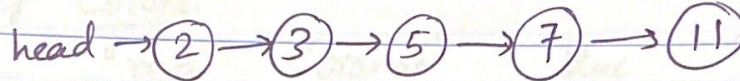


# Lecture 18: Hashing and Hash Tables

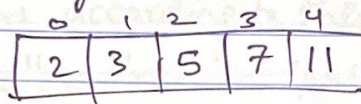
Previously: (sorted) sets

1. (sorted) linked list



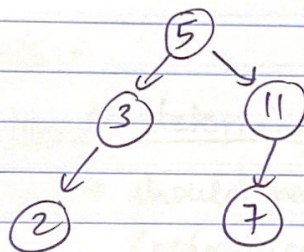
} all ops  $O(n)$

2. (Sorted) array



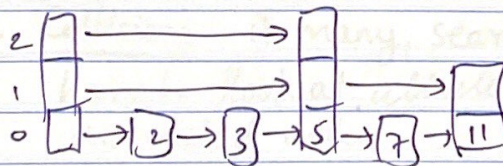
} find in  $O(\log n)$   
binary search  
 other ops  $O(n)$

3. (balanced) BST:



} All ops in  $O(h)$  ( $h = \text{height}$ )  
 $\Rightarrow$  unbalanced  $\rightarrow h = O(n)$   
 $\Rightarrow$  balanced  $\rightarrow h = O(\log n)$   
 (AVL)

4. Skiplists



probabilistic guarantees

Expected running times for ops  $O(\log n)$

} average over random choices of algorithm, not user.



⇒ Idea: use randomness to assign random numerical values to elements, "sort" according to values.

eg Colors:

"red"	"orange"	"blue"	"pink"
↓	↓	↓	↓
2	6	1	3

Sort colors according to these values!

"blue" < "red" < "pink" < "orange"

Use sorting for Sorted Set implementation

### ISSUES:-

1. How to determine order from colors?

→ should be deterministic so get same # (reproducible) from each color name.

→ assign values to chars in name of color?

2. Collisions - if many, searching is inefficient b/c have to look at all elements with same associated value.

3. Semantic Equivalence - need to associate same value to o & p if  $(o.equals(p) == true)$ .

Analogy: Balls in Bins (Assignment 5)

objects	↑	assigned values
"red"	→	2



(3)

Suppose: associate a "random looking" value (an int) to each object wrt property:

- 1)  $\text{getVal}(\text{obj}) \rightarrow$  gives value associated to obj;
- 2) values appear random - hard to find correlation between  $\text{getVal}(o)$  and  $\text{getVal}(p)$
- 3) Values reproducible: multiple calls to  $\text{getVal}(o)$  will be same.
- 4)  $\text{getVals}$  respects semantic equivalence.  
 $\rightarrow o.\text{equals}(p)$  then  $\text{getVal}(o) == \text{getVal}(p)$

Q Given  $\text{getVal}$  method,  
how to implement Unsorted set as efficiently as sorted set?

Possibility: Store pair of  $\langle \text{obj}, \text{getVal}(\text{obj}) \rangle$   
make comparable using this  
store pairs, in say, AVL tree.

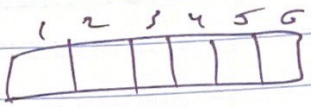
$\rightarrow$  add, find, remove in  $O(\log n)$  time  
assuming  $\text{getVal}$  is  $O(1)$ .

Q Can we do better?  
- faster?  
 $\rightarrow$  simpler?



(4)

⇒ Suggestion: Have array, use `getVal(o)` to determine index at which `o` should be stored.

Example	values	
"red"	5	
"orange"	1	
"yellow"	4	
"green"	2	
"blue"	3	
"violet"	3	

`add("red") → getVal("red") = 5`  
`add("blue") → ... = 3`  
`find("green") → 2`  
`add("violet") → oops!!`

⇒ Dealing with collisions:

"Chaining" each index of array refers to linked list of elements

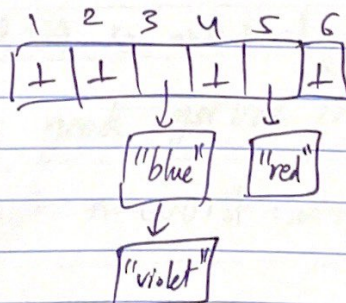
→ add: search for elt:

1. go to associated index
2. search list
3. append element if not found

→ find: do 1 & 2 above.

→ remove: 1 and 2; remove elt from list if found.

\* prev. example:



Each index of array refers to head of linked list.



- to do  $op(x)$ .

↑ find, add, remove

- perform that operation on list at index `getVal(x)`

$\Rightarrow \text{add("green")} \{- \text{go to index 2}\}$

$\Rightarrow$  remove("blue") { - go to index 3.  
                                - remove node from list }

★ Question: If array of size  $(n)$  storing  $n$  elements, what is the expected time to find

$$am_r$$

01 - i - n-1

--	--	--	--	--	--

occupancy of  $\text{arr}[i] = n_i$

$$E(\text{occupancy}) = \sum_{i=0}^n (\text{occupancy of } i) P(\text{get } (i))$$

$$= (n_0 + n_1 + \dots + n_{n-1}) \frac{1}{n} = n \cdot \frac{1}{n} = \underline{\underline{1}}$$

Conclusion: expected time find/add/remove is  $O(1)$

Issue: `getVal()` is not truly random.

→ defining good `getVal` methods is an art form!

→ don't get n-worst case guarantees.