

## Lecture 4: Data Structures

Harvard SEAS - Fall 2023

2023-09-14

## 1 Announcements

- Please fill out PS0 feedback survey.
- Next SRE: Thursday 9/21
- Participation is not a contest! Many ways to contribute!

## 2 Recommended Reading

- CLRS Chapter 10
- Roughgarden II, Sec. 10.0–10.1, 11.1
- CS50 Week 5: <https://cs50.harvard.edu/x/2022/weeks/5/>

## 3 Use of Reductions

The use of reductions is mostly described by the following lemma, which we'll return to many times in CS 120:

**Lemma 3.1.** *Let  $\Pi$  and  $\Gamma$  be computational problems such that  $\Pi \leq \Gamma$ . Then:*

1. *If there exists an algorithm solving  $\Gamma$ , then there exists an algorithm solving  $\Pi$ .*
2. *If there does not exist an algorithm solving  $\Pi$ , then there does not exist an algorithm solving  $\Gamma$ .*
3. *If there exists an algorithm solving  $\Gamma$  with runtime  $g(n)$ , and  $\Pi \leq_{T,h} \Gamma$ , then there exists an algorithm solving  $\Pi$  with runtime  $O(T(n) + g(h(n)))$ .*
4. *If there does not exist an algorithm solving  $\Pi$  with runtime  $O(T(n) + g(h(n)))$ , and  $\Pi \leq_{T,h} \Gamma$ , then there does not exist an algorithm solving  $\Gamma$  with runtime  $O(g(n))$ .*

## 4 Interval Scheduling

The lecture notes for the previous class (Lecture 3) introduced the IntervalScheduling-Decision computational problem, motivated by checking whether a collection of time reservations (e.g. radio airtime) have any conflicts:

<b>Input</b>	: A collection of intervals $[a_0, b_0], \dots, [a_{n-1}, b_{n-1}]$ , where each $a_i, b_i \in \mathbb{R}$ and $a_i \leq b_i$
<b>Output</b>	: YES if the intervals are disjoint from each other (for all $i \neq j$ , $[a_i, b_i] \cap [a_j, b_j] = \emptyset$ ) NO otherwise

**Computational Problem** IntervalScheduling-Decision

The notes from Lecture 3 show that IntervalScheduling-Decision can be solved efficiently by a reduction to *Sorting*:

**Proposition 4.1.** *IntervalScheduling-Decision on inputs that consist of  $n$  intervals can be reduced to Sorting an array of  $n$  key-value pairs in time  $O(n)$ . That is, IntervalScheduling-Decision  $\leq_{O(n),n}$  Sorting.*

As a corollary, IntervalScheduling-Decision can be solved in time  $O(n \log n)$ .

**Q:** Suppose we have already solved IntervalScheduling in this way, and another interval  $[a', b']$  is given to us (e.g. another listener tries to buy some airtime). Do we need to spend time  $O(n \log n)$  again to decide whether we can fit that interval in?

**A:** We can use binary search to find where  $[a', b']$  would be in the sorted array and check for conflicts with the adjacent time slots. Using binary search on a sorted array takes  $O(\log n)$  time.

## 5 Static Data Structure

The sorted array in the above solution is an example of a (static) data structure. Let's abstract what data structures are supposed to do.

**Definition 5.1.** A *static data structure problem* is a quadruple  $\Pi = (\mathcal{I}, \mathcal{O}, \mathcal{Q}, f)$  where:

- $\mathcal{I}$  is a (typically infinite) set of possible inputs  $x$ , and  $\mathcal{O}$  is a (sometimes infinite) set of possible outputs  $y$ .
- $\mathcal{Q}$  is a set of *queries*, and
- for every  $x \in \mathcal{I}$  and  $q \in \mathcal{Q}$ ,  $f(x, q) \subseteq \mathcal{O}$  is a set of *valid answers*.

