

Lecture 5

- Treaps
- Find, insert, delete, split, and join in treaps
- Randomized search trees
- Randomized search tree time costs

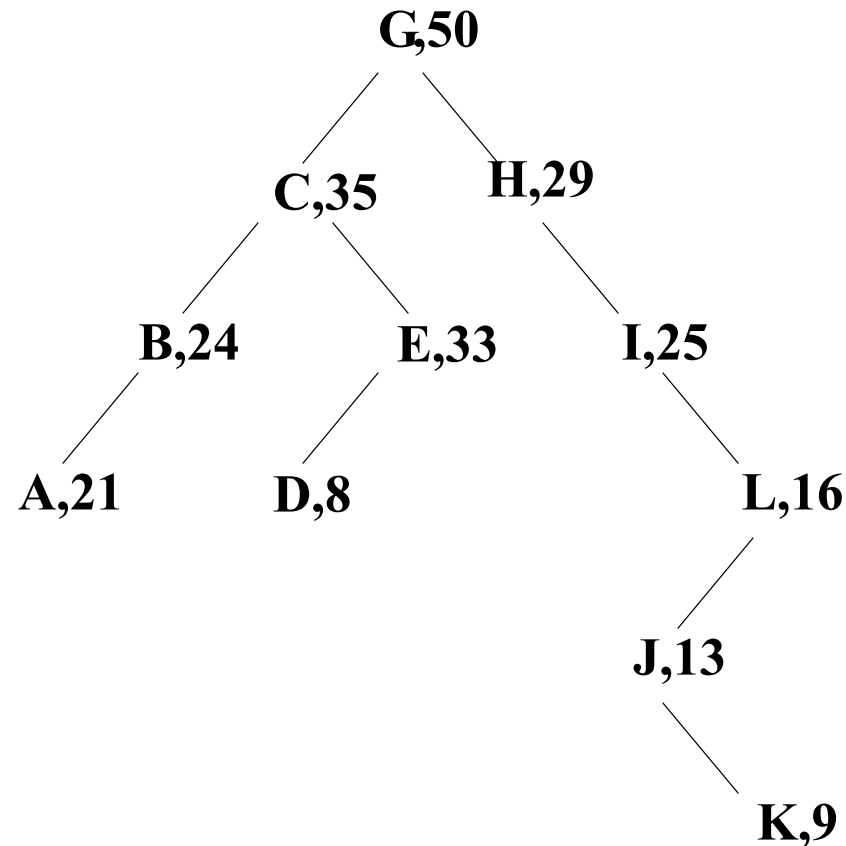
Reading: “Randomized Search Trees” by Aragon & Seidel, *Algorithmica* 1996,
<http://sims.berkeley.edu/~aragon/pubs/rst96.pdf>; Weiss, Chapter 12
section 5

Trees, heaps, and treaps

- A binary search tree (BST) is a binary tree:
 - Each Node in a BST contains a key; key values are comparable to each other
 - A BST has the BST ordering property: For every node X , the key in X is greater than all keys in the left subtree of X , and less than all keys in the right subtree of X
- A heap is a binary tree:
 - Each Node in a heap contains a priority; priority values are comparable to each other
 - A heap has the heap ordering property: For every node X , the priority in X is greater than or equal to all priorities in the left and right subtrees of X
- A treap is a binary tree:
 - Nodes in a treap contain both a key, and a priority
 - A treap has the BST ordering property with respect to its keys, and the heap ordering property with respect to its priorities

Treaps: an example

- Suppose keys are letters, with alphabetic ordering; priorities are integers, with numeric ordering. This tree is a treap:



Uniqueness of treaps

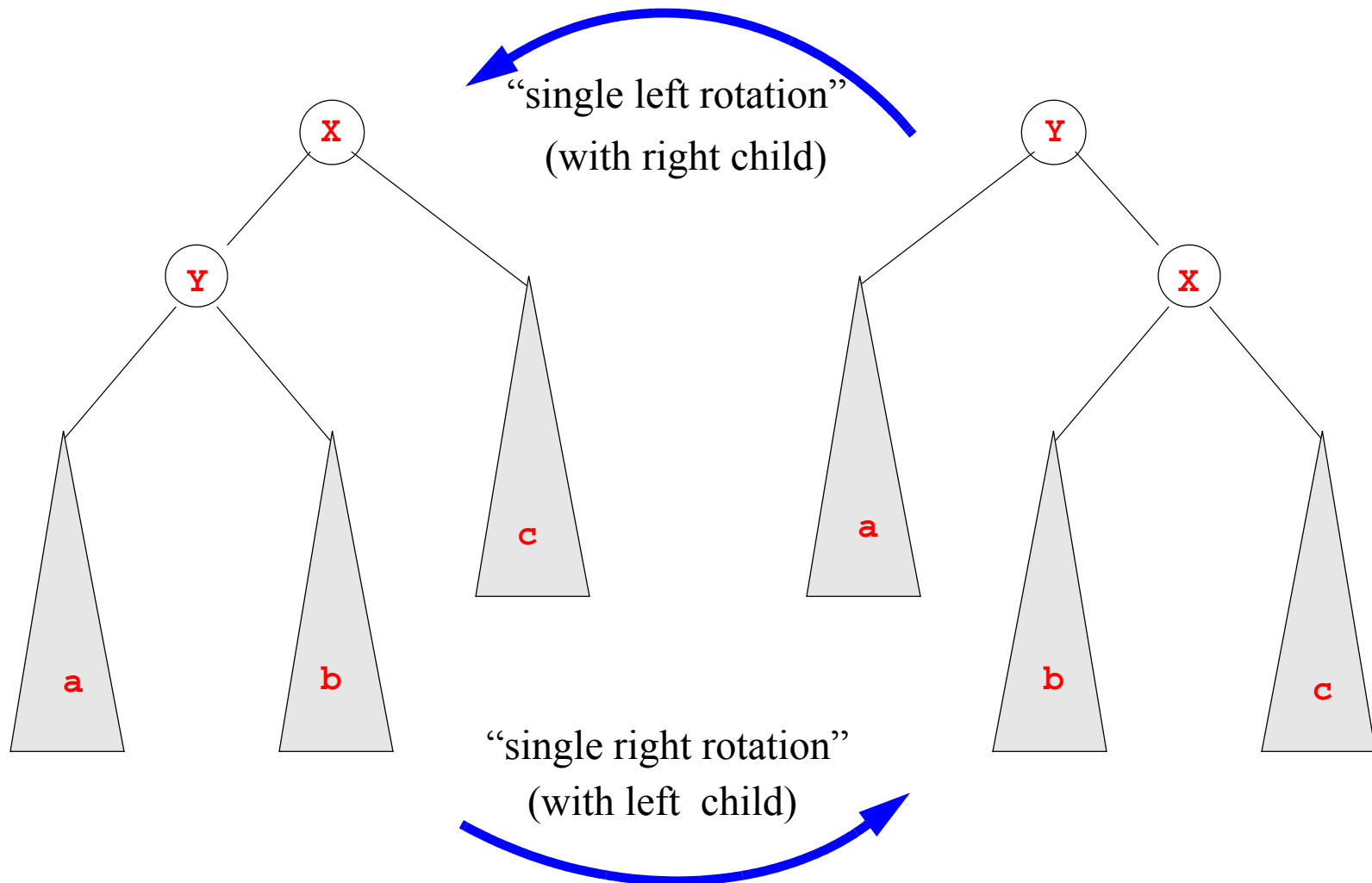
- Given a set of (key,priority) pairs, with all the key values unique, you could always construct a treap containing those (key,priority) pairs:
 - Start with an empty treap
 - Insert the (key,priority) pairs in decreasing order of priority, using the usual binary search tree insert algorithm that pays attention to the key values only
 - The result is a treap: BST ordering of keys is enforced by the BST insert, heap ordering of priorities is enforced by inserting in priority sorted order
- If the priority values as well as the key values are unique, the treap containing the (key,priority) pairs is unique
- For example, the treap on the previous page is the unique treap containing these pairs:
(G,50),(C,35),(E,33),(H,29),(I,25),(B,24),(A,21),(L,16),(J,13),(K,9),(D,8)
- Of course we will really be interested in algorithms that create a treap from (key,priority) pairs no matter what order they are inserted in.

