

Lesson 7

A Review of Data Structures: *Fully Developing the Container of Knowledge*

Wholeness of the Lesson

An analysis of the average-case and worst-case running times of many familiar data structures (for instance, array lists, linked lists, stacks, queues, hashtables, binary search trees) highlights their strengths and potential weaknesses; clarifies which data structures should be used for different purposes; and points to aspects of their performance that could potentially be improved. Likewise, finer levels of intelligence are more expanded but at the same time more discriminating. For this reason, action that arises from a higher level of consciousness spontaneously computes the best path for success and fulfillment.

Overview of the Lesson

- ◆ The List Data Type
- ◆ Stacks
- ◆ Queues
- ◆ Hashtables
- ◆ Binary Search Trees
- ◆ Balanced Binary Search Trees

The List Data Type

- List Types with Running Times

| Operation | ArrayList Running Time | LinkedList Running Time |
|-----------------------|------------------------|-------------------------|
| get(index) | | |
| find(object) | | |
| addFirst(object) | | |
| addLast(object) | | |
| insert(index, object) | | |
| remove(index) | | |
| remove(object) | | |
| Print | | |

The List Data Type

- List Types with Running Times

| Operation | ArrayList Running Time | LinkedList Running Time |
|-----------------------|-------------------------------|---------------------------------------------------------|
| get(index) | $O(1)$ | $O(n)$ (avg and worst) |
| find(object) | $O(n)$ | $O(n)$ |
| addFirst(object) | $O(n)$ (restructure needed) | $O(1)$ |
| addLast(object) | $O(1)$ (unless resize needed) | $O(1)$ (if last pointer maintained) |
| insert(index, object) | $O(n)$ (restructure needed) | $O(n)$ (find insertion point) |
| remove(index) | $O(n)$ (restructure needed) | $O(n)$ (find insertion point; doubly linked preferable) |
| remove(object) | $O(n)$ (find & restructure) | $O(n)$ (find insertion point) |
| Print | $O(n)$ | $O(n)$ |

The STACK ADT

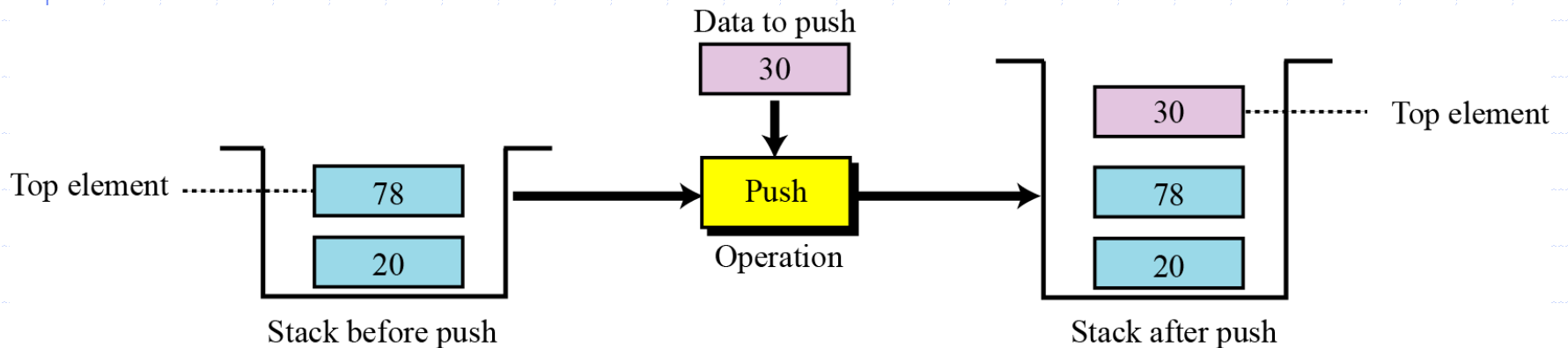
◆ **Definition:** A STACK is a LIST in which insertions and deletions can occur relative to just one designated position (called the *top of the stack*).

