

Algorithms

(Trees)

Pramod Ganapathi

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

April 3, 2021



Contents

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

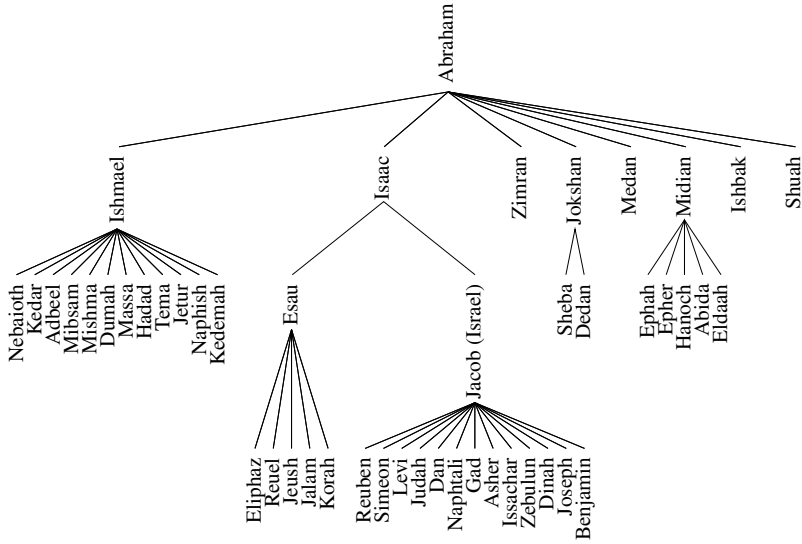
Dictionary operations (for unique keys)

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

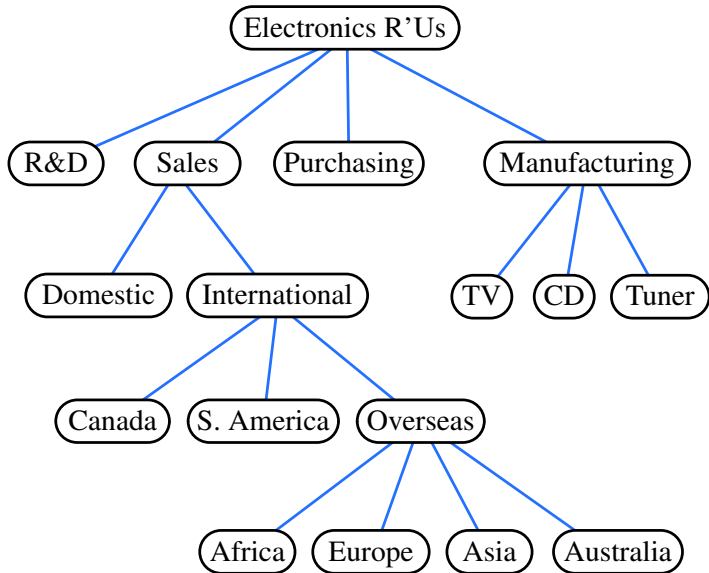
* = Amortized

General Trees and Binary Trees

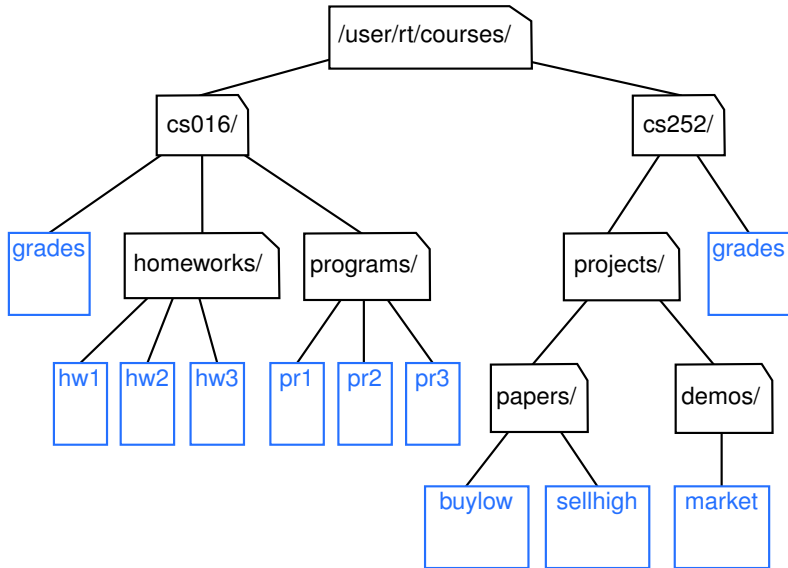
Family tree



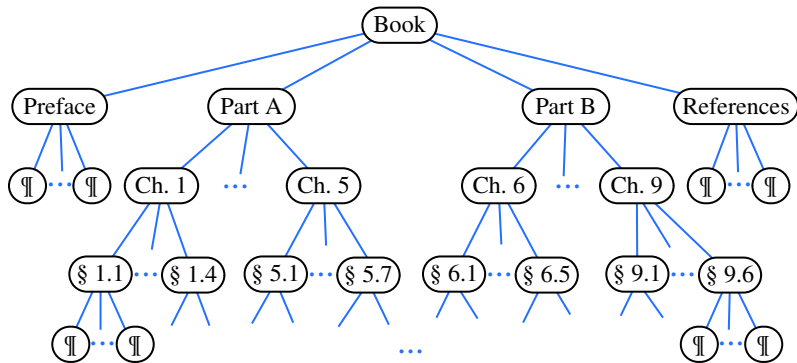
Company organization tree



File system tree



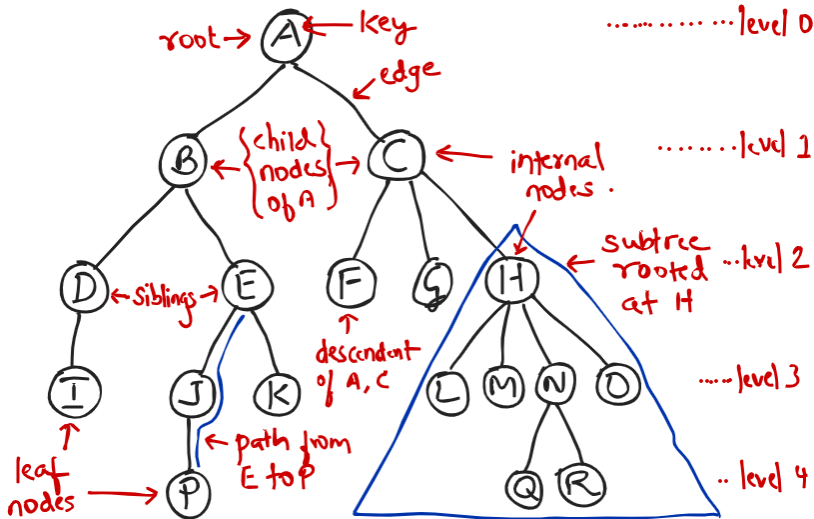
Book organization tree



Terminology

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



Binary tree

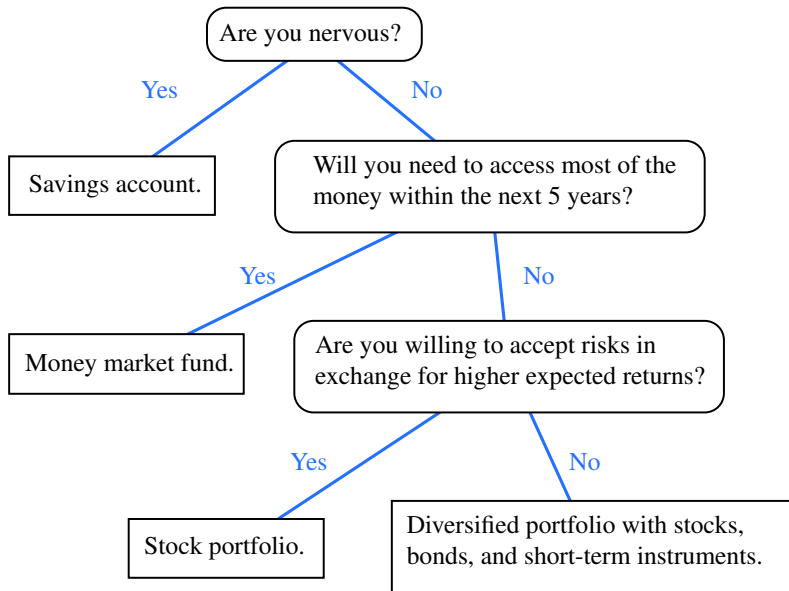
A **binary tree** is an ordered tree with the following properties:

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.

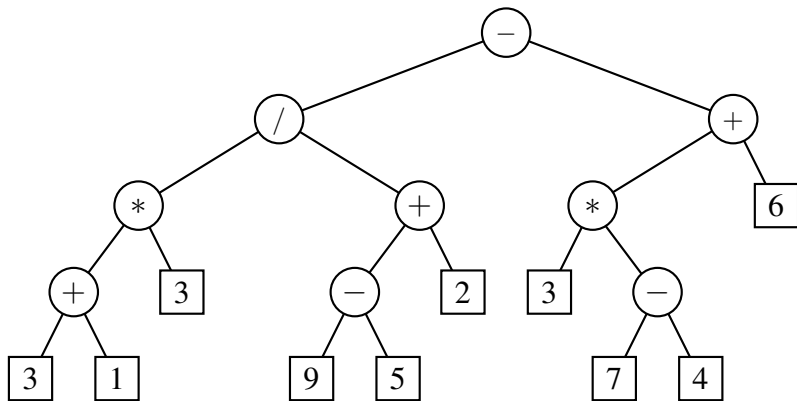
A **recursive definition** of the binary tree:

- 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 .

