# Advanced Data Structures

## Medians and Order Statistics

# Order Statistics

- The *i*th *order statistic* in a set of *n* elements is the *i*th smallest element
- The *minimum* is thus the 1st order statistic
- The *maximum* is the *n*th order statistic
- The *median* is the *n*/2 order statistic
    - If *n* is even, there are 2 medians
- *How can we calculate order statistics?*
- *What is the running time?*

# Order Statistics

- *Ho man comparisons are needed to find the minimum element in a set? The ma imum?*

- *Can e find the minimum and ma imum ith less than t ice the cost?*

- Yes:
  - Walk through elements by pairs
    - Compare each element in pair to the other
    - Compare the largest to maximum, smallest to minimum
  - Total cost: 3 comparisons per 2 elements = O(3n /2)

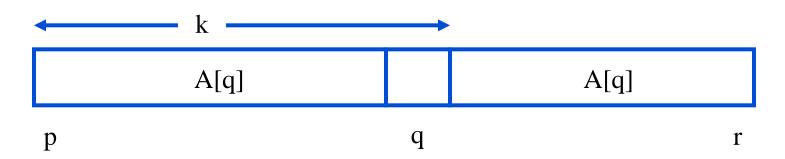# Finding Order Statistics: The Selection Problem

- A more interesting problem is *selection*: finding the $i$th smallest element of a set

- We will show:
  - A practical randomized algorithm with O(n) expected running time
  - A cool algorithm of theoretical interest only with O(n) worst-case running time

# Randomized Selection

- Key idea: use partition() from quicksort
  - But, only need to examine one subarray
  - This savings shows up in running time: O(n)

  q = RandomizedPartition(A, p, r)

| A[q] | | A[q] |
|------|--|------|

p                                q                                r

# Randomized Selection

```
RandomizedSelect(A, p, r, i)
    if (p == r) then return A[p];
    q = RandomizedPartition(A, p, r)
    k = q - p + 1;
    if (i == k) then return A[q];
    if (i < k) then
        return RandomizedSelect(A, p, q-1, i);
    else
        return RandomizedSelect(A, q+1, r, i-k);
```

# Randomized Selection

- Analyzing `RandomizedSelect()`
  - Worst case: partition always 0:n-1

    $T(n) = T(n-1) + O(n)$  = **???**

    $\quad\quad\quad = O(n^2)$     (arithmetic series)

    - No better than sorting!

  - "Best" case: suppose a 9:1 partition

    $T(n) = T(9n/10) + O(n)$  = **???**

    $\quad\quad\quad = O(n)$     (Master Theorem, case 3)

    - Better than sorting!

    - *What if this had been a 99:1 split?*

# Randomized Selection

- Average case
  - For upper bound, assume $i$

# Randomized Selection

● Assume $T(n) \le cn$ for sufficiently large $c$:

$$T(n) \le \frac{}{n} \sum_{k=n/}^{n-1} T(k) + \Theta(n)$$

$$\le \frac{}{n} \sum_{k=n/}^{n-1} ck + \Theta(n) \qquad \textcolor{blue}{\textit{W}} \qquad \textcolor{blue}{\textbf{?}}$$

$$= \frac{c}{n} \left( \sum_{k=1}^{n-1} k - \sum_{k=1}^{n/-1} k \right) + \Theta(n) \qquad \textcolor{red}{\textbf{\textit{"S}}} \quad \textbf{"}$$

$$= \frac{c}{n} \left( \frac{1}{}(n-1)n - \frac{1}{} \left( \frac{n}{} - 1 \right) \frac{n}{} \right) + \Theta(n) \quad \textcolor{blue}{\textit{W}} \qquad \textcolor{blue}{\textbf{?}}$$

$$= c(n-1) - \frac{c}{} \left( \frac{n}{} - 1 \right) + \Theta(n) \qquad \textcolor{blue}{\textit{W}} \qquad \textcolor{blue}{\textbf{?}}$$

# Randomized Selection

● Assume $T(n) \leq cn$ for sufficiently large $c$:

$$T(n) \leq c(n-1) - \frac{c}{2}\left(\frac{n}{2} - 1\right) + \Theta(n) \qquad \textbf{\textit{T}}$$

$$= cn - c - \frac{cn}{4} + \frac{c}{2} + \Theta(n) \qquad \textbf{\textit{M}}$$

$$= cn - \frac{cn}{4} - \frac{c}{2} + \Theta(n) \qquad \textbf{\textit{S}} \qquad \textbf{\textit{/2}}$$

$$= cn - \left(\frac{cn}{4} + \frac{c}{2} - \Theta(n)\right) \qquad \textbf{\textit{R}}$$

$$\leq cn \quad \text{(if c is big enough)} \qquad \textbf{\textit{W}}$$

# Worst-Case Linear-Time Selection

- Randomized algorithm works well in practice
- What follows is a worst-case linear time algorithm, really of theoretical interest only
- Basic idea:
  - Generate a good partitioning element
  - Call this element

