

ADVANCED PROGRAMMING TECHNIQUES

PART IV

Data Structures

Benjamin BOGOSEL

Ecole Polytechnique

Department of Applied Mathematics

- 1 Introduction
- 2 File and Queue
- 3 List, Array, Sequence
- 4 Trees
- 5 Priority Queues
- 6 Dictionaries

- A data structure is a way of organizing and storing information: facilitate access or other purposes
- A data structure has an **interface**: a series of procedure to add, delete, access, re-organize the data
- a data structure can also store additional information ★ for example: the heap size for a heap
- An **abstract data type** (ADT) defines the properties of the structure and of the interface

In the following

- **Dynamic sets**: the number of elements may grow or decrease
- the objects may have multiple attributes, but we can concentrate only on the **key** with respect to which the object is identified
- Some data sets assume that there exists a **total order** on the keys: every two keys can be compared!

Two types: searching/access and modifications

Searching:

- $\text{SEARCH}(S, k)$: returns a pointer x to an element in S such that $x.\text{key} = k$ or NONE if the element is not in S
- $\text{MINIMUM}(S)$, $\text{MAXIMUM}(S)$: returns a pointer for the smallest/largest key
- $\text{SUCCESSOR}(S, x)$, $\text{PREDECESSOR}(S, x)$: returns a pointer for the immediately greater/smaller than x in S , NONE if x is the maximum/minimum

Modifications:

- $\text{INSERT}(S, x)$: insert the element x in S
- $\text{DELETE}(S, x)$: delete the element x in S

- Generally, for an ADT multiple implementations are possible
- The complexity of operations depends on the **implementation**, not on the ADT
- Basic implementation bricks depend on the programming language
- A data structure can be implemented using another existing data structure

- Pile: collection of objects – last in first out (LIFO)
- Queue: collection of objects – first in first out (FIFO)
- Double file: combine the two
- List: ordered collection of objects accessible by their position
- Vector: collection of objects accessible depending on their rank
- Tree: collection of objects structured like a tree
- Priority queue: access to the element of maximal key
- Dictionary: structure implementing the three operations: search, insert, delete

- 1 Introduction
- 2 File and Queue**
- 3 List, Array, Sequence
- 4 Trees
- 5 Priority Queues
- 6 Dictionaries

- Dynamic set of objects: LIFO (last in first out)
- Interface
 - $\text{STACK-EMPTY}(S)$: true if the pile is empty
 - $\text{PUSH}(S, x)$: pushes the value x on the pile S
 - $\text{POP}(S)$: extracts and returns the value at the top of the pile

Applications:

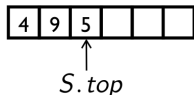
- The **undo** option
- calling functions in a compiler, etc

Implementation:

- with an array: fixed size *a priori*
- linked list (dynamically allocated)

Implementation through an array

- S is an array containing the elements of the pile
- $S.top$ is the position of the top of the pile



```
PUSH( $S, x$ )  
1  if  $S.top == S.length$   
2      error "overflow"  
3   $S.top = S.top + 1$   
4   $S[S.top] = x$ 
```

```
STACK-EMPTY( $S$ )  
1  return  $S.top == 0$ 
```

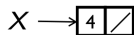
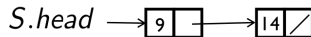
```
POP( $S$ )  
1  if STACK-EMPTY( $S$ )  
2      error "underflow"  
3  else  $S.top = S.top - 1$   
4      return  $S[S.top + 1]$ 
```

★ Complexity in time and space $O(1)$: inconvenient – the space occupied does not depend on the number of objects

- Data structure containing a sequence of elements containing
 - *x.data*: the data
 - *x.next* pointer to the next element
 - (if doubly linked) *x.prev* pointer to the previous element
- *L.head*: pointer to the first element in the list
- (if doubly linked) *L.tail*: pointer to the last element in the list
- void pointers at end of the list (in the corresponding direction)

Implementing a pile using a list

★ S is a simple list



PUSH(S, x)

```
1  $x.next = S.head$   
2  $S.head = x$ 
```

STACK-EMPTY(S)

```
1 if  $S.head == NIL$   
2   return TRUE  
3 else return FALSE
```

