# Lecture 6: Data Structures III

**COMP526: Efficient Algorithms**

Updated: October 22, 2024

Will Rosenbaum
University of Liverpool

AC : 787201

# Announcements

1. Third Quiz, due Friday
   - Similar format to before
   - Covers fundamental data structures (Lectures 4–6)
   - Quiz is **closed resource**
     - No books, notes, internet, etc.
     - Do not discuss until after submission deadline (Friday night, after midnight)
2. Programming Assignment (Draft) Posted *Today*
   - Due Wednesday, 13 November
3. Attendance Code:

787201

# Meeting Goals

- Finish up heaps
  - Give an efficient array-backed PRIORITYQUEUE
- Introduce two more ADTs:
  - ORDEREDSET
  - MAP
- Introduce binary search trees
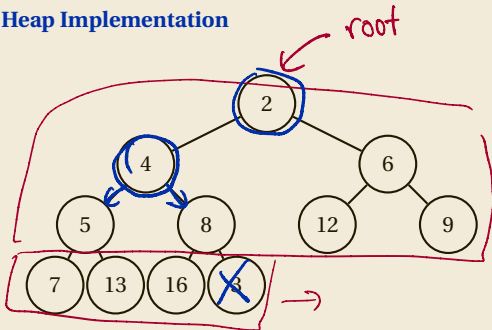- Discuss balanced binary search trees

# Heaps

# Last Time: Priority Queues and Heaps

**Priority Queues, Formally**

- $S$ is the state of the queue, initially $S = \varnothing$
- $S.\text{INSERT}(x, p(x)) : S = x_0 x_1 \cdots x_i x_{i+1} \cdots x_{n-1} \mapsto x_0 x_1 \cdots x_i\, x\, x_{i+1} \cdots x_{n-1}$
  - where $p(x_i) \le p(x) < p(x_{i+1})$
- $S.\text{MIN}() :$ returns $x_0$ where $S = x_0 x_1 \cdots x_{n-1}$
- $S.\text{REMOVEMIN}() : xS \mapsto S$, returns $x$
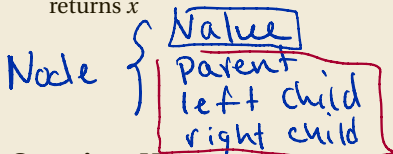
**Heap Implementation**



root

- INSERT via BUBBLEUP procedure
- REMOVEMIN via TRICKLEDOWN procedure
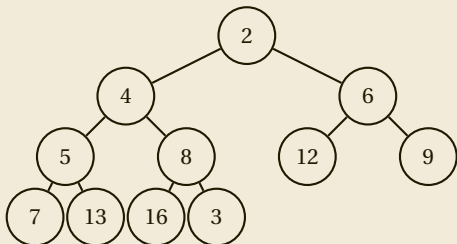
$O(\log n)$ steps

# Last Time: Priority Queues and Heaps

**Priority Queues, Formally**

- $S$ is the state of the queue, initially $S = \varnothing$
- $S.\text{INSERT}(x, p(x)) : S = x_0 x_1 \cdots x_i x_{i+1} \cdots x_{n-1} \mapsto x_0 x_1 \cdots x_i \, x \, x_{i+1} \cdots x_{n-1}$
  - where $p(x_i) \leq p(x) < p(x_{i+1})$
- $S.\text{MIN}() :$ returns $x_0$ where $S = x_0 x_1 \cdots x_{n-1}$
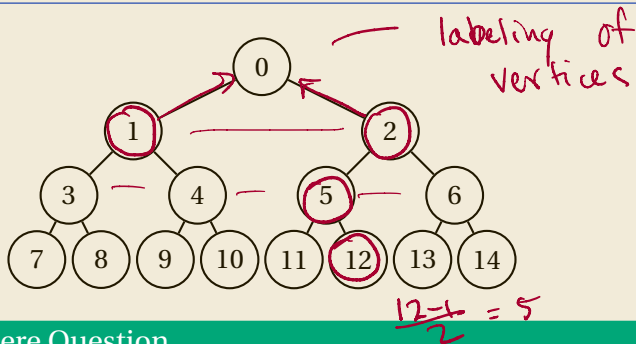- $S.\text{REMOVEMIN} : xS \mapsto S$, returns $x$

**Heap Implementation**



- INSERT via BUBBLEUP procedure
- REMOVEMIN via TRICKLEDOWN procedure
- **Issue:** using NODEs incurs overhead
  - locality of reference
  - storing additional references

Node $\left\{ \begin{array}{l} \text{Value} \\ \text{Parent} \\ \text{left child} \\ \text{right child} \end{array} \right.$