# Worst-Case Linear-Time Selection

● The algorithm in words:

1. Divide $n$ elements into groups of 5

2. Find median of each group (*Ho  ?  Ho   long?*)

3. Use Select() recursively to find median  of the $\lfloor n/5 \rfloor$ medians

4. Partition the $n$ elements around  .  Let $k = \text{rank}(\ )$

5. **if** (i == k) **then** return x

   **if** (i < k) **then** use Select() recursively to find $i$th smallest element in first partition

   **else** (i > k) use Select() recursively to find ($i$-$k$)th smallest element in last partition

# Worst-Case Linear-Time Selection

- *How many of the 5-element medians are ≤ ?*
  - At least 1/2 of the medians $= \lfloor \lfloor n/5 \rfloor / 2 \rfloor = \lfloor n/10 \rfloor$
- *How many elements are ≤ ?*
  - At least $3 \lfloor n/10 \rfloor$ elements
- For large $n$, $3 \lfloor n/10 \rfloor \geq n/4$  *(How large?)*
- So at least $n/4$ elements $\leq$
- Similarly: at least $n/4$ elements $\geq$

# Worst-Case Linear-Time Selection

- Thus after partitioning around  , step 5 will call Select() on at most $3n/4$ elements

- The recurrence is therefore:

$$T(n) \quad T(\lfloor n/5 \rfloor) \quad T(3n/4) \qquad (n)$$

$$T(n/5) \quad T(3n/4) \qquad (n)$$

$$cn/5 \quad 3cn/4 \quad \Theta(n)$$

$$= 19cn/20 + \Theta(n)$$

$$= cn - (cn/20 - \Theta(n))$$

$$\leq cn \quad \text{if } c \text{ is big enough}$$

# Worst-Case Linear-Time Selection

● Intuitively:

- Work at each level is a constant fraction (19/20) smaller

  - Geometric progression!

- Thus the O(n) work at the root dominates

# Linear-Time Median Selection

● Given a "black box" O(n) median algorithm, what can we do?

■ *i*th order statistic:

◆ Find median

◆ Partition input around

◆ if (*i*   (n+1)/2) recursively find *i*th element of first half

◆ else find (*i* - (n+1)/2)th element in second half

◆ T(n) = T(n/2) + O(n) = O(n)

# Linear-Time Median Selection

- Worst-case O(n lg n) quicksort
  - Find median   and partition around it
  - Recursively quicksort two halves
  - T(n) = 2T(n/2) + O(n) = O(n lg n)

# Dynamic Order Statistics

- We've seen algorithms for finding the $i$th element of an unordered set in O($n$) time

- Next, a structure to support finding the $i$th element of a dynamic set in O(lg $n$) time

  - *What operations do d namic sets usuall  support?*

  - *What structure   orks    ell for these?*

  - *Ho    could    e use this str -0 (t) 0 (r)-0 (uc) 0R55379cm*

# Order Statistic Trees

- OS Trees augment red-black trees:
  - Associate a *si e* field with each node in the tree
  - **x->size** records the size of subtree rooted at **x**, including **x** itself:
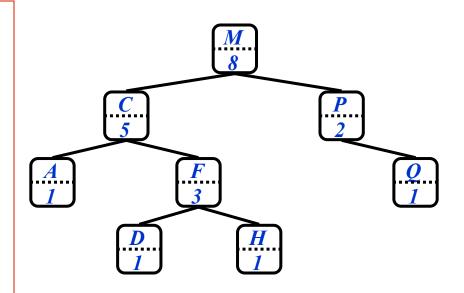
# Selection On OS Trees



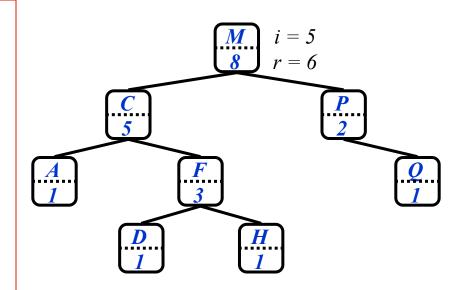*Ho    can    e use this propert
to select the ith element of the set?*

# OS-Select

```
OS-Select(x, i)

{

    r = x->left->size + 1;
    if (i == r)
        return x;
    else if (i < r)
        return OS-Select(x->left, i);
    else
        return OS-Select(x->right, i-r);

}
```
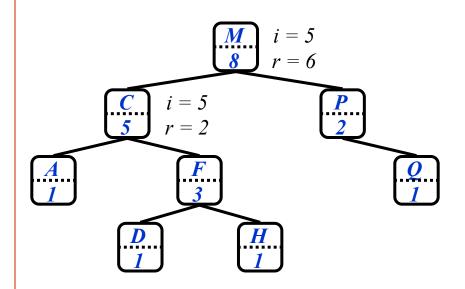
# OS-Select Example

- Example: show OS-Select(*root*, 5):

```
OS-Select(x, i)
{
  r = x->left->size + 1;
  if (i == r)
    return x;
  else if (i < r)
    return OS-Select(x->left, i);
  else
    return OS-Select(x->right, i-r);
}
```