Decision tree



Arithmetic expression tree

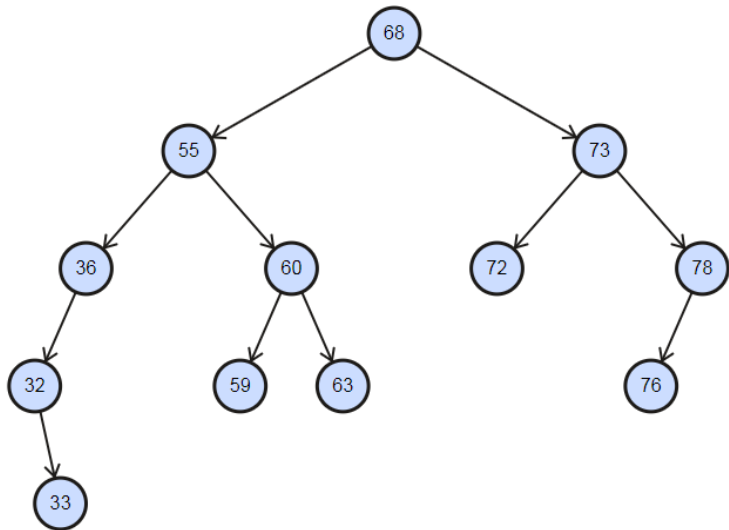


Tree represents $((((3 + 1) * 3) / ((9 - 5) + 2)) - ((3 * (7 - 4)) + 6))$.

Terminology

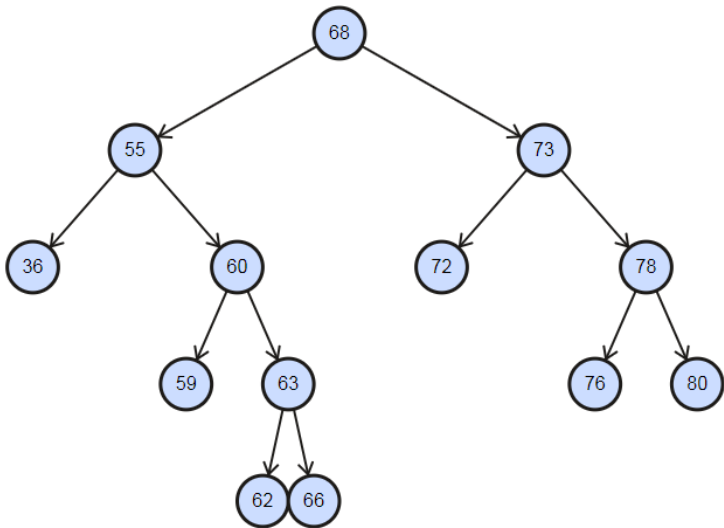
Term	Meaning
Left subtree	Subtree rooted at the left child of an internal node
Right subtree	Subtree rooted at the right child of an internal node
Proper/full tree	A tree in which every node has either 0 or 2 children
Complete tree	Tree in which all except possibly the last level is completely filled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

Tree example



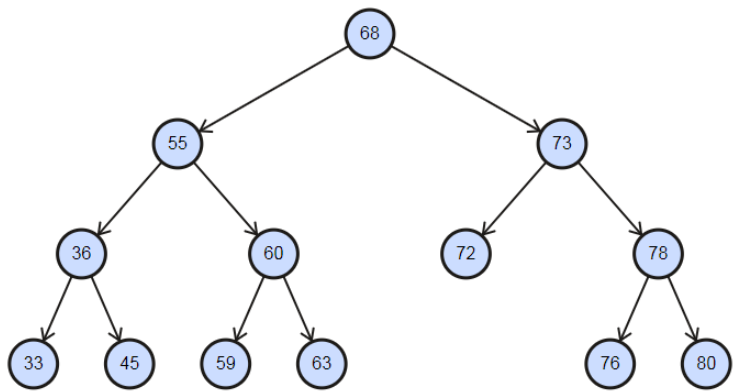
Binary ✓, Proper ✗, Complete ✗, Perfect ✗

Tree example



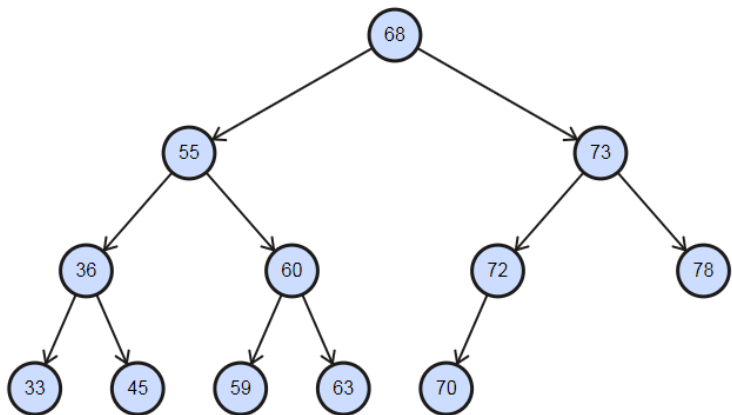
Binary ✓, Proper ✓, Complete ✗, Perfect ✗

Tree example



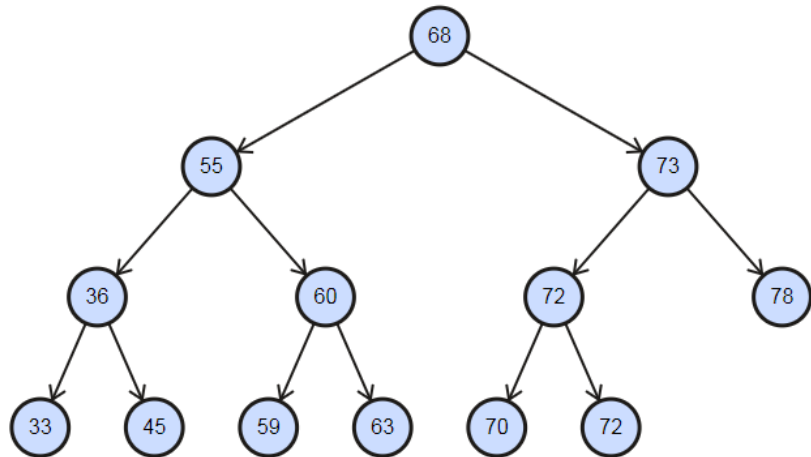
Binary ✓, Proper ✓, Complete ✗, Perfect ✗

Tree example



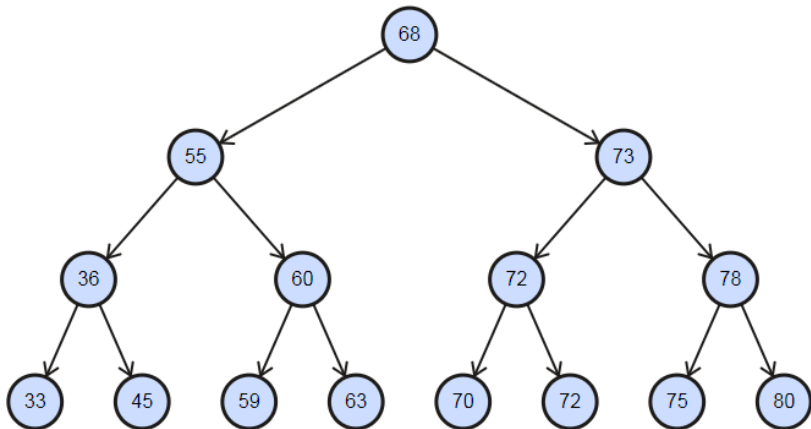
Binary ✓, Proper ✗, Complete ✓, Perfect ✗

Tree example



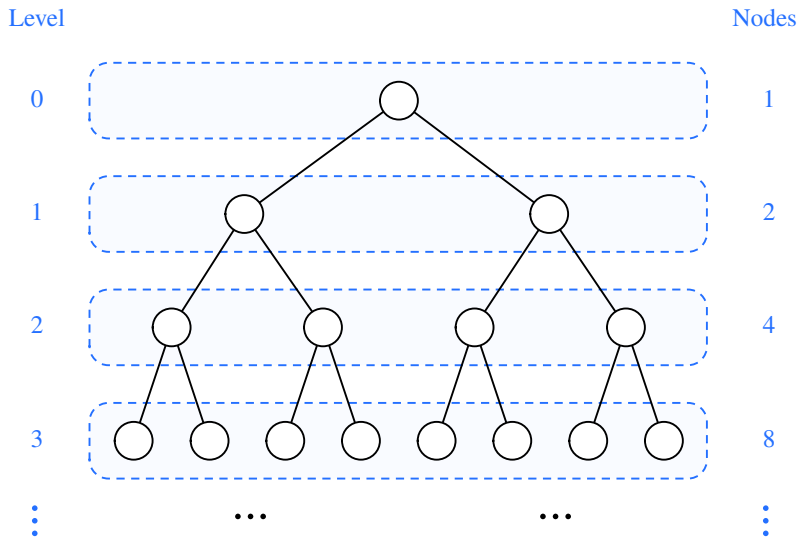
Binary ✓, Proper ✓, Complete ✓, Perfect ✗

