Algorithms (Trees)

Pramod Ganapathi

Department of Computer Science State University of New York at Stony Brook

September 21, 2021



- Dictionary ADT
- General Trees and Binary Trees
- Binary Search Trees
- Balanced Search Trees
 - (2,4)-Trees
 - B Trees
- Tries and Suffix Trees

Dictionary ADT represents a collection of items, where, each item can be a key or a (key, value) pair.

ADT	ltem	Ordered?	Duplicates?	Implementation
Set	key	X	×	Hash table
Sorted set	key	✓	×	Balanced tree
Multiset	key	×	1	Hash table
Sorted multiset	key	✓	1	Balanced tree
Мар	(key, value)	X	×	Hash table
Sorted map	(key, value)	✓	×	Balanced tree
Multimap	(key, value)	×	1	Hash table
Sorted multimap	(key, value)	✓	1	Balanced tree

Map

• A map is a collection of key-value pairs (k, v), where, keys are unique.

Key	Value
Dictionary word	Word meaning
User ID	User record
Employee ID	Employee record
Student ID	Student record
Patient ID	Patient record
Profile ID	Person details
Order ID	Order details
Transaction ID	Transaction details
URL	Web page
Full file name	File

Method	Functionality
add(e)	Adds the element e to S (if not already present).
remove(e)	Removes the element e from S (if it is present).
<pre>contains(e)</pre>	Returns whether e is an element of S .
iterator()	Returns an iterator of the elements of S .
addAll(T)	Updates S to also include all elements of set T ,
	effectively replacing S by $S \cup T$.
retainAll(T)	Updates \boldsymbol{S} so that it only keeps those elements
	that are also elements of set T , effectively replac-
	ing S by $S \cap T$.
removeAll(T)	Updates \boldsymbol{S} by removing any of its elements that
	also occur in set T , effectively replacing S by $S-T$.

 Set = unordered set; Map = unordered map. java.util.HashSet is an implementation of the set ADT. java.util.HashMap is an implementation of the map ADT.

Method	Functionality
first()	Returns the smallest element in S .
last()	Returns the largest element in S .
<pre>ceiling(e)</pre>	Returns the smallest element $\geq e$.
floor(e)	Returns the largest element $\leq e$.
lower(e)	Returns the largest element $< e$.
higher(e)	Returns the smallest element $> e$.
<pre>subSet(e1,e2)</pre>	Returns an iteration of all elements greater than
	or equal to $e1$, but strictly less than $e2$.
<pre>pollFirst()</pre>	Returns and removes the smallest element in S .
pollLast()	Returns and removes the largest element in S .

- java.util.TreeSet is an implementation of the sorted set ADT. java.util.TreeMap is an implementation of the sorted map ADT.
- TreeSet and TreeMap use balanced search tree

Method	Functionality
add(e)	Adds a single occurrences of e to the multiset.
<pre>contains(e)</pre>	Returns true if the multiset contains an element $= e$.
count(e)	Returns the number of occurrences of e in the multiset.
remove(e)	Removes a single occurrence of e from the multiset.
<pre>remove(e, n)</pre>	Removes n occurrences of e from the multiset.
size()	Returns the number of elements of the multiset
	(including duplicates).
iterator()	Returns an iteration of all elements of the multiset
	(repeating those with multiplicity greater than one).

- Java does not include any form of a multiset. Guava = Google Core Libraries for Java. Guava's Multiset is an implementation of the multiset ADT. Guava's Multimap is an implementation of the multimap ADT.
- Similarly, one can define sorted multiset ADT

Dictionary operations (for unique keys)