# OS-Select Example

- Example: show OS-Select(*root*, 5):
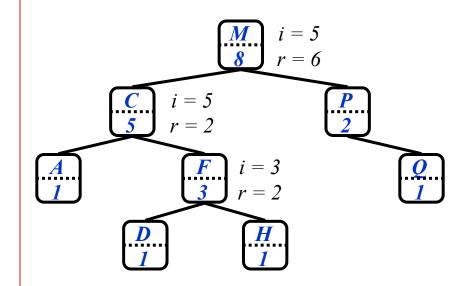
```
OS-Select(x, i)
{
  r = x->left->size + 1;
  if (i == r)
    return x;
  else if (i < r)
    return OS-Select(x->left, i);
  else
    return OS-Select(x->right, i-r);
}
```

# OS-Select Example

- Example: show OS-Select(*root*, 5):
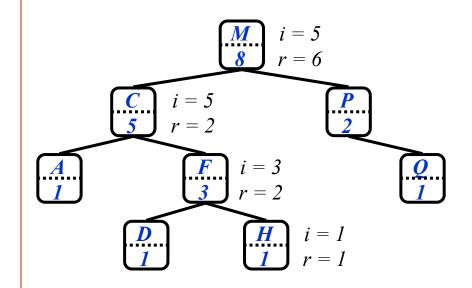
```
OS-Select(x, i)
{
  r = x->left->size + 1;
  if (i == r)
    return x;
  else if (i < r)
    return OS-Select(x->left, i);
  else
    return OS-Select(x->right, i-r);
}
```

# OS-Select Example

- Example: show OS-Select(*root*, 5):

```
OS-Select(x, i)
{
  r = x->left->size + 1;
  if (i == r)
    return x;
  else if (i < r)
    return OS-Select(x->left, i);
  else
    return OS-Select(x->right, i-r);
}
```

# OS-Select Example

- Example: show OS-Select(*root*, 5):

```
OS-Select(x, i)
{
  r = x->left->size + 1;
  if (i == r)
    return x;
  else if (i < r)
    return OS-Select(x->left, i);
  else
    return OS-Select(x->right, i-r);
}
```

# OS-Select: A Subtlety

```
OS-Select(x, i)
{
    r = x->left->size + 1;
    if (i == r)
        return x;
    else if (i < r)
        return OS-Select(x->left, i);
    else
        return OS-Select(x->right, i-r);
}
```

*Oops…*

- *What happens at the leaves?*

- *Ho   can   e deal elegantl    ith this?*

# OS-Select

```
OS-Select(x, i)
{
    r = x->left->size + 1;
    if (i == r)
        return x;
    else if (i < r)
        return OS-Select(x->left, i);
    else
        return OS-Select(x->right, i-r);
}
```
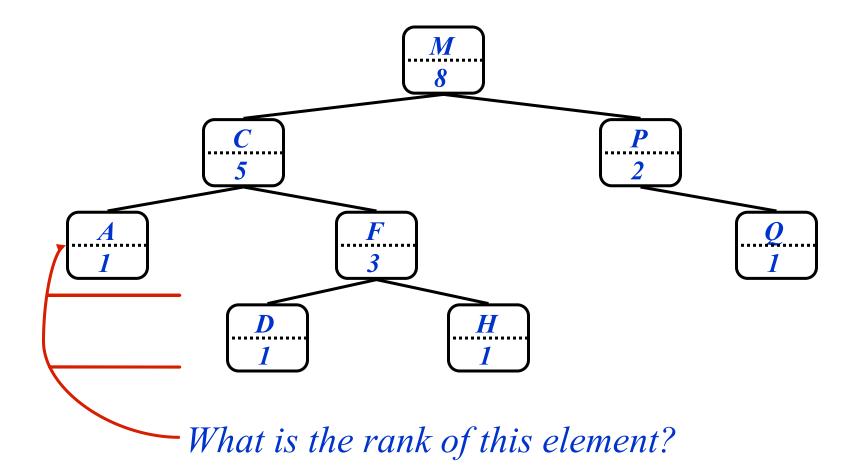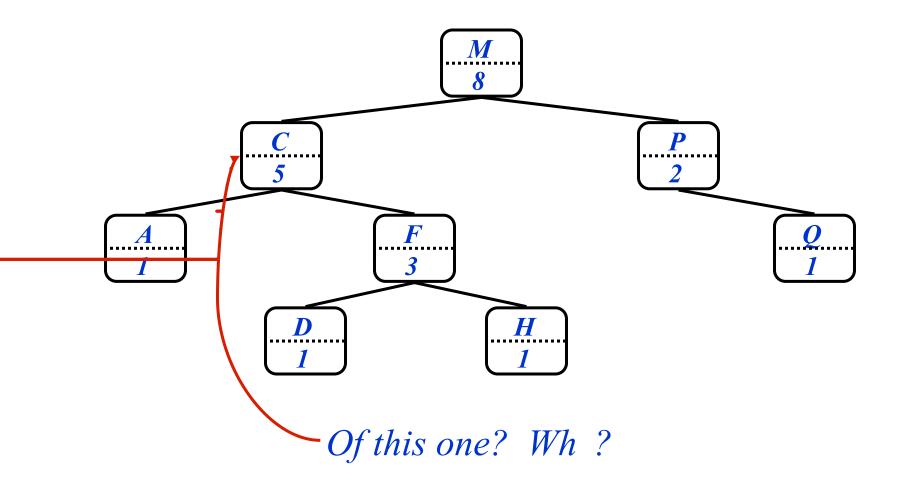
- *What   ill be the running time?*

