Carleton University
Alina Shaikhet
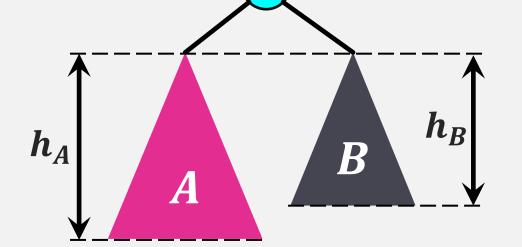
**COMP 2402**
# Scapegoat Trees

# Balanced Search Trees

A node in a tree is **height-balanced** if the heights of its subtrees differ by no more than **1**

$$\text{If } |h_A - h_B| \leq 1 \text{ then } u \text{ is } \textbf{height-balanced.}$$
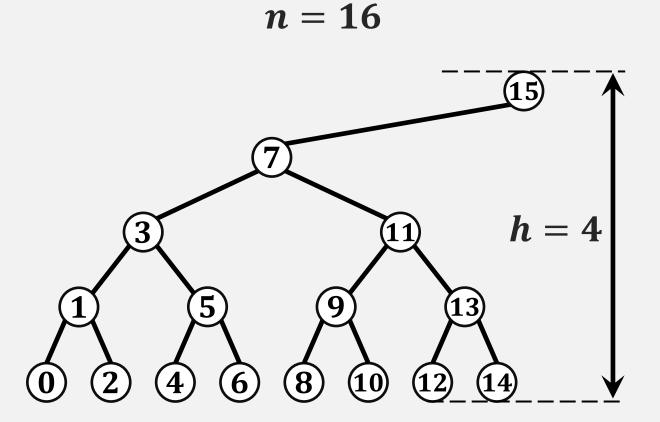
A tree is **height-balanced** iff every node in it is **height-balanced.**

When we say that a tree is **balanced,** we usually mean that its height is $c \log n$, where $c$ is some constant and $n$ is the number of elements in the tree.

# Balanced Search Trees

Is it possible to make the height of this tree smaller?

$$n = 16$$



$$h = 4$$

A **perfect** binary tree is a binary tree in which all interior nodes have two children, and all leaves have the same depth (or same level).

# Scapegoat Trees

A **ScapegoatTree** keeps itself balanced by performing partial rebuilding operations (where an entire subtree is deconstructed and rebuilt into a perfectly balanced subtree).

A **ScapegoatTree** is a **BinarySearchTree** that, in addition to keeping track of the number $n$ of nodes in the tree also keeps a counter $q$ that maintains an **upper-bound on the number of nodes**.

$$n \leq q \leq 2n$$

$$q/2 \leq n \leq q$$

A **ScapegoatTree** can look surprisingly unbalanced, however, it always maintains logarithmic height.

The height of the **ScapegoatTree** does not exceed

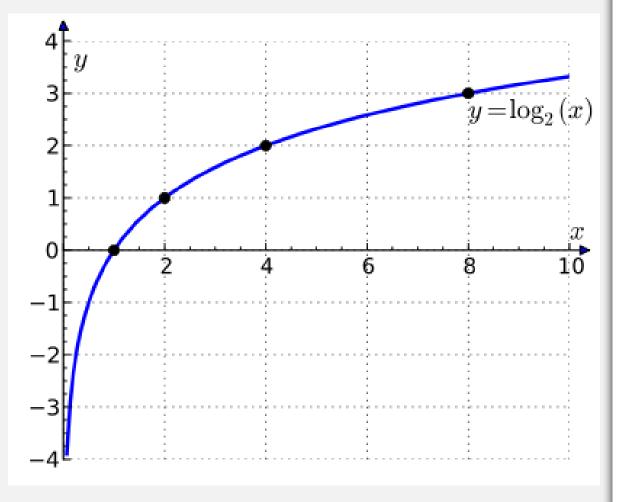$$\log_{3/2} q \leq \log_{3/2} 2n < \log_{3/2} n + 2$$

# Logarithms

$$y = \log_b n$$

$\log_b n$ is the number of times we divide $n$ by $b$ before it becomes $\leq 1$.

$$b^y = b^{\log_b n} = n$$

$$2^8 = 256$$
$$8 = \log_2 256$$

$$\log_b xy = \log_b x + \log_b y$$



$y = \log_2(x)$

$$\log_b n = \frac{\log_a n}{\log_a b} = \frac{\log_a n}{c} = \frac{1}{c} \log_a n$$

That is why the logarithm base is irrelevant when using Big $O$ notation.