	Worst case			Average case		
Data structure	Search	Insert	Delete	Search	Insert	Delete
Sorted array	$\mathcal{O}\left(\log n\right)$	$\mathcal{O}\left(n ight)$	$\mathcal{O}\left(n ight)$	$\mathcal{O}\left(\log n\right)$	$\mathcal{O}\left(n ight)$	$\mathcal{O}\left(n ight)$
Unsorted list	$\mathcal{O}\left(n ight)$	$\mathcal{O}(1)$	$\mathcal{O}\left(n ight)$	$\mathcal{O}\left(n ight)$	$\mathcal{O}(1)$	$\mathcal{O}\left(n ight)$
Hashing	$\mathcal{O}\left(n ight)$	$\mathcal{O}\left(n ight)$	$\mathcal{O}\left(n ight)$	$\mathcal{O}\left(1\right)^{*}$	$\mathcal{O}\left(1\right)^{*}$	$\mathcal{O}\left(1\right)^{*}$
BST	$\mathcal{O}\left(n ight)$	$\mathcal{O}\left(n ight)$	$\mathcal{O}\left(n ight)$	$\mathcal{O}\left(\log n\right)$	$\mathcal{O}\left(\log n\right)$	$\mathcal{O}\left(\log n\right)$
Splay tree	$\mathcal{O}\left(\log n\right)^*$					
Scapegoat tree	$\mathcal{O}\left(\log n\right)$	$\mathcal{O}\left(\log n\right)^*$	$\mathcal{O}\left(\log n\right)^*$	$\mathcal{O}\left(\log n\right)$	$\mathcal{O}\left(\log n\right)$	$\mathcal{O}\left(\log n\right)$
AVL tree	$\mathcal{O}\left(\log n\right)$					
Red-black tree	$\mathcal{O}\left(\log n\right)$					
AA tree	$\mathcal{O}\left(\log n\right)$					
(a, b)-tree	$\mathcal{O}\left(\log n\right)$					
B-tree	$\mathcal{O}(\log n)$	$\mathcal{O}\left(\log n\right)$	$\mathcal{O}(\log n)$	$\mathcal{O}(\log n)$	$\mathcal{O}(\log n)$	$\mathcal{O}\left(\log n\right)$

* = Amortized

General Trees and Binary Trees (HOME)

Family tree



Company organization tree



File system tree



Book organization tree



Term	Meaning
Tree	ADT that stores elements hierarchically
Parent node	Immediate previous-level node
Child nodes	Immediate next-level nodes
Root node	Top node of the tree
Sibling nodes	Nodes that are children of the same parent
External nodes	Nodes without children
Internal nodes	Nodes with one or more children
Ancestor node	Parent node or ancestor of parent node
Descendent node	Child node or descendent of child node
Subtree	Tree consisting of the node and its descendants
Edge	Pair of nodes denoting a parent-child relation
Path	Pair of nodes denoting an ancestor-descendant relation
Ordered tree	Tree with a meaningful linear order among child nodes

Terminology



- A binary tree is an ordered tree with the following properties:
- 1. Every node has at most two children.
- 2. Each child node is labeled as a left child or a right child.
- 3. A left child precedes a right child in the order of children. A recursive definition of the binary tree:
- An empty tree.
- A nonempty tree having a root node r, which stores an element, and two binary trees that are respectively the left and right subtrees of r.



Arithmetic expression tree



Tree represents ((((3+1)*3)/((9-5)+2)) - ((3*(7-4))+6)).

Term	Meaning
Left subtree	Subtree rooted at the left child of an internal node
Right subtree	Subtree rooted at the right child of an internal node
Proper/full tree	A tree in which every node has either 0 or 2 children
Complete tree	Tree in which all except possibly the last level is com- pletely filled and the nodes in the last level are as far left as possible
Perfect tree	Complete tree in which the last level is completely filled



Binary ✓, Proper ✗, Complete ✗, Perfect ✗



Binary ✓, Proper ✓, Complete ✗, Perfect ✗



Binary ✓, Proper ✓, Complete ✗, Perfect ✗



Binary ✓, Proper ✗, Complete ✓, Perfect ✗



Binary ✓, Proper ✓, Complete ✓, Perfect ✗



Binary ✓, Proper ✓, Complete ✓, Perfect ✓

Levels and maximum number of nodes



Let

- T = nonempty binary tree
- $n_{\text{external}} = \text{number of external nodes}$
- $n_{\text{internal}} = \text{number of internal nodes}$

