

Efficient Data Structures

Data Structures and Algorithms (094224)

Yuval Emek

Winter 2022/23

Data structures revisited

- **Focus:** a dynamic set of objects
- Objects are typically organized in some **well structured digraph**
- A designated node attribute: **key** = the object's identifier
 - Keys are totally ordered
 - For simplicity, assume **uniqueness**
- Supporting (a subset of the) operations:
insert, delete, search, minimum, maximum, successor, predecessor
 - **FAST!!!**
- Actual data stored in satellite attributes
 - Sometimes a pointer to the actual data (storage considerations)

1 Binary search trees

- Binary search tree operations

2 2-3 trees

- 2-3 tree operations
 - Dynamic updates
- Discussion

3 Binary heaps

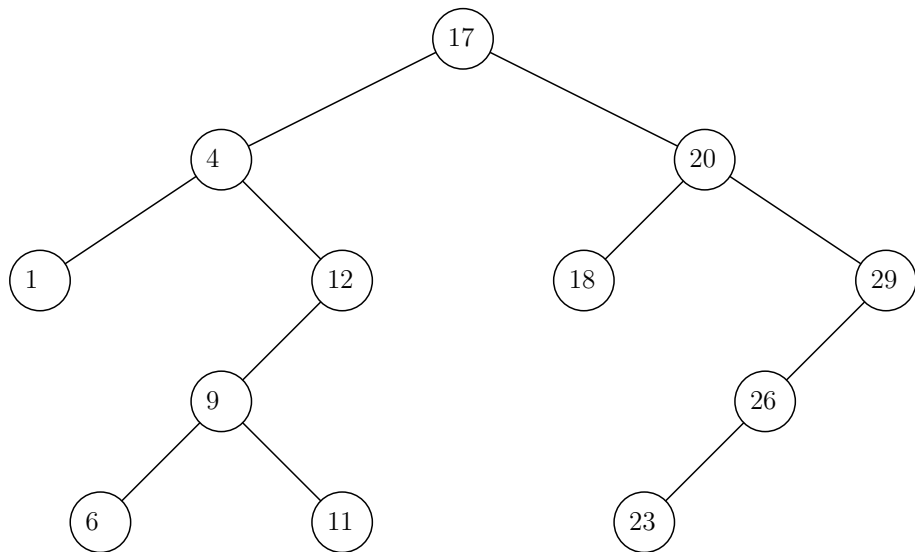
- Binary heap operations
- Discussion

A (rooted) binary tree organization

Binary search tree (עץ חיפוש בינארי)

- Objects associated with nodes of a rooted tree T with **degrees ≤ 2**
 - a.k.a. **binary tree**
- **Data structure's attribute:**
 - **root** = pointer to root of T
- **Node attributes (additional):**
 - **left** = pointer to left child in T
 - **right** = pointer to right child in T
 - **p** = pointer to parent in T
- **The binary search tree property:**
if y is a node in the left (resp., right) subtree of node x , then $y.key < x.key$ (resp., $y.key > x.key$)
 - Binary search tree is **not unique** for a given set of keys

Example



1 Binary search trees

- Binary search tree operations

2 2-3 trees

- 2-3 tree operations
 - Dynamic updates
- Discussion

3 Binary heaps

- Binary heap operations
- Discussion

Search in the binary search tree rooted at x for a node whose key is k

`Tree_Search(x, k)`

- 1: **if** $x == nil$ or $x.key == k$ **then**
- 2: return x
- 3: **if** $k < x.key$ **then**
- 4: return `Tree_Search($x.left, k$)`
- 5: **else** return `Tree_Search($x.right, k$)`

- Run-time: $O(\text{height}(T))$

Finding the minimum/maximum element

Find the node with the smallest key in the binary search tree rooted at x

`Tree_Minimum(x)`

```
1: while  $x.left \neq NIL$  do  
2:    $x = x.left$   
3: return  $x$ 
```

- Undefined if the tree is empty
- Run-time: $O(\text{height}(T))$
- How do we find the maximum?