Tree example



Binary ✓, Proper ✓, Complete ✓, Perfect ✓

Levels and maximum number of nodes



Properties of binary tree

Let

- T = nonempty binary tree
- n_{external} = number of external nodes
- n_{internal} = number of internal nodes
- $n = n_{\text{external}} + n_{\text{internal}}$
- d_{max} = maximum depth of the tree

Then

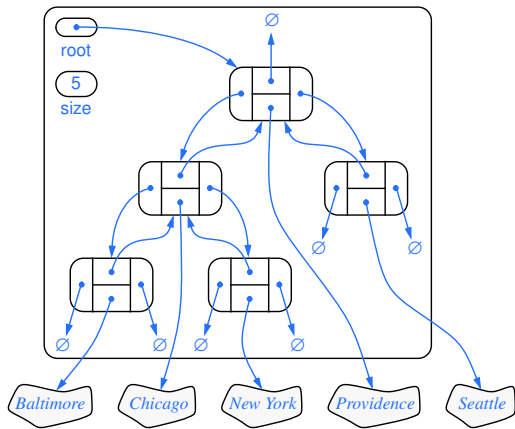
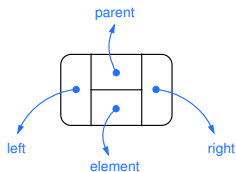
- $d_{\text{max}} + 1 \leq n \leq 2^{d_{\text{max}}+1} - 1$
- $1 \leq n_{\text{external}} \leq 2^{d_{\text{max}}}$
- $d_{\text{max}} \leq n_{\text{internal}} \leq 2^{d_{\text{max}}} - 1$
- $\log(n + 1) - 1 \leq d_{\text{max}} \leq n - 1$

Properties of proper binary tree

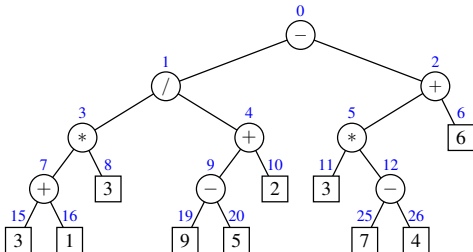
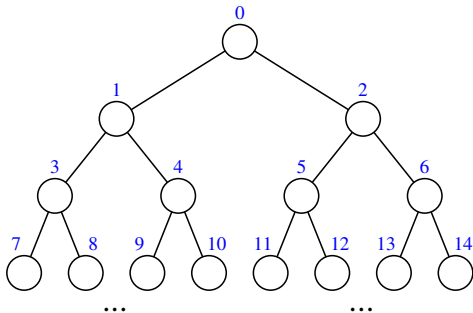
If T is a proper nonempty binary tree,

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

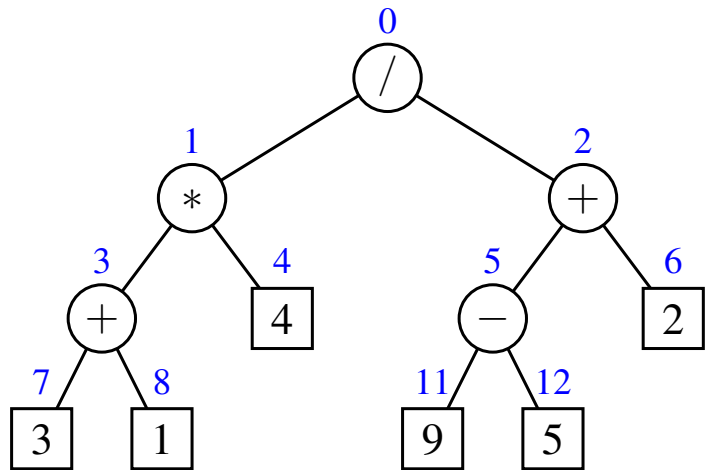
Implementing a binary tree using linked structure



Implementing a binary tree using array



Implementing a binary tree using array



/	*	+	+	4	-	2	3	1			9	5		
0	1	2	3	4	5	6	7	8	9	10	11	12	13	14

Implementing a binary tree using array

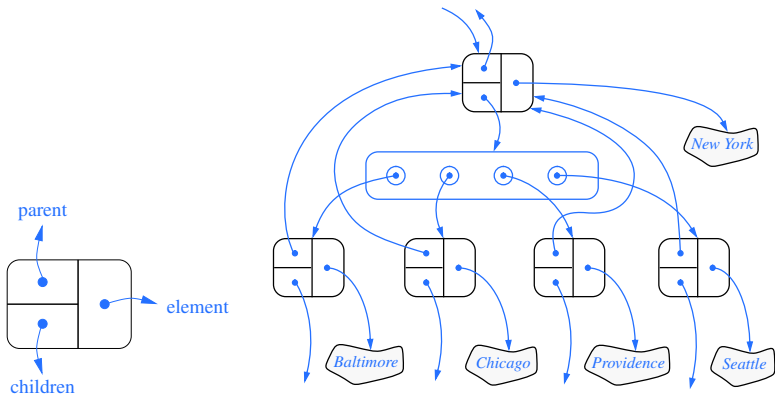
- **Level numbering** or **level ordering**

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.
- $0 \leq f(p) \leq 2^n - 1$, where n = number of nodes in T

Implementing a general tree using linked structure



Tree traversals

- 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	✓	✓
Inorder traversal	✓	✗
Postorder traversal	✓	✓
Breadth-first traversal	✓	✓

Preorder/inorder/postorder traversals

PREORDERTRAVERSAL(*root*)

1. **if** *root* \neq *null* **then**
2. VISIT(*root*)
3. PREORDERTRAVERSAL(*root.left*)
4. PREORDERTRAVERSAL(*root.right*)

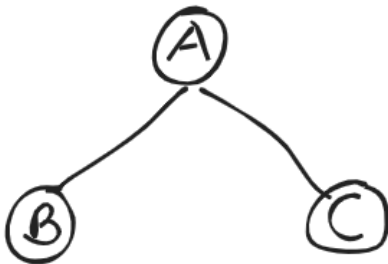
INORDERTRAVERSAL(*root*)

1. **if** *root* \neq *null* **then**
2. INORDERTRAVERSAL(*root.left*)
3. VISIT(*root*)
4. INORDERTRAVERSAL(*root.right*)

POSTORDERTRAVERSAL(*root*)

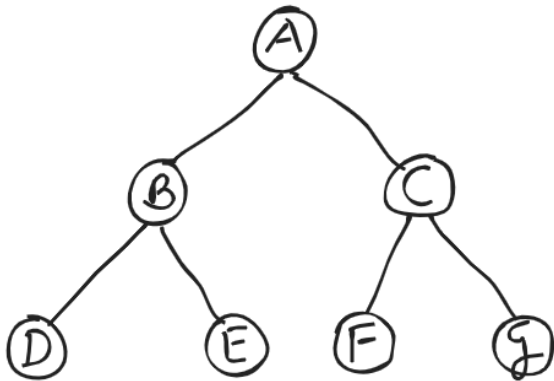
1. **if** *root* \neq *null* **then**
2. POSTORDERTRAVERSAL(*root.left*)
3. POSTORDERTRAVERSAL(*root.right*)
4. VISIT(*root*)

Preorder/inorder/postorder traversals



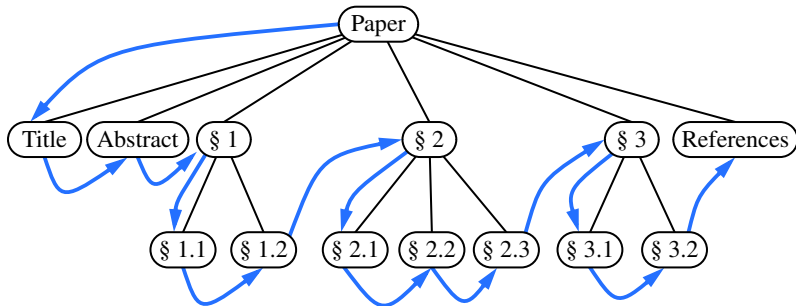
- Preorder traversal = A B C
- Inorder traversal = B A C
- Postorder traversal = B C A

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

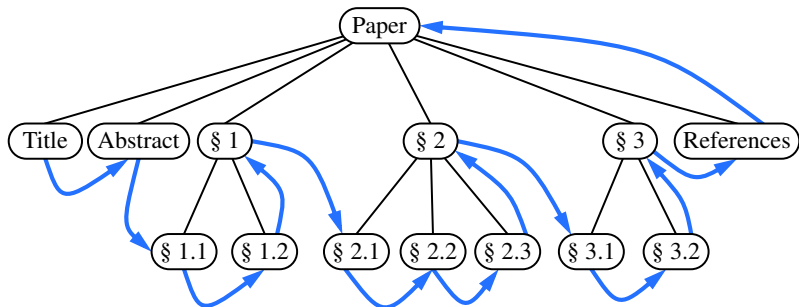


Preorder traversal: Table of contents

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

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

Postorder traversal

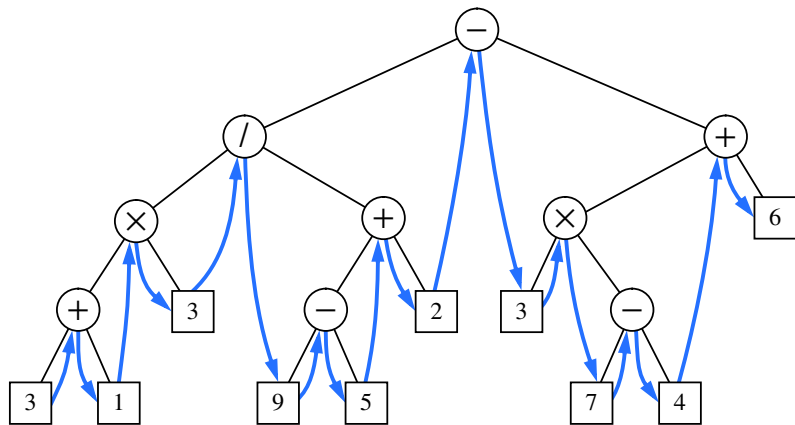


Postorder traversal: Compute disk space

COMPUTEDISKSPACE(*root*)

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

Inorder traversal: Arithmetic expression



$$(((3 + 1) \times 3) / ((9 - 5) + 2)) - ((3 \times (7 - 4)) + 6)$$

Breadth-first traversal

General tree.