Operations on treaps

- Treaps permit insert, delete, and find operations (also others but these are basic)
- Finding a key in a treap is very easy: just use the usual BST search algorithm
- Insert and delete of keys are slightly more complicated, since the operations must respect *both* the BST and heap ordering properties as invariants
- Recall that insert and delete in heaps use “bubble up” and “trickle down” exchanges to restore the heap ordering property
- Those same operations won’t work in treaps, because they can destroy the BST ordering property
- The trick is to use AVL rotations instead of exchanges:
 - An AVL rotation always preserves BST ordering, and so it can be used to move a (key,priority) pair up or down the tree to correct a failure of heap ordering without disturbing BST ordering

AVL rotations

- Recall the AVL single rotations:

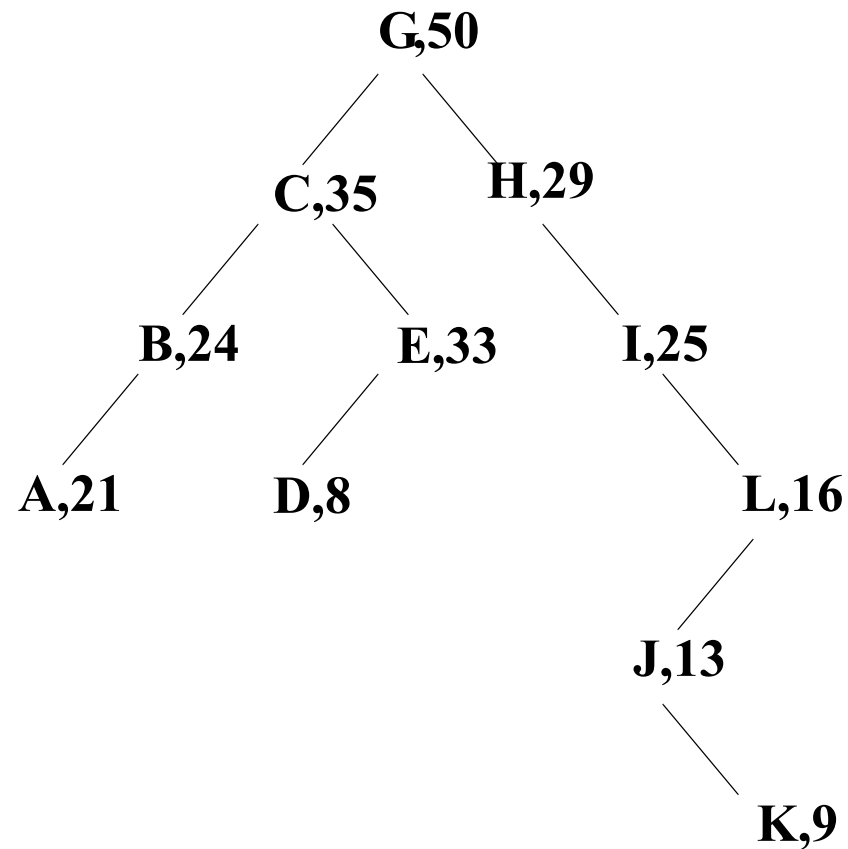


Insert in treaps

- To insert a key-priority pair (K,P) into a treap, do the following:
 - Insert the pair as a new leaf, using the usual BST insert algorithm which pays attention to the key value K
 - Then rotate the node up using AVL rotations as necessary, until the priority of its parent is greater than or equal to P , or the node becomes the root
 - (To rotate up, use a left rotation if the node is a right child of its parent, a right rotation if it is a left child)
- Since AVL rotations are constant-time operations, insert in a treap can be performed in time $O(H)$, where H is the height of the treap

