

Data structures are

FUNDAMENTAL!

- All fields of CS involve storing, retrieving and processing data

- Information retrieval
- Geographic Inf. Systems
- Machine Learning
- Text/String processing
- Computer graphics
- ...

Basic elements in study of data structures

- **Modeling**: How real-world objects are encoded
- **Operations**: Allowed functions to access + modify structure
- **Representation**: Mapping to memory
- **Algorithms**: How are ops. performed?

Course Overview:

- Fundamental data structures + algorithms
- Mathematical techniques for analyzing them
- Implementation

Introduction to Data Structures

- Elements of data structures
- Our approach
- Short review of asymptotics

Our approach:

- **Theoretical**: Algorithms + Asymptotic Analysis
- **Practical**: Implementation + practical efficiency

Common:

- $O(1)$: **constant time** 😊
[Hash map]
- $O(\log n)$: **log time** (very good!)
[Binary search]
- $O(n^p)$: ($p = \text{constant}$) **Poly time**
e.g. $O(\sqrt{n})$ [Geometric search]

Asymptotic: "Big-O"

- Ignore constants
- Focus on large n

$$T(n) = 34n^2 + 15n \cdot \log n + 143$$

$$T(n) = O(n^2)$$

Asymptotic Analysis:

- Run time as a function of $n \leftarrow$ no. of items
- Worst-case, average-case, randomized
- **Amortized**: Average over a series of ops.

Linear List ADT:

Stores a sequence of elements $\langle a_1, a_2, \dots, a_n \rangle$. Operations:

init() - create an empty list

get(i) - returns a_i

set(i, x) - sets i^{th} element to x

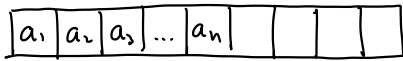
insert(i, x) - inserts x prior to i^{th} (moving others back)

delete(i) - deletes i^{th} item (moving others up)

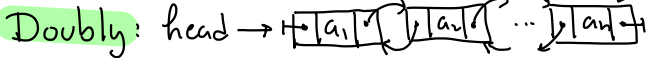
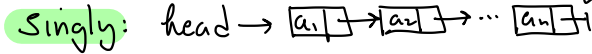
length() - returns num. of items

Implementations:

Sequential: Store items in an array



Linked allocation: linked list



Performance varies with implementation

Abstract Data Type (ADT)

- Abstracts the functional elements of a data structure (math) from its implementation (algorithm/programming)

Basic Data Structures I

- ADTs
- Lists, Stacks, Queues
- Sequential Allocation

Doubling Reallocation:

When array of size n overflows

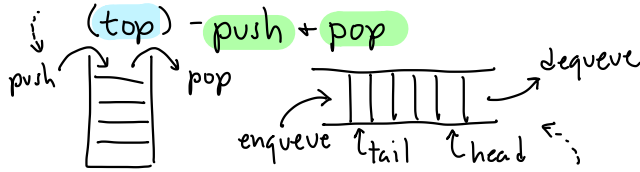
- allocate new array size $2n$
- copy old to new
- remove old array

Dynamic Lists + Sequential Allocation

Allocation: What to do when your array runs out of space?

Deque ("deck"): Can insert or delete from either end

Stack: All access from one side

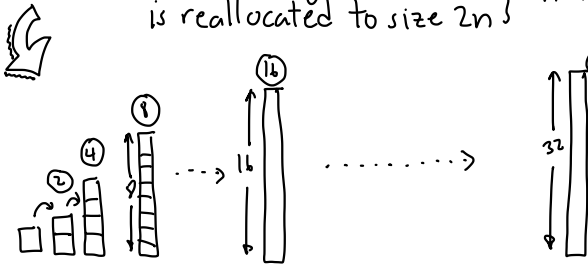


Queue: FIFO list: **enqueue** inserts at **tail** and **dequeue** deletes from **head**

Cost model (Actual cost)

Cheap: No reallocation \rightarrow 1 unit

Expensive: Array of size n is reallocated to size $2n$



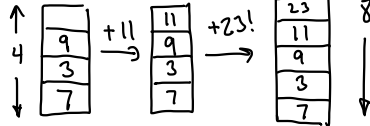
1 2 3 4 5 6 7 8 9 16 17
 $+1 +1 +1 +1 +1 +1 +1 +1 +1 +1 +1 +1 +1 +1 +1$

Total = $17 + (2+4+8+16+32) = 79$

Dynamic (Sequential) Allocation

- When we overflow, double

Eg. Stack



Basic Data Structures II

- Amortized analysis of dynamic stack

Amortized Cost: Starting from an empty structure, suppose that any sequence of m ops takes time $T(m)$. The **amortized cost** is $T(m)/m$.

Thm: Starting from an empty stack, the amortized cost of our stack operations is at most 5.