POP(S)

```
1 if STACK-EMPTY( $S$ )  
2   error "underflow"  
3 else  $x = S.head$   
4    $S.head = S.head.next$   
5   return  $x$ 
```

★ Complexity in time $O(1)$, complexity in space $O(n)$ for n operations

Application: paranthesis mismatch

★ Check the pairing of parantheses in a string of characters:

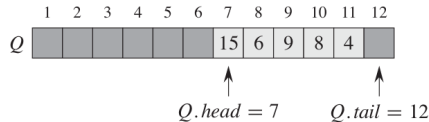
Example: $((x) + (y))/2 \rightarrow no$, $[-b + sqrt(4 * (a) * c)]/(2 * a) \leftarrow yes$

★ Solution based on a pile:

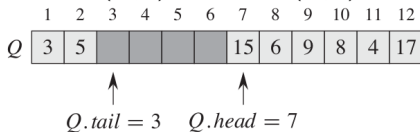
- Iterate through the characters of the string
- If we encounter a left paranthesis: push it on the pile
- If we encounter a right paranthesis pop the last element of the pile and **check if it has the same type**
- If a mismatch occurs return FALSE
- If the pile is empty at the end of the string, without any mismatches, return TRUE

- data structure of the type FIFO (first in first out)
- Interface:
 - $\text{ENQUEUE}(Q, s)$: put the element s at the end of the queue Q
 - $\text{DEQUEUE}(S)$: pop the element at the top of the queue Q
- Implementation using a circular table
 - Q is an array of fixed length $Q.length$ (putting more than $Q.length$ elements in the file gives an error)
 - $Q.head$ is the position of the top of the queue
 - $Q.tail$ is the void position at the end of the queue
 - initially: $Q.head = Q.tail = 1$.

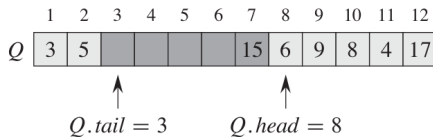
ENQUEUE and DEQUEUE



ENQUEUE(Q , 17), ENQUEUE(Q , 3), ENQUEUE(Q , 5)



DEQUEUE(Q) \rightarrow 15



ENQUEUE and DEQUEUE

ENQUEUE(Q, x)

```
1   $Q[Q.tail] = x$   
2  if  $Q.tail == Q.length$   
3       $Q.tail = 1$   
4  else  $Q.tail = Q.tail + 1$ 
```

DEQUEUE(Q)

```
1   $x = Q[Q.head]$   
2  if  $Q.head == Q.length$   
3       $Q.head = 1$   
4  else  $Q.head = Q.head + 1$   
5  return  $x$ 
```

- complexity in time $O(1)$, complexity in space $O(1)$

Using a list: straightforward implementation, space complexity depends on the length of the list

- 1 Introduction
- 2 File and Queue
- 3 List, Array, Sequence**
- 4 Trees
- 5 Priority Queues
- 6 Dictionaries

- Dynamic set of ordered objects based on their positions
- Generalizes the structures seen previously
- Interface:
 - functions for a double list (insertion and removal at beginning and end)
 - insert **before or after** a position
 - remove the element at a given position
 - replace an element
 - first and last elements of the list
 - previous or next elements in the list
- similarly to a double list: doubly linked list

- Dynamic set of objects occupying ranks given by consecutive integers
- Interface:
 - element at a given rank
 - replace at given rank
 - insert at given rank
 - remove at given rank
 - vector size

Implementation: list, extensible array

- insertion
 - $O(n)$ for an individual operation: n is the size of the vector
 - $O(n)$ for n inserts at the end of the vector
 - $O(n^2)$ for n inserts at the beginning of the vector (shift all elements right...)
- removal: similar

- 1 Introduction
- 2 File and Queue
- 3 List, Array, Sequence
- 4 Trees**
- 5 Priority Queues
- 6 Dictionaries

Abstract data structure for a tree

- Main idea: data is associated to the nodes of a tree
- The nodes are accessible depending on their relative position in the tree