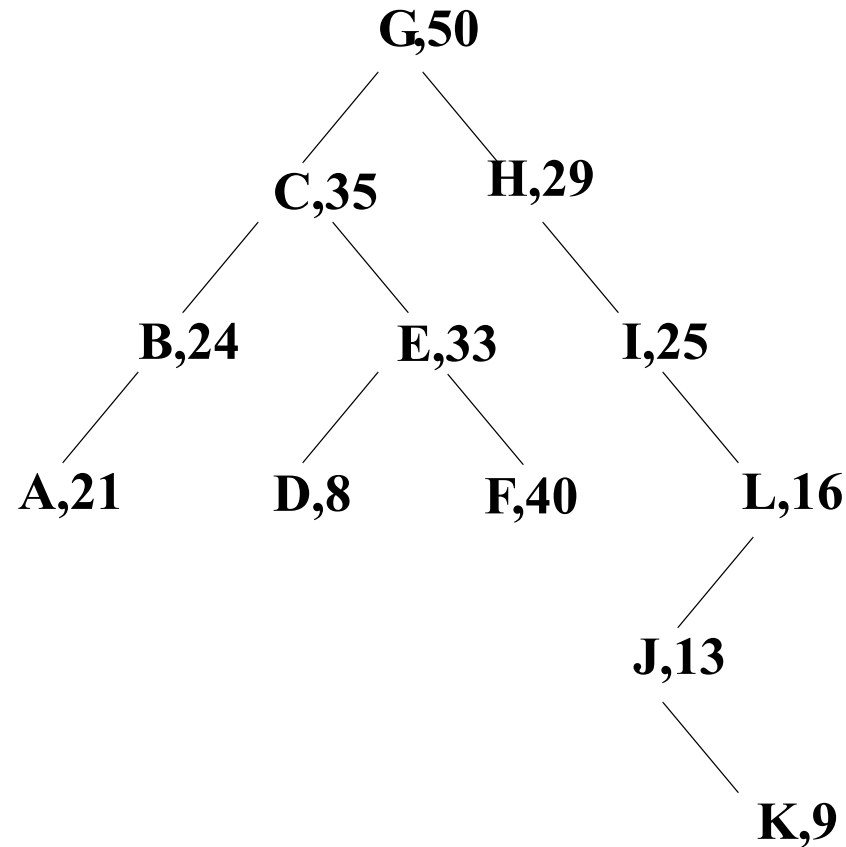
Insert in treaps: an example

- Insert the (key,priority) pair (F,40) in this treap:



Insert in treaps: an example, step 1

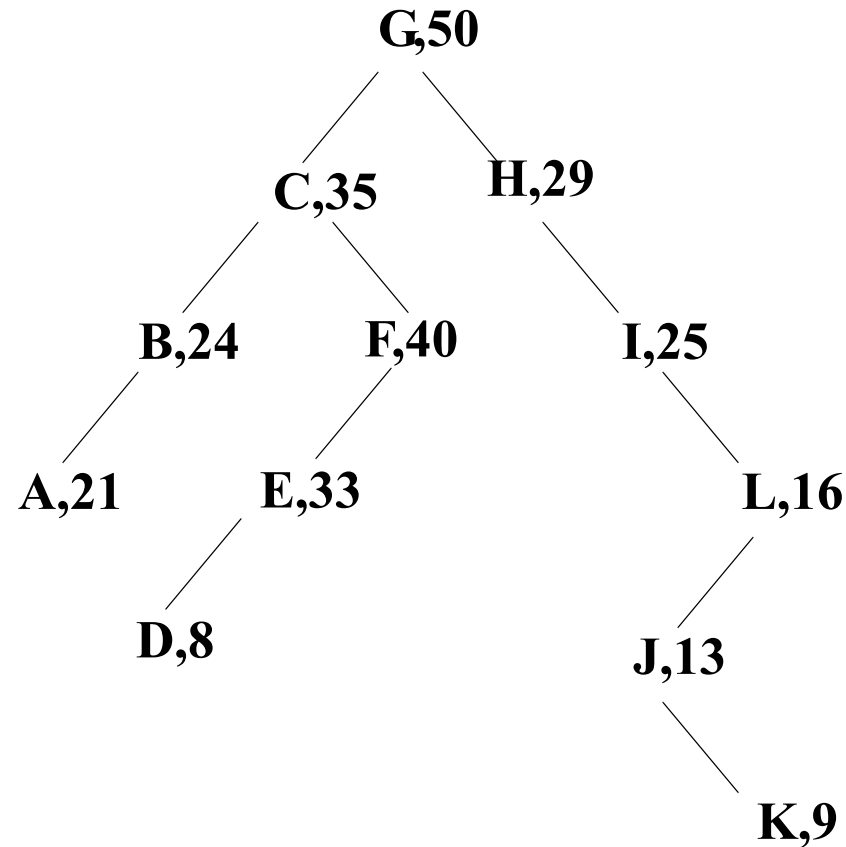
- Ordinary BST insert gives this:



- But the heap ordering property is violated. Need to rotate up to correct it

Insert in treaps: an example, step 2

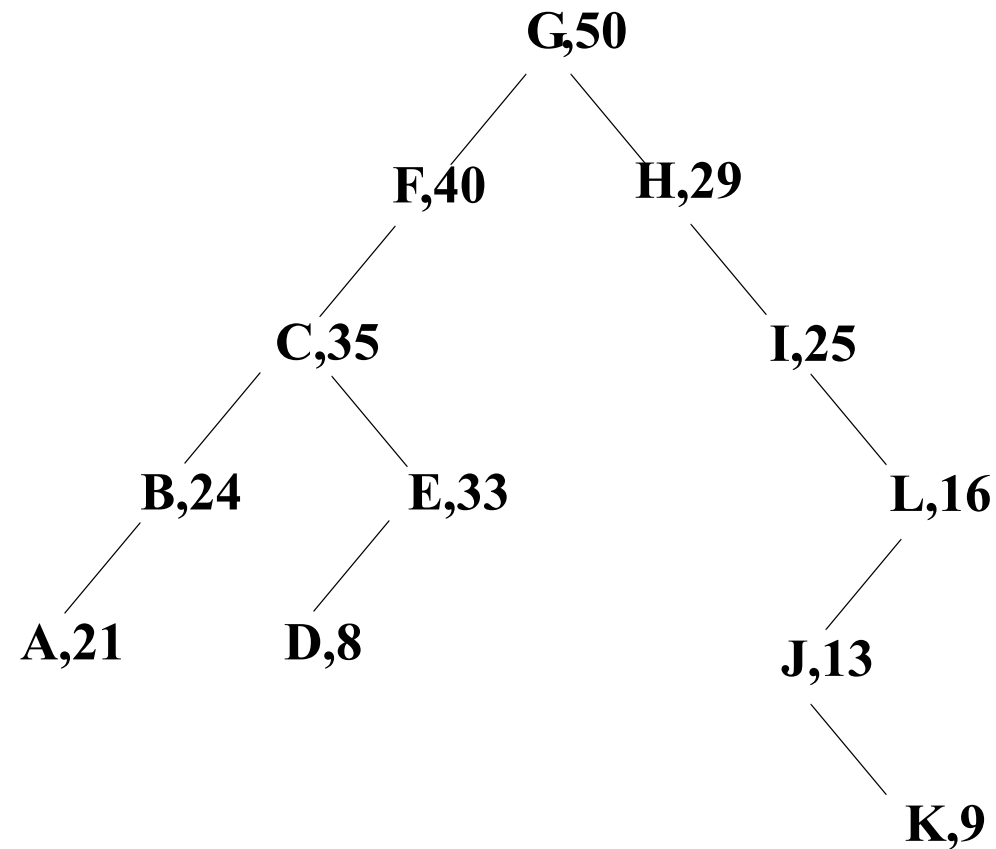
- One AVL left rotation gives this:



- But the heap ordering property is still violated. Need to rotate again to correct it

Insert in treaps: an example, step 3

- Now the tree is again a treap:

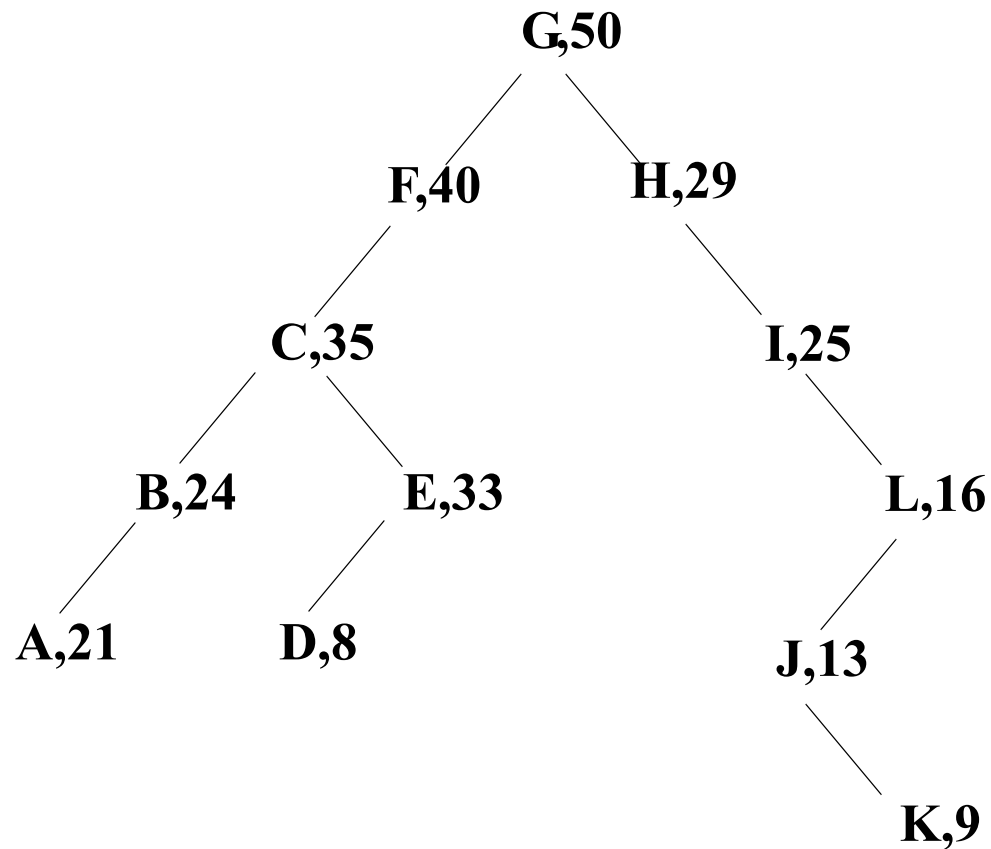


Delete in treaps

- To delete a key K , do the following:
 - Search for the node X containing K using the usual BST find algorithm
 - If the node X is a leaf, just delete the node (unlink it from its parent)
 - Otherwise, use AVL rotations to rotate the node down until it becomes a leaf; then delete it
 - (If there are 2 children, always rotate with the child that has the larger priority, to preserve heap ordering: use a left rotation if the right child has larger priority, right rotation otherwise)
- Since AVL rotations are constant-time operations, delete in a treap can be performed in time $O(H)$, where H is the height of the treap
- (Note that the AVL-rotation-to-leaf trick also works for delete in ordinary BST's)