◆ **Operations:**

| | |
|------|-----------------------------------------------------|
| pop | remove top of the stack and return this object) |
| push | insert object as new top of stack |
| peek | view object at top of the stack without removing it |

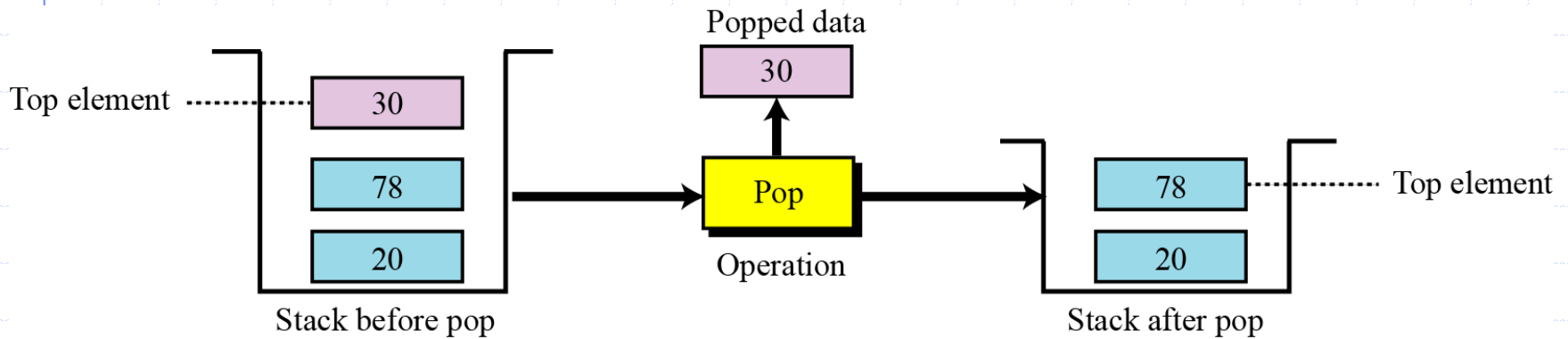
The STACK ADT

◆ push operation:



The STACK ADT

◆ pop operation:



The STACK ADT

◆ Running Times:

- All operations, under a reasonable implementation, run in constant time – that is, $O(1)$.

◆ Implementation:

- Stacks can be implemented using arrays (rightmost element is top) or linked lists.

The QUEUE ADT

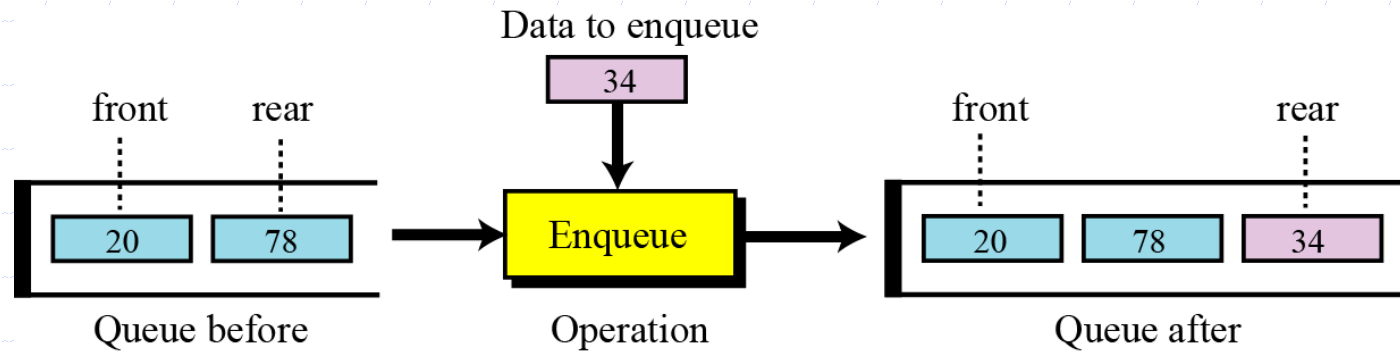
◆ **Definition.** Like a STACK, a QUEUE is a specialized LIST in which insertions may occur only at a designated position (the *back*) and deletions may occur only at a designated position (the *front*).