•
$$n = n_{\text{external}} + n_{\text{internal}}$$

• $d_{\max} = \max$ imum depth of the tree Then

•
$$dmax + 1 \le n \le 2^{d_{\max}+1} - 1$$

- $1 \le n_{\text{external}} \le 2^{d_{\max}}$
- $d_{\max} \le n_{\text{internal}} \le 2^{d_{\max}} 1$
- $\log(n+1) 1 \le d_{\max} \le n-1$

If \boldsymbol{T} is a proper nonempty binary tree,

•
$$2d_{\max} + 1 \le n \le 2^{d_{\max} + 1} - 1$$

•
$$d_{\max} + 1 \le n_{\text{external}} \le 2^{d_{\max}}$$

•
$$d_{\max} \le n_{\text{internal}} \le 2^{d_{\max}} - 1$$

•
$$\log(n+1) - 1 \le d_{\max} \le (n-1)/2$$

•
$$n_{\text{external}} = n_{\text{internal}} + 1$$

Implementing a binary tree using linked structure



Implementing a binary tree using array



Implementing a binary tree using array





- Level numbering or level ordering For every node p of T, let f(p) be the whole number defined as: $f(p) = \begin{cases} 0 & \text{if } p \text{ is the root,} \\ 2f(q) + 1 & \text{if } p \text{ is the left child of position } q, \\ 2f(q) + 2 & \text{if } p \text{ is the right child of position } q. \end{cases}$
- Then, node p will be stored at index f(p) in the array.
- $0 \le f(p) \le 2^n 1$, where n = number of nodes in T

Implementing a general tree using linked structure



• A traversal of a tree T is a systematic way of accessing or visiting all the nodes of T.

Traversal	Binary tree?	General tree?
Preorder traversal	~	1
Inorder traversal	✓	×
Postorder traversal	✓	1
Breadth-first traversal	✓	 ✓

Preorder/inorder/postorder traversals

PREORDERTRAVERSAL(root)
 if root ≠ null then VISIT(root) PREORDERTRAVERSAL(root.left) PREORDERTRAVERSAL(root.right)
INORDERTRAVERSAL(root)
 if root ≠ null then INORDERTRAVERSAL(root.left) VISIT(root) INORDERTRAVERSAL(root.right)
PostorderTraversal(root)
 if root ≠ null then POSTORDERTRAVERSAL(root.left) POSTORDERTRAVERSAL(root.right) VISIT(root)

Preorder/inorder/postorder traversals



- Preorder traversal = A B C
- Inorder traversal = B A C
- Postorder traversal = B C A
Preorder/inorder/postorder traversals



- Preorder traversal = A [left] [right] = A B D E C F G
- Inorder traversal = [left] A [right] = D B E A F C G
- Postorder traversal = [left] [right] A = D E B F G C A

Preorder traversal



Preorder traversal: Table of contents

Paper	Paper
Title	Title
Abstract	Abstract
§1	§1
§1.1	§1.1
§1.2	§1.2
§2	§2
§2.1	§2.1

Postorder traversal



COMPUTEDISKSPACE(root)

- 1. $space \leftarrow root.value$
- 2. for each child child of root node do
- 3. $space \leftarrow space + COMPUTEDISKSPACE(root.child)$
- 4. return space

Inorder traversal: Arithmetic expression



General tree.

BREADTHFIRSTTRAVERSAL()

- 1. Q.enqueue(root)
- 2. while Q is not empty do
- $\textbf{3.} \quad curr \leftarrow Q.\mathsf{dequeue}()$
- 4. VISIT(curr)
- 5. for each child child of curr node do
- $6. \qquad Q. {\tt enqueue}(curr.child)$

Binary tree.