For such a data-structure problem, we want to design efficient algorithms that preprocess the input  $x$  into a data structure that allows for quickly answering queries  $q$  that come later. For example, to be able to determine whether a new interval conflicts with one of the original ones, it suffices to solve the following data-structure problem.

To formalize these two types of queries using Definition 5.1, we can take

$$\mathcal{Q} = \{(\text{search}, K) : K \in \mathbb{R}\} \cup \{(\text{next-smaller}, K) : K \in \mathbb{R}\}$$

Note that both types of queries may have no valid solution, or may have multiple solutions. (Why?) If we removed the Next-smaller queries and only kept Search queries, we would have the (static) *Dictionary* data structure problem, which we will study in a couple of weeks.

**Input** : An array of key-value pairs  $x = ((K_0, V_0), \dots, (K_{n-1}, V_{n-1}))$ , with each  $K_i \in \mathbb{R}$   
**Queries** :

- Search( $K$ ) for  $K \in \mathbb{R}$ : output  $(K_i, V_i)$  such that  $K_i = K$ .
- Next-smaller( $K$ ) for  $K \in \mathbb{R}$ : output  $(K_i, V_i)$  such that  $K_i = \max\{K_j : K_j < K\}$ .

### Data-Structure Problem Static Predecessors

#### Using Static Predecessors+Successors for Interval Scheduling:

Note that a solution to the Static Predecessors problem is not immediately a sufficient solution to determining whether a new interval conflicts with one of the original ones. (We assume here that we have already checked that none of the original ones conflict.) While we can use the Next-smaller query to check whether an interval with an earlier start time conflicts with a given interval, we would also need an analogous Next-larger (successor) query to check if a later interval conflicts. We can define a Static Predecessors and Successors problem by including the input query:

- Next-larger( $K$ ) for  $K \in \mathbb{R}$ : output  $(K_i, V_i)$  such that  $K_i = \min\{K_j : K_j > K\}$ .

Then, given an interval  $[a^*, b^*]$  we can check whether any conflicts exist by

- Call Search( $a^*$ ) and check that it returns  $\perp$ .
- Call Next-smaller( $a^*$ ), which returns  $(a_i, b_i)$ . Check that  $a^* > b_i$
- Call Next-larger( $a^*$ ), which returns  $(a_j, b_j)$ . Check that  $b^* < a_j$

Note that a Successor Data Structure can be constructed from a Predecessor Data Structure. (How?)

**Remark.** The terminology *predecessor* comes from the fact that when  $K$  is a key in the input array  $x$ , then Next-smaller( $K$ ) should return the *predecessor* of  $K$ . CLRS only defines the predecessor problem where the query is a key of already in the set. However, the more general formulation we have given (where  $K$  can be any real number) is more standard and more useful.

Predecessor Data Structures (and the equivalent Successor Data Structures) have many applications. They enable one to perform a “range select” — selecting all of the elements of a dataset whose keys fall within a given range. This is a fundamental operation across many application and systems, including relational databases (e.g. a university database selecting all CS alumni who graduated in the 1990’s), NoSQL data stores (e.g. selecting all users of a social network within a given age range), and ML systems (e.g. filtering intermediate results during neural network training sessions).

**Definition 5.2.** For a (static) data structure problem  $\Pi = (\mathcal{I}, \mathcal{O}, \mathcal{Q}, f)$ , an *implementation* is a pair of algorithms (Preprocess, Eval) such that

- for all  $x \in \mathcal{I}$  and  $q \in \mathcal{Q}$  such that  $f(x, q) \neq \emptyset$ , we have  $\text{Eval}(\text{Preprocess}(x), q) \in f(x, q)$ , and
- there is a special output  $\perp \notin \mathcal{O}$  such that for all  $x \in \mathcal{I}$  and  $q \in \mathcal{Q}$  such that  $f(x, q) = \emptyset$ , we have  $\text{Eval}(\text{Preprocess}(x), q) = \perp$ .