Interface: For a tree T and a node n

- $\text{PARENT}(T, n)$: give the parent of the node n
- $\text{ISEMPTY}(T)$: true if the tree is void
- $\text{CHILDREN}(T, n)$: give a data structure containing the children of node n (ordered or not)
- $\text{ISROOT}(T, n)$: true if n is the root
- $\text{ISINTERNAL}(T, n)$: true if the node is internal
- $\text{ISEXTERNAL}(T, n)$: true if the node is external
- $\text{GETDATA}(T, n)$: get data from node n
- $\text{ROOT}(T)$: give the root node
- $\text{SIZE}(T)$: the number of nodes in the tree

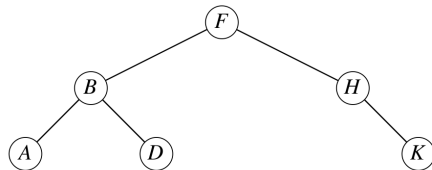
- Depth of a node: count the number of parents till reaching the root node $O(n)$
- Compute the height of the tree
- etc...

First solution: level numerotation

- ★ Root at position 1
- ★ from parent r reach children at positions $2r$ and $2r + 1$
- ★ may contain void elements if the tree is not complete
- ★ complexity in space $O(2^n)$, complexity in time $O(1)$

Other options: linked structures - for each node keep pointers to its parent and its left/right children

- A way of ordering the nodes in the tree, to iterate through them
- In depth
 - Infix (in order)
 - Prefix (in preorder)
 - Suffix (in postorder)
- In breadth

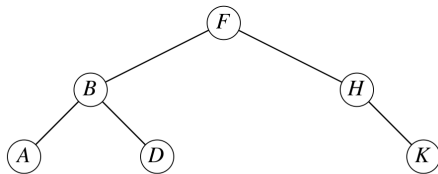


$\Rightarrow \langle A, B, D, F, H, K \rangle$

★ each node is visited after its left child and before its right child

```
INORDER-TREE-WALK( $T, x$ )  
1  if HASLEFT( $T, x$ )  
2      INORDER-TREE-WALK( $T, \text{LEFT}(x)$ )  
3  print GETDATA( $T, x$ )  
4  if HASRIGHT( $T, x$ )  
5      INORDER-TREE-WALK( $T, \text{RIGHT}(x)$ )
```

Prefix traversal

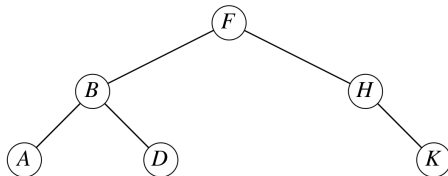


$\Rightarrow \langle F, B, A, D, H, K \rangle$

★ each node is visited before its children

```
PREORDER-TREE-WALK( $T, x$ )  
1  print GETDATA( $T, x$ )  
2  if HASLEFT( $T, x$ )  
3      PREORDER-TREE-WALK( $T, \text{LEFT}(x)$ )  
4  if HASRIGHT( $T, x$ )  
5      PREORDER-TREE-WALK( $T, \text{RIGHT}(x)$ )
```

Postfix traversal



$\Rightarrow \langle A, D, B, K, H, F \rangle$

★ each node is visited after its children

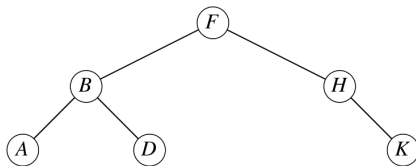
```
POSTORDER-TREE-WALK( $T, x$ )  
1  if HASLEFT( $T, x$ )  
2      POSTORDER-TREE-WALK( $T, \text{LEFT}(x)$ )  
3  if HASRIGHT( $T, x$ )  
4      POSTORDER-TREE-WALK( $T, \text{RIGHT}(x)$ )  
5  print GETDATA( $T, x$ )
```

All depth traversals are $\Theta(n)$ in time

- We need to pass through each node: $T(n) = \Omega(n)$
- By definition and recurrence of the traversals we find $T(n) = O(n)$

Breadth walks

- ★ visit the closest node to the root which was not yet visited
- ★ Using a queue in $\Theta(n)$