Finding a successor/predecessor

Find the node y with the smallest key among those with $y.key > x.key$

`Tree_Successor(x)`

```
1: if  $x.right \neq NIL$  then  
2:   return Tree_Minimum(x.right)  
3:  $y = x.p$   
4: while  $y \neq NIL$  and  $x == y.right$  do  
5:    $x = y$   
6:    $y = y.p$   
7: return  $y$ 
```

- Run-time: $O(\text{height}(T))$
- How do we find a predecessor?

the opposite

Inserting a new node

Insert the new node z (DS attributes initialized to NIL) into T

$\text{Tree_Insert}(T, z)$

```
1: if  $T.\text{root} == NIL$  then
2:    $T.\text{root} = z$ 
3: else
4:    $y = T.\text{root}$ 
5:    $x = NIL$ 
6:   while  $y \neq NIL$  do
7:      $x = y$ 
8:     if  $z.\text{key} < y.\text{key}$  then
9:        $y = y.\text{left}$ 
10:    else  $y = y.\text{right}$ 
11:    $z.p = x$ 
12:   if  $z.\text{key} < x.\text{key}$  then
13:      $x.\text{left} = z$ 
14:   else  $x.\text{right} = z$ 
```

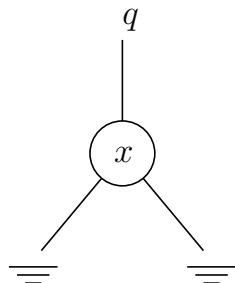
• Run-time: $O(\text{height}(T))$

Deleting a node

Deleting node x

case 1: x is a leaf

Not much to do



q

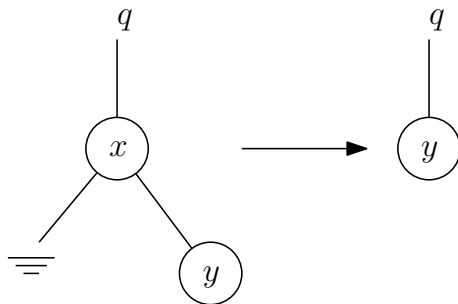
- Run-time: $O(1)$

Deleting a node — cont.

Deleting node x

case 2: x has a right child y and no left child

Replace x with y



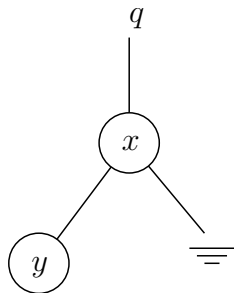
• Run-time: $O(1)$

Deleting a node — cont.

Deleting node x

case 3: x has a left child y and no right child

Replace x with y



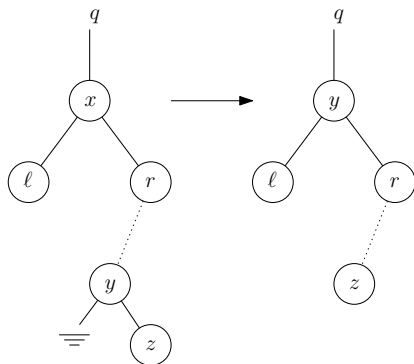
- Run-time: $O(1)$

Deleting a node — cont.

Deleting node x

case 4: x has two children

- 1 Find x 's successor y (in x 's right subtree)
- 2 Swap x and y
- 3 Remove x
 - We know how to do it because now x has ≤ 1 child