Our goal is for Eval to run as fast as possible. (As we'll see in examples below, sometimes there are multiple types of queries, in which case we often separately measure the running time of each type.) Secondarily, we would like Preprocess to also be reasonably efficient and to minimize the memory usage of the data structure  $\text{Preprocess}(x)$ .

Note that there is no pre-specified problem that the algorithms Preprocess and Eval are required to solve individually; we only care that *together* they correctly answer queries. Thus, a big part of the creativity in solving data structure problems is figuring out what the form of  $\text{Preprocess}(x)$  should be. Our first example (from today) takes it to be a sorted array, but we will see other possibilities in the upcoming classes (like binary search trees and hash tables).

Let's formalize our solution to the Static Predecessor Problem:

- $\text{Preprocess}(x)$ : Sort the array  $x$  of key-value pairs in  $O(n \log n)$  time to obtain a sorted array  $x' = ((K'_0, V'_0), \dots, (K'_{n-1}, V'_{n-1}))$ .
- $\text{Eval}(x', (\text{search}, K))$ : Use binary search to find  $i \in [n]$  such that  $K'_i = K$  in time  $O(\log n)$  and return  $(K'_i, V'_i)$ . If no such  $i$  exists, return  $\perp$ .
- $\text{Eval}(x', (\text{next-smaller}, K))$ : Use binary search to find  $i = \max\{j : K'_j < K\}$  in time  $O(\log n)$  and return  $(K'_i, V'_i)$ . If no such  $i$  exists (namely, if  $K'_0 \geq K$ ), return  $\perp$ .

Thus we obtain:

**Theorem 5.3.** *Static Predecessors has an implementation in which*

- $\text{Eval}(x', (\text{search}, K))$  and  $\text{Eval}(x', (\text{next-smaller}, K))$  both take time  $O(\log n)$
- $\text{Preprocess}(x)$  takes time  $O(n \log n)$
- $\text{Preprocess}(x)$  has size  $O(n)$

## 6 Implementing Data Structures

*This will not be covered in class, but you may find it useful to see how the mathematical descriptions above translate into concrete code.*

As you saw on Problem Set 0, we can represent data structures in Python by using Python **classes** similarly to C **structs**. However, we can attach *methods* that implement the Preprocess and Eval functions of a data-structure solution. The implementation details of these methods (as well as the private attributes) can be hidden from a user of the class, creating an “abstraction barrier” that allows the user to focus only on the data-structure problem that it solves, without concern for the particular solution. (This is the notion of an “abstract data type” that some of you may have seen in CS51.) For example:

```
class StaticPredecessor:
    def __init__(self, x: list): # preprocessing method
        # Python's built-in sorting doesn't take keys as explicit inputs,
        # but asks the user to specify a function that defines the keys
        self.sortedarray = MergeSort(x)

    def search(self, K: float):
        left = 0
```

```

right = len(self.sortedarray) - 1
mid = 0

while left <= right:

    mid = (left + right) // 2

    # If K is greater, ignore left half
    if self.sortedarray[mid][0] < K:
        left = mid + 1

    # If K is smaller, ignore right half
    elif self.sortedarray[mid][0] > K:
        right = mid - 1

    else:
        return self.sortedarray[i]

# if no solution is found
return None

def nextsmaller(self, K: float):
    # ...

MyDS = StaticPredecessor([(5,"a"),(3,"b"),(7,"c"),(2,"d")]) # runs __init

MyDS.nextsmaller(4) # should return (3,"b")

```

Note: Despite the name, a Python `list` is implemented as an array rather than as a linked list.

## 7 Dynamic Data Structures

As you might have been wondering for the Interval Scheduling Problem, it is often the case that we do not get all of our input data at once, but rather it comes in through incremental updates over time. This gives rise to *dynamic* data-structure problems.