Delete in treaps: an example

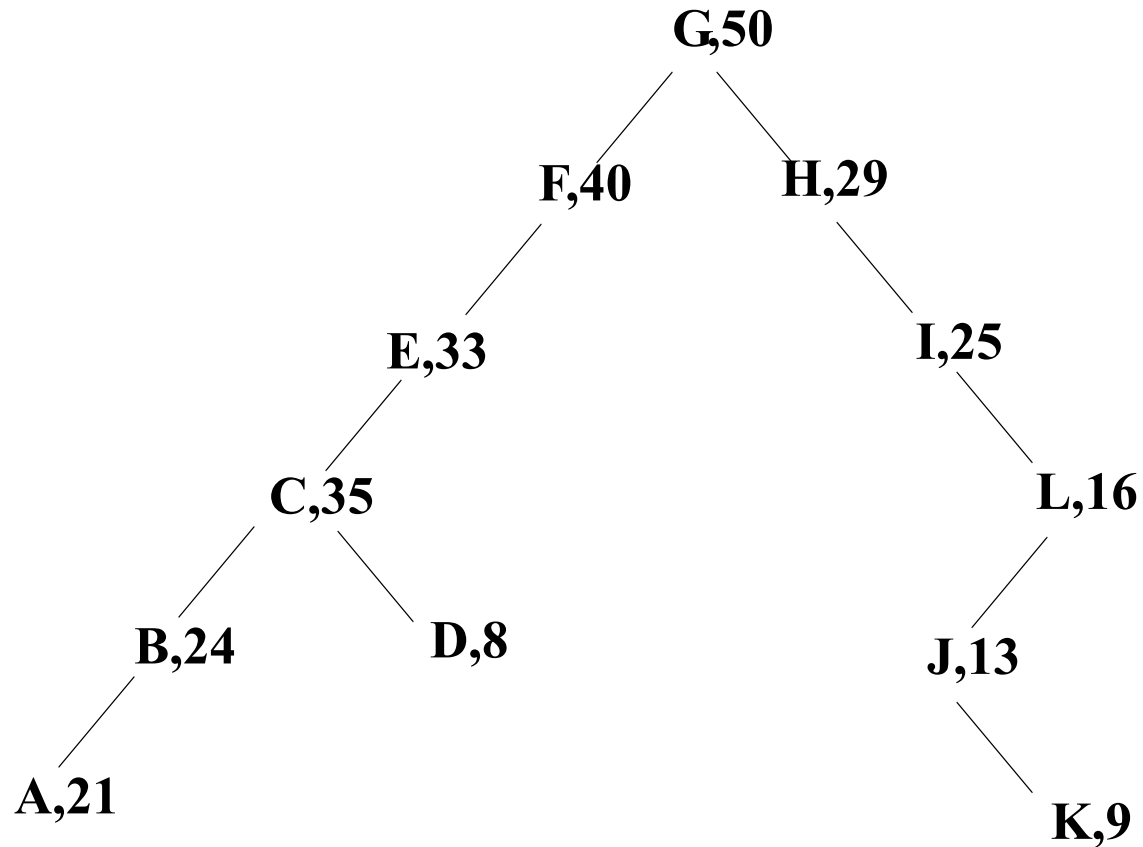
- Delete the key C from this treap:



- Find the node containing C. It is not a leaf, so need to rotate it down

Delete in treaps: an example, step 2

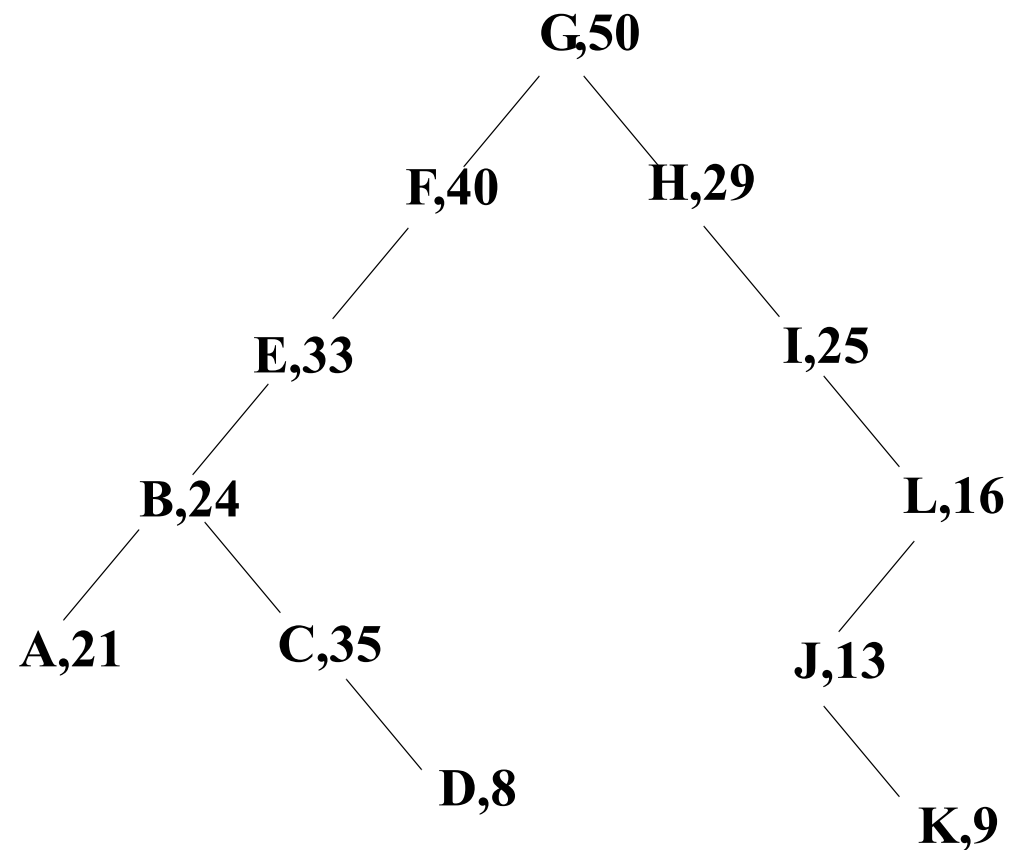
- Rotating the node with its larger priority child gives this:



- The node containing C is still not a leaf, so need to rotate down again

Delete in treaps: an example, step 3

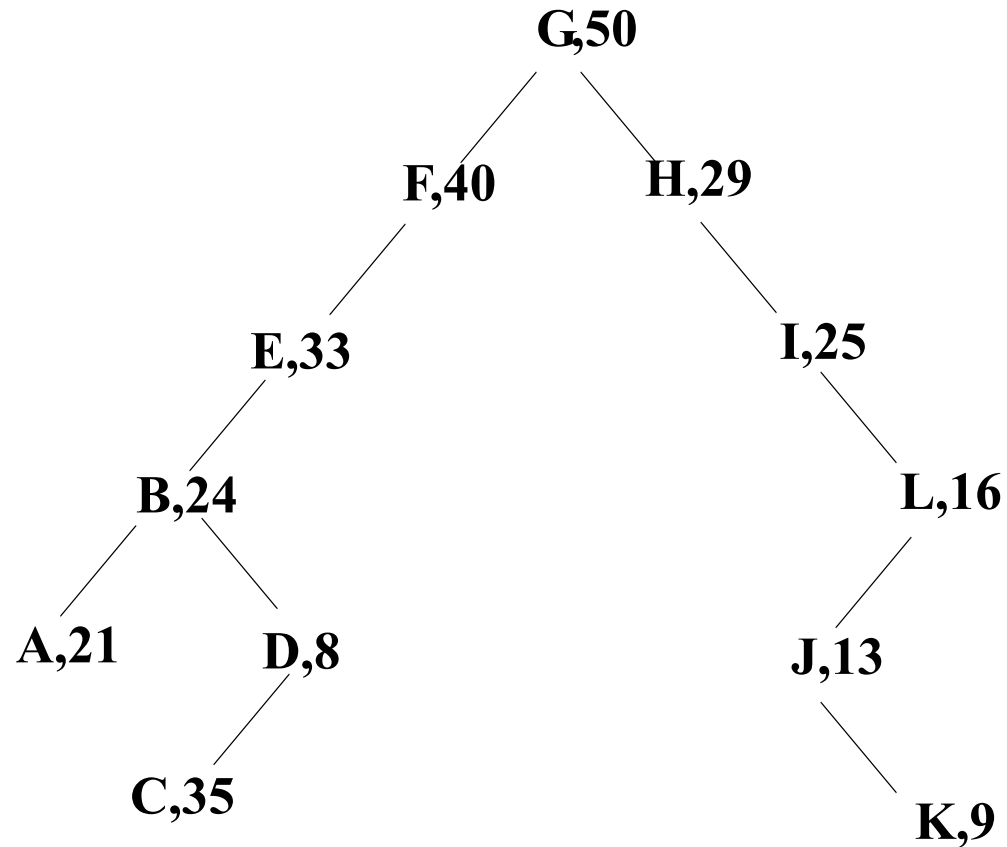
- Rotating the node with its larger priority child now gives this:



- The node containing C is still not a leaf, so need to rotate down yet again

Delete in treaps: an example, step 4

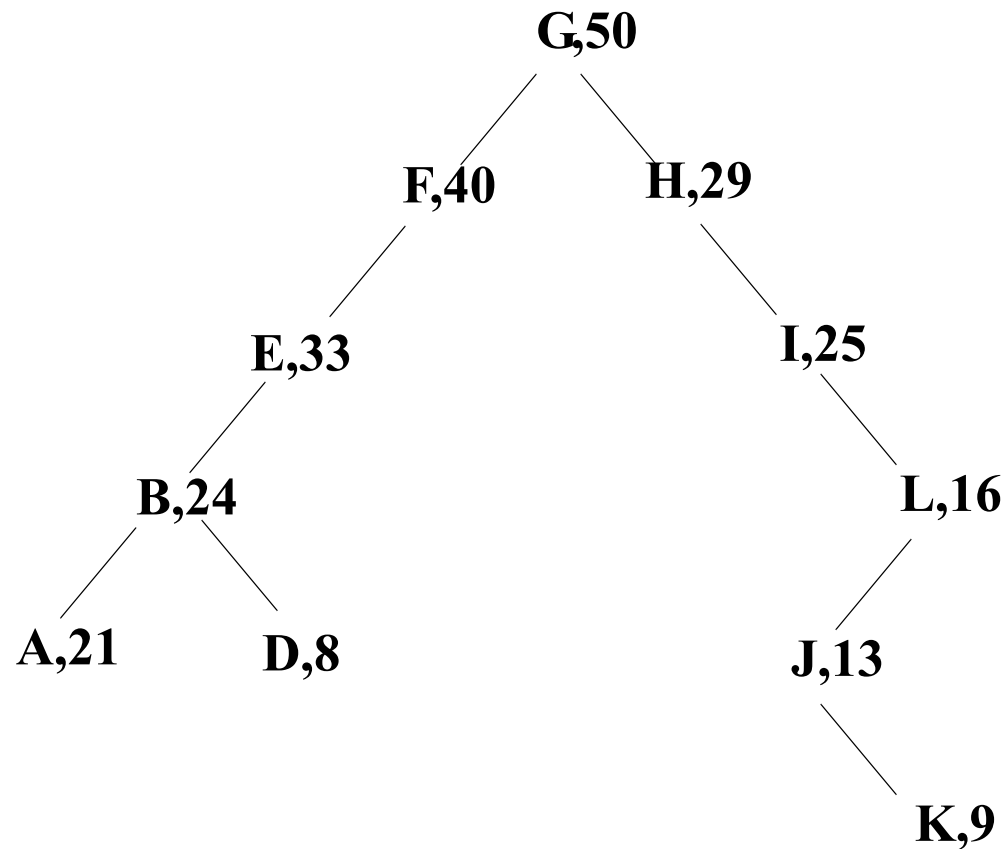
- Rotating the node with its child now gives this:



- Now the node containing C is a leaf, so just delete it

Delete in treaps: an example, step 5

- After clipping off the leaf, the result is again a treap:



Why treaps?