BREADTHFIRSTTRAVERSAL()

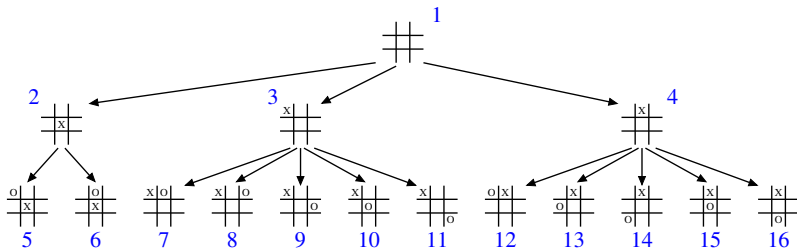
1. $Q.enqueue(root)$
2. **while** Q is not empty **do**
3. $curr \leftarrow Q.dequeue()$
4. VISIT($curr$)
5. **for** each child $child$ of $curr$ node **do**
6. $Q.enqueue(curr.child)$

Binary tree.

BREADTHFIRSTTRAVERSAL()

1. $Q.enqueue(root)$
2. **while** Q is not empty **do**
3. $curr \leftarrow Q.dequeue()$
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

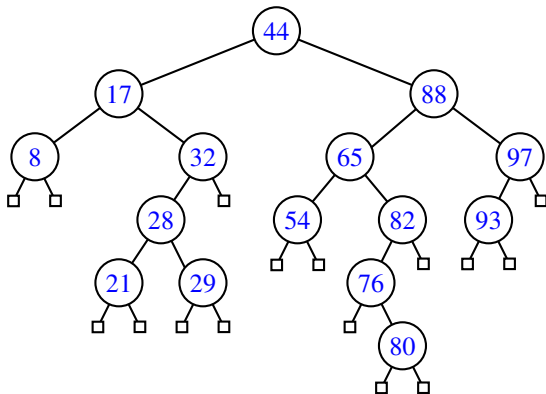


Binary Search Trees (BST)

Binary search tree (BST)

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

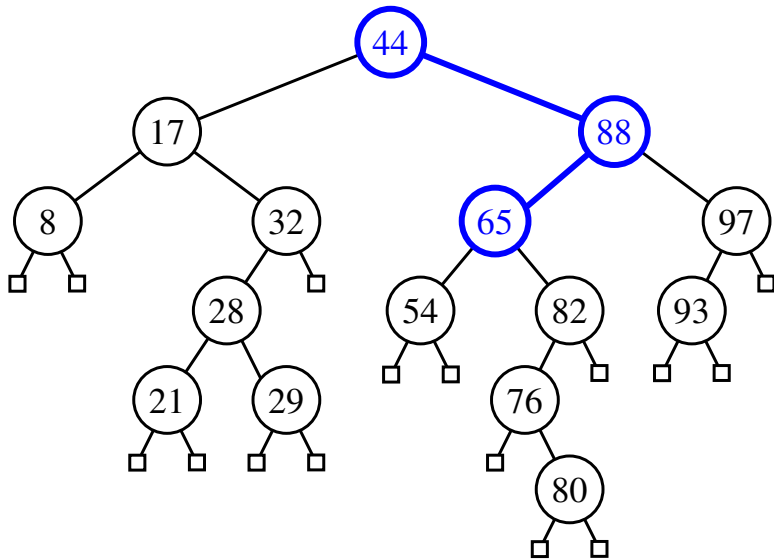
- Node p stores an element, say $p.key$.
- Keys stored in the left subtree of p are less than $p.key$.
- Keys stored in the right subtree of p are greater than $p.key$.



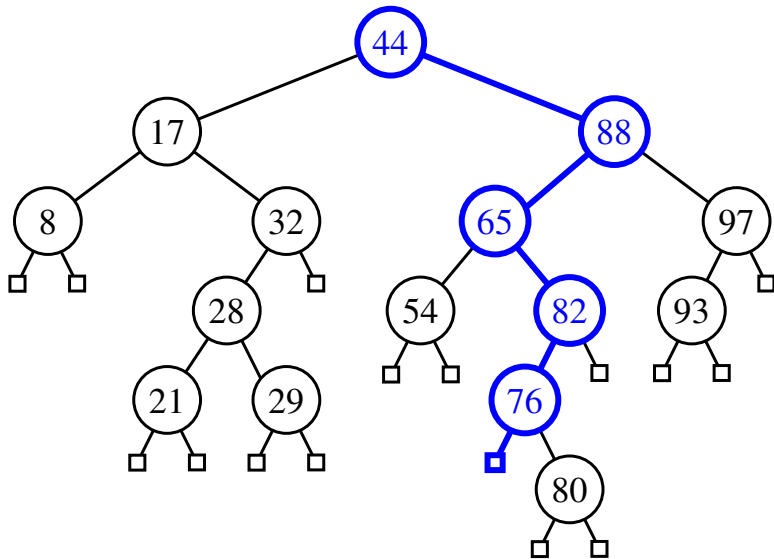
Binary search tree node

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

Search: 65 exists



Search: 68 does not exist



Search: Recursive algorithm

SEARCH(*curr*, *target*)

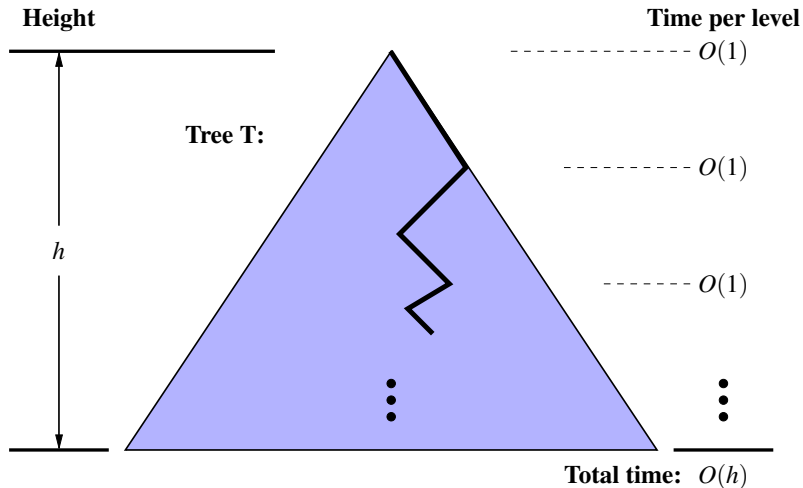
1. **if** *curr* = *null* **then**
2. **return** *curr* ▷ unsuccessful search
3. **else if** *target* < *curr.key* **then**
4. **return** SEARCH(*curr.left*, *target*) ▷ recur on left subtree
5. **else if** *target* > *curr.key* **then**
6. **return** SEARCH(*curr.right*, *target*) ▷ recur on right subtree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

Search: Iterative algorithm

SEARCH(*curr*, *target*)

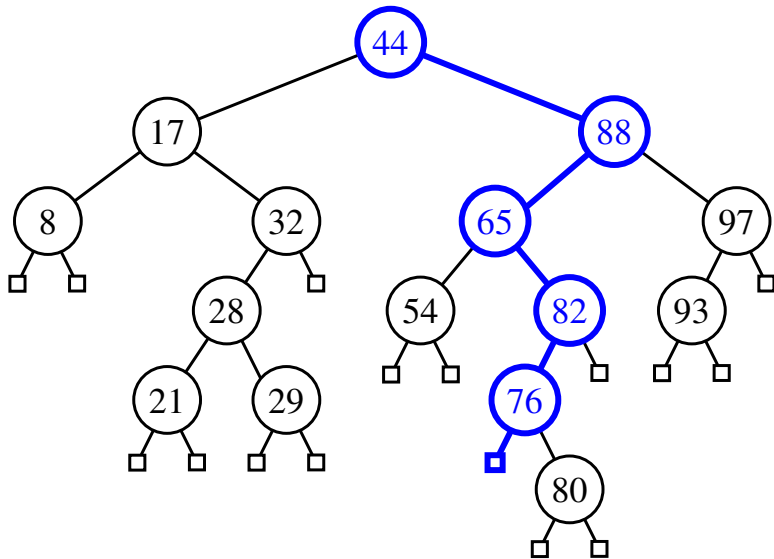
1. **while** *curr* \neq *null* **do**
2. **if** *target* < *curr.key* **then**
3. *curr* \leftarrow *curr.left* ▷ recur on left subtree
4. **else if** *target* > *curr.key* **then**
5. *curr* \leftarrow *curr.right* ▷ recur on right subtree
6. **else if** *target* = *curr.key* **then**
7. **return**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