$\Rightarrow \langle F, B, H, A, D, K \rangle$

BREADTH-TREE-WALK(T)

```
1  Q = "Empty queue"
2  if not ISEMPTY( $T$ )
3      ENQUEUE( $Q$ , root( $T$ ))
4  while not QUEUE-EMPTY( $Q$ )
5       $y$  = DEQUEUE( $Q$ )
6      print GETDATA( $T$ ,  $y$ )
7      if HASLEFT( $T$ ,  $y$ )
8          ENQUEUE( $Q$ , LEFT( $y$ ))
9      if HASRIGHT( $T$ ,  $y$ )
10         ENQUEUE( $Q$ , RIGHT( $y$ ))
```

- 1 Introduction
- 2 File and Queue
- 3 List, Array, Sequence
- 4 Trees
- 5 Priority Queues**
- 6 Dictionaries

- Dynamical set of objects classified by a priority order
 - extract object with maximal priority
 - priority given by a key in a set having total order
- Interface:
 - $\text{INSERT}(S, x)$: insert element x in S
 - $\text{MAXIMUM}(S)$: return the element in S with largest key
 - $\text{EXTRACT-MAX}(S)$: delete the element of S having the largest key

- Static array
 - Q is a static array with fixed length
 - Elements in Q are sorted in increasing order according to their **keys**
 - Complexity of extraction: $O(1)$
 - Complexity of insertion: $O(n)$
 - Space complexity: $O(1)$
- Linked list
 - Q is a linked list, sorted according to its keys.
 - Complexity of extraction: $O(1)$
 - Complexity of insertion: $O(n)$
 - Complexity in space: $O(n)$

Using a heap

- the top priority element is on top of the heap: $O(1)$
- insertion/extraction costs $O(\log n)$

HEAP-MAXIMUM(A)

```
1  return  $A[1]$ 
```

HEAP-EXTRACT-MAX(A)

```
1  if  $A.heap-size < 1$   
2      error "heap underflow"  
3   $max = A[1]$   
4   $A[1] = A[A.heap-size]$   
5   $A.heap-size = A.heap-size - 1$   
6  MAX-HEAPIFY( $A, 1$ ) // reconstruit le tas  
7  return  $max$ 
```

- 1 Introduction
- 2 File and Queue
- 3 List, Array, Sequence
- 4 Trees
- 5 Priority Queues
- 6 Dictionaries**

Dictionaries are data dynamical data structures using comparable keys, supporting the following operations

- $\text{SEARCH}(S, k)$ returns a pointer x towards an element in S such that $x.\text{key} = k$ and *NONE* if k does not belong to S
- $\text{SEARCH}(S, x)$ inserts the element x in the dictionary S . If an element with the same key exists, the value is replaced
- $\text{DELETE}(S, x)$ removes the element x from S and does nothing if x is not in the dictionary.

We always assume that any two keys are comparable.

Two general objectives:

- minimize the insert and access cost
- minimize the storage cost

Many implementations are possible.

First solution:

- store the dictionary entries into a linked list
- to search an element, loop through the list
- Insertion (similarly for Deletion):
 - search the key
 - if found, then replace its value
 - if not found, then put it in the top of the list
- Complexity (worst case):
 - Insertion: $\Theta(N)$
 - Search: $\Theta(N)$
 - Delete: $\Theta(N)$

Second solution:

- store elements in a **vector** which we maintain sorted
- use binary search for searching: $\Theta(\log n)$
- Insertion: search the position and insert at given rank - shift all elements towards the right
- Deletion: delete and shift remaining elements towards the left
- Complexity:
 - Insertion $\Theta(N)$
 - Deletion $\Theta(N)$
 - Search $\Theta(\log N)$

Up to this point

Implementation	Worst case			Average		
	SEARCH	INSERT	DELETE	SEARCH	INSERT	DELETE
List	$\Theta(n)$	$\Theta(n)$	$\Theta(n)$	$\Theta(n)$	$\Theta(n)$	$\Theta(n)$
Sorted vector	$\Theta(\log n)$	$\Theta(n)$	$\Theta(n)$	$\Theta(\log n)$	$\Theta(n)$	$\Theta(n)$

Can we do better?