◆ **Operations:**

| | |
|---------|-----------------------------------------------------------------------|
| dequeue | remove the element at the front (usually also returns this object) |
| enqueue | insert object at the back |
| peek | view object at front of queue without removing it |

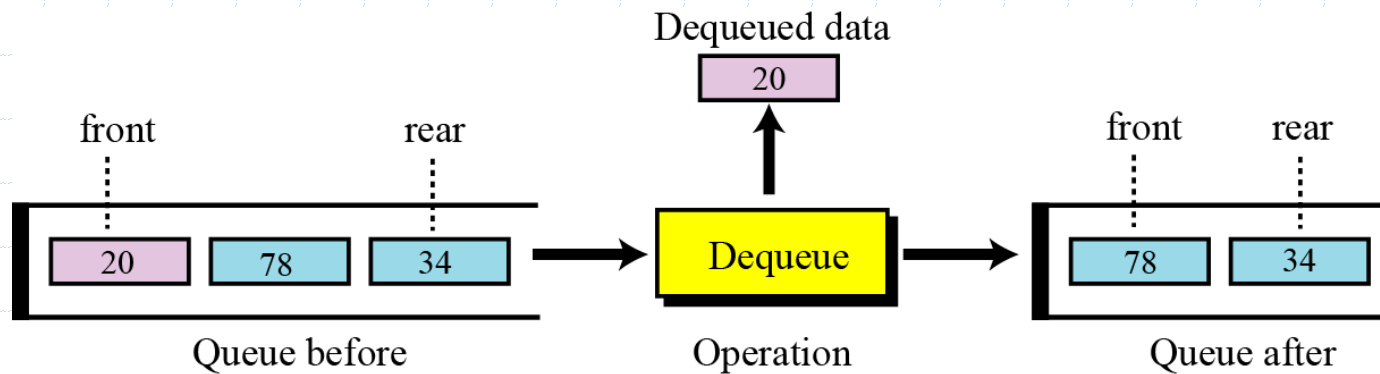
The QUEUE ADT

◆ enqueue operation:



The QUEUE ADT

◆ dequeue operation:



The QUEUE ADT

- ◆ Running Times:

- All operations, under a reasonable implementation, run in constant time.

- ◆ Implementation:

- Queues can be implemented using an array or linked list.

Main Point

Stacks and queues achieve $O(1)$ performance of their main operations, which involve either reading/removing the top element or inserting a new element either at the top or the end. Stacks and queues achieve their high level of efficiency by concentrating on a single point of input (top of stack or end of queue) and a single point of output (top of stack or front of queue).

Stacks and queues make use of the principle from Maharishi Vedic Science that the dynamism of creation arises in the concentration of dynamic intelligence to a point value ("collapse of infinity to a point").

The Hashtable ADT

- ◆ A Hashtable is a generalization of an array in which any object can be used as a key instead of just integers. A Hashtable has *keys* which are used to look up corresponding *values*.
- ◆ Operations:

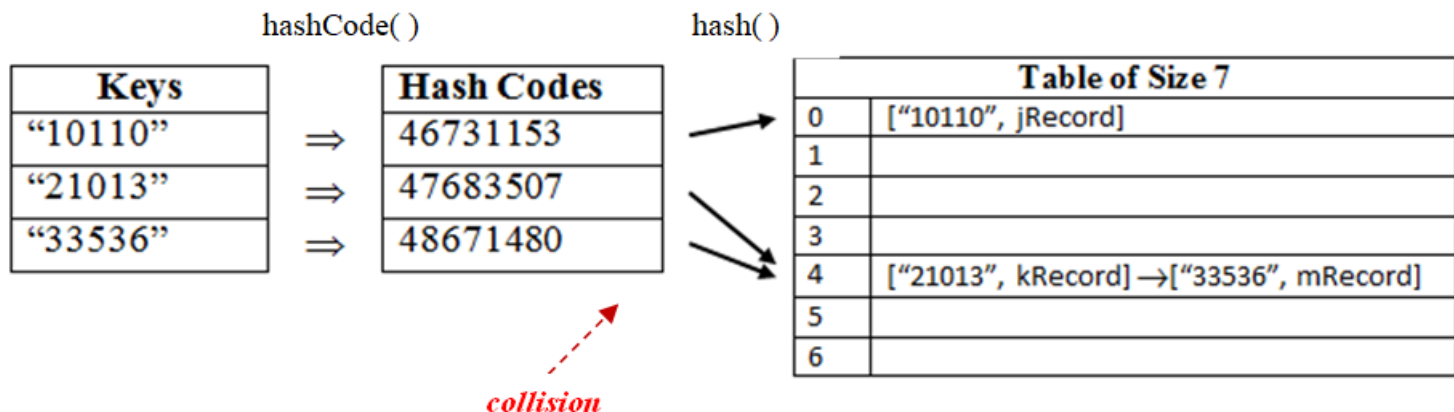
| | |
|--------|--------------------------------------------------------------------------------------------------------------|
| get | returns the value to which the specified key is mapped, or null if this map contains no mapping for the key. |
| put | maps the specified key to the specified value in this hashtable. |
| remove | removes the key (and its corresponding value) from this hashtable. |

