span on each round, we expect a logarithmic number of rounds with a geometrically decreasing input. We've seen such recurrences before; they solve to expected linear work and log-squared span.

**Algorithm 1.8.** *Removing duplicates with hashing.*

```
1 fun removeDuplicates S =  
2   if |S| = 0 then ⟨⟩ else  
3   let  
4     val E = Seq.enum S  
5  
6     val m =  $\frac{3|S|}{2}$   
7     val base = ⟨NONE : 0 ≤ i < m⟩  
8     val updates = ⟨(h(k,m), SOME(i,k)) : (i,k) ∈ E⟩  
9     val T = Seq.inject (base, updates)  
10  
11    fun accept (i,k) = (T[h(k,m)] = SOME (i,k))  
12    val A = ⟨k : (i,k) ∈ E | accept(i,k)⟩  
13  
14    fun retry k = let val SOME(−, k′) = T[h(k,m)] in k ≠ k′ end  
15  in  
16    Seq.append (A, removeDuplicates ⟨k ∈ S | retry(k)⟩)  
17  end
```

## 1.4 Additional Exercises

### Exercise 1.9.

**Task 1.10.** Design a data structure which supports the following operations:

	<i>Work</i>	<i>Span</i>	<i>Description</i>
<i>fromSeq S</i>	$O( S )$	$O(\log^2  S )$	Constructs a dynamic median data structure from the collection of keys in $S$
<i>median M</i>	$O(1)$	$O(1)$	Returns the median of all keys stored in $M$
<i>insert (M, k)</i>	$O(\log  M )$	$O(\log  M )$	Inserts $k$ into $M$

For simplicity, you may assume that all elements inserted into such a structure are distinct.

**Exercise 1.11.** Prove a lower bound of  $\Omega(\log n)$  for `deleteMin` in comparison-based meldable priority queues. That is, prove that any meldable priority queue implementation which has a logarithmic `meld` cannot support `deleteMin` in faster than logarithmic time.