BREADTHFIRSTTRAVERSAL()

- 1. Q.enqueue(root)
- 2. while Q is not empty do
- $\textbf{3.} \quad curr \leftarrow Q.\mathsf{dequeue}()$
- 4. VISIT(*curr*)
- 5. **if** left child exists **then** Q.enqueue(curr.left)
- 6. if right child exists then Q.enqueue(curr.right)

Breadth-first traversal: Game trees



Binary Search Trees (BST) HOME

Binary search tree (BST)

A binary search tree is a proper binary tree T such that, for each internal node p of T:

- Node *p* stores an element, say *p.key*.
- Keys stored in the left subtree of p are less than p.key.
- Keys stored in the right subtree of p are greater than p.key.



```
class Node<T>
1.
     {
2.
3.
         T key;
         Node<T> left;
4.
5.
         Node<T> right;
6.
7.
         Node(T item, Node<T> lchild, Node<T> rchild)
         { key = item; left = lchild; right = rchild; }
8.
9
         Node(T item)
10.
         { this(item, null, null); }
11.
12.
```

Search: 65 exists



Search: 68 does not exist



SEARCH(curr, target)	
1. if $curr = null$ then	
2. return <i>curr</i>	ho unsuccessful search
3. else if $target < curr.key$ then	
4. return SEARCH(curr.left, target)	ho recur on left subtree
5. else if $target > curr.key$ then	
6. return SEARCH(curr.right, target)	▷ recur on right subtree
7. else if $target = curr.key$ then	
8. return <i>curr</i>	▷ successful search

SE	EARCH $(curr, target)$	
1.	while $curr \neq null$ do	
2.	if $target < curr.key$ then	
3.	$curr \leftarrow curr.left$	ho recur on left subtree
4.	else if $target > curr.key$ then	
5.	$curr \leftarrow curr.right$	▷ recur on right subtree
6.	else if $target = curr.key$ then	
7.	return curr	ho successful search
8.	return null	▷ unsuccessful search

Search: Analysis



Runtime $\in \Theta(h) \in \mathcal{O}(n)$

Insert 68



Insert 68



INSERT(curr, item)	
Input: Root of tree and item to be inserted Output: New root after item insertion	
1. if $curr = null$ then2. $curr \leftarrow NODE(item)$ \triangleright item doe	es not exist
3. else if $curr \neq null$ then 4. if $item < curr.key$ then	
5. $curr.left \leftarrow \text{INSERT}(curr.left, item) $ \triangleright recur on l	eft subtree
6. else if $item > curr.key$ then	
7. $curr.right \leftarrow \text{INSERT}(curr.right, item) > \text{recur on right}$	ght subtree
8. else if $item = curr.key$ then	
9. do nothing \triangleright	item exists
10. return <i>curr</i>	

In	INSERT(curr, item)			
In	Input: Root of tree and item to be inserted			
Οι	itput: Inserted node			
1.	1. $prev \leftarrow null$			
2.	while $curr \neq null$ do			
3.	$prev \leftarrow curr$			
4.	if $item < curr.key$ then			
5.	$curr \leftarrow curr.left$	▷ recur on left subtree		
6.	else if $item > curr.key$ then			
7.	$curr \leftarrow curr.right$	\triangleright recur on right subtree		
8.	else if $item = curr.key$ then			
9.	return curr	\triangleright item exists		
10.	$curr \leftarrow \text{NODE}(item)$	\triangleright item does not exist		
11.	if $prev \neq null$ then			
12.	if $item < prev.key$ then $prev.left \leftarrow curr$			
13.	if $item > prev.key$ then $prev.right \leftarrow curr$			
14.	return curr			

Insert: Analysis



Runtime $\in \Theta(h) \in \mathcal{O}(n)$

Delete 32: Node 32 has one child



Delete 32: Node 32 has one child



Delete 88: Node 88 has two children



Delete 88: Node 88 has two children



Deleting a node (with a particular key) has four cases:

1. Node is not found.

Do nothing.

- 2. Node is found and it has 0 nonempty children. Delete the node.
- 3. Node is found and it has 1 nonempty child. Delete the node.

Its nonempty child will take the location of the node.

4. Node is found and it has 2 nonempty children.

Locate the predecessor of the node.

 ${\sf Predecessor} = {\sf curr.left.right.right....right}$

Predecessor will take the location of the node.

Predecessor's left child will take the location of the predecessor.