Hashing Strategy

- ◆ In a Hashtable, values lying in a large range are compressed into a (typically) much smaller range consisting of array indices.
- ◆ This transformation is accomplished in two Steps:
 - create a *hashcode* for each key – typically, this is a big integer
 - create a *hash value* for each hashcode – the hash value is an index in the underlying table for the Hashtable.

Hashing Strategy

| User's View of Hashtable | |
|--------------------------|---------------------|
| Key (= Employee ID) | Value (= Record) |
| "10110" | jRecord |
| "21013" | kRecord |
| "33536" | mRecord |



$\text{hashcode}(\text{key}) = \text{key.hashCode}()$

$\text{hash}(\text{hashcode}) = \text{hashcode} \% \text{tablesize}$

Running Times for HashMap

- ◆ There are several factors that contribute to the running time of hashtable operations:
 - Cost of computing hashCode().
 - Cost of computing hash values.
 - The degree to which these two computations minimize collisions

Cost of Computing hashCode()

- ◆ Good hashCode functions will create a hashCode based on the instance variables of the object. Suppose your class has instance variables `u`, `v`, `w` and corresponding hashCodes `hash_u`, `hash_v`, `hash_w`, Then:

```
@Override
public int hashCode() {
    int result = 17;
    result += 31 * result + hash_u;
    result += 31 * result + hash_v;
    result += 31 * result + hash_w;
    return result;
}
```

If computing each of the hashCodes required only $O(1)$ time, then running time to execute `hashCode()` for an instance of the object that has `u`, `v`, `w` as its instance variables is also $O(1)$.

Cost of Computing hashCode()

- ◆ *Hashing Strings.* Here is how `hashCode` is overridden in the Java `String` class: (in the `String` class, the variable `value` stores `this.toCharArray()`).

```
public int hashCode() {  
    int h = hash;  
    if (h == 0 && value.length > 0) {  
        char val[] = value;  
  
        for (int i = 0; i < value.length; i++) {  
            h = 31 * h + val[i];  
        }  
        hash = h;  
    }  
    return h;  
}
```

If all strings used as keys in a hashtable have length $< m$ (for some fixed integer m), then the cost of computing each `hashCode` is $O(m)$. This should be viewed as a fixed constant running time, therefore $O(1)$.

Cost of Computing hashCode()

- ◆ *Cases in which hashCode is expensive.* Whenever keys for a hashtable have variable length without some clear upper bound, the cost of performing hashCode may end up being proportional to (or grow even faster than) the table size itself. In such cases, we cannot claim hashCode runs in $O(1)$ time.

Example: You could use Java Lists as keys in a hashtable. If you know for sure that the Lists that will be used as keys always have size less than a fixed number m , then hashCodes of these lists can still be computed in $O(1)$ time. But if the Lists may end up containing arbitrarily many elements, then this is no longer the case.

Cost of Computing Hash Values

- ◆ *Traditional computation.* Suppose you have an underlying table of size n , so that keys for your hashtable need to be mapped into one of the slots $0, 1, 2, \dots, n - 1$. A generally effective method to compute hash values from hash codes is:

$$\text{hashCode} \bmod n$$

Cost of Handling Collisions

- ◆ When the hash value function h maps two *hashCodes* to the same slot, we have a *collision*. When there are too many collisions, the hash table is said to be *skewed*. In the worst case, all *hashCodes* are mapped to the same slot. In that case, all values are stacked up (typically in a list) in a single slot. This results in $O(n)$ performance for *find*, *insert*, and *remove* operations.