**Question.** How can we represent heaps as arrays?

# A Clue: Number the Vertices



labeling of vertices

0

1    2

3    4    5    6

7    8    9    10    11    12    13    14
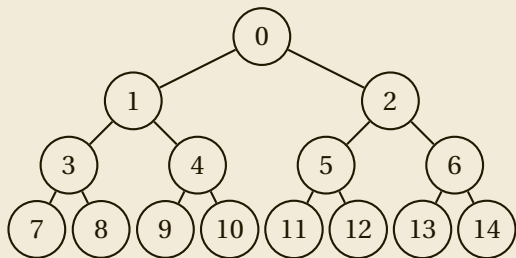
$\frac{12-1}{2} = 5$

## PollEverywhere Question

Suppose a vertex is assigned a label $i > 0$ in this numbering of the vertices. What is the label of $i$'s parent in the labeling?
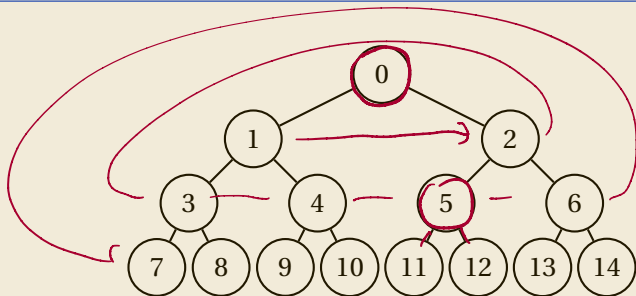
pollev.com/comp526

# A Clue: Number the Vertices



**Relationships:**

- If $i > 0$, then $i$'s parent has index $\lfloor (i-1)/2 \rfloor$ ← round down

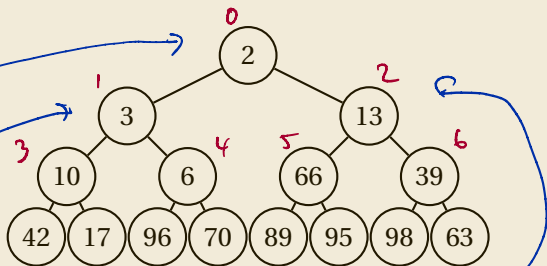# A Clue: Number the Vertices



**Relationships:**

- If $i > 0$, then $i$'s parent has index $\lfloor (i-1)/2 \rfloor$
- $i$'s left child has index $2i+1$
- $i$'s right child has index $2i+2$

# Arrays as Heaps

Associate numbering of tree vertices as array indexes!

## Complete binary tree representation

- If $i > 0$, then $i$'s parent has index $\lfloor (i-1)/2 \rfloor$
- $i$'s left child has index $2i + 1$
- $i$'s right child has index $2i + 2$



## Array representation

| 0 | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9 | 10 | 11 | 12 | 13 | 14 |
|---|---|---|---|---|---|---|---|---|---|----|----|----|----|----|
| 2 | 3 | 13 | 10 | 6 | 66 | 39 | 42 | 17 | 96 | 70 | 89 | 95 | 98 | 63 |

# Example: Array BUBBLEUP

We can apply heap procedures directly to the array without reference to the tree itself!

- If $i > 0$, then $i$'s parent has index $\lfloor (i-1)/2 \rfloor$
- $i$'s left child has index $2i+1$
- $i$'s right child has index $2i+2$

*(handwritten annotation: first index w/out element (array view))*

```
1:  procedure INSERT(p)
2:      v ← new vertex storing p
3:      u ← first vtx with < 2 children
4:      add v as u's child
5:      PARENT(v) ← u
6:      while value(v) < value(u) and u ≠ ⊥ do
7:          SWAP(value(v), value(u))
8:          v ← u
9:          u ← PARENT(v)
10:     end while
11: end procedure
```

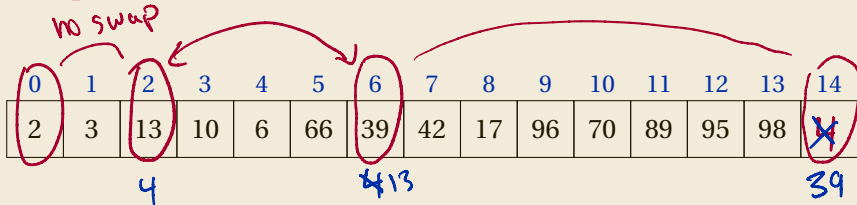*(handwritten annotation: || heap size)*

# Example: Array BUBBLEUP

We can apply heap procedures directly to the array without reference to the tree itself!