(Can we use successor instead of predecessor?)

Delete: Recursive algorithm

_		
DI	ELETE(curr, item)	
In	out: Root of tree and item to be deleted	
Οι	itput: New root after item deletion	
1.	if $curr = null$ then	
2.	do nothing	⊳ item does not exist
3.	else if $item < curr.key$ then	
4.	$curr.left \leftarrow \text{DELETE}(curr.left, item)$	⊳ recur on left
5.	else if $item > curr.key$ then	
6.	$curr.right \leftarrow \text{Delete}(curr.right, item)$	⊳ recur on right
7.	else if $item = curr.key$ then	\triangleright item exists
8.	if $curr.left = null$ then	ho 0 or 1 child
9.	$curr \leftarrow curr.right$	
10.	else if $curr.right = null$ then	ho 1 child
11.	$curr \leftarrow curr.left$	
12.	else	⊳ 2 children
13.	$curr.key \leftarrow FindMax(curr.left).key$	\triangleright find predecessor
14.	$curr.left \leftarrow \text{Delete}(curr.left, curr.key)$	\triangleright delete predecessor
15.	return curr	

Problem

How do you write an iterative algorithm for deleting an item?

Delete: Analysis



Runtime $\in \Theta(h) \in \mathcal{O}(n)$

Balanced Search Trees (HOME)

Balanced search trees: Motivation

Data structure	Search	Insert	Delete
Binary search tree	$\mathcal{O}\left(n ight)$	$\mathcal{O}\left(n ight)$	$\mathcal{O}\left(n ight)$
Balanced search tree	$\mathcal{O}\left(\log n\right)$	$\mathcal{O}\left(\log n\right)$	$\mathcal{O}\left(\log n\right)$

(2,4)-trees

- A (2,4)-tree or 2-3-4 tree is a balanced search tree.
- A (2,4)-tree satisfies two properties:
 - 1. Size property. Every non-empty node has 2, 3, or 4 children.
 - 2. Depth property. All empty nodes have the same depth.





There are three types of non-empty nodes:

- 2-nodes have 2 children and 1 key. e.g.: [11], [12], [15], [17]
- 3-nodes have 3 children and 2 keys. e.g.: [3 4], [5 10], [13 14]
- 4-nodes have 4 children and 3 keys. e.g.: [6 7 8]

Search: 24 exists



Search: 12 does not exist





Size and depth properties are satisfied.


Overflow: Size property is violated at [13 14 15 17].



Size property at [13 14 15 17] will be fixed via split operation.



Overflow: Size property is violated at [5 10 12 15].



Size property at [5 10 12 15] will be fixed via split operation.



Size and depth properties are satisfied.

Insert: Node split







Insert 10, 8





Underflow: Size property is violated is [4].



Size property will be fixed via transfer operation.





Underflow: Size property is violated is [12], which has non-empty children. It will be fixed via swap with predecessor. Underflow: Size property is violated is [11].



Size property will be fixed via fusion operation.







Delete

 $\begin{array}{l} n_e = \mbox{node with empty children} \\ n_{\neq e} = \mbox{node with non-empty children} \\ s_{3,4} = \mbox{immediate sibling of } n_e \mbox{ is a 3-node or a 4-node} \\ s_2 = \mbox{immediate sibling of } n_e \mbox{ is a 2-node} \\ p = \mbox{parent of } n_e \end{array}$

- Deletion of $n_{\neq e}$ can always be reduced to n_e
- Suppose deleted node is:
 - 1. $n_{\neq e}$.

Swap with the n_e predecessor

2. n_e and $s_{3,4}$ exists.

Transfer a child and key of $s_{3,4}$ to p and a key of p to n_e .

3. n_e and $s_{3,4}$ does not exist.

Fuse/merge n_e with s_2 to get n'_e . Move key from p to n'_e .

Method	Running time
Search	$\mathcal{O}\left(\log n\right)$
Insert	$\mathcal{O}\left(\log n ight)$
Delete	$\mathcal{O}\left(\log n ight)$



Computer memory



Cache-efficient algorithms: Example

Problem

How do you efficiently sort a 1 GB file of natural numbers?