Remember: T is a tree, with a given root

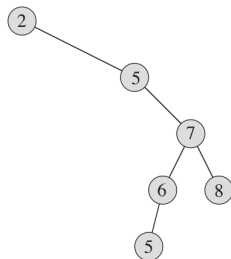
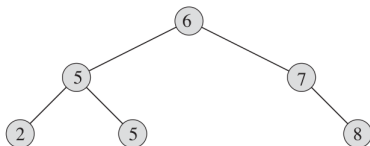
★ for every node x we have access to

- $x.parent$: parent
- $x.key$
- $x.left$: left child
- $x.right$: right child

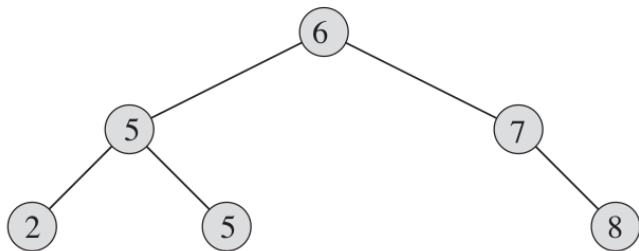
Binary search trees

A **binary search tree** is a binary tree implementing a dictionary with operation costing $O(h)$ where h is the height of the tree!

- every node has an associated key
- the tree verifies the following property for two nodes x, y
 - If y is in the left **sub-tree** at x then $y.key < x.key$
 - If y is in the right **sub-tree** at x then $y.key \geq x.key$



Going through a binary search tree



$\Rightarrow \langle 2, 5, 5, 6, 7, 8 \rangle$

The infix walk in a binary tree allows to see the keys in increasing order

★ Binary search

```
TREE-SEARCH( $x, k$ )  
1  if  $x == \text{NIL}$  or  $k == x.\text{key}$   
2      return  $x$   
3  if  $k < x.\text{key}$   
4      return TREE-SEARCH( $x.\text{left}, k$ )  
5  else return TREE-SEARCH( $x.\text{right}, k$ )
```

★ Complexity: $T(n) = O(h)$ where h is the height of the tree

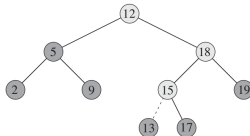
★ Worst case: $h = n$

Other operations

- ★ Minimal/maximal key: leftmost-rightmost node, complexity $O(h)$
- ★ Successor, Predecessor: minimum/maximum in the right/left sub-trees, complexity $O(h)$
- ★ Insertion: search the key $x.key$ in the tree, insert x where the search stopped: $O(h)$

TREE-INSERT(T, z)

```
1   $y = \text{NIL}$ 
2   $x = T.\text{root}$ 
3  while  $x \neq \text{NIL}$ 
4       $y = x$ 
5      if  $z.\text{key} < x.\text{key}$ 
6           $x = x.\text{left}$ 
7      else  $x = x.\text{right}$ 
8   $z.\text{parent} = y$ 
9  if  $y == \text{NIL}$ 
10     // Tree  $T$  was empty
11      $T.\text{root} = z$ 
12 elseif  $z.\text{key} < y.\text{key}$ 
13      $y.\text{left} = z$ 
14 else  $y.\text{right} = z$ 
```



- ★ Deletion: $O(h)$

- ★ all operations are $O(h)$ where h is the height of the tree
- ★ inserting keys in arbitrary order may lead to different heights
- ★ The average height of a tree obtained is $\Theta(\log n)$

Summary

Implementation	Worst case			Average		
	SEARCH	INSERT	DELETE	SEARCH	INSERT	DELETE
List	$\Theta(n)$	$\Theta(n)$	$\Theta(n)$	$\Theta(n)$	$\Theta(n)$	$\Theta(n)$
Sorted vector	$\Theta(\log n)$	$\Theta(n)$	$\Theta(n)$	$\Theta(\log n)$	$\Theta(n)$	$\Theta(n)$
Binary search tree	$\Theta(n)$	$\Theta(n)$	$\Theta(n)$	$\Theta(\log n)$	$\Theta(\log n)$	$\Theta(\log n)$

- ★ Can we get $O(\log n)$ in the worst case?
- ★ Yes: **Weight-Balanced Trees**: trees with guaranteed height of order $\Theta(\log n)$
- ★ Can we do better (on average) ? Yes: hash tables