- If $i > 0$, then $i$'s parent has index $\lfloor (i-1)/2 \rfloor$
- $i$'s left child has index $2i + 1$
- $i$'s right child has index $2i + 2$

```
1: procedure INSERT(p)
2:     v ← new vertex storing p
3:     u ← first vtx with < 2 children
4:     add v as u's child
5:     PARENT(v) ← u
6:     while value(v) < value(u) and u ≠ ⊥ do
7:         SWAP(value(v), value(u))
8:         v ← u
9:         u ← PARENT(v)
10:    end while
11: end procedure
```

**Example.** INSERT(4)

no swap

| 0 | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9 | 10 | 11 | 12 | 13 | 14 |
|---|---|---|---|---|---|---|---|---|---|----|----|----|----|----|
| 2 | 3 | 13 | 10 | 6 | 66 | 39 | 42 | 17 | 96 | 70 | 89 | 95 | 98 | ✗ |

4          4 13                                    39

# Array Backed Operations

Using arrays, we can define INSERT and REMOVEMIN much more cleanly!

```
 1: procedure INSERT(p)
 2:     i ← n                    ▷ n is heap size
 3:     a[i] ← p
 4:     n ← n + 1
 5:     j ← ⌊(i − 1)/2⌋          ▷ j is i's parent
 6:     while i > 0 and a[i] < a[j] do
 7:         SWAP(a, i, j)
 8:         i ← j
 9:         j ← ⌊(i − 1)/2⌋
10:     end while
11: end procedure
```

```
 1: procedure REMOVEMIN
 2:     m ← a[0]
 3:     a[0] ← a[n − 1]
 4:     n ← n − 1
 5:     i ← 0
 6:     j ← argmin{a[2i + 1], a[2i + 2]}
 7:     while j < n and a[i] > a[j] do
 8:         SWAP(a, i, j)
 9:         i ← j
10:         j ← argmin{a[2i + 1], a[2i + 2]}
11:     end while
12:     return m
13: end procedure
```

Both of these operations still complete after $O(\log n)$ iterations

- very little overhead, since only array operations are used!

# Ordered Sets and Maps

# Adding Order to Elements

*priorities*

**Question.** What made our operations on heaps efficient?

- **Answer:** Order! We can order/compare priorities.

Two more ADT with **ordered** elements:

**Ordered Sets** store a collection (set) of *distinct* elements from an ordered universe.

- CONTAINS($x$) check if the set contains $x' = x$ and return $x'$
- ADD($x$) add $x$ to the set if $x$ was not present
- REMOVE($x$) remove $x$ if $x$ was present

# Adding Order to Elements

**Question.** What made our operations on heaps efficient?

- **Answer:** Order! We can order/compare priorities.

Two more ADT with **ordered** elements:

**Ordered Sets** store a collection (set) of *distinct* elements from an ordered universe.

- CONTAINS($x$) check if the set contains $x' = x$ and return $x'$
- ADD($x$) add $x$ to the set if $x$ was not present
- REMOVE($x$) remove $x$ if $x$ was present

$$a[k] \leftarrow v$$
$$a[k]$$

**Maps**[a] store a collection of *values* with associated ordered *keys* with array-like access.

- PUT($k, v$) set the value associated with key $k$ to $v$
- GET($k$) return the value associated with key $k$
- REMOVE($k$) remove the pair associated with $k$
- CONTAINS($k$) check if the map contains a value associated with $k$

---

[a]Aka: associative arrays, dictionaries (Python `dict`), symbol table

# Ordered Sets vs Maps

**Ordered Sets**

- CONTAINS($x$) check if the set contains $x' = x$ and return $x'$
- ADD($x$) add $x$ to the set if $x$ was not present
- REMOVE($x$) remove $x$ if $x$ was present

**Maps**

- PUT($k, v$) set the value associated with key $k$ to $v$
- GET($k$) return the value associated with key $k$
- REMOVE($k$) remove the pair associated with $k$
- CONTAINS($k$) check if the map contains a value associated with $k$

### PollEverywhere Question

If we are given an ORDEREDSET implementation, how could we use it to implement a MAP?

pollev.com/comp526

# Ordered Sets vs Maps

**Ordered Sets**

- CONTAINS($x$) check if the set contains $x' = x$ and return $x'$
- ADD($x$) add $x$ to the set if $x$ was not present
- REMOVE($x$) remove $x$ if $x$ was present

**Maps**

- PUT($k, v$) set the value associated with key $k$ to $v$
- GET($k$) return the value associated with key $k$
- REMOVE($k$) remove the pair associated with $k$
- CONTAINS($k$) check if the map contains a value associated with $k$

## Maps via Ordered Sets

- Create an ordered set that stores pairs $(k, v)$ (tuple)
- Compare $(k, v) \le (k', v') \iff k \le k'$
- CONTAINS, REMOVE are same