5

# SSet Implementation

- add($x$),
- remove($x$),

$$O(\log n)$$
**ignoring rebuild()**

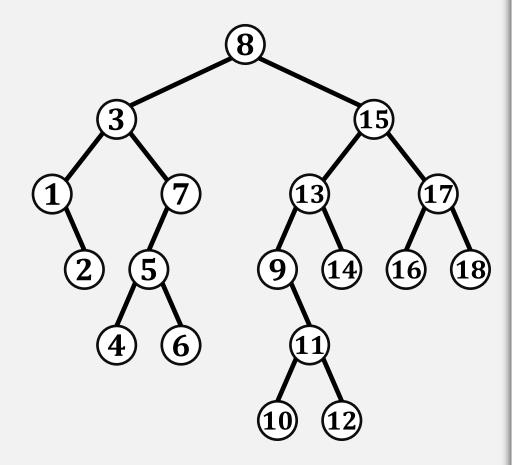- find($x$) – find the **smallest** value that is $\geq x$. ← $$O(\log n)$$

find($x$) is implemented using the standard algorithm for searching in a **BinarySearchTree.** This takes time proportional to the height of the tree which is

$$< \log_{3/2} n + 2$$

# SSet: remove($x$)

Use the usual algorithm for deleting the value $x$ stored in **BinarySearchTree:**

- Find node $u$ that contains value $x$.
- If $u$ is a leaf, then detach $u$ from its parent.
- If $u$ has only one child, then splice $u$ from the tree by having $u$.parent adopt $u$'s child.
- If $u$ has two children, then find a node $w$, that has less than two children, such that $w.x$ can replace $u.x$. Then remove $w$. Option **1**: Find the largest element in the left subtree of $u$, Option **2**: Find the smallest element in the right subtree of $u$.
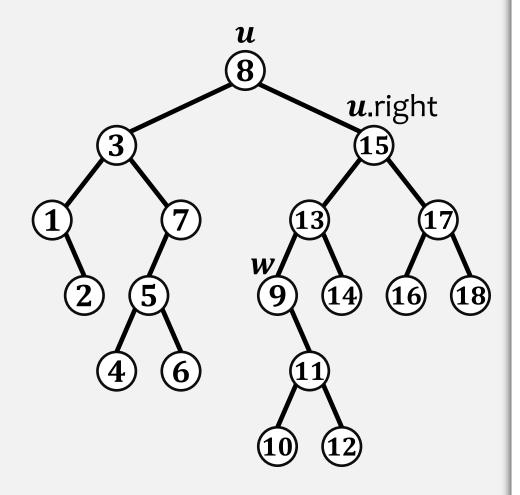
# SSet: remove($x$)

Let's choose option **2**.

Start at $u$.right and keep going to the left child until there is no left child. The node you stopped at is $w$.

- $w.x$ is the smallest value in the subtree rooted at $u$.right.
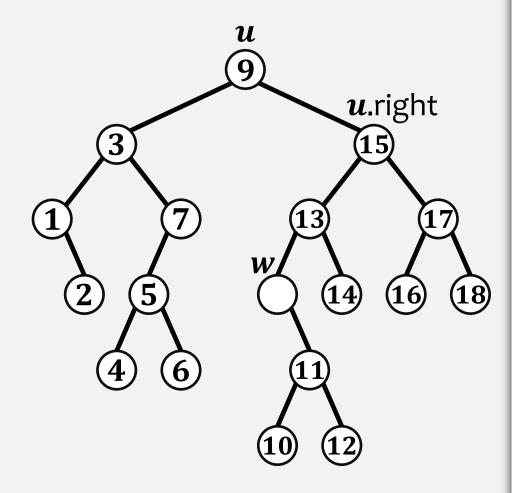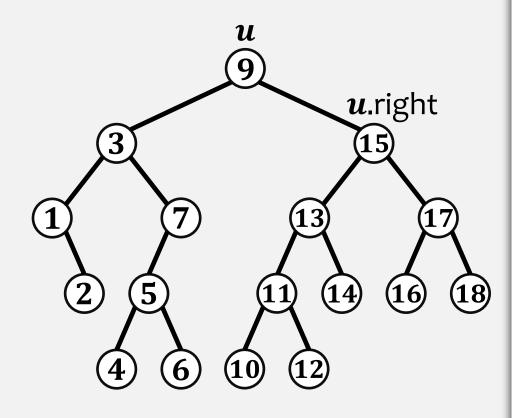
- $w$ has no left child.

remove(**8**)

# SSet: remove($x$)

remove($\mathbf{8}$)

Let's choose option $\mathbf{2}$.

Start at $\boldsymbol{u}$.right and keep going to the left child until there is no left child. The node you stopped at is $\boldsymbol{w}$.