Search: Analysis

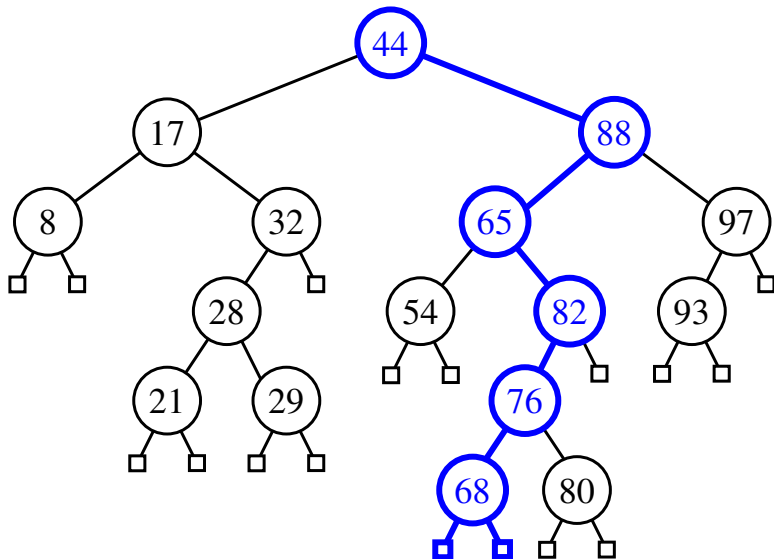


$$\text{Runtime} \in \Theta(h) \in \mathcal{O}(n)$$

Insert 68



Insert 68



Insert: Recursive algorithm

INSERT(*curr*, *item*)

Input: Root of tree and item to be inserted

Output: New root after item insertion

1. **if** *curr* = *null* **then**
2. *curr* \leftarrow NODE(*item*) ▷ item does not exist
.....
3. **else if** *curr* \neq *null* **then**
4. **if** *item* < *curr.key* **then**
5. *curr.left* \leftarrow INSERT(*curr.left*, *item*) ▷ recur on left subtree
6. **else if** *item* > *curr.key* **then**
7. *curr.right* \leftarrow INSERT(*curr.right*, *item*) ▷ recur on right subtree
8. **else if** *item* = *curr.key* **then**
9. do nothing ▷ item exists
.....
10. **return** *curr*

Insert: Iterative algorithm

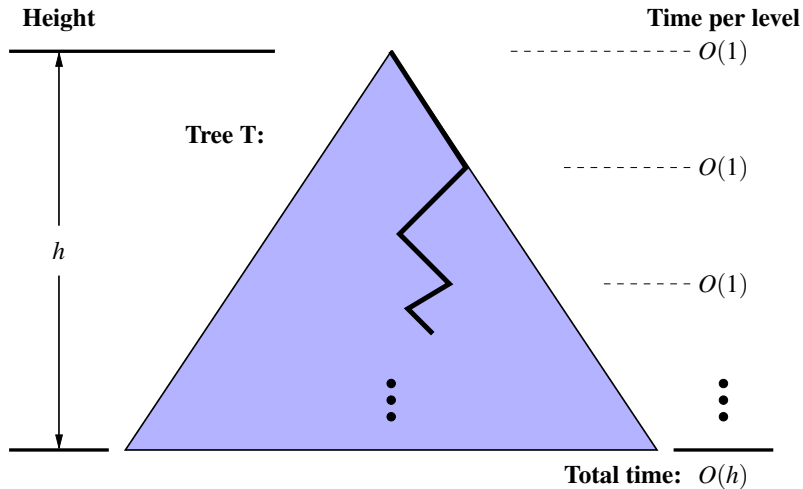
INSERT(*curr*, *item*)

Input: Root of tree and item to be inserted

Output: Inserted node

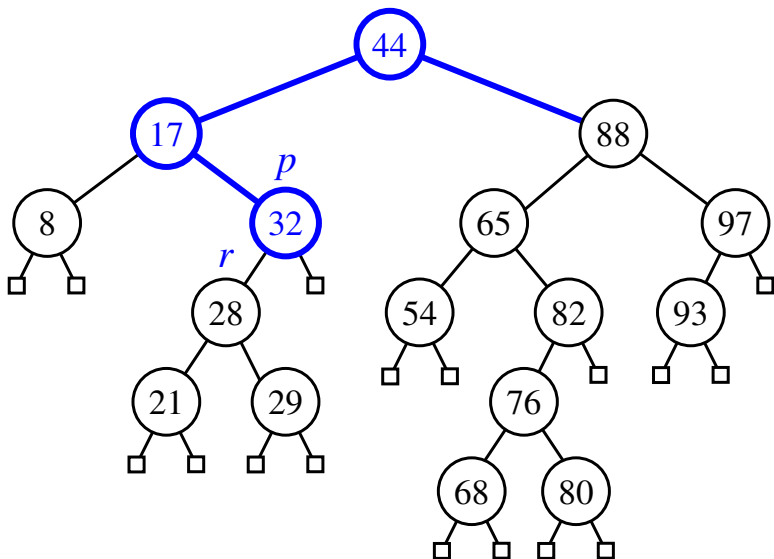
1. *prev* \leftarrow *null*
2. **while** *curr* \neq *null* **do**
3. *prev* \leftarrow *curr*
4. **if** *item* < *curr.key* **then**
5. *curr* \leftarrow *curr.left* ▷ recur on left subtree
6. **else if** *item* > *curr.key* **then**
7. *curr* \leftarrow *curr.right* ▷ recur on right subtree
8. **else if** *item* = *curr.key* **then**
9. **return** *curr* ▷ item exists
-
10. *curr* \leftarrow NODE(*item*) ▷ 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

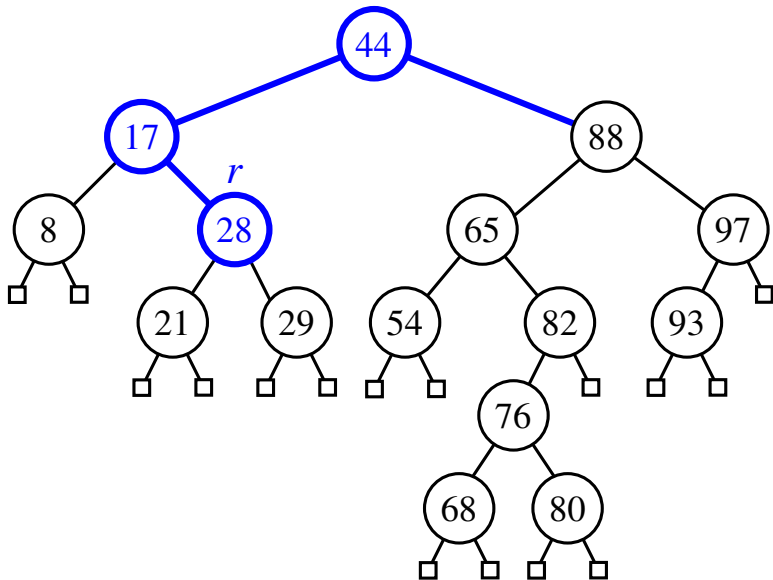


$$\text{Runtime} \in \Theta(h) \in \mathcal{O}(n)$$

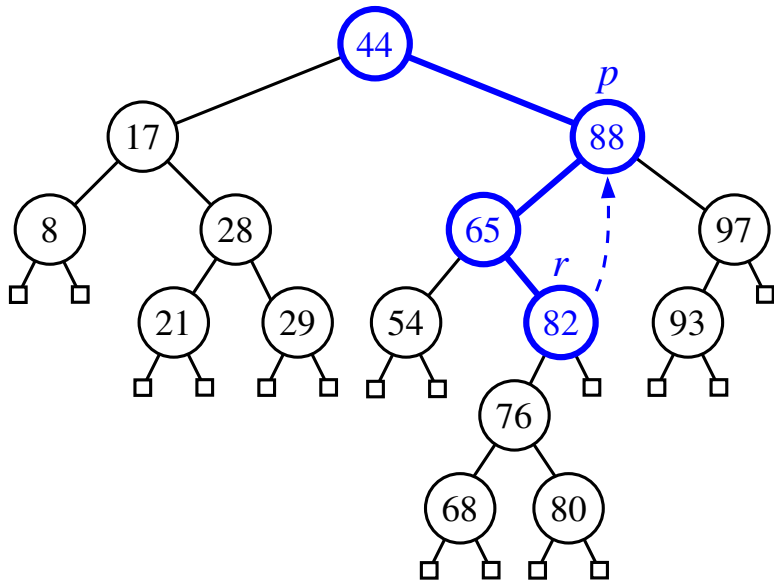
Delete 32: Node 32 has one child



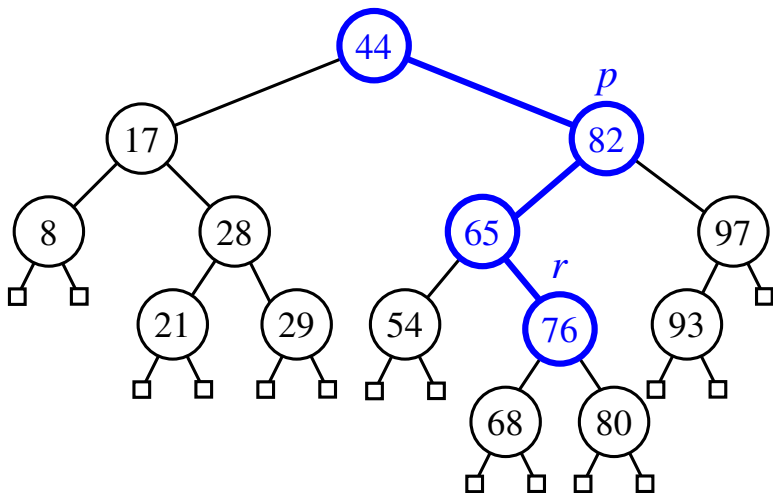
Delete 32: Node 32 has one child



Delete 88: Node 88 has two children



Delete 88: Node 88 has two children



Delete

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.