- Run-time: $O(\text{height}(T))$

Balanced trees

- **Supported operations:**
insert, delete, search, minimum, maximum, successor, predecessor
- Run-time of each operation in a binary search tree T is $O(\text{height}(T))$
- $\text{height}(T)$ can be $\Omega(n)$ in the worst case
 - How?
 - It is always $\Omega(\log n)$
 - #nodes increases from one level to the next by factor ≤ 2
- Looking for a data structure that supports these operations in time $O(\log n)$ in the **worst case**
 - Referred to as a **balanced tree** (עץ מאוזן)
- There exist balanced tree DSs based on binary search trees
 - AVL trees, red-black trees
- **In this lecture:** balanced tree DS that follows a different approach
- **Note:** the height of a binary search tree is $O(\log n)$ on average
 - What does “on average” mean for a DS?
 - Usually efficient in practice

- 1 Binary search trees
 - Binary search tree operations

- 2 2-3 trees
 - 2-3 tree operations
 - Dynamic updates
 - Discussion

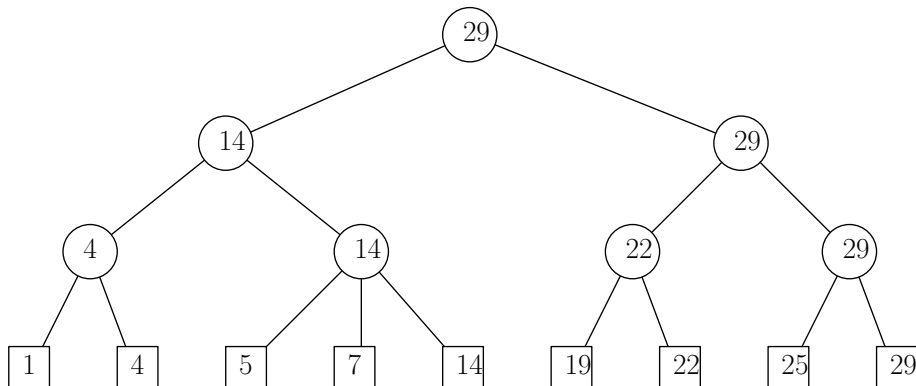
- 3 Binary heaps
 - Binary heap operations
 - Discussion

A different kind of rooted tree

2-3 tree $\langle 2-3 \text{ רע} \rangle$

- Rooted tree T satisfying:
 - Each internal node has **degree 2 or 3**
 - All leaves are at the **same level**
- Objects are stored **only** at the leaves
 - Internal nodes are there for organization purposes
- **Data structure's attribute:**
 - **root** = pointer to root of T
- **Node attributes:**
 - **left** = pointer to left child in T
 - **middle** = pointer to middle child in T
 - **right** = pointer to right child in T
 - **p** = pointer to parent in T
 - **key** = maximum key in its subtree
 - Leaves have additional satellite attributes (storing the actual data)
- When degree is 2, **left** and **middle** children are in use
- **The 2-3 tree property:**
keys in left subtree $<$ keys in middle subtree $<$ keys in right subtree

Example



Structural properties

- Consider a 2-3 tree T with n leaves
- **Observation 1:** in the leaf level, keys are sorted from left to right
 - By induction on height
- **Observation 2:** $\# \text{nodes in level } i + 1 \geq 2 \times \# \text{nodes in level } i$
 - Each internal node has ≥ 2 children
- **Corollary 1:** $\text{height}(T) = O(\log n)$
 - Will see: run-time of each supported operation is $O(\text{height}(T))$
- **Corollary 2:** $\# \text{internal nodes} < \# \text{leaves}$
 - By induction on height
 - Don't waste space (asymptotically)