- To PUT($k, v$), use REMOVE($(k, \cdot)$) then ADD($(k, v)$)
- To GET($k$), use $(k, v) \leftarrow$ CONTAINS($(k, \cdot)$) and return $v$
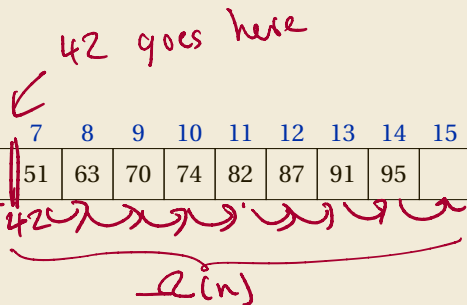
# Ordered Sets via Arrays
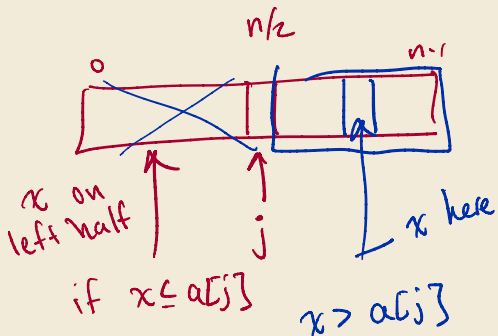
ORDEREDSETS can be implemented by arrays:

- Maintain a sorted array $a = [x_0, x_1, \ldots, x_n]$ with each $x_i \le x_{i+1}$.
- ADD($x$) and REMOVE($x$) implemented in $\Theta(n)$ worst case time
  - To ADD find index $i$ such that $x_i \le x < x_{i+1}$
  - Shift elements $x_j$ with $j \ge i+1$ to next index
    - This uses $\Theta(n)$ time
  - Set $a[i+1] \leftarrow x$

**Example.** How to ADD(42)?

42 goes here

| 0 | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9 | 10 | 11 | 12 | 13 | 14 | 15 |
|---|---|---|---|---|---|---|---|---|---|----|----|----|----|----|----|
| 2 | 3 | 10 | 28 | 31 | 34 | 39 | 51 | 63 | 70 | 74 | 82 | 87 | 91 | 95 | |

42

$a(n)$

# Ordered Sets via Arrays

ORDEREDSETs can be implemented by arrays:

- Maintain a sorted array $a = [x_0, x_1, \ldots, x_n]$ with each $x_i \leq x_{i+1}$.
- ADD($x$) and REMOVE($x$) implemented in $\Theta(n)$ worst case time
  - To ADD find index $i$ such that $x_i \leq x < x_{i+1}$
  - Shift elements $x_j$ with $j \geq i+1$ to next index
    - This uses $\Theta(n)$ time
  - Set $a[i+1] \leftarrow x$

**Question.** How can we implement CONTAINS($x$) more quickly?

# Efficient Search

**Idea.** Binary Search:

- Start at the *middle index j*
  - $x \le a[j] \implies$ index of $x$ must be $i \le j$
  - otherwise $i > j$
- Apply procedure to remaining interval with half excluded
  - compare $x$ to midpoint of remaining interval
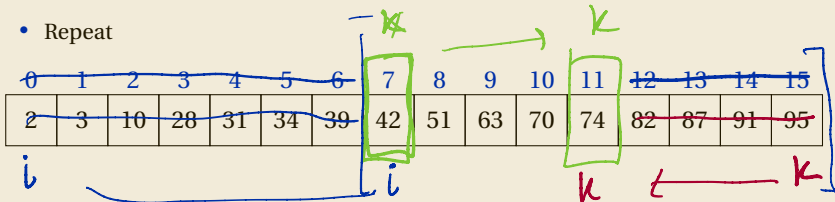  - eliminate half of the interval
- Repeat



n/2

n-1

0

x on left half

j

x here

if $x \le a[j]$

$x > a[j]$

# Efficient Search

**Idea.** Binary Search:

- Start at the *middle index j*

  - $x \le a[j] \implies$ index of $x$ must be $i \le j$
  - otherwise $i > j$

- Apply procedure to remaining interval with half excluded

  - compare $x$ to midpoint of remaining interval
  - eliminate half of the interval

- Repeat

*left endpt of active interval*

1: **procedure** BINARYSEARCH(x)
2:    $i \leftarrow 0, k \leftarrow n-1$    ← *right endpt.*
3:    $j \leftarrow \lfloor (i+k)/2 \rfloor$
4:    **while** $i < j$ **do**
5:       **if** $x \le a[j]$ **then**
6:          $k \leftarrow j$
7:       **else**
8:          $i \leftarrow j$
9:       **end if**
10:   **end while**
11:   **return** $i$
12: **end procedure**

*find(72)*

| 0 | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9 | 10 | 11 | 12 | 13 | 14 | 15 |
|---|---|---|---|---|---|---|---|---|---|----|----|----|----|----|----|
| 2 | 3 | 10 | 28 | 31 | 34 | 39 | 42 | 51 | 63 | 70 | 74 | 82 | 87 | 91 | 95 |

# Efficiency of Binary Search

## PollEverywhere

What is the (worst case) running time of BINARYSEARCH on an array of length $n$?