# Determining The Rank Of An Element



*What is the rank of this element?*

# Determining The Rank Of An Element



*Of this one?  Wh  ?*

# Determining The Rank Of An Element

$M$

# Determining The
# Rank Of An Element



*What about the rank of this element?*

# Determining The Rank Of An Element



*This one? What's the pattern here?*

# OS-Rank

```
OS-Rank(T, x)
{
    r = x->left->size + 1;
    y = x;
    while (y != T->root)
        if (y == y->p->right)
            r = r + y->p->left->size + 1;
        y = y->p;
    return r;
}
```

- *What will be the running time?*

*E ample 1:*
*find rank of element   ith ke   H*

```
OS-Rank(T, x)

{

    r = x->left->size + 1;

    y = x;

    while (y != T->root)

        if (y == y->p->right)

            r = r + y->p->left->size + 1;

        y = y->p;

    return r;

}
```

*M*
*8*

*C*
*5*

*P*
*2*

*A*
*1*

*F*
*3*

*O*
*1*

*D*
*1*

*H*
*1*

*= 1*

# Determining The Rank Of An Element

*E ample 1:*

*find rank of element  ith ke  H*
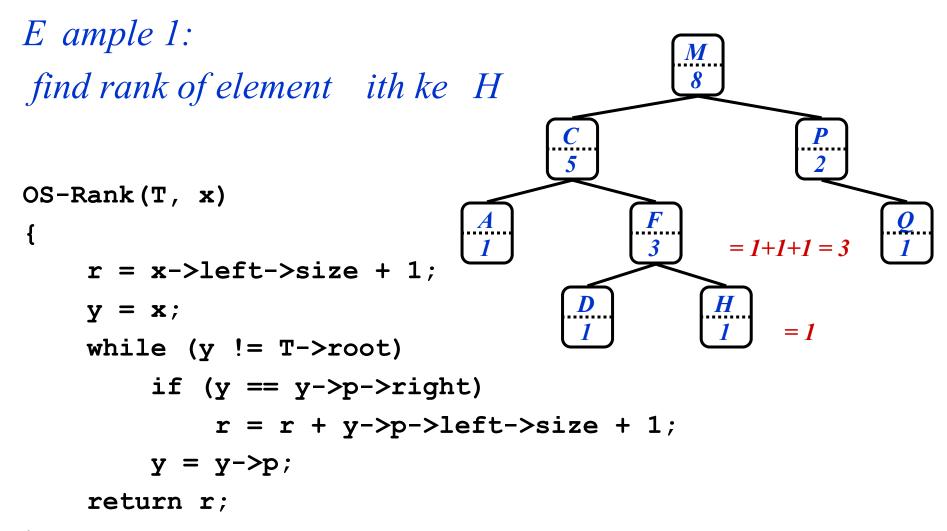


```
OS-Rank(T, x)
{
    r = x->left->size + 1;
    y = x;
    while (y != T->root)
        if (y == y->p->right)
            r = r + y->p->left->size + 1;
        y = y->p;
    return r;
}
```

# Determining The Rank Of An Element

*E ample 1:*

*find rank of element  ith ke  H*



```
OS-Rank(T, x)
{
    r = x->left->size + 1;
    y = x;
    while (y != T->root)
        if (y == y->p->right)
            r = r + y->p->left->size + 1;
        y = y->p;
    return r;
}
```

# Determining The Rank Of An Element

*E ample 1:*

*find rank of element  ith ke  H*

```
OS-Rank(T, x)
{
    r = x->left->size + 1;
    y = x;
    while (y != T->root)
        if (y == y->p->right)
            r = r + y->p->left->size + 1;
        y = y->p;
    return r;
}
```