- $\boldsymbol{w}.\boldsymbol{x}$ is the smallest value in the subtree rooted at $\boldsymbol{u}$.right.

- $\boldsymbol{w}$ has no left child.

delete node $\boldsymbol{w}$

9

# SSet: remove($x$)

remove($\mathbf{8}$)

Let's choose option $\mathbf{2}$.

Start at $\boldsymbol{u}$.right and keep going to the left child until there is no left child. The node you stopped at is $\boldsymbol{w}$.

- $\boldsymbol{w}.\boldsymbol{x}$ is the smallest value in the subtree rooted at $\boldsymbol{u}$.right.

- $\boldsymbol{w}$ has no left child.

Note that remove($\boldsymbol{x}$) can never increase the height of the tree.

$\boldsymbol{u}$

$\boldsymbol{u}$.right

delete node $\boldsymbol{w}$

10

# SSet: remove($x$)

- remove $x$ using standard BST algorithm for removal.
- $n--$
- if $q > 2n$ then <u>rebuild(root)</u> and set $q = n$

$$2^0 + 2^1 + \cdots + 2^k = 2^{k+1} - 1$$

**this function rebuilds the tree into a perfectly balanced tree**

For Complete Binary Tree
$$h = \lfloor \log_2 n \rfloor$$

$h = \log_2(n+1) - 1$

There are $2^{h+1} - 1$ nodes in a **full** binary tree of height $h$:

$$2^{h+1} - 1 = n$$
$$2^{h+1} = n + 1$$
$$h + 1 = \log_2 2^{h+1} = \log_2(n+1)$$

11

# SSet:  remove($x$)

- remove  $x$  using standard BST algorithm for removal.
- $n--$
- if  $q > 2n$  then  <u>rebuild(root)</u>  and set  $q = n$

$O(n)$

**this function rebuilds the tree into a perfectly balanced tree**

After rebuild() for every node $u$ :

$$|\text{size}(u.\,\text{left}) - \text{size}(u.\,\text{right})| \leq 1$$

If we ignore the cost of rebuilding, the running time of the  remove($x$)  operation is proportional to the height of the tree and is therefore  $O(\log n)$.

Alina Shaikhet – COMP 2402 – Carleton University

# SSet: add($x$)

- Use the standard algorithm for adding $x$ to **BinarySearchTree**:

- $n++$

- $q++$

- If $(\text{depth}(u) > \log_{3/2} q)$ then walk from $u$
  back up to the root looking for a scapegoat $w$:

$$\frac{\text{size}(w.\text{child})}{\text{size}(w)} > \frac{2}{3}$$

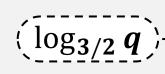$w.$ child is a child
of $w$ on the path
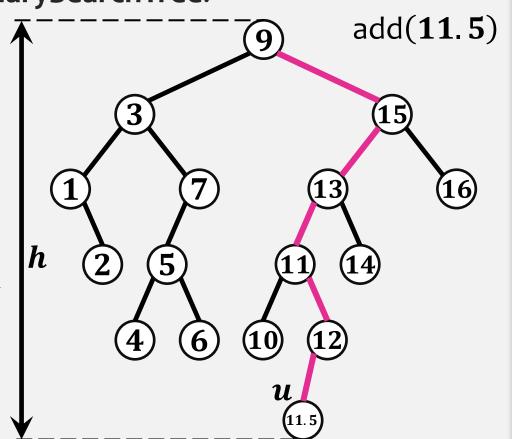from the root to $u$

Then rebuild($w$)  $\textcolor{magenta}{O(\text{size}(w))}$

If we ignore the cost of finding
the scapegoat and rebuilding:  $O(\log n)$

add($x$) can increase
the height of a tree

$\log_{3/2} q$

add($\mathbf{11.5}$)

# SSet:  add($x$) – scapegoat $w$

$$\frac{\text{size}(A + \bigcirc)}{\text{size}(w)} > \frac{2}{3}$$

$$\text{size}(A + \bigcirc) > \frac{2}{3}\text{size}(w)$$

$\log_{3/2} q$

$w$

$A$

$B$

rebuild($w$)

$w$

14

# Does the scapegoat exist?

Suppose, that there does not exist a scapegoat $w$ on the path from $u$ (with value $x$) to the root.

Then for all the nodes $w$ on the path

$$\frac{\text{size}(w)}{\text{size}(\text{parent of } w)} \le \frac{2}{3}$$

Thus, the length of the path from the root to $u$ is at most $\log_{3/2} n \le \log_{3/2} q$

So, we didn't cross the $\boxed{\log_{3/2} q}$ line.