pollev.com/comp526

```
1: procedure BINARYSEARCH(x)
2:     i ← 0, k ← n − 1
3:     j ← ⌊(i + k)/2⌋
4:     while i < j do
5:         if x ≤ a[j] then
6:             k ← j
7:         else
8:             i ← j
9:         end if
10:    end while
11:    return i
12: end procedure
```

# Efficiency of Binary Search

### Proposition

The worst-case running time of BINARYSEARCH is $\Theta(\log n)$.

```
1: procedure BINARYSEARCH(x)
2:     i ← 0, k ← n − 1
3:     j ← ⌊(i + k)/2⌋
4:     while i < j do
5:         if x ≤ a[j] then
6:             k ← j
7:         else
8:             i ← j
9:         end if
10:     end while
11:     return i
12: end procedure
```

# Efficiency of Binary Search

_size of active interval_

### Proposition

The worst-case running time of
BINARYSEARCH is $\Theta(\log n)$.

### Proof.

- Consider the value of $k - i$.

- After $\ell$ iterations of the loop,
  have $k - i \leq \frac{n}{2^\ell}$ (induction)

- Termination when $k - i \leq 1$

- $\ell = \lceil \log n \rceil + 1 \implies \frac{n}{2^\ell} \leq 1$
  □

_↑_
_log n_
_rounded up_

_process terminates._

```
1: procedure BINARYSEARCH(x)
2:     i ← 0, k ← n − 1
3:     j ← ⌊(i + k)/2⌋
4:     while i < j do
5:         if x ≤ a[j] then
6:             k ← j
7:         else
8:             i ← j
9:         end if
10:    end while
11:    return i
12: end procedure
```

# Making All Operations Efficient?

## A Nagging Question

For ORDEREDSETs, we can perform all operations in $o(n)$ time?

- Array implementation only gives CONTAINS in $O(\log n)$ time
- Other operations are $\Theta(n)$
- This seems harder than efficient PRIORITYQUEUE as elements can be added *and* removed from anywhere in the data structure

# Making All Operations Efficient?

## A Nagging Question

For ORDEREDSETS, we can perform all operations in $o(n)$ time?

- Array implementation only gives CONTAINS in $O(\log n)$ time
- Other operations are $\Theta(n)$
- This seems harder than efficient PRIORITYQUEUE as elements can be added *and* removed from anywhere in the data structure

**Up next:** A solution in two parts

1. Binary Search Trees
2. Balancing Binary Trees

# Binary Search Trees

# Binary Search Tree Definition

## Definition

Suppose $T$ is a binary tree and every vertex $v$ in $T$ has an associated value. We say $T$ is a **binary search tree** (**BST**) if for every vertex (value) $v$:

*(handwritten: $\leq 2$ children)*

1. every *left descendant* $u$ satisfies $u \leq v$,
2. every *right descendant* $w$ satisfies $w \geq v$.

*(handwritten annotations: "≤ 15" over left subtree, "≥ 15" over right subtree)*

# BST Search

Given a BST $T$, how can we search for a value $x$ in $T$?

CONTAINS(19)?

Start at root



- 19 must be to right of 18 ⎫ ⟹ 19
- 18 doesn't have right child ⎭ not in BST

# BST Search

Given a BST $T$, how can we search for a value $x$ in $T$?

$\perp$ "perp" to indicate non-existent node

1: **procedure** CONTAINS($x$)
2:      $v$ = tree root
3:      **while** $v \neq x$ and $v \neq \perp$ **do**
4:          **if** $x < v$ **then**
5:              $v \leftarrow$ LEFTCHILD($v$)    go left
6:          **else**
7:              $v \leftarrow$ RIGHTCHILD($v$)   go right
8:          **end if**
9:      **end while**
10:     **return** $v$
11: **end procedure**

# BST Search

Given a BST $T$, how can we search for a value $x$ in $T$?

```
 1: procedure CONTAINS(x)
 2:     v = tree root
 3:     while v ≠ x and v ≠ ⊥ do
 4:         if x < v then
 5:             v ← LEFTCHILD(v)
 6:         else
 7:             v ← RIGHTCHILD(v)
 8:         end if
 9:     end while
10:     return v
11: end procedure
```

## PollEverywhere

What is the (worst case) running time of CONTAINS on a tree with $n$ vertices?

pollev.com/comp526

# **BST** CONTAINS **Efficiency**

## Observation

The (worst-case) running time of
CONTAINS on $T$ is $\Theta(h)$ where $h$ is
the **height** of $T$

- $h$ is the length of the longest
  path from root to any leaf in $T$

The height of $T$ can be:

- As small as $\log n$
- As large as $n - 1$

## The Moral

The efficiency of CONTAINS
depends on the structure of $T$.