- 1 Binary search trees
 - Binary search tree operations

- 2 2-3 trees
 - 2-3 tree operations
 - Dynamic updates
 - Discussion

- 3 Binary heaps
 - Binary heap operations
 - Discussion

Initialization

Create an (empty) 2-3 tree

`2-3_Init(T)`

1: new internal node x

2: new leaves ℓ, m

3: $\ell.key = -\infty$

4: $m.key = +\infty$

5: $\ell.p = m.p = x$

6: $x.key = +\infty$

7: $x.left = \ell$

8: $x.middle = m$

9: $T.root = x$

▷ DS attributes initialized to *NIL*

▷ DS attributes initialized to *NIL*

- Run-time: $O(1)$

- ℓ and m are called *sentinel nodes*

Search in the 2-3 tree rooted at x for a node whose key is k

$2_3_Search(x, k)$

```
1: if  $x$  is a leaf then  
2:   if  $x.key == k$  then  
3:     return  $x$   
4:   else return  $NIL$   
5: if  $k \leq x.left.key$  then  
6:   return  $2\_3\_Search(x.left, k)$   
7: else if  $k \leq x.middle.key$  then  
8:   return  $2\_3\_Search(x.middle, k)$   
9: else return  $2\_3\_Search(x.right, k)$ 
```

- $+\infty$ sentinel node ensures existence of right child in line 9

Finding the minimum/maximum element

Find the leaf with the smallest key in the 2-3 tree T

$2_3_Minimum(T)$

```
1:  $x = T.root$ 
2: while  $x$  is not a leaf do
3:    $x = x.left$ 
4:  $x = x.p.middle$ 
5: if  $x.key \neq +\infty$  then
6:   return  $x$ 
7: else error:  $T$  is empty
```

→ sentinal node פ"י ופ"י
final node $(-\infty)$ ופ"י

- How do we find the maximum?

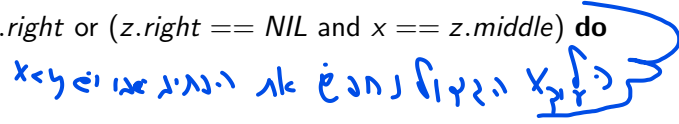
הי"ו -> זהר"ם אם הי"ו..

Finding a successor/predecessor

Find the leaf y with the smallest key among those with $y.key > x.key$

2.3_Successor(x)

```
1:  $z = x.p$ 
2: while  $x == z.right$  or ( $z.right == NIL$  and  $x == z.middle$ ) do
3:    $x = z$ 
4:    $z = z.p$ 
5: if  $x == z.left$  then
6:    $y = z.middle$ 
7: else  $y = z.right$ 
8: while  $y$  is not a leaf do
9:    $y = y.left$ 
10: if  $y.key < +\infty$  then
11:   return  $y$ 
12: else return  $NIL$ 
```



- How do we find a predecessor?

- 1 Binary search trees
 - Binary search tree operations

- 2 2-3 trees
 - 2-3 tree operations
 - Dynamic updates
 - Discussion

- 3 Binary heaps
 - Binary heap operations
 - Discussion

Procedure Update_Key

Update the key of x to the maximum key in its subtree;
(only) $x.middle$ and $x.right$ may be NIL

Update_Key(x)

- 1: $x.key = x.left.key$
- 2: **if** $x.middle \neq NIL$ **then**
- 3: $x.key = x.middle.key$
- 4: **if** $x.right \neq NIL$ **then**
- 5: $x.key = x.right.key$

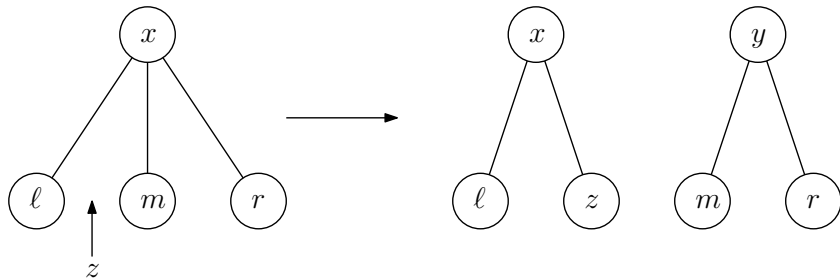
Procedure Set_Children

Set ℓ , m , and r to be the left, middle, and right children, respectively, of x ;
(only) m and r may be NIL

Set_Children(x, ℓ, m, r)