Cost of Handling Collisions

- ◆ In order for hashtable operations to be $O(1)$, collisions have to be handled well.
- ◆ *Separate Chaining.* Separate chaining provides a way to handle collisions by placing in a list all key/value pairs for which the hash values associated with the key are equal: All key/value pairs belonging to the same slot are lined up in a list that lives in that slot.

Cost of Handling Collisions

- ◆ *Load Factor.* Suppose the underlying table of a hashtable has n slots and we insert m elements. Then the *load factor* for the hashtable is given by
$$\alpha = m/n.$$

- Assuming input data has no patterns that would cause the table to be skewed, and assuming hashing distributes keys reasonably evenly, then

The expected size of the list in each slot in the table is α .

- Under these assumptions, the running time of put, get and remove is always $O(\alpha)$.
 - ◆ put: $O(1)$ to locate the slot (via hashing), $O(\alpha)$ to put the key in the list
 - ◆ get: $O(1)$ to locate the slot (via hashing), $O(\alpha)$ to find the requested key in the list
 - ◆ remove: $O(1)$ to locate the slot, $O(\alpha)$ to find the requested key in the list, $O(1)$ to remove from the list (recall the list is a *linked* list)

Cost of Handling Collisions

- ◆ When m is $O(n)$ (that is, number of elements is not “too much” larger than tableSize), then α is $O(1)$.
- ◆ *Rehashing.* If the number m of elements becomes significantly larger than the table size n , then the hashtable is no longer efficient. Java handles this by rehashing the hashtable whenever the load factor gets too big (the threshold can be set by the user but defaults to 0.75).

Summary on Hashtables

- ◆ It is accurate to assert that the hashtable operations `get`, `put`, `remove`, `containsKey` perform in $O(1)$ time on average, under the following assumptions:
 - Both `hashCode` and `hash` functions perform in $O(1)$. This will be the case if keys are typical (numeric objects, Strings of bounded length, objects whose instance variables are not too complex)
 - The number of elements in the table is $O(n)$, where n is the table size
- ◆ These assumptions are realistic as long as:
 - The developer overrides `hashCode()` in a way that mixes up values adequately
 - The keys used in a hashtable are not too complex (avoiding using potentially long collection classes or strings whose lengths are not uniformly bounded)

Main Point

Hashtables are a generalization of the concept of an array. They support (nearly) random access of table elements by looking up with a (possibly) non-integer key, and therefore their main operations have an average-case running time of $O(1)$. Hashtables illustrate the principle of *Do less and accomplish more* by providing extremely fast implementation of the main List operations.

Binary Search Trees

- ◆ Question: What is the best data structure to keep data in sorted order in memory?
- ◆ How hard is it to implement this strategy: maintain sorted order to optimize searches ?
 - If we are using an ArrayList, then maintaining the list in sorted order is expensive because each time we add a new element, it must be inserted into the correct spot, and this requires array copy routines.
 - If we are using a LinkedList, searches are not very efficient.
- ◆ *The Need.* We need a data structure that performs insertions efficiently in order to maintain sorted order (like a linked list) but that also performs finds efficiently (as in an array list where binary search is highly efficient)

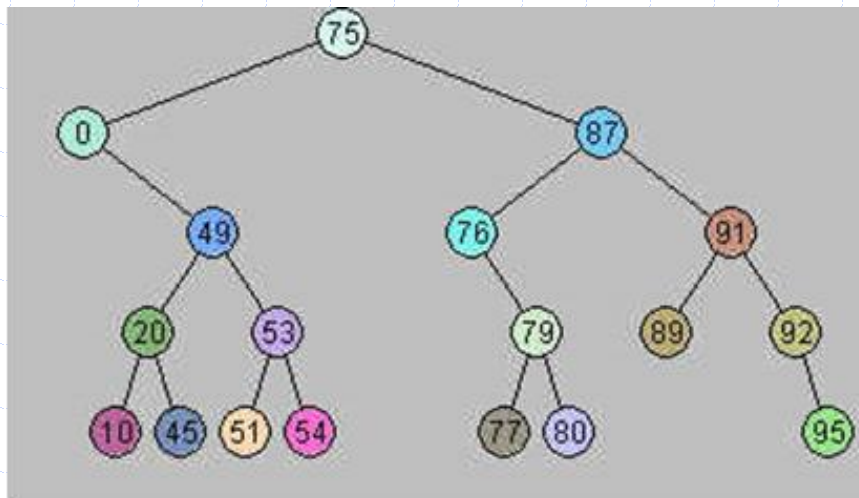
A Solution: Binary Search Trees

- ◆ A *binary tree* is a tree that each node has a reference to a left and right child node (though some references may be null).
- ◆ A *binary search tree* (BST) is a binary tree in which the BST Rule is satisfied:

BST Rule:

At each node N , every value in the left subtree of N is less than the value at N , and every value in the right subtree of N is greater than the value at N .

For the moment, we assume all values are of type Integer.



A Solution: Binary Search Trees

- ◆ The fundamental operations on a BST are:

```
public boolean find(Integer val)
public void insert(Integer val)
public boolean remove(Integer val)
public void print()
```

- ◆ When implemented properly, BSTs perform insertions and deletions faster than can be done on Linked Lists and performs any `find` with as much efficiency as `BinarySearch` on a sorted array.
- ◆ In addition, because of the BST Rule, the BST keeps all data in sorted order, and the algorithm for displaying all data in its sorted order is very efficient.

Insertion into a BST

Recursive algorithm for insertion of Integer x [an iterative implementation is given in the demo code]

- ❖ If the root of the tree being examined is null, create a new root having value x
- ❖ Else:
 - if x is less than the value in the root, insert to the left subtree of the root
 - if x is bigger than the value in the root, insert to the right subtree

Exercise: Insert the following into an initially empty BST:
1,8,2,3,9,5,4

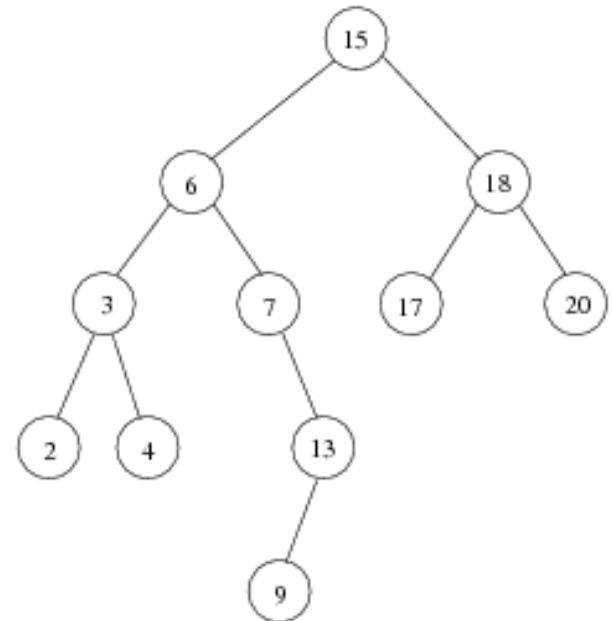
Searching a BST

Recursive algorithm for finding an Integer x

The recursive *find* operation for BSTs is in reality the binary search algorithm in the context of BSTs.

1. If the root is null, return false
2. If the value in the root equals x , return true
3. If x is less than the value in the root, return the result of searching the left subtree
4. If x is greater than the value in the root, return the result of searching the right subtree

Example: search
for 13



BST In-Order Traversal

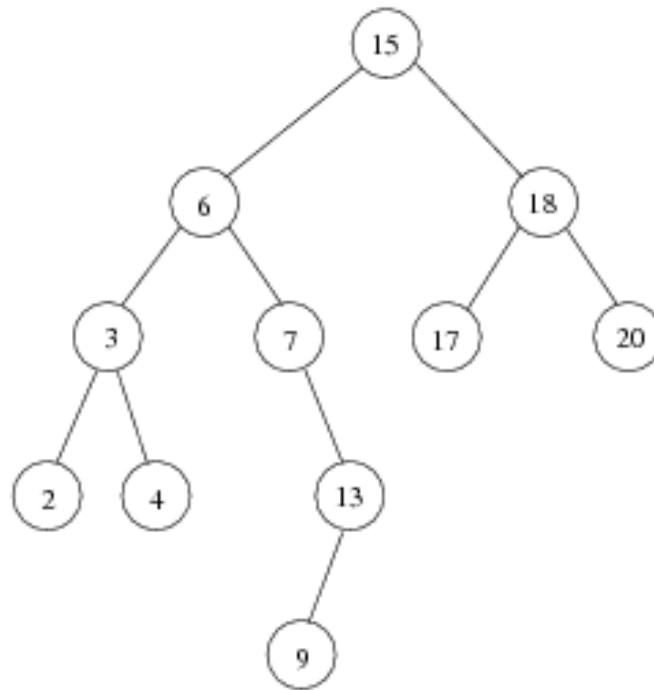
Recursive print algorithm – outputs values in every node in sorted order