balanced tree

val looking for!
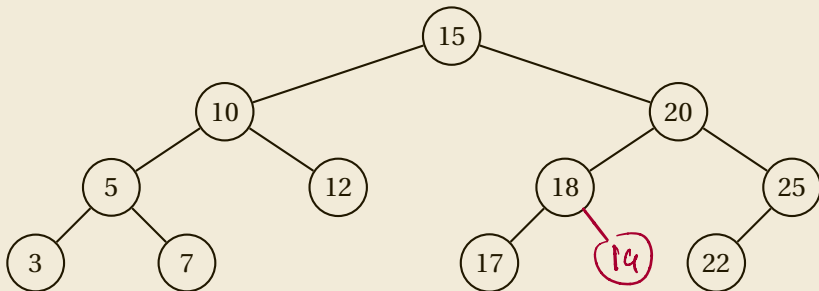
# BST Add

## Question

How could we ADD(19) to the following BST so it remains a BST?



start by search for 19

maintains BST properties

# BST Add

**Question**

How could we ADD(19) to the following BST so it remains a BST?



**Observation.** To ADD($x$), we should add a new vertex wherever the CONTAINS($x$) execution fails to find $x$.

# Adding in Pseudocode

```
1: procedure ADD(x)
2:     v, u ← root
3:     while v ≠ ⊥ do
4:         if x = v then
5:             return
6:         else if x < v then
7:             u ← v
8:             v ← LEFTCHILD(v)
9:         else
10:            u ← v
11:            v ← RIGHTCHILD(v)
12:        end if
13:    end while
14:    if x < u then
15:        set x as u's left child
16:    else
17:        set x as u's right child
18:    end if
19: end procedure
```

**Example.** ADD(8)

# Adding in Pseudocode

```
1: procedure ADD(x)
2:     v, u ← root
3:     while v ≠ ⊥ do
4:         if x = v then
5:             return
6:         else if x < v then
7:             u ← v
8:             v ← LEFTCHILD(v)
9:         else
10:            u ← v
11:            v ← RIGHTCHILD(v)
12:        end if
13:    end while
14:    if x < v then
15:        set x as v's left child
16:    else
17:        set x as v's right child
18:    end if
19: end procedure
```

### PollEverywhere Question

Describe a sequence of ADD(x) operations starting from an empty BST such that every operation takes $\Omega(n)$ time.
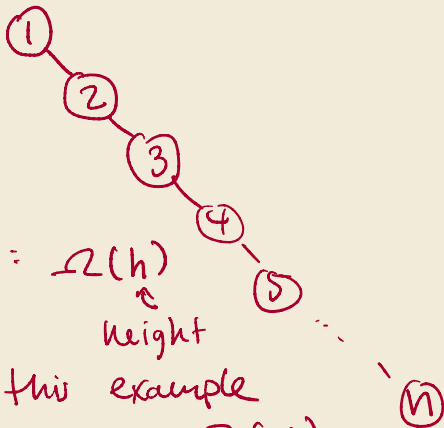
pollev.com/comp526

# Adding in Pseudocode

```
1: procedure ADD(x)
2:     v, u ← root
3:     while v ≠ ⊥ do
4:         if x = v then
5:             return
6:         else if x < v then
7:             u ← v
8:             v ← LEFTCHILD(v)
9:         else
10:             u ← v
11:             v ← RIGHTCHILD(v)
12:         end if
13:     end while
14:     if x < v then
15:         set x as v's left child
16:     else
17:         set x as v's right child
18:     end if
19: end procedure
```
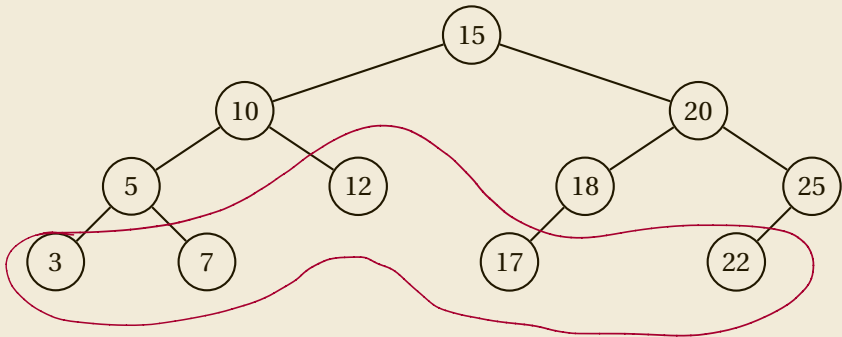
**A Bad Sequence:** 1 2 3 ... n



ops: $\Omega(h)$

$\uparrow$ height

For this example
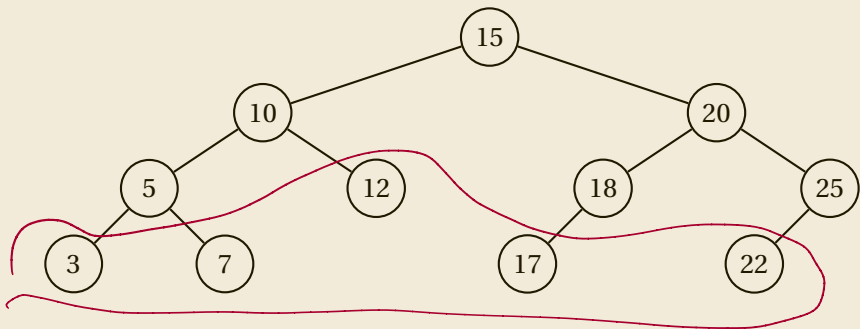
h = n-1 ⟹ $\Omega(n)$ ops.