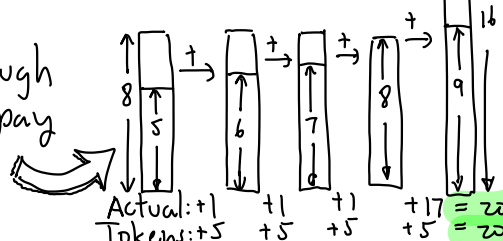
[i.e. any seq. of m ops has cost $\leq 5 \cdot m$]

Charging Argument:

- Each request of push/pop we charge user 5 work tokens
- We use 1 token to pay for the operation + put other 4 in bank account.
- Will show there is enough in bank account to pay actual costs.

Proof:

- Break the full sequence after each reallocation \rightarrow run **1 2 | 3 | 4 5 | 6 7 8 9 | 10 11 ... 16 | 17**
- At start of a run there are $n+1$ items in stack and array size is $2n$
- There are at least n ops before the end of run
- During this time we collect at least $5n$ tokens
 - \rightarrow 1 for each op
 - \rightarrow 4 for deposit
- Next reallocation costs $4n$, but we have enough saved! \square



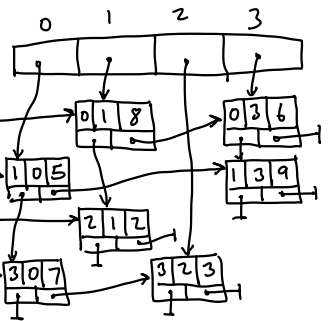
Fixed Increment: Increase by a fixed constant
 $n \rightarrow n + 100$

Fixed factor: Increase by a fixed constant factor (not nec. 2)
 $n \rightarrow 5 \cdot n$

Squaring: Square the size (or some other power)
 $n \rightarrow n^2$ or $n \rightarrow \lceil n^{1.5} \rceil$

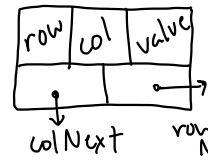
Dynamic Stack:
 - Showed doubling \Rightarrow Amortized $O(1)$
 - Other strategies?

0	8	0	6
5	0	0	9
0	2	0	0
7	0	3	0



Basic Data Structures III
 - Dynamic Stack - Wrap-up
 - Multilists + Sparse Matrices

Node:

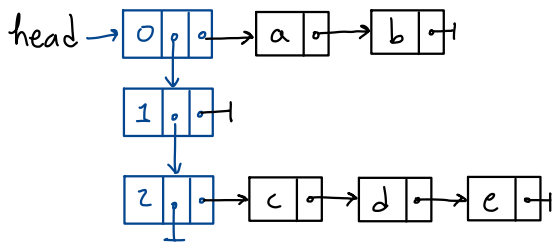


Idea: Store only non-zero entries linked by row and column

Which of these provide $O(1)$ amortized cost per operation?

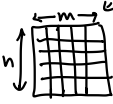
Leave as exercise ☹️
 (spoiler alert!)
 Fixed increment \rightarrow no
 Fixed factor \rightarrow yes
 Squaring \rightarrow ?? (depends on cost model)

Multilists: Lists of lists



Sparse Matrices:

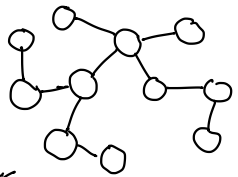
An $n \times m$ matrix has $n \cdot m$ entries and takes (naively) $O(n \cdot m)$ space



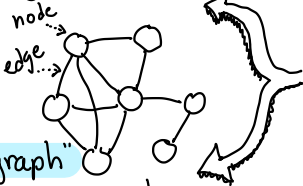
Sparse matrix: Most entries are zero

Tree (or "Free Tree")

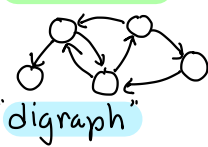
- undirected
- connected
- acyclic graph



Undirected



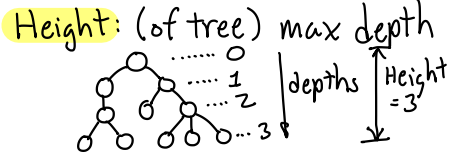
Directed



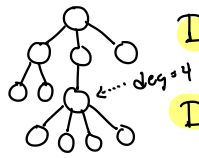
Graph: $G=(V,E)$
 V = finite set of **vertices** (nodes)
 E = set of **edges** (pairs of vertices)

Trees: Basic Concepts and Definitions

Depth: path length from root

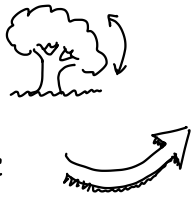
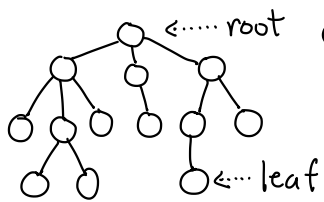