- 1: $\langle x.left, x.middle, x.right \rangle = \langle \ell, m, r \rangle$
- 2: $\ell.p = x$
- 3: **if** $m \neq NIL$ **then**
- 4: $m.p = x$
- 5: **if** $r \neq NIL$ **then**
- 6: $r.p = x$
- 7: **Update_Key**(x)

New leaf insertion



New leaf insertion — procedure `Insert_And_Split`

Insert node z as a child of node x ;
split x if necessary and return the new node

`Insert_And_Split(x, z)`

```
1:  $\langle \ell, m, r \rangle = \langle x.left, x.middle, x.right \rangle$ 
2: if  $r == NIL$  then
3:   if  $z.key < \ell.key$  then
4:     Set_Children( $x, z, \ell, m$ )
5:   else if  $z.key < m.key$  then
6:     Set_Children( $x, \ell, z, m$ )
7:   else Set_Children( $x, \ell, m, z$ )
8:   return  $NIL$ 
9: new internal node  $y$ 
10: ...
```

New leaf insertion — procedure `Insert_And_Split`

`Insert_And_Split(x, z)` — cont.

```
10: if  $z.key < \ell.key$  then  
11:   Set_Children(x, z,  $\ell$ , NIL)  
12:   Set_Children(y, m, r, NIL)  
13: else if  $z.key < m.key$  then  
14:   Set_Children(x,  $\ell$ , z, NIL)  
15:   Set_Children(y, m, r, NIL)  
16: else if  $z.key < r.key$  then  
17:   Set_Children(x,  $\ell$ , m, NIL)  
18:   Set_Children(y, z, r, NIL)  
19: else Set_Children(x,  $\ell$ , m, NIL)  
20:   Set_Children(y, r, z, NIL)  
21: return  $y$ 
```

Inserting a new leaf

Insert the new leaf z (DS attributes initialized to NIL) into T

2.3_Insert(T, z)

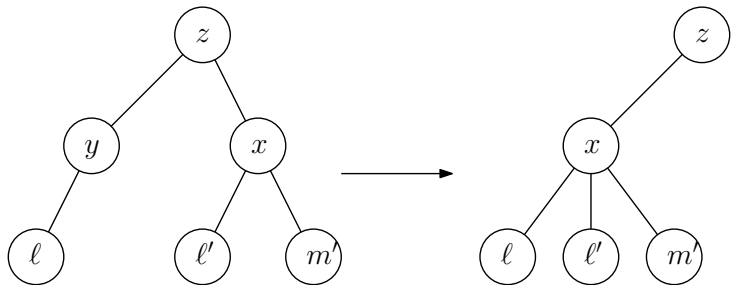
```
1:  $y = T.root$ 
2: while  $y$  is not a leaf do
3:   if  $z.key < y.left.key$  then  $y = y.left$ 
4:   else if  $z.key < y.middle.key$  then  $y = y.middle$ 
5:   else  $y = y.right$ 
6:  $x = y.p$ 
7:  $z = \text{Insert\_And\_Split}(x, z)$ 
8: while  $x \neq T.root$  do
9:    $x = x.p$ 
10:  if  $z \neq NIL$  then
11:     $z = \text{Insert\_And\_Split}(x, z)$ 
12:  else  $\text{Update\_Key}(x)$ 
13: ...
```

Inserting a new leaf

2_3_Insert(T, z) — cont.

```
13: if  $z \neq NIL$  then  
14:   new internal node  $w$   
15:   Set_Children( $w, x, z, NIL$ )  
16:    $T.root = w$ 
```


Leaf deletion



Leaf deletion — procedure Borrow_Or_Merge

Borrow a child from a sibling x of y or merge x and y ;
return a pointer to the parent of y (and x)