1. If the root is null, return
2. Print the left subtree of the root, in sorted order
3. Print the root
4. Print the right subtree of the root, in sorted order

Note: There are several commonly used procedures for visiting all nodes in a binary tree: *in-order traversal*, *pre-order traversal*, *post-order traversal*, and others. In this course we will only discuss in-order traversal since it is the only one that outputs elements in sorted order.

BST In-Order Traversal

◆ Example:

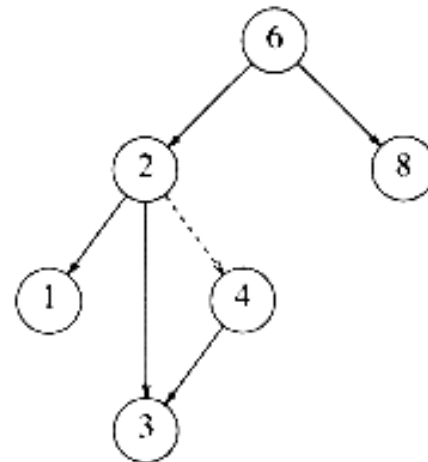
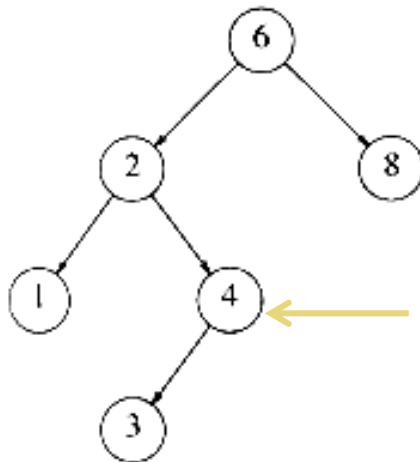


In-Order: 2, 3, 4, 6, 7, 9, 13, 15, 17, 18, 20

Deletions in a BST

Algorithm to remove Integer x

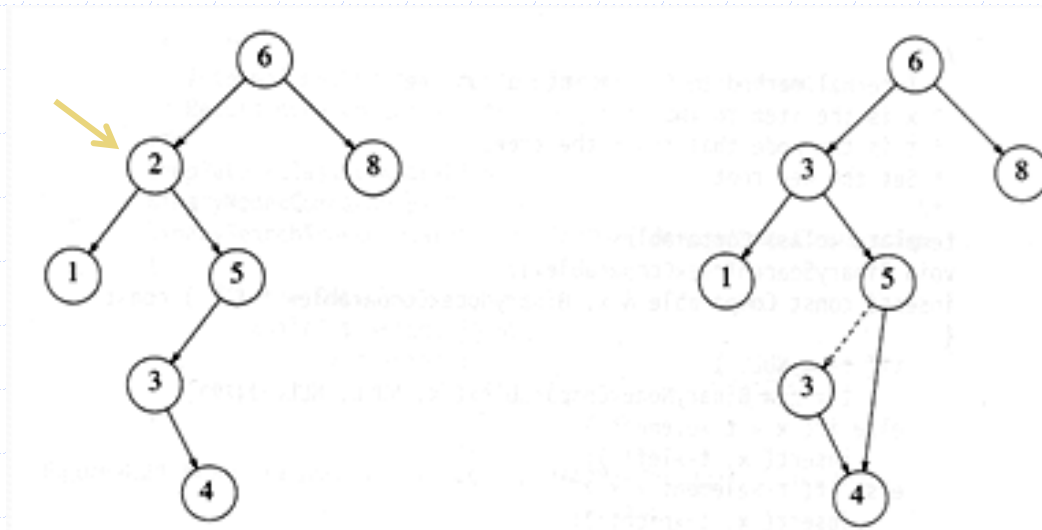
- Case I: Node to remove is a leaf node
 - ◆ Find it and set to null
- Case II: Node to remove has one child
 - ◆ Create new link from parent to child, and set node to be removed to null