Predecessor = 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

DELETE(*curr*, *item*)

Input: Root of tree and item to be deleted

Output: New root after item deletion

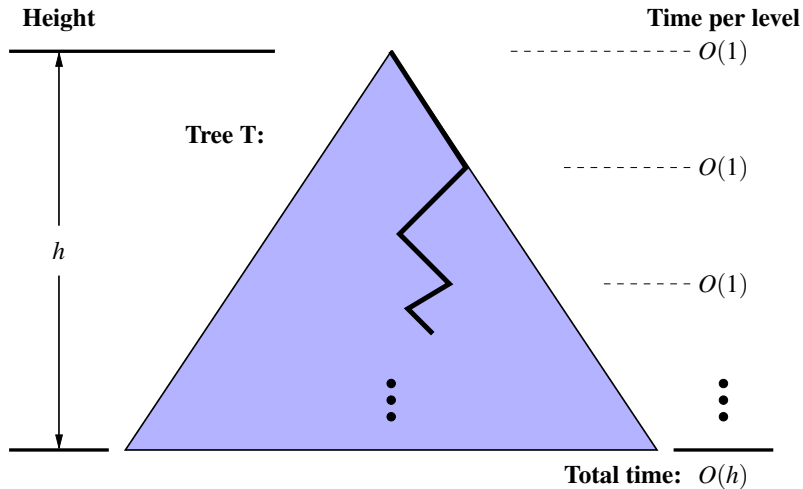
1. **if** *curr* = *null* **then**
2. do nothing ▷ item does not exist
3. **else if** *item* < *curr.key* **then**
4. *curr.left* ← DELETE(*curr.left*, *item*) ▷ recur on left
5. **else if** *item* > *curr.key* **then**
6. *curr.right* ← DELETE(*curr.right*, *item*) ▷ recur on right
7. **else if** *item* = *curr.key* **then** ▷ item exists
-
8. **if** *curr.left* = *null* **then** ▷ 0 or 1 child
9. *curr* ← *curr.right*
10. **else if** *curr.right* = *null* **then** ▷ 1 child
11. *curr* ← *curr.left*
12. **else** ▷ 2 children
13. *curr.key* ← FINDMAX(*curr.left*).*key* ▷ find predecessor
14. *curr.left* ← DELETE(*curr.left*, *curr.key*) ▷ delete predecessor
-
15. **return** *curr*

Delete: Iterative algorithm

Problem

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

Delete: Analysis



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

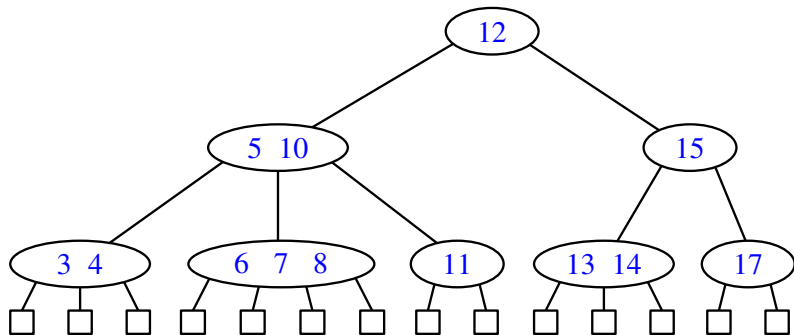
Balanced Search Trees

Balanced search trees: Motivation

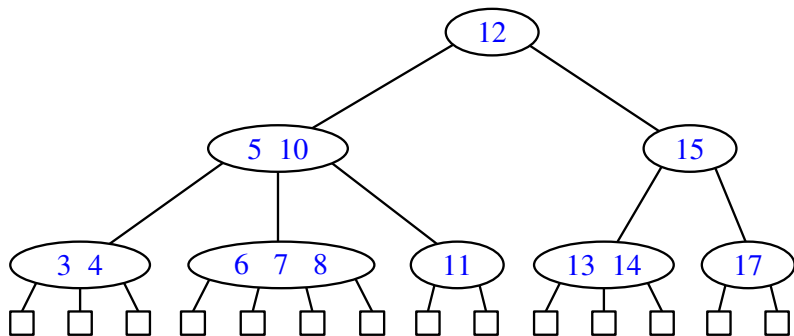
Data structure	Search	Insert	Delete
Binary search tree	$\mathcal{O}(n)$	$\mathcal{O}(n)$	$\mathcal{O}(n)$
Balanced search tree	$\mathcal{O}(\log n)$	$\mathcal{O}(\log n)$	$\mathcal{O}(\log n)$

(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.



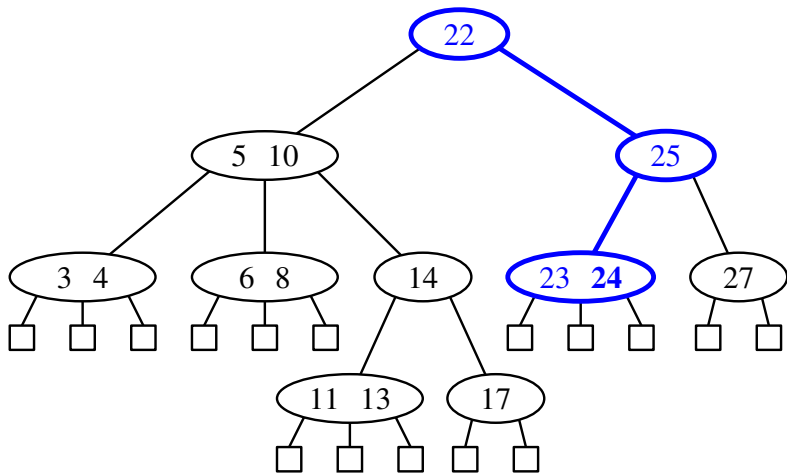
(2,4)-trees



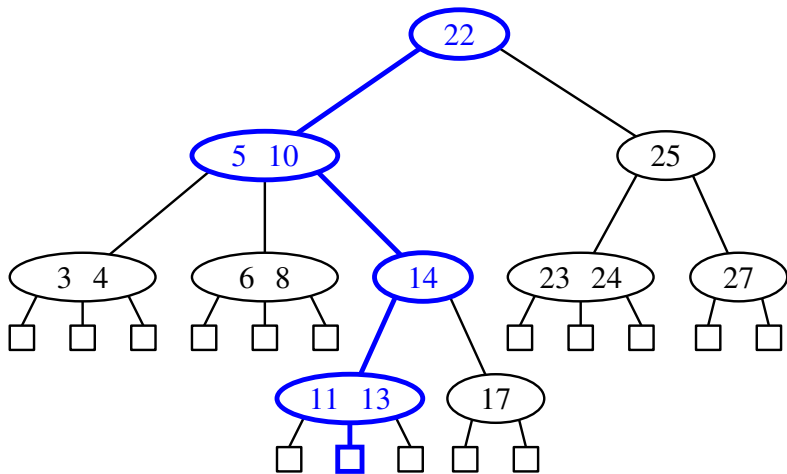
There are three types of non-empty nodes:

- **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]

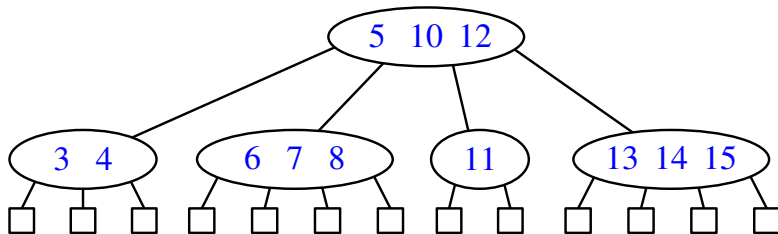
Search: 24 exists



Search: 12 does not exist

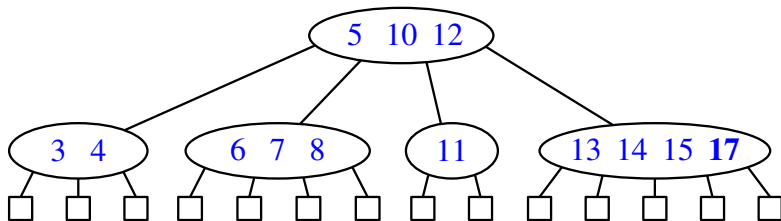


Insert 17



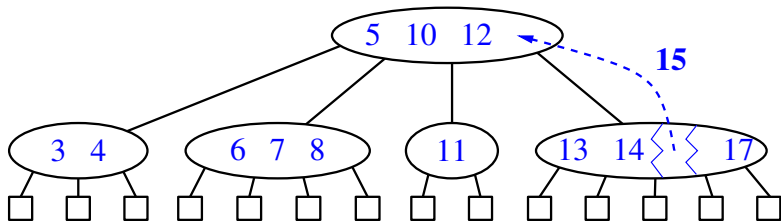
Size and depth properties are satisfied.