# BST Remove

How could we remove an element from a BST?

# BST Remove
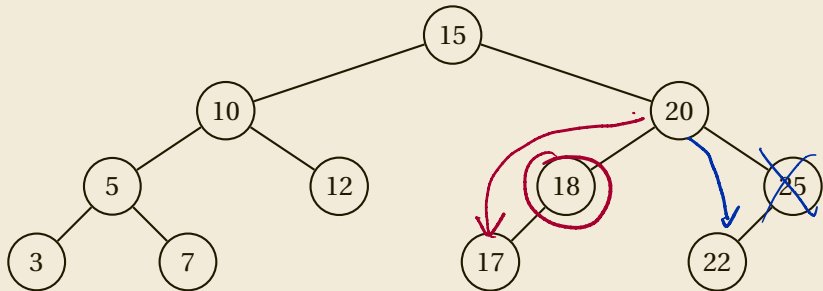
## Question

How could we remove an element from a BST?



**Case 1: A leaf**. Just remove it!

# BST Remove

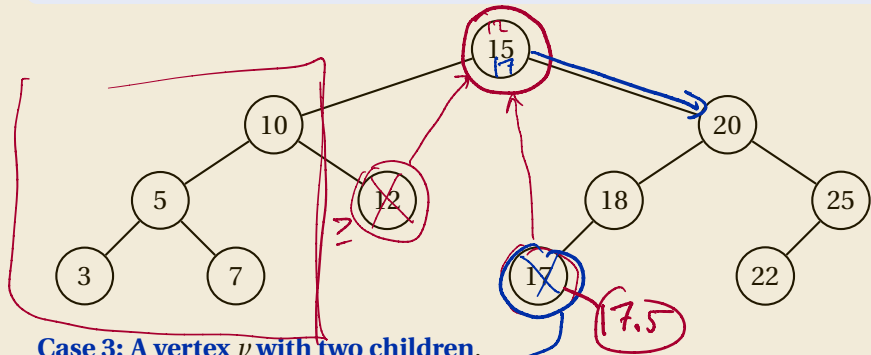**Question**

How could we remove an element from a BST?



**Case 2: A vertex *v* with single child**. Splice! Set *v*'s child to be its parent's child.

# BST Remove

How could we remove an element from a BST?



**Case 3: A vertex $v$ with two children**.

1. Find *next smallest* value $w$.
2. Copy $w$'s value to $v$.
3. Remove $w$

must have 0 or 1 children, so remove accordingly

# So Far…

…we've implemented

- CONTAINS($x$)
- ADD($x$)
- REMOVE($x$)

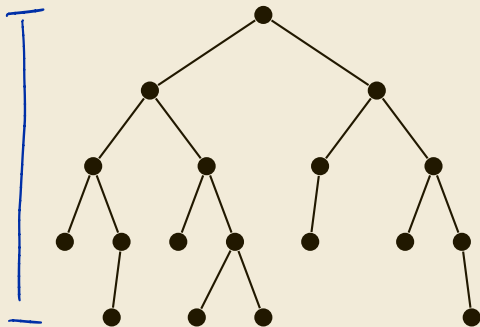for ORDEREDSETs.

**But** we haven't improved *efficiency*

- All of these operations can cost as much as $\Theta(n)$
  - efficiency depends on previous operations performed!

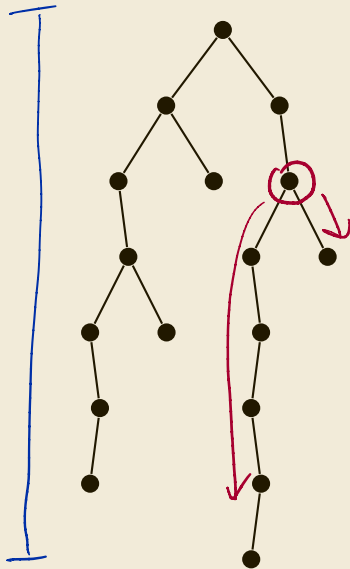**Idea.** We can *restructure* BSTs.