*M 8* = 5

*C 5* = 5

*P 2*

*A 1*

*F 3* = 3

*Q 1*

*D 1*

*H 1* = 1

# Determining The Rank Of An Element

$$\frac{M}{8}$$

$$\frac{C}{5}$$

$$\frac{P}{2}$$

# Determining The Rhank Of An Element

*E ample 2:*

*find rank of element  ith ke  P*

```
OS-Rank(T, x)
{
    r = x->left->size + 1;
    y = x;
    while (y != T->root)
        if (y == y->p->right)
            r = r + y->p->left->size + 1;
        y = y->p;
    return r;
}
```
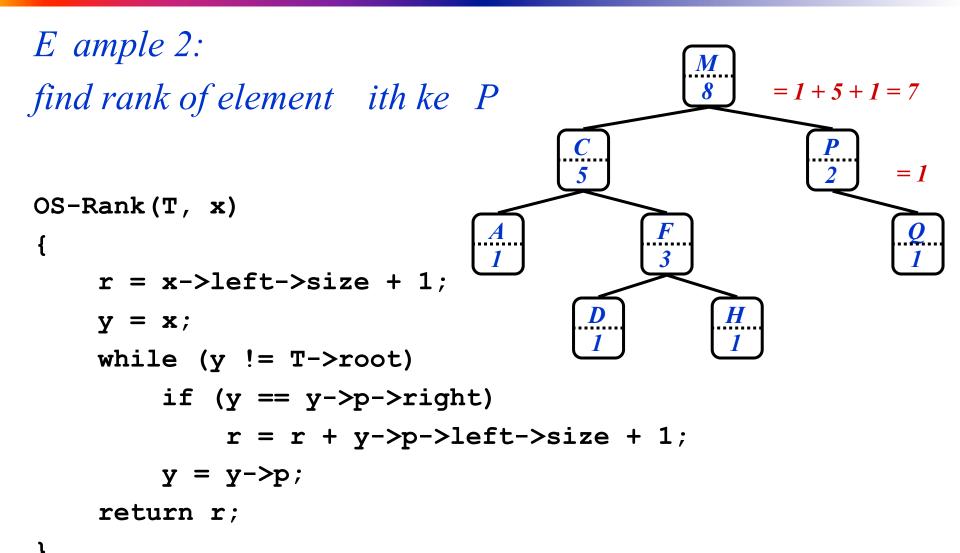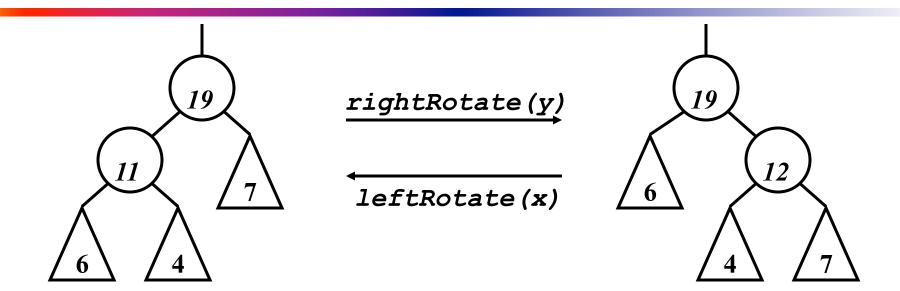


= 1 + 5 + 1 = 7

= 1

# Maintaining Subtree Sizes

- So by keeping subtree sizes, order statistic operations can be done in $O(\lg n)$ time

- Maintain sizes during Insert() and Delete() operations
  - Insert(): Increment size fields of nodes traversed during search down the tree
  - Delete(): Decrement sizes along a path from the deleted node to the root
  - Both: Update sizes correctly during rotations

# Maintaining Size Through Rotation



- Salient point: rotation invalidates only    and
- Can recalculate their sizes in constant time
  - *Wh  ?*

# Augmenting Data Structures: Methodology

- Choose underlying data structure
  - E.g., red-black trees
- Determine additional information to maintain
  - E.g., subtree sizes
- Verify that information can be maintained for operations that modify the structure
  - E.g., Insert(), Delete()   (don't forget rotations!)
- Develop new operations
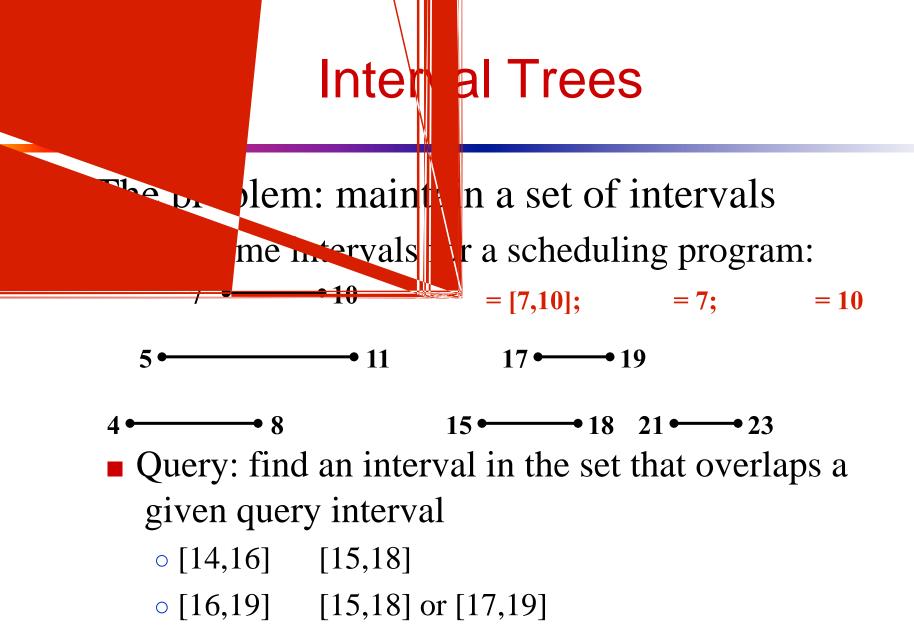  - E.g., OS-Rank(), OS-Select()