Insert 17



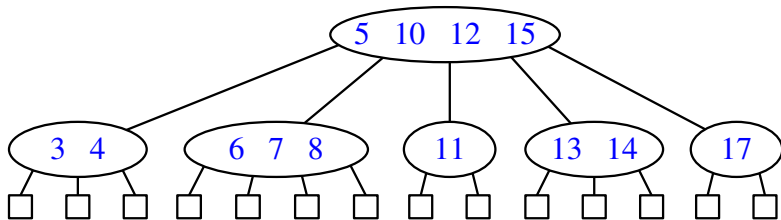
Overflow: Size property is violated at [13 14 15 17].

Insert 17



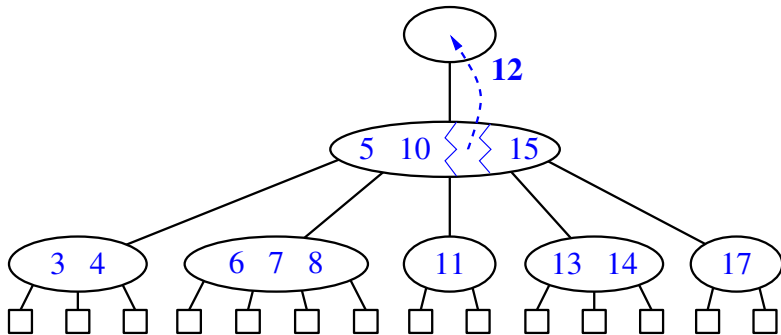
Size property at [13 14 15 17] will be fixed via **split** operation.

Insert 17



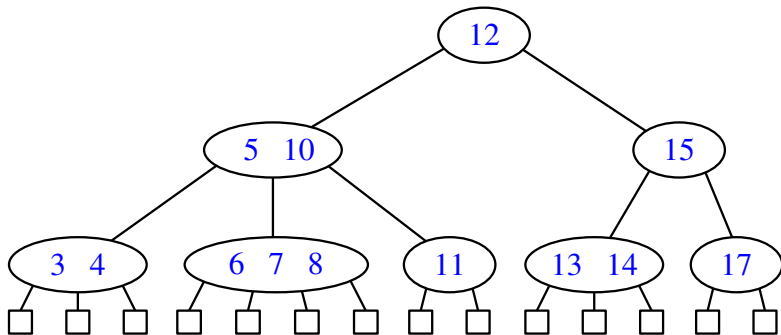
Overflow: Size property is violated at [5 10 12 15].

Insert 17



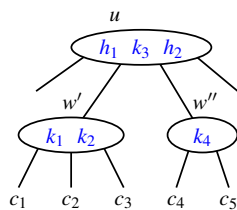
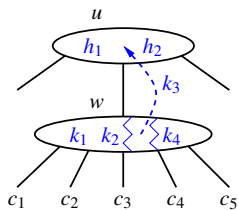
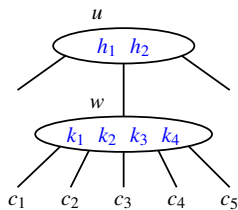
Size property at [5 10 12 15] will be fixed via **split** operation.

Insert 17

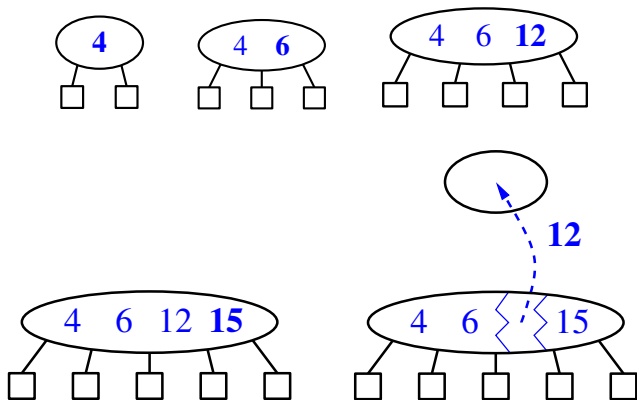


Size and depth properties are satisfied.

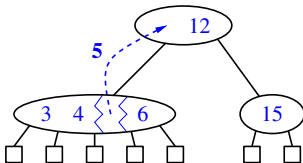
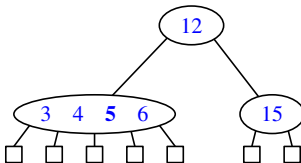
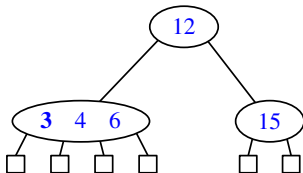
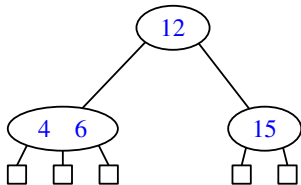
Insert: Node split



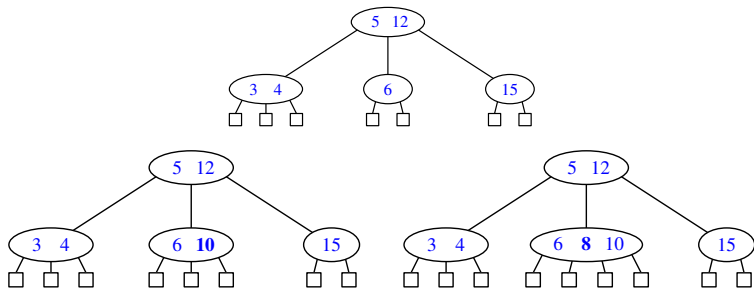
Insert 4, 6, 12, 15



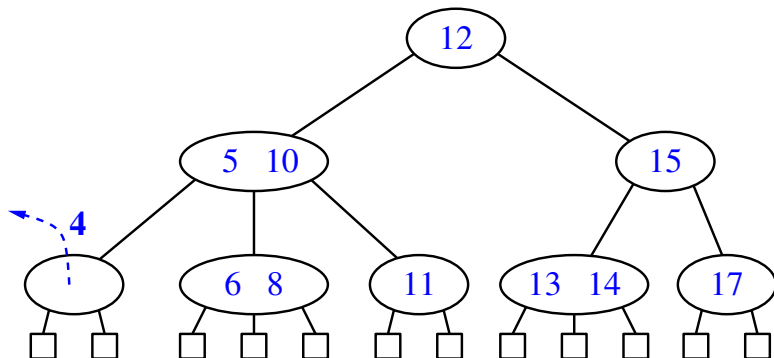
Insert 3, 5



Insert 10, 8

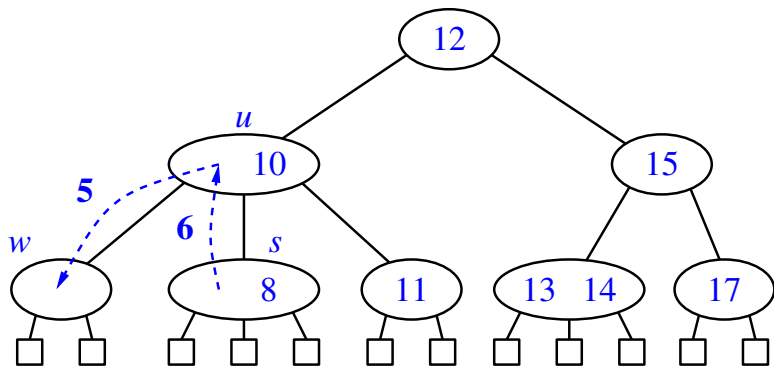


Delete 4



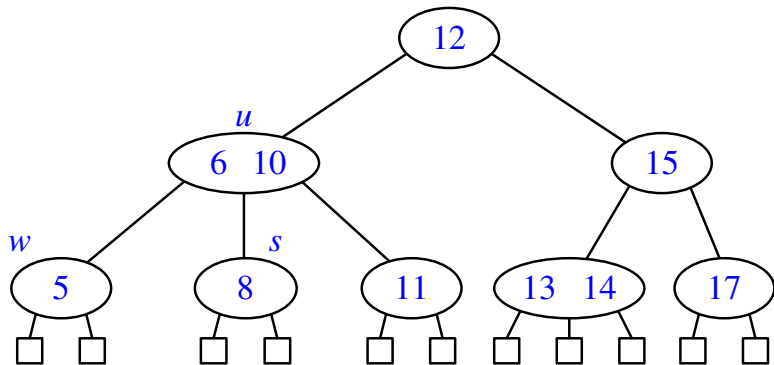
Underflow: Size property is violated is [4].