**Definition 7.1.** A *dynamic data structure problem* is a quintuple  $\Pi = (\mathcal{I}, \mathcal{O}, \mathcal{U}, \mathcal{Q}, f)$  where:

- $\mathcal{I}$  is a (sometimes infinite) set of possible initial inputs  $x$ , and  $\mathcal{O}$  is a (sometimes infinite) set of possible outputs  $y$ .
- $\mathcal{U}$  is a set of *updates*,
- $\mathcal{Q}$  is a set of *queries*, and
- for every  $x \in \mathcal{I}$ ,  $u_0, u_1, \dots, u_{m-1} \in \mathcal{U}$ , and  $q \in \mathcal{Q}$ ,  $f(x, u_0, \dots, u_{m-1}, q) \subseteq \mathcal{O}$  is a set of *valid answers*.

Often we take  $\mathcal{I} = \emptyset$ , since the inputs can usually be constructed through a sequence of updates. For example:

**Updates :**

- $\text{Insert}(K, V)$  for  $K \in \mathbb{R}$ : add one copy of  $(K, V)$  to the multiset  $S$  of key-value pairs being stored. ( $S$  is initially empty.)
- $\text{Delete}(K)$  for  $K \in \mathbb{R}$ : delete one key-value pair of the form  $(K, V)$  from the multiset  $S$  (if there are any remaining).

**Queries :**

- $\text{Search}(K)$  for  $K \in \mathbb{R}$ : output  $(K', V') \in S$  such that  $K' = K$ .
- $\text{Next-smaller}(K)$  for  $K \in \mathbb{R}$ : output  $(K', V') \in S$  such that  $K' = \max\{K'' : \exists V'' (K'', V'') \in S, K'' < K\}$ .

**Data-Structure Problem** Dynamic Predecessors

A *multiset* is like a set but can contain more than one copy of an element. The multiset  $S$  appearing in the above definition is only used to define the functionality of the data structure, namely how queries should be answered. How this set is actually maintained is up to the particular solution (i.e. data structure implementation).

To solve a dynamic data structure problem, we need to now come up with algorithms that also implement the updates.

The definition of what it means to implement a dynamic data structure with algorithms gets a bit cumbersome (and we don't expect you to remember it), but we include it here in the notes for completeness:

**Definition 7.2.** For a dynamic data structure problem  $\Pi = (\mathcal{I}, \mathcal{O}, \mathcal{Q}, \mathcal{U}, f)$ , an *implementation* is a triple of algorithms (Preprocess, EvalQ, EvalU) such that for all  $x \in \mathcal{I}$ ,  $u_0, u_1, \dots, u_{m-1} \in \mathcal{U}$  and  $q \in \mathcal{Q}$ , if  $f(x, u_0, u_1, \dots, u_{m-1}) \neq \emptyset$ , we have:

$$\text{EvalQ}(\text{EvalU}(\dots \text{EvalU}(\text{EvalU}(\text{Preprocess}(x), u_0), u_1), \dots, u_{m-1}), q) \in f(x, u_0, u_1, \dots, u_{m-1}, q),$$

else

$$\text{EvalQ}(\text{EvalU}(\dots \text{EvalU}(\text{EvalU}(\text{Preprocess}(x), u_0), u_1), \dots, u_{m-1}), q) = \perp.$$

Now, we would like both EvalU and EvalQ to be extremely fast. Dynamic data structures can also be conveniently implemented using `classes` in Python, similarly to Section 6, with the update operations also implemented as methods.

## 8 Exercise

*Left as food for thought prior to the next class!*

Suppose we use a sorted array to implement a data structure storing a dynamically changing multiset  $S$  of key-value pairs with insertions and deletions. How efficiently can you perform each of the following operations (in the worst case), as a function of the current number  $n$  of elements of  $S$ ? Possible answers are  $O(1)$ ,  $O(\log n)$ ,  $O(n)$ , and “other”.

1.  $\text{Insert}(K, V)$
2.  $\text{Delete}(K)$
3.  $\text{Min}()$ : return the smallest key  $K$  in  $S$ .
4.  $\text{Rank}(K)$ : return the *number* of pairs  $(K', V') \in S$  such that  $K' < K$ .