# Advanced Data Structures

Augmenting Data Structures:
Interval Trees

# Review: Methodology For Augmenting Data Structures

- Choose underlying data structure

- Determine additional information to maintain

- Verify that information can be maintained for operations that modify the structure

- Develop new operations

# Interval Trees

- The problem: maintain a set of intervals
  - E.g., time intervals for a scheduling program:

        7               10

# Interval Trees

The problem: maintain a set of intervals

...me intervals for a scheduling program:

7 ———— 10        = [7,10];        = 7;        = 10

5 ———— 11        17 ———— 19

4 ———— 8        15 ———— 18    21 ———— 23

- Query: find an interval in the set that overlaps a given query interval
  - [14,16]      [15,18]
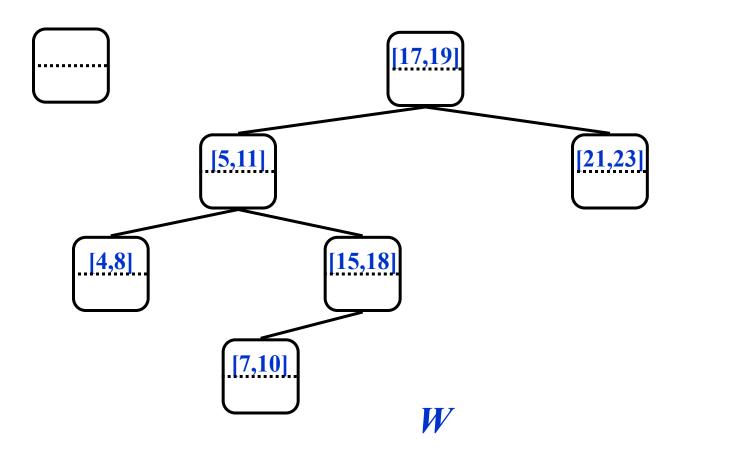  - [16,19]      [15,18] or [17,19]
  - [12,14]      NULL
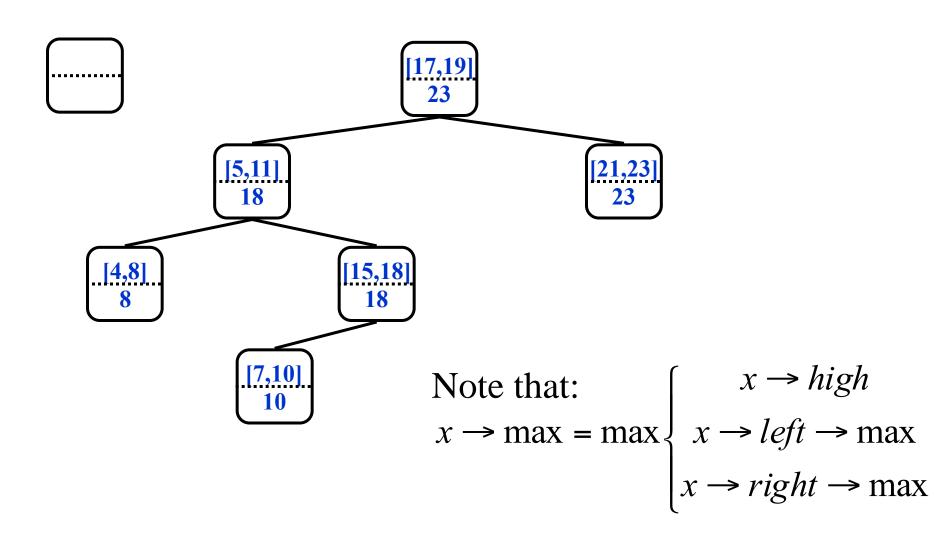
# Interval Trees

- Following the methodology:
    - Pick underlying data structure
    - Decide what additional information to store
    - Figure out how to maintain the information
    - Develop the desired new operations

# Interval Trees

- Following the methodology:
  - *Pick underl ing data structure*
    - Red-black trees will store intervals, keyed on $i \quad lo$
  - Decide what additional information to store
  - Figure out how to maintain the information
  - Develop the desired new operations

# Interval Trees

- Following the methodology:
  - Pick underlying data structure
    - Red-black trees will store intervals, keyed on $i \cdot lo$
  - *Decide what additional information to store*
    - We will store $max$, the maximum endpoint in the subtree rooted at $i$
  - Figure out how to maintain the information
  - Develop the desired new operations

# Interval Trees



[17,19]

[5,11]    [21,23]