Deletions in a BST - continued

Algorithm to remove Integer x

- Case III: Node to remove has two children
 - Find smallest node in right subtree – say it stores an Integer y . This node has at most one child
 - Replace x with y in node to be removed. Delete node that used to store y – this is done as in one of the two cases above



Average Case Analysis of BST Operations

- ◆ What does "average-case" mean in the case of BSTs? There is more than one interpretation, though in each case the result is the same. A BST with n nodes is said to be *randomly built* if n distinct integers are randomly chosen and inserted successively into an initially empty BST.
- ◆ The main result is:

The average depth of a node in a randomly built BST having n nodes is $O(\log n)$.
- ◆ Therefore, the operations of insert, delete and search perform in $O(\log n)$ on average.

Running Time of In-Order Traversal

Algorithm: InOrderTraversal(*node*)

Input: A node *node* of a binary tree B with *n* nodes and with left child *left* and right child *right*

Output: Every node in the subtree of B at *node* is printed

if *node* \neq null **then**

InOrderTraversal(*left*)

print(*node*)

InOrderTraversal(*right*)

Using the technique of counting self calls, we can show that the running time of In-order Traversal is $O(n)$. First, exactly one self-call is made for each node in the tree, so there are exactly n self-calls. Second, the processing time for each self-call, apart from other self-calls, is $O(1)$. Therefore, running time $O(n)$.

Using BSTs for Sorting

- ◆ Consider the following procedure for sorting a list of Integers:
 - Insert them into a BST
 - Print the results (or modify "print" so that it puts values in a list)
- ◆ *Exercise:* How good is this new sorting algorithm?

Handling Duplicates in a BST

- ◆ Running times of find and insert given above, and their implementations, make the assumption that there are no duplicate values.
- ◆ To handle duplicate values, store in each Node a List (instead of a value); then when a value is added to a Node, it is instead added to the List.
- ◆ Example: Suppose we are storing Employees in a BST, ordered by Name. Our BST will be created so that the value in each node is a List<Employee> instance. Then, after insertions, all Employees with the same name will be found in a single List located in a single Node.

Introduction to Balanced BSTs

- ◆ In order to avoid worst case scenarios for BSTs, several kinds of *balanced* BSTs have been devised. Typically, these work by formulating a *balance condition* ensures that, as long as the condition is satisfied, the BST (having n nodes) must have height $O(\log n)$. The balance condition is forced to be true by performing balancing steps after every insertion and deletion.

Introduction to Balanced BSTs: Red-Black Trees

- ◆ A BST is *red-black* if it has the following 4 properties:
 - Every node is colored either red or black
 - The root is colored black.
 - If a node is red, its children are black.
 - For each node n , every path from n to a NULL reference has the same number of black nodes.

BSTs in the Java Libraries: Java uses red-black trees as a background data structure for several of their data structures – namely `TreeSet` and `TreeMap`.

Connecting the Parts of Knowledge With The Wholeness of Knowledge

1. A Binary Search Tree can be used to maintain data in sorted order more efficiently than is possible using any kind of list. Average case running time for insertions and searches is $O(\log n)$.
2. In a Binary Search Tree that does not incorporate procedures to maintain balance, insertions, deletions and searches all have a worst-case running time of $\Omega(n)$. By incorporating balance conditions, the worst case can be improved to $O(\log n)$.
3. *Transcendental Consciousness* is the field of perfect balance. All differences have Transcendental Consciousness as their common source.
4. *Impulses Within The Transcendental Field*. The sequential unfoldment that occurs within pure consciousness and that lies at the basis of creation proceeds in such a way that each new expression remains fully connected to its source. In this way, the balance between the competing emerging forces is maintained.
5. *Wholeness Moving Within Itself*. In Unity Consciousness, balance between inner and outer has reached such a state of completion that the two are recognized as alternative viewpoints of a single unified wholeness.