Borrow_Or_Merge(y)

```
1:  $z = y.p$ 
2: if  $y == z.left$  then
3:    $x = z.middle$ 
4:   if  $x.right \neq NIL$  then
5:     Set_Children( $y, y.left, x.left, NIL$ )
6:     Set_Children( $x, x.middle, x.right, NIL$ )
7:   else Set_Children( $x, y.left, x.left, x.middle$ )
8:     delete  $y$ 
9:     Set_Children( $z, x, z.right, NIL$ )
10:  return  $z$ 
11: ...
```

Leaf deletion — procedure Borrow_Or_Merge

Borrow_Or_Merge(*y*) — cont.

```
11: if y == z.middle then  
12:   x = z.left  
13:   if x.right ≠ NIL then  
14:     Set_Children(y, x.right, y.left, NIL)  
15:     Set_Children(x, x.left, x.middle, NIL)  
16:   else Set_Children(x, x.left, x.middle, y.left)  
17:     delete y  
18:     Set_Children(z, x, z.right, NIL)  
19:   return z  
20: ...
```

Leaf deletion — procedure Borrow_Or_Merge

Borrow_Or_Merge(*y*) — cont.

20: $x = z.middle$

21: **if** $x.right \neq NIL$ **then**

22: Set_Children($y, x.right, y.left, NIL$)

23: Set_Children($x, x.left, x.middle, NIL$)

24: **else** Set_Children($x, x.left, x.middle, y.left$)

25: delete y

26: Set_Children($z, z.left, x, NIL$)

27: return z

Deleting a leaf

Delete leaf x from T

`2_3_Delete(T, x)`

```
1:  $y = x.p$ 
2: if  $x == y.left$  then
3:   Set_Children( $y, y.middle, y.right, NIL$ )
4: else if  $x == y.middle$  then
5:   Set_Children( $y, y.left, y.right, NIL$ )
6: else Set_Children( $y, y.left, y.middle, NIL$ )
7: delete  $x$ 
8: ...
```

▷ $\deg(y)$ may be < 2

Deleting a leaf

2.3_Delete(T, x) — cont.

```
8: while  $y \neq NIL$  do
9:   if  $y.middle == NIL$  then
10:    if  $y \neq T.root$  then
11:       $y = \text{Borrow\_Or\_Merge}(y)$ 
12:    else  $T.root = y.left$ 
13:       $y.left.p = NIL$ 
14:      delete  $y$ 
15:      return
16:    else  $\text{Update\_Key}(y)$ 
17:       $y = y.p$ 
```

- 1 Binary search trees
 - Binary search tree operations

- 2 2-3 trees
 - 2-3 tree operations
 - Dynamic updates
 - Discussion

- 3 Binary heaps
 - Binary heap operations
 - Discussion

More on 2-3 trees

- All operations run in time $O(\text{height}(T)) = O(\log n)$
 - A balanced tree DS
- Sometimes the actual keys are stored only at the leaves
 - Internal nodes store **pointers** to keys in the corresponding leaves
 - Saving storage if keys themselves are large
- 2-3 trees are a special case of **B^+ trees**
 - Parameterized by a **degree parameter d**
 - The degree of each internal node x satisfies $\lceil d/2 \rceil \leq \deg(x) \leq d$
 - The degree of the root r satisfies $2 \leq \deg(r) \leq d$
 - Height proportional to $\log_d n$
- **Often:** x stores maximum key in **each child subtree**
 - $\deg(x)$ keys in total
- Very efficient for **external memories**
 - Access is very expensive in comparison to internal memory
 - Read/write a whole **page** rather than a single address
 - Adjust the parameter d so that size of a node \approx size of page
 - Save on #page read/write operations

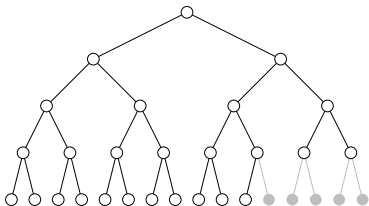
- 1 Binary search trees
 - Binary search tree operations