[4,8]    [15,18]

[7,10]

*W*    *?*

# Interval Trees



Note that:
$$x \to \max = \max \begin{cases} x \to high \\ x \to left \to \max \\ x \to right \to \max \end{cases}$$
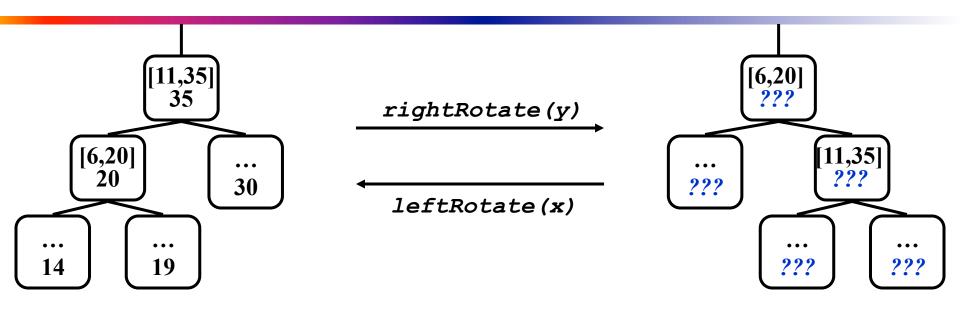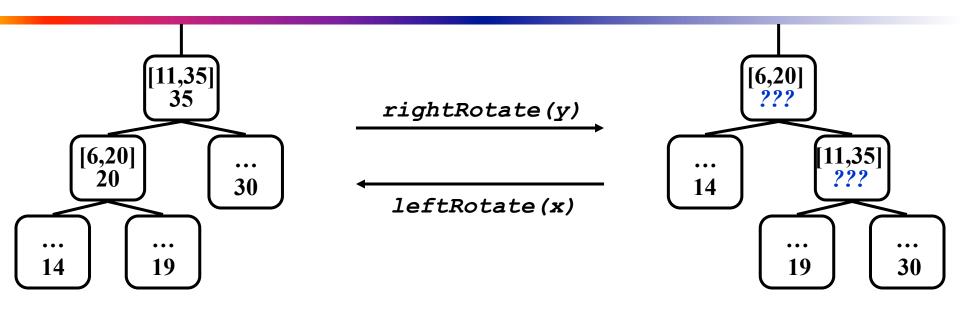
# Interval Trees

- Following the methodology:
  - Pick underlying data structure
    - Red-black trees will store intervals, keyed on $i$ $lo$
  - Decide what additional information to store
    - Store the maximum endpoint in the subtree rooted at $i$
  - *Figure out ho to maintain the information*
    - *Ho ould e maintain ma field for a BST?*
    - *What's different?*
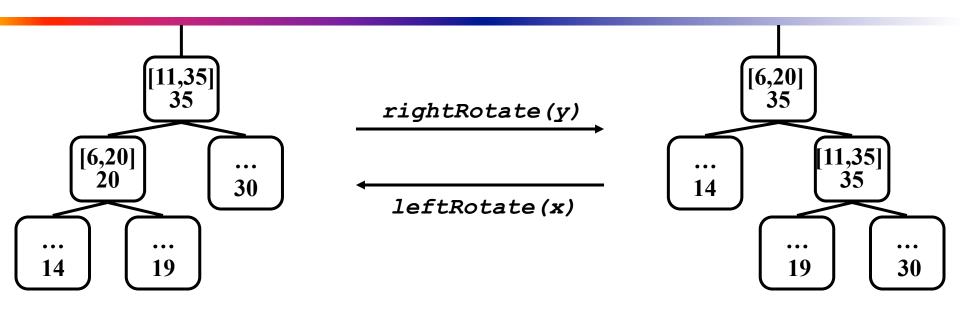  - Develop the desired new operations

# Interval Trees



- *What are the new max values for the subtrees?*

# Interval Trees



- *What are the new max values for the subtrees?*
- A: Unchanged
- *What are the new max values for and ?*

# Interval Trees



- *What are the new max values for the subtrees?*
- A: Unchanged
- *What are the new max values for and ?*
- A: root value unchanged, recompute other