Problem

How do you efficiently sort a 1 GB file of natural numbers?

Workout

Do you want to use quicksort or merge sort, usually implemented in a standard library's sorting algorithm? Your computer program might still take hours to run. Reason? Your algorithm is computation-efficient but not communication-efficient and communication is more expensive than computation. Reducing communication (via good use of cache) leads to reduced running time. An algorithm that makes good use of cache is called cache-efficient. A cache-efficient sorting algorithm might take just a few minutes to sort a 1 GB file of numbers. Example: External-memory merge sort. An algorithm must have the following two features in order to make good use of cache.

- 1. Spatial data locality
- 2. Temporal data locality

• Meaning?

Whenever a cache block is brought into the cache, it contains as much useful data as possible.

• How to exploit?

Group data in blocks (or pages). Move data in blocks.

• Meaning?

Whenever a cache block is brought into the cache, as much useful work as possible is performed on this data before removing the block from the cache.

• Necessary condition?

Total computations is asymptotically greater than space i.e., $T(n)\in\omega\left(S(n)\right)$

• How to exploit?

Design recursive divide-and-conquer algorithms

Cache complexity

- Cache complexity is the asymptotic number of cache misses or page faults incurred by an algorithm.
- Cache-efficient algorithms incur fewer cache misses.
- Cache-efficient algorithms try to exploit both spatial and temporal data locality.
- Terminology: B = data block size, M = cache size



Problem	Cache-inefficient algo	Cache-efficient algo	
Sorting	Merge sort	Ext-memory merge sort	
	$\mathcal{O}\left(n\log n\right)$	$\mathcal{O}\left(\frac{n}{B}\log_{\frac{M}{B}}\frac{n}{B}\right)$	
Balanced tree	(2,4)-tree	B tree	
	$\mathcal{O}\left(\log n\right)$	$\mathcal{O}\left(\log_B n\right)$	
Matrix product	Iterative	Recursive D&C	
	$\mathcal{O}\left(n^{3} ight)$	$\mathcal{O}\left(\frac{n^3}{B\sqrt{M}}\right)$	

- (*a*, *b*)-tree is a straightforward generalization of (2,4)-tree in which the complexities depend on *a* and *b*
- By choosing proper values for *a* and *b*, we get a balanced search tree that has excellent external-memory performance
- (a, b)-tree is a multiway search tree such that each node has between a and b children and stores between a − 1 and b − 1 entries

- An (a, b)-tree is a balanced multiway search tree.
- An (a, b)-tree satisfies three properties:
 - 1. $2 \le a \le (b+1)/2$
 - 2. Size property. Every non-empty node has children in the range [a, b].
 - 3. Depth property. All empty nodes have the same depth.

- B tree of order d is an (a, b) tree with $a = \lfloor d/2 \rfloor$ and b = d.
- B trees are analyzed for cache complexity.
- B trees are cache-efficient, when d = B, as they exploit spatial data locality.



	(2,4)-tree		B tree	
Method	Communication	Computation	Communication	Computation
Search	$\mathcal{O}\left(\log n\right)$	$\mathcal{O}\left(\log n\right)$	$\mathcal{O}\left(\log_B n\right)$	$\mathcal{O}\left(\log n\right)$
Insert	$\mathcal{O}\left(\log n\right)$	$\mathcal{O}\left(\log n\right)$	$\mathcal{O}\left(\log_B n\right)$	$\mathcal{O}\left(\log n\right)$
Delete	$\mathcal{O}\left(\log n\right)$	$\mathcal{O}\left(\log n ight)$	$\mathcal{O}\left(\log_B n\right)$	$\mathcal{O}\left(\log n\right)$

 B trees (and variants such as B+ trees, B* trees, B# trees) are used for file systems and databases.
 Microsoft: NTFS
 Mac: HFS, HFS+
 Linux: BTRFS, EXT4, JFS2
 Databases: Oracle, DB2, Ingres, SQL, PostgreSQL