Contradiction!

$n$

$\le \frac{2}{3}n$

$\le \left(\frac{2}{3}\right)^2 n$

$\le \left(\frac{2}{3}\right)^3 n$

**# nodes on the path**

$\le \log_{3/2} n$

$x$

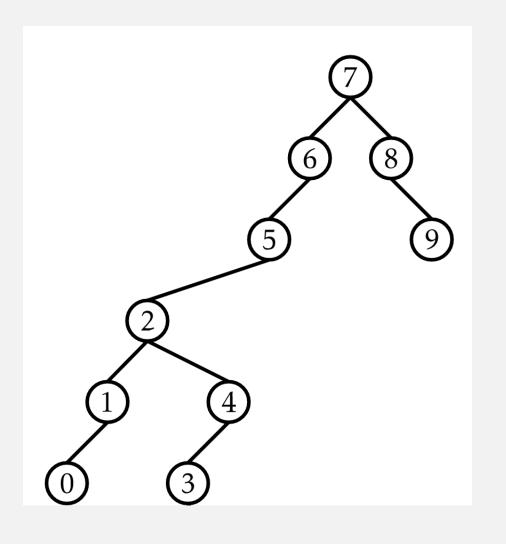Alina Shaikhet – COMP 2402 – Carleton University

15

# Scapegoat Trees

A **ScapegoatTree** can look surprisingly unbalanced, however, it always maintains logarithmic height:

$$h \leq \log_{3/2} q$$

For the tree in the example:

$$n = 10 = q$$

$$\text{height} = 5$$

$$5 < \log_{3/2} 10 \approx 5.679$$



from the ods textbook

# Scapegoat Trees

A **ScapegoatTree** can look surprisingly unbalanced, however, it always maintains logarithmic height:
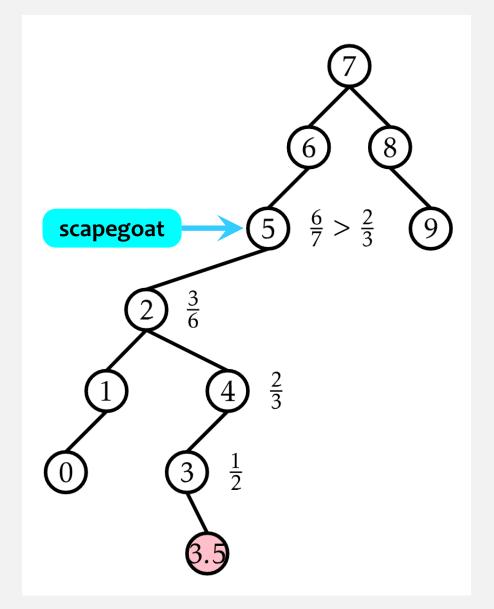
$$h \leq \log_{3/2} q$$

For the tree in the example:

$$n = \cancel{10}\;^{11} = q$$

height = $\cancel{5}$ 6

$$5 < \log_{3/2} 10 \approx 5.679$$

$$6 > \log_{3/2} 11 \approx 5.914$$



from the ods textbook

# SSet: add($x$)

```java
boolean add(T x) {
    // first do basic insertion keeping track of depth
    Node<T> u = newNode(x);
    int d = addWithDepth(u);
    if (d > log32(q)) {
        // depth exceeded, find scapegoat
        Node<T> w = u.parent;
        while (3*size(w) <= 2*size(w.parent))
            w = w.parent;
        rebuild(w.parent);
    }
    return d >= 0;
}
```

# rebuild($u$)

- traverse (in-order) $u$'s subtree and collect all its nodes into an array $a$
- recursively build a balanced subtree using $a$:

$$m = \frac{a.\,\mathrm{length}}{2}$$

$a[m]$ becomes the root of the new subtree,
$a[0], \dots, a[m-1]$ get stored recursively in the **left** subtree, and
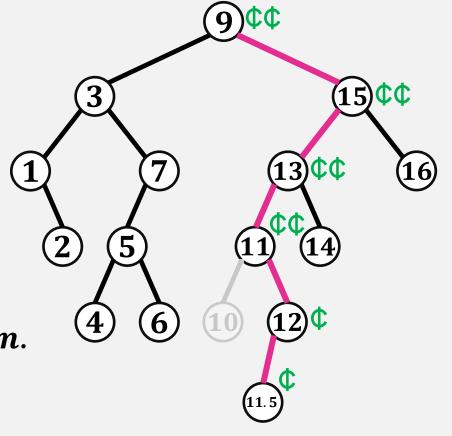$a[m+1], \dots, a[a.\,\mathrm{length}-1]$ get stored recursively in the right subtree.