# Interval Trees

- Following the methodology:
  - Pick underlying data structure
    - Red-black trees will store intervals, keyed on $i$  $lo$
  - Decide what additional information to store
    - Store the maximum endpoint in the subtree rooted at $i$
  - Figure out how to maintain the information
    - Insert: update max on way down, during rotations
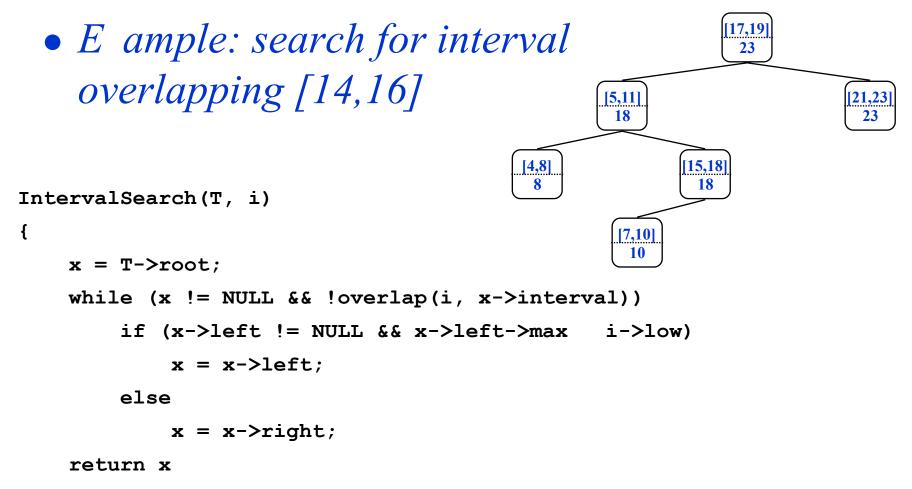    - Delete: similar
  - *Develop the desired ne   operations*

# Searching Interval Trees

```
IntervalSearch(T, i)
{
    x = T->root;
    while (x != NULL && !overlap(i, x->interval))
        if (x->left != NULL && x->left->max   i->low)
            x = x->left;
        else
            x = x->right;
    return x
}
```

- *What   ill be the running time?*

# IntervalSearch() Example

- *E ample: search for interval overlapping [14,16]*

```
IntervalSearch(T, i)

{

    x = T->root;

    while (x != NULL && !overlap(i, x->interval))

        if (x->left != NULL && x->left->max   i->low)

            x = x->left;

        else

            x = x->right;

    return x

}
```

Tree structure:

- [17,19] 23 (root)
  - [5,11] 18
    - [4,8] 8
    - [15,18] 18
      - [7,10] 10
  - [21,23] 23

# IntervalSearch() Example

- *E ample: search for interval overlapping [12,14]*

Tree nodes:
- [17,19] 23
- [5,11] 18
- [21,23] 23
- [4,8] 8
- [15,18] 18
- [7,10] 10

```
IntervalSearch(T, i)

{

    x = T->root;

    while (x != NULL && !overlap(i, x->interval))

        if (x->left != NULL && x->left->max   i->low)

            x = x->left;

        else

            x = x->right;

    return x

}
```
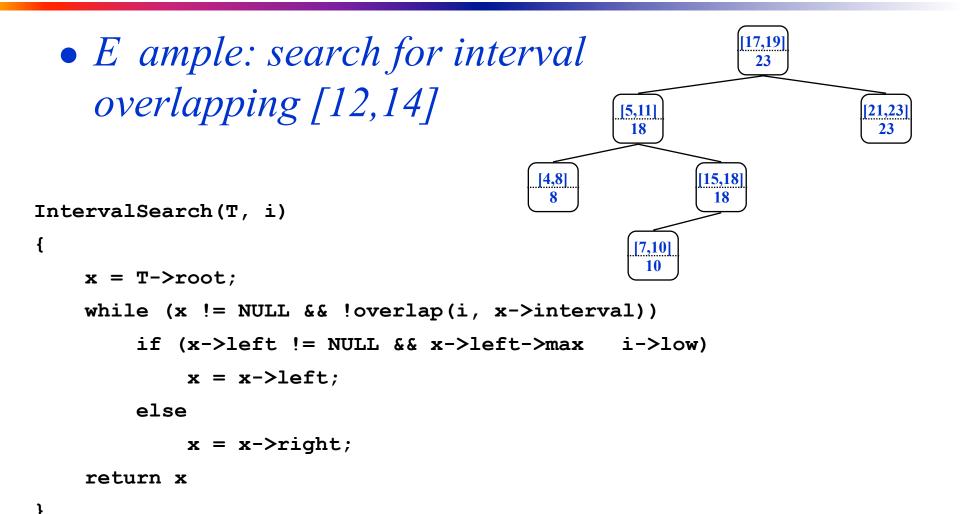
# Correctness of IntervalSearch()

- Key idea: need to check only 1 of node's 2 children
  - Case 1: search goes right
    - Show that ∃ overlap in right subtree, or no overlap at all
  - Case 2: search goes left
    - Show that ∃ overlap in left subtree, or no overlap at all

# Correctness of IntervalSearch()

- Case 1: if search goes right, $\exists$ overlap in the right subtree or no overlap in either subtree
  - If $\exists$ overlap in right subtree, we're done
  - Otherwise:
    - $x \to$ left $=$ NULL, or $x \to$ left $\to$ max $<$ i $\to$ low (*Why ?*)
    - Thus, no overlap in left subtree!

```
while (x != NULL && !overlap(i, x->interval))
        if (x->left != NULL && x->left->max   i->low)
            x = x->left;
        else
            x = x->right;
    return x;
```

# Correctness of IntervalSearch()

- Case 2: if search goes left, ∃ overlap in the left subtree or no overlap in either subtree
  - If ∃ overlap in left subtree, we'