- Treaps are worth studying because...
 - they permit very easy implementations of *split* and *join* operations, as well as pretty simple implementations of *insert*, *delete*, and *find*
 - they are the basis of randomized search trees, which have performance comparable to balanced search trees but are simpler to implement
 - they also lend themselves well to more advanced tree concepts, such as weighted trees, interval trees, etc.
- We will look at the first two of these points

Tree splitting

- The tree splitting problem is this:
 - Given a tree and a key value K not in the tree, create two trees: One with keys less than K , and one with keys greater than K
- This is easy to solve with a treap, once the insert operation has been implemented:
 - Insert $(K, \text{INFINITY})$ in the treap
 - Since this has a higher priority than any node in the heap, it will become the root of the treap after insertion
 - Because of the BST ordering property, the left subtree of the root will be a treap with keys less than K , and the right subtree of the root will be a treap with keys greater than K . These subtrees then are the desired result of the split
- Since insert can be done in time $O(H)$ where H is the height of the treap, splitting can also be done in time $O(H)$
- (yes, this same idea could be used in an ordinary BST as well...)

Tree joining

- The tree joining or merging problem is this:
 - Given two trees T_1 , T_2 , such that each key in T_1 is less than all keys in T_2 , create a new tree T that contains all and only the keys from T_1 and T_2
- This is easy to do with a treap, once the delete operation has been implemented:
 - Create a “dummy” node with any key value and any priority
 - Make the root of T_1 be the left child, and the root of T_2 be the right child, of this dummy node
 - Perform a delete operation on the dummy node
- Since delete can be done in time $O(H)$ where H is the height of the treap, joining can also be done in time $O(H)$
- (yes, this same idea could be used in an ordinary BST as well...)

Disadvantages of treaps

- Treaps permit easy implementations of find, insert, delete, split, and join operations
- All these operations take worst case time $O(H)$, where H is the height of the treap
- However, treaps (like BST's) can become very unbalanced, so that $H = O(N)$, and that's bad
- Maintaining a strict balance condition (like the AVL or red-black property) in a treap would be impossible in general, if the user supplies both key values and priorities: remember that treaps with unique keys and priorities are unique
- However, if priorities are generated randomly by the insert algorithm, balance can be maintained with high probability and the operations stay very simple
 - that is the idea behind randomized search trees

Randomized search trees

- Randomized search trees were invented by Cecilia Aragon and Raimund Seidel, in early 1990's
- RST's are treaps in which priorities are assigned randomly by the insert algorithm when keys are inserted
- To implement a randomized search tree:
 - Adapt a treap implementation and its insert method that takes a (key,priority) pair as argument
 - To implement the RST insert method that takes a key as argument:
 - call a random number generator to generate a uniformly distributed random priority (a 32-bit random int is more than enough in typical applications; fewer bits can also be made to work well) that is independent of the key
 - call the treap insert method with the key value and that priority
- That's all there is to it: none of the other treap operations need to be changed at all
- (The RST implementation should take care to hide the random priorities, however)

Analysis of randomized search trees

- We will do an average case analysis of the “successful find” operation: how many steps are required, on average, to find that the key you’re looking for is in the tree?
- Suppose you have a RST with N nodes x_1, \dots, x_N , holding keys k_1, \dots, k_N and priorities p_1, \dots, p_N , such that x_i is the node holding key k_i and priority p_i
- As always when doing average-case analysis, you have to be clear about your probabilistic assumptions. We will make 2:
 - Assumption #1: Each key k_1, \dots, k_N in the tree is equally likely to be searched for
 - Assumption #2: The priorities p_1, \dots, p_N are randomly uniformly generated independently of each other and of the keys
- For convenience we will assume that:
 - keys are listed in sorted order: $k_i < k_{i+1}$ for all $0 < i < N$ (though keys can be inserted in any order),
 - all priorities are distinct

Expected node depth

- Let the depth of node x_i be $d(x_i)$, so the number of comparisons required to find key k_i in this tree is $d(x_i)$
- First, we will find the expected value (i.e. average) of $d(x_i)$, the depth of the node containing the i^{th} smallest key
- For this average-case analysis, we will average (not over all key insertion sequences, but) over all ways of generating random priorities during key insertion
- Let $Pr(p_1, \dots, p_N)$ be the probability of generating the N priority values p_1, \dots, p_N
 - Note that under the assumption that keys k_1, \dots, k_N are listed in sorted order, the priority values p_1, \dots, p_N determine the shape of the treap and the location of every key in it, so in particular they determine the depth of node x_i
- Then we can write the expected value (i.e. average) of $d(x_i)$ as:

$$E[d(x_i)] = \sum_{p_1, \dots, p_N} Pr(p_1, \dots, p_N) d(x_i)$$

Expected depth and ancestors

- Define A_{ij} to be the “indicator function” for the RST’s ancestor relation:

$$A_{ij} = \begin{cases} 1, & \text{if } x_i \text{ is an ancestor of } x_j \\ 0, & \text{otherwise} \end{cases}$$

- (We will take $A_{ii} = 1$ for all i , i.e. we consider a node to be an ancestor of itself.)
- Now note that the depth of a node is just the number of ancestors it has; so we can write

$$d(x_i) = \sum_{m=1}^N A_{mi}$$

- ... and so the expected value of the depth of node x_i is (using the fact that the expectation of a sum is the sum of expectations):

$$E[d(x_i)] = \sum_{p_1, \dots, p_N} Pr(p_1, \dots, p_N) \sum_{m=1}^N A_{mi} = \sum_{m=1}^N E[A_{mi}]$$

- So now what is the expected value of A_{mi} ?