- 2 2-3 trees
 - 2-3 tree operations
 - Dynamic updates
 - Discussion

- 3 Binary heaps
 - Binary heap operations
 - Discussion

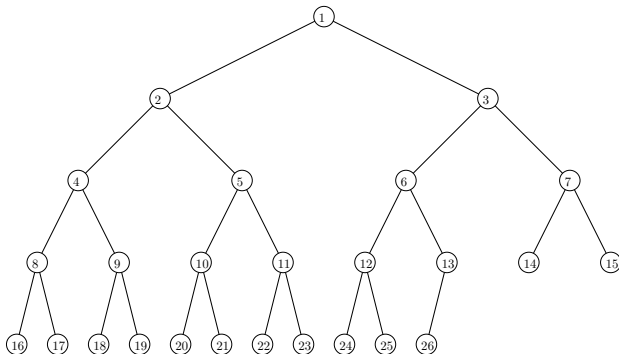
Binary trees revisited

- A **complete** (שלם) binary tree:
 - the degree of each internal node is exactly 2; and
 - all leaves are at the same depth
- Height h complete binary tree has 2^h leaves and $2^h - 1$ internal nodes
- A **nearly complete** (כמעט שלם) binary tree:
 - obtained from a height h complete binary tree by deleting the k **rightmost** leaves for some $0 \leq k < 2^h$
- **Properties of nearly complete binary trees:**
 - A generalization of complete binary trees
 - All leaves are in the last two levels
 - $\# \text{internal nodes} \leq \# \text{leaves} \leq \# \text{internal nodes} + 1$



Array representation

- Representing an n -node nearly complete binary tree in **array** $A[1 \dots n]$
 - The root is identified with $A[1]$
 - Left child of internal node $A[i]$ is identified with $A[2i]$
 - Right child of internal node $A[i]$ is identified with $A[2i + 1]$
 - Parent of non-root node $A[i]$ is identified with $A[\lfloor i/2 \rfloor]$
 - Define $\text{Left}(i) = 2i$, $\text{Right}(i) = 2i + 1$, $\text{Parent}(i) = \lfloor i/2 \rfloor$



A nearly complete binary tree data structure

Binary heap (ערמה בינארית)

(sometimes simply heap)

- Objects associated with the nodes of a nearly complete binary tree T
- T represented in (prefix of) array $A[1 \dots A.length]$
 - $A.length$ is an upper bound on #objects in the dynamic set
 - Known in advance
 - Exists in many applications
 - n objects stored in $A[1 \dots n]$
- **Data structure's attribute** (on top of the standard array's attributes):
 - $heap\text{-}size = \# \text{currently stored objects } (n \leq A.length)$
- **The binary heap property:**
 $A[\text{Parent}(i)].key < A[i].key$ for every non-root $A[i]$ each almost complete binary tree that have this is a binary heap
- Sometimes referred to as a *priority queue* (תור עדיפויות) or *minimum heap* (ערמת מינימום)
 - Switching the relevant operations/inequalities yields a *maximum heap* (ערמת מקסימום)

- 1 Binary search trees
 - Binary search tree operations

- 2 2-3 trees
 - 2-3 tree operations
 - Dynamic updates
 - Discussion

- 3 Binary heaps
 - Binary heap operations
 - Discussion

Building a heap — Procedure Heapify

Precondition: subtrees rooted at $A[\text{Left}(i)]$ and $A[\text{Right}(i)]$ are heaps;
modify A so that the subtree rooted at $A[i]$ is a heap

Heapify(A, i)

```
1:  $\ell = \text{Left}(i)$ 
2: if  $\ell \leq A.\text{heap-size}$  and  $A[\ell].\text{key} < A[i].\text{key}$  then
3:    $\text{smallest} = \ell$ 
4: else  $\text{smallest} = i$ 
5:  $r = \text{Right}(i)$ 
6: if  $r \leq A.\text{heap-size}$  and  $A[r].\text{key} < A[\text{smallest}].\text{key}$  then
7:    $\text{smallest} = r$ 
8: if  $\text{smallest} \neq i$  then
9:   swap  $A[i]$  and  $A[\text{smallest}]$ 
10:  Heapify( $A, \text{smallest}$ )
```