- Goal: ensure that the BST has small **height**.
- After each update, check and update tree structure.
  - maintain BST property
  - updates performed efficiently

# Balanced Binary Trees

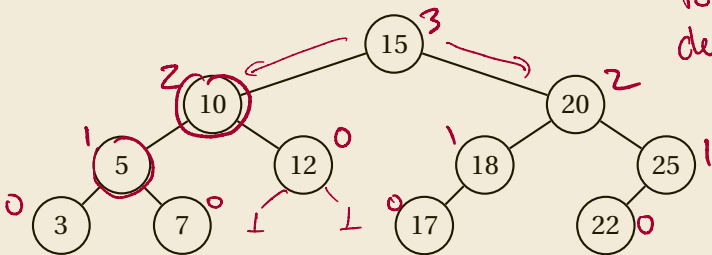# Distinguishing the Good from the Bad



Good

# Height Balanced Trees

## Definition (Left and Right Height)

Let $v$ be a vertex in a tree. We define:

- $h(\bot) = -1$
- $h(v) = 1 + \max(h(\text{LeftChild}(v)), h(\text{RightChild}(v)))$
- $h_\ell(v) = h(\text{LeftChild}(v))$
- $h_r(v) = h(\text{RightChild}(v))$

= dist to farthest descendent leaf
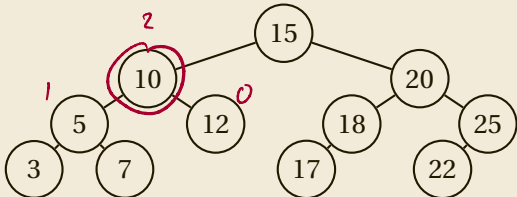
# Height Balanced Trees

## Definition (Left and Right Height)

Let $v$ be a vertex in a tree. We define:

- $h(\bot) = -1$
- $h(v) = 1 + \max(h(\text{LEFTCHILD}(v)), h(\text{RIGHTCHILD}(v)))$
- $h_\ell(v) = h(\text{LEFTCHILD}(v))$
- $h_r(v) = h(\text{RIGHTCHILD}(v))$

## Def. (Height Balanced)

We call a tree **height balanced** if for every vertex $v$, $|h_\ell(v) - h_r(v)| \leq 1$.

# Properties of Height Balanced Trees

## Proposition

Suppose $T$ is a height balanced tree of height $h$. Then $T$ has $n \geq 2^{h/2}$ vertices.

# Properties of Height Balanced Trees

## Proposition

Suppose $T$ is a height balanced tree of height $h$. Then $T$ has $n \geq 2^{h/2}$ vertices.

## Proof.

Let $M(h)$ denote the minimum size of a height balanced tree of height $h$.

- Observe that $M(0) = 1$, $M(1) = 2$.
- In general $M(h) \geq 1 + M(h-1) + M(h-2)$
  - one subtree of the root is a height balanced tree of height $h-1$
  - other subtree is height balanced with height at least $h-2$
- So $M(h) \geq 2M(h-2)$
- Inductive argument $\implies M(h) \geq 2^{h/2}$.

$\square$

# Properties of Height Balanced Trees

## Proposition

Suppose $T$ is a height balanced tree of height $h$. Then $T$ has $n \geq 2^{h/2}$ vertices.

## Consequences.

If $T$ is a height balanced tree with $n$ vertices, then its height $h$ satisfies $h \leq 2 \log n$

$\implies$ CONTAINS($x$) takes time $O(\log n)$

$\implies$ ADD($x$) takes time $O(\log n)$

$\implies$ REMOVE($x$) takes time $O(\log n)$

# Maintaining Height Balance

**Our Strategy.** Maintain a BST that is height balanced **for any sequence of operations performed**.

- No one is *forcing* us to keep the tree structure determined by our ADD/REMOVE operations
    - there are many valid BSTs that store the same collection of elements!

# Maintaining Height Balance

**Our Strategy.** Maintain a BST that is height balanced **for any sequence of operations performed**.

- No one is *forcing* us to keep the tree structure determined by our ADD/REMOVE operations
  - there are many valid BSTs that store the same collection of elements!
- Starting from a balanced tree, ADD($x$) may introduce imbalance.
- If imbalance is introduced try to fix it:
  - find closest unbalanced vertex to $x$ and correct its balance
  - look for other imbalance and correct it

# Maintaining Height Balance

**Our Strategy.** Maintain a BST that is height balanced **for any sequence of operations performed**.

- No one is *forcing* us to keep the tree structure determined by our ADD/REMOVE operations
    - there are many valid BSTs that store the same collection of elements!
- Starting from a balanced tree, ADD($x$) may introduce imbalance.
- If imbalance is introduced try to fix it:
    - find closest unbalanced vertex to $x$ and correct its balance
    - look for other imbalance and correct it

**For next time.** Think about how you could implement this strategy.

- *Where* could imbalance occur? And how much?
- What *local* operations can fix the imbalance?
- What is the worst-case running time of restoring balance?

# Next Time: Sorting

- Finishing Balanced BSTs
- The Sorting Task
- Efficient Sorting by Divide and Conquer

# Scratch Notes