Probability of being an ancestor

- A_{mi} is an indicator function: it is a random variable that takes only values 0,1
- The expected value of any indicator function is just equal to the probability that that indicator function has value 1
- So, $E[A_{mi}] = Pr(A_{mi} = 1)$, i.e., the probability that node x_m is an ancestor of x_i
- Seidel & Aragon 1996 prove the following lemma (recall we are assuming that if $i < j$, the key in node x_i is smaller than the key in x_j , and priorities are distinct):

Lemma: x_m is an ancestor of x_i if and only if among all priorities p_h such that h lies between the indices m and i inclusive, p_m is the largest

- So, probability that node x_m is an ancestor of x_i is just the probability that the random priority generated for x_m is higher than the other $|m - i|$ priorities generated for nodes with indexes between m and i inclusive
- But since the priorities are generated randomly and independently, each of those $|m - i| + 1$ nodes have equal probability of having the highest priority! So,

$$E[A_{mi}] = \frac{1}{|m - i| + 1}$$

Expected depth of node x_i

- This lets us get what we were seeking first, the expected depth of the node with key k_i in a RST:

$$E[d(x_i)] = \sum_{m=1}^N E[A_{mi}] = \sum_{m=1}^N \frac{1}{|m-i|+1}$$

- It is interesting to see that this depends on i , the position of the key in the ordered set of keys in the RST. Keys in the middle of the ordering will on average be somewhat deeper than smaller or larger keys. This is true of randomly constructed BST's as well
- For example, if the keys are integers $1,2,\dots,1000$, the expected depth of the node with key 500 is 12.59, while the expected depth of the node with key 1 or 1000 is 7.49
- But we are really interested in the average number of comparisons needed to find a key, assuming that all keys are equally likely to be searched. This is just the expected node depth, averaged over all nodes:

$$D_{avg}(N) = \sum_{i=1}^N \frac{1}{N} (E[d(x_i)]) = \frac{1}{N} \sum_{i=1}^N \sum_{m=1}^N \frac{1}{|m-i|+1}$$

Simplifying the double summation

- Let's simplify that double summation

$$\sum_{i=1}^N \sum_{m=1}^N \frac{1}{|m-i|+1}$$

- Note that there are N^2 terms. We can write the N^2 different values of the denominator $|m-i|+1$ as entries in an $N \times N$ matrix, where rows are indexed by m , columns by i . For example, for $N=9$:

1	2	3	4	5	6	7	8	9
2	1	2	3	4	5	6	7	8
3	2	1	2	3	4	5	6	7
4	3	2	1	2	3	4	5	6
5	4	3	2	1	2	3	4	5
6	5	4	3	2	1	2	3	4
7	6	5	4	3	2	1	2	3
8	7	6	5	4	3	2	1	2
9	8	7	6	5	4	3	2	1

- You can see that in general there will be N 1's, $2(N-1)$ 2's, $2(N-2)$ 3's, ..., $4(N-1)$'s, and $2N$'s
- That lets us rewrite the double summation, as shown next

The solution

- So we can write

$$\sum_{i=1}^N \sum_{m=1}^N \frac{1}{|m-i|+1} = 2 \sum_{i=1}^N \frac{N-i+1}{i} - N = 2(N+1) \sum_{i=1}^N \frac{1}{i} - 3N$$

- And thus we have that the average number of comparisons for a successful find (i.e., the average node depth) in a randomized search tree is

$$D_{avg}(N) = \frac{2(N+1)}{N} \sum_{i=1}^N \frac{1}{i} - 3$$

- This is exactly the same as the average node depth for a binary search tree, under the assumption that all key insertion sequences are equally likely!
- Think about it, that seems right:
 - The RST's treap structure is identical to a BST with keys inserted in order of their priority
 - With random priorities, all priority orderings in an RST are equally likely
- So what is the difference in practice between a RST and a BST?

Comparing RST's and vanilla BST's

- We have seen that in the average depth of a node in a N-node RST is the same as in a N-node BST: for large N it is approximately $2 \ln(N) = 1.386 \log_2 N$, which is $O(\log N)$
- This seems very good, but the analysis for the BST depended on the assumption that all sequences of key insertions were equally likely
- Often in practice “bad” sequences of BST key insertions (in which the keys are somewhat sorted) can in fact be more likely than others
- Also, if a particular sequence is “bad”, it will be bad (leading to much worse than $2 \ln N$ average depth) *every* time a BST is built with that sequence
- However, in a randomized search tree, the average case analysis is *independent of the sequence of key insertions*
- If a good random number generator is used to generate the treap priorities, the probability of constructing a “bad” randomized search tree is very low, no matter what the sequence of key insertions is
- That's why RST's are better than vanilla BST's!

More properties of randomized search trees

- The expected value of node depth in a RST is $O(\log N)$, and thus average time cost for successful find is $O(\log N)$
- Similar considerations show that unsuccessful find, insert, delete, split, and join in a RST all have average time cost $O(\log N)$
- It is *possible* for a randomized search tree to be badly unbalanced, with height significantly worse than $\log_2 N$, where N is the number of nodes in the treap; however this is unlikely to happen
 - To be badly unbalanced, the random priorities have to be correlated in a certain way with key values, and with a good random number generator this will be unlikely to occur
- An interesting question is: *How* unlikely is that to happen?

More properties of randomized search trees, cont'd

- The expected time costs are like average time costs, averaged over many constructions of a treap with the same N keys, but different random priorities
- Question: For a single randomized search tree with N keys, how likely is it that its height is much greater than a deterministically balanced search tree such as AVL?
- Aragon and Seidel analyzed this question and showed that the answer is:
 - It is possible, but it can be considered extremely unlikely
- They derive this formula (here e is the natural logarithm base, \ln is log base e , and c is any positive constant):

$$Pr(H > 2c \ln N) < N \left(\frac{N}{e} \right)^{-c \ln(c/e)}$$

- Example: $N=10000$, so $\ln N = 9.210\dots$. Pick $c = 5.429\dots$ so $2c \ln N = 100$. Plugging in numbers, we get $Pr(H > 100) < 4.081 \cdot 10^{-10}$
- That is: the chance that the height of a randomized search tree with 10,000 nodes will be greater than 100 is less than one in two billion
- So, although randomized search trees do not absolutely guarantee good performance, they will almost certainly provide good performance in practice

Next time

- Randomized data structures
- Random number generation
- Skip lists: ideas and implementation
- Skip list time costs

Reading:

“Skip Lists: A Probabilistic Alternative to Balanced Trees” (author William Pugh, available online)

Weiss, Chapter 10 section 4