Degree (of node): number of children



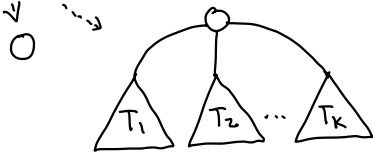
Degree (of tree): max. degree of any node

Rooted tree: A free tree with **root node**

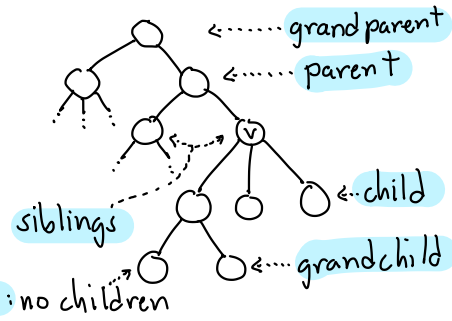


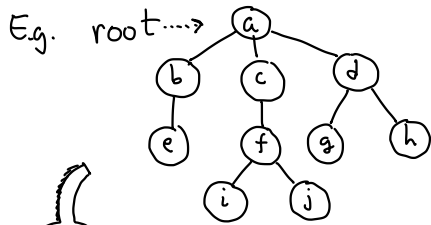
Formal definition:

Rooted tree: is either
 - single node (root)
 - set of one or more rooted trees (subtrees) joined to a common root



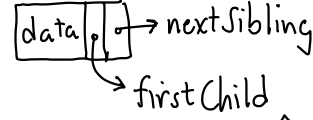
"Family" Relations





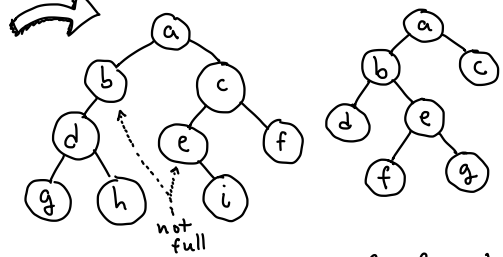
Representing rooted trees:
Each node stores a (linked) list of its children

Node structure:



Trees Representation + Binary Trees (I)

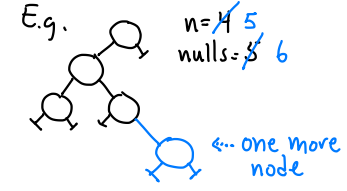
(Not full) Full:



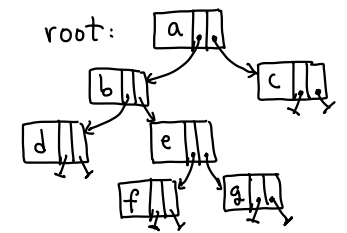
Full: Every non-leaf node has 2 children

Wasted space?

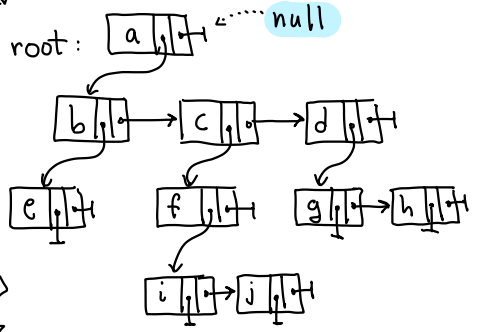
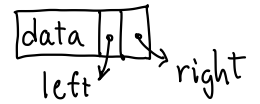
Theorem: A binary tree with n nodes has $n+1$ null links



In Java: class BTNode<E> {
E data;
BTNode<E> left;
BTNode<E> right;
....
}



Node structure:



called the **Binary representation**



Binary tree: A rooted tree of degree 2, where each node has two children (possibly null) **left + right**



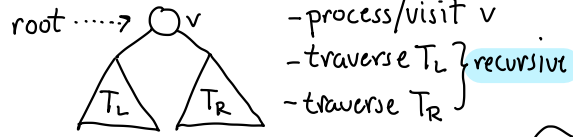
```

traverse(BTNode v) {
  if (v == null) return;
  visit/process v ← Preorder
  traverse (v.left)
  visit/process v ← Inorder
  traverse (v.right)
  visit/process v ← Postorder
}

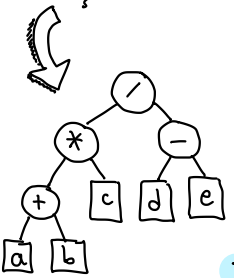
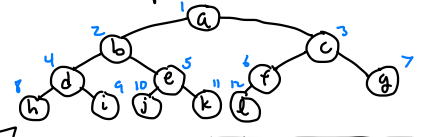
```

Traversals: How to (systematically) visit the nodes of a rooted tree?