Delete 4

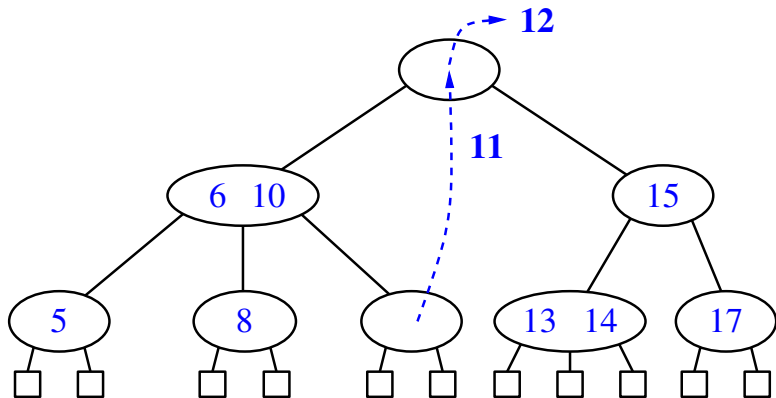


Size property will be fixed via **transfer** operation.

Delete 12



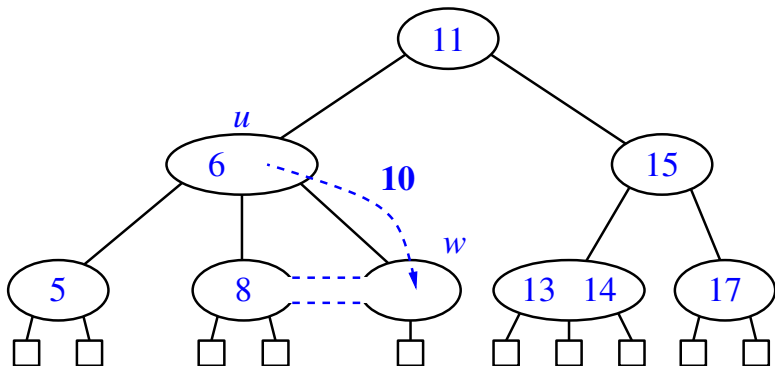
Delete 12



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

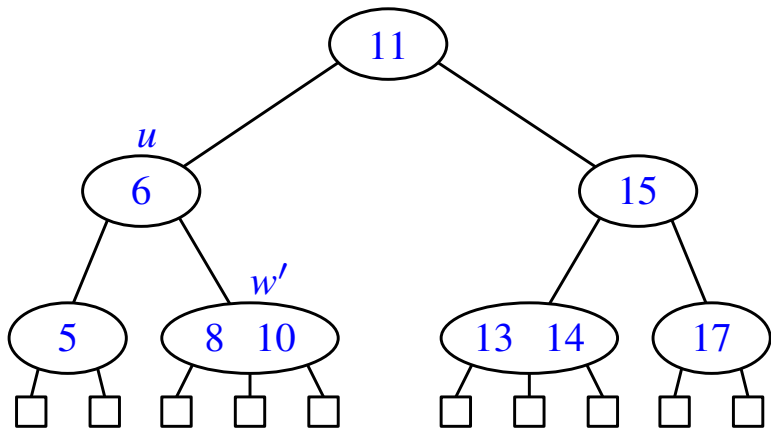
Underflow: Size property is violated is [11].

Delete 12

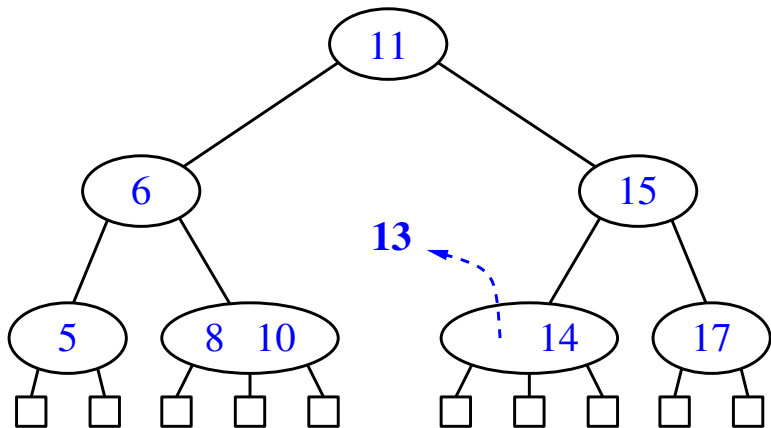


Size property will be fixed via **fusion** operation.

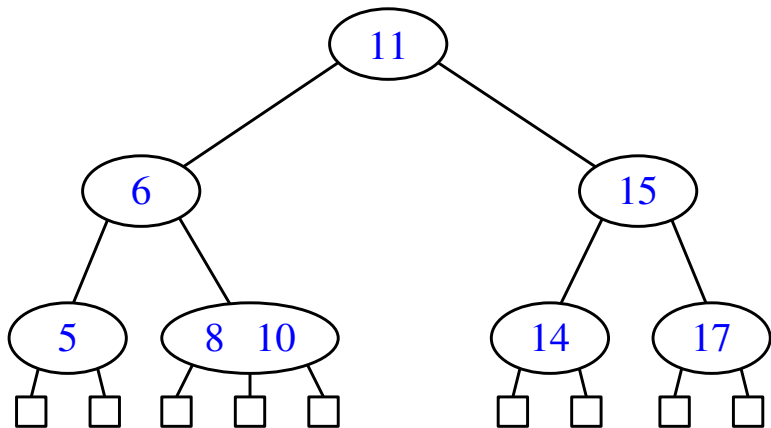
Delete 12



Delete 13



Delete 13



Delete

n_e = node with empty children

$n_{\neq e}$ = node with non-empty children

$s_{3,4}$ = immediate sibling of n_e is a 3-node or a 4-node

s_2 = immediate sibling of n_e is a 2-node

p = parent of n_e

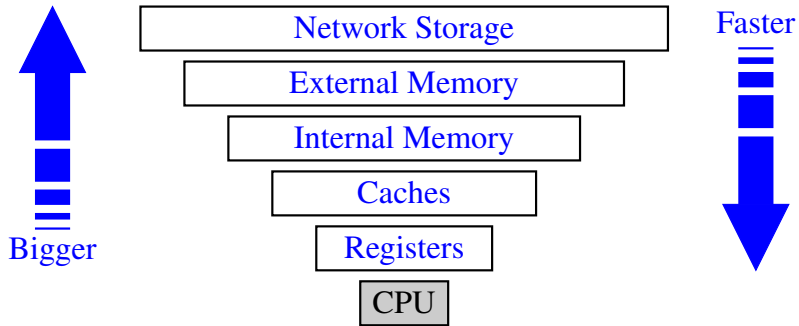
- Deletion of $n_{\neq e}$ can always be reduced to n_e
- Suppose deleted node is:
 1. $n_{\neq e}$.
Swap with the n_e predecessor
 2. n_e and $s_{3,4}$ exists.
Transfer a child and key of $s_{3,4}$ to p and a key of p to n_e .
 3. n_e and $s_{3,4}$ does not exist.
Fuse/merge n_e with s_2 to get n'_e . Move key from p to n'_e .

(2,4)-trees: Complexity

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

B Trees

Computer memory



Cache-efficient algorithms: Example

Problem

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

Cache-efficient algorithms: Example

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**.

Cache data locality

An algorithm must have the following two features in order to make good use of cache.

1. Spatial data locality
2. Temporal data locality

Spatial 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.

Temporal data locality

- 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?

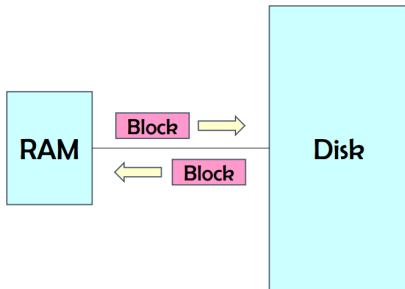
Total computations is asymptotically greater than space
i.e., $T(n) \in \omega(S(n))$

- 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



Cache-efficient algorithms

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

(a, b) -trees

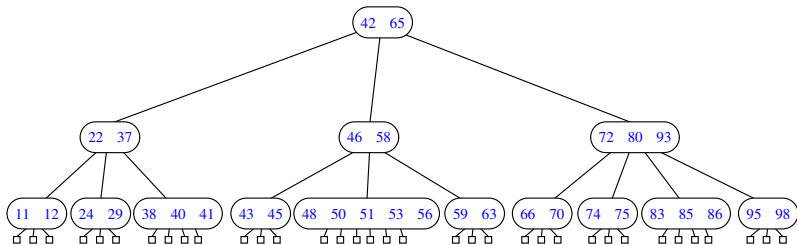
- (a, b) -tree is a straightforward generalization of $(2,4)$ -tree in which the complexities depend on a and b
- By choosing proper values for a and b , we get a balanced search tree that has excellent **external-memory performance**
- (a, b) -tree is a multiway search tree such that each node has between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

(a, b) -trees

- An (a, b) -tree is a balanced multiway search tree.
- An (a, b) -tree satisfies three properties:
 1. $2 \leq a \leq (b + 1)/2$
 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 trees

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



B trees: Complexity

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

- 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