- $A[i]$ **seeps down** the tree
- Run-time: $O(\text{height}(T_i)) = O(\log n)$
 - T_i = subtree rooted at $A[i]$

Building a heap from an arbitrary array

Build a heap from the objects in $A[1 \dots A.length]$

Build_Heap(A)

- 1: $A.heap\text{-}size = A.length$
- 2: **for** $i = A.length, \dots, 2, 1$ **do**
- 3: *Heapify*(A, i)

- Naive run-time analysis:
 - n iterations
 - $O(\log n)$ time per iteration
 - $\implies O(n \log n)$
- Can do better

Building a heap — run-time analysis

- Recall: run-time of $\text{Heapify}(A, i)$ is $O(\text{height}(T_i))$
- Height of a nearly complete binary tree T with n nodes is $\lfloor \lg n \rfloor$
- #nodes of height j in $T \leq 2^{\lfloor \lg n \rfloor - j} = O(n/2^j)$
- \implies run-time of Build_Heap is

$$O\left(\sum_{j=0}^{\lfloor \lg n \rfloor} j \cdot \frac{n}{2^j}\right) \leq O(n) \cdot \sum_{j=0}^{\infty} \frac{j}{2^j} = O(n)$$

Extracting the minimum

Delete the minimum object from heap A and return it

$\text{Heap_Extract_Min}(A)$

- 1: **if** $A.\text{heap-size} < 1$ **then**
- 2: error “the heap is empty”
- 3: $\text{min} = A[1]$
- 4: $A[1] = A[A.\text{heap-size}]$
- 5: $A.\text{heap-size} = A.\text{heap-size} - 1$
- 6: $\text{Heapify}(A, 1)$
- 7: **return** min

- Run-time: $O(\log n)$
- What's the run-time of returning (a pointer to) the minimum object without deleting it? $O(1)$

Decreasing a key

Precondition: $k < A[i].key$;

change the key of $A[i]$ to k

$\text{Heap_Decrease_Key}(A, i, k)$

- 1: **if** $k > A[i].key$ **then**
- 2: error “new key is larger than current key”
- 3: $A[i].key = k$
- 4: **while** $i > 1$ and $A[i].key < A[\text{Parent}(i)].key$ **do**
- 5: swap $A[i]$ and $A[\text{Parent}(i)]$
- 6: $i = \text{Parent}(i)$

- The updated $A[i]$ **seeps up** the tree
- Run-time: $O(\log n)$

Inserting a new node

Insert the new heap node x into heap A

$\text{Heap_Insert}(A, x)$

1: $s = A.\text{heap-size} + 1$

2: $A[s] = x$

3: $A[s].\text{key} = \infty$

4: $A.\text{heap-size} = s$

5: $\text{Heap_Decrease_Key}(A, s, x.\text{key})$

• Run-time: $O(\log n)$

▷ copy all satellite attributes

▷ a valid heap of size s

- 1 Binary search trees
 - Binary search tree operations

- 2 2-3 trees
 - 2-3 tree operations
 - Dynamic updates
 - Discussion

- 3 Binary heaps
 - Binary heap operations
 - Discussion

Comparison to balanced trees

- **Balanced trees** can be adapted to support:
 - Returning a pointer to the minimum object in time $O(1)$
 - Extracting the minimum object in time $O(\log n)$
 - Decreasing the key of a given node in time $O(\log n)$
 - How?
- **Advantages of heaps over balanced trees:**
 - Building a new heap with n objects is (asymptotically) faster
 - Speed-up (non-asymptotic) due to smaller constants
 - Speed-up (non-asymptotic) due to array representation
 - No pointers, whole DS stored in the same memory page