# Analysis

Credit scheme:

- each node stores a number of credits

- during an **insertion** or **deletion** of $u$, we give one credit to each node on the path to $u$.

- during a deletion we also store an additional credit



Invariant:

1. The number of $q$-credits saved up is at least $q - n$.

2. The number of credits on any node $u$ is at least
$$|\text{size}(u.\text{left}) - \text{size}(u.\text{right})| - 1$$

$$n \leq q \leq 2n$$

# Analysis

Invariant:
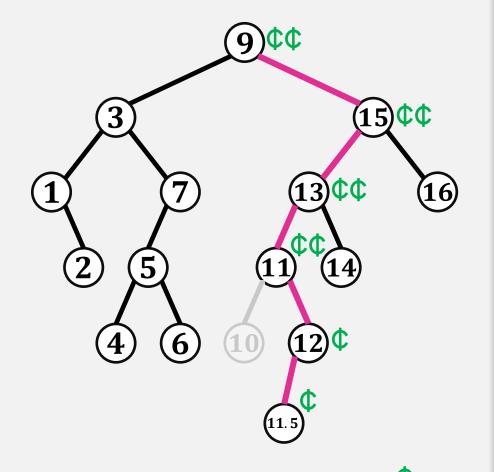1.  The number of $q$-credits saved up is at least $q - n$.
2.  The number of credits on any node $u$ is at least
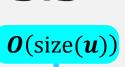$$|\text{size}(u.\,\text{left}) - \text{size}(u.\,\text{right})| - 1$$

If we call rebuild($u$) during a **deletion**, it is because we decremented $n$ and now $q > 2n$.
But we have $q - n$ $q$-credits stored.

$$q - n > 2n - n = n$$

We use these to pay for the $O(n)$ time it takes to rebuild the root.



$$n \leq q \leq 2n$$

# Analysis

$O(\text{size}(u))$

**Invariant:**
1. The number of $q$-credits saved up is at least $q - n$.
2. The number of credits on any node $u$ is at least

$$|\text{size}(u.\text{left}) - \text{size}(u.\text{right})| - 1$$

If we call rebuild$(u)$ during an **insertion**, it is because $u$ is a scapegoat. Assume that path goes to the left:

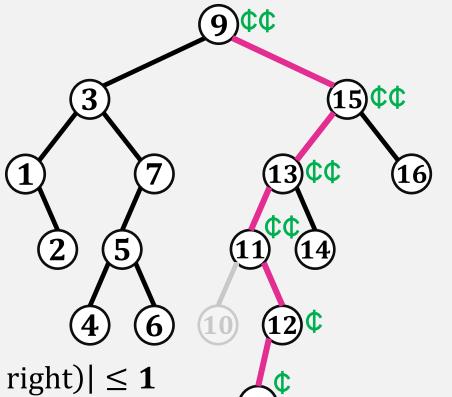$$\frac{\text{size}(u.\text{left})}{\text{size}(u)} > \frac{2}{3} \qquad \text{size}(u.\text{left}) > 2\,\text{size}(u.\text{right})$$

$$\text{size}(u) = 1 + \text{size}(u.\text{left}) + \text{size}(u.\text{right})$$

$$\text{size}(u.\text{left}) - \text{size}(u.\text{right}) > \frac{1}{3}\,\text{size}(u)$$

The last time a subtree containing $u$ was rebuilt (or when $u$ was inserted), we had $|\text{size}(u.\text{left}) - \text{size}(u.\text{right})| \leq 1$

The number of add/remove operations that have affected $u$'s subtrees since then is at least $\dfrac{1}{3}\text{size}(u) - 1 \leq \#$ credits stored at $u$

$$n \leq q \leq 2n$$

23

# Analysis

- During an insertion or deletion, we give one credit to each node on the path to the inserted (or deleted) node $u$.

  So we hand out at most $\log_{3/2} q \leq \log_{3/2} m$ credits per operation.

- During a deletion we also store an additional credit with $q$ .

  Thus, in total we give out at most $O(m \log m)$ credits.

# Theorem 8.1

A **ScapegoatTree** implements the **SSet** interface. Ignoring the cost of rebuild($u$) operations, a **ScapegoatTree** supports the operations add($x$), remove($x$), and find($x$) in $O(\log n)$ time per operation.

Furthermore, beginning with an empty **ScapegoatTree**, any sequence of $m$ add($x$) and remove($x$) operations results in a total of $O(m \log m)$ time spent during all calls to rebuild($u$).