Binary Tree Traversals (can be generalized)



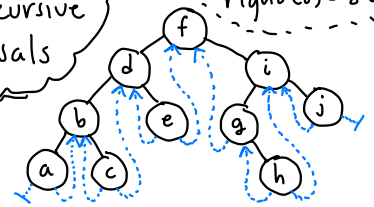
Complete Binary Tree: All levels full (except last)



Preorder: / * + a b c - d e
Postorder: a b + c * d e - /
Inorder: (a + b) * c / d - e

Binary Trees:
Traversals, Extension,
and More

Challenge:
Nonrecursive
traversals

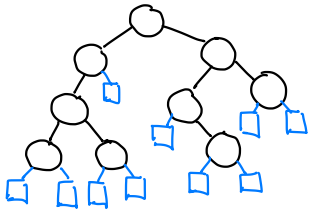


Those wasteful null links...

Thm: An extended binary tree with n internal nodes (black) has $n+1$ external nodes (blue)

Another way to save space...
Threaded binary tree:

Extended binary tree: Replace each null link with a special leaf node: **external node**



Observation: Every extended binary tree is full

Store (useful) links in the null links. (Use a mark bit to distinguish link types.)
 Eg. **Inorder Threads:**
 Null left → inorder predecessor
 Null right → " successor

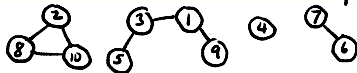
Examples:

- Given prime p , $a \equiv b \pmod p$

Eg. $p=5$; Partition: $\{0,5,10,\dots\}, \{1,6,11,\dots\}, \{2,7,12,\dots\}, \dots$

- Given graph G , vertices u, v ,

$u \equiv v$ if in same connected component



Partition: $\{2,8,10\}, \{1,3,5,9\}, \{4\}, \{6,7\}$

Union-Find:

Given set $S = \{1, 2, \dots, n\}$ maintain a partition supporting ops:

Init: Each element in its own set $\{1\}, \{2\}, \dots, \{n\}$

Union(s, t): Merge two sets s & t , and replace with their union

Find(x): Return the set containing x

Example: Suppose: $S_1 = \{1,5\}, S_2 = \{2,6,8\}, S_3 = \{3,4,7\}$

Union(S_1, S_2) $\rightarrow \{1,3,4,5,7\}, \{2,6,8\}$

Find(5) $\rightarrow S_1$, Find(8) $\rightarrow S_2$

Equivalence Relation:

Binary relation over set S such that $\forall a, b, c \in S$:

reflexive: $a \equiv a$

symmetric: $a \equiv b \Rightarrow b \equiv a$

transitive: $a \equiv b \wedge b \equiv c \Rightarrow a \equiv c$

Any equivalence relation defines a partition over S .

Disjoint Set Union-Find I

A simple approach to finding is to trace the path to the root -

Set = Element = int $1 \leq x \leq n$

```

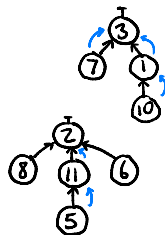
Set Simple-Find(Element x) {
    while (parent[x] != null)
        x ← parent[x]
    return x
}
    
```

Set Identifiers are indices of root nodes

Eg. Find(7) = 3

Find(10) = 3

Find(5) = 2



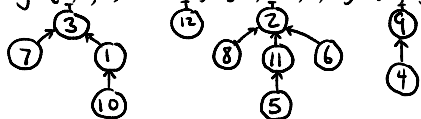
Two items in same set iff Find(x) = Find(y)

Inverted-Tree Approach:

- Store elements of each set in tree with links to parent

- Root node is set identifier

Eg. $\{1,3,7,10\}$ $\{12\}$ $\{2,5,6,8,11\}$ $\{4,9\}$



Array-Based Implementation:

Assume: $S = \{1, 2, \dots, n\}$

parent[1..n], where parent[i] is its parent index or 0 = null if root

1	2	3	4	5	6	7	8	9	10	11	12
3	0	0	9	11	2	3	2	0	1	2	0



Set Union (Set s, Set t) {

```

if (rank[s] > rank[t])
    swap s + t
parent[s] ← t
rank[t] ← max(rank[t], 1 + rank[s])
return t
    
```

Recall: These are just array indices of roots

How to Union?

- Just link one tree under the other
- How to maintain low heights?
- Rank: Based on height of tree. Link lower rank as child



Lemma: Assuming rank-based merging a tree of height h has at least 2^h nodes.

Proof: By induction on num. of unions
 Basis: Single node. $h=0$, $2^0=1$ nodes
 Step: Consider the last of series of unions. Let $T' + T''$ be trees to merge: Heights: $h' + h''$
 Sizes: $n' + n''$
 By induction: $n' \geq 2^{h'}$ $n'' \geq 2^{h''}$

Disjoint Set Union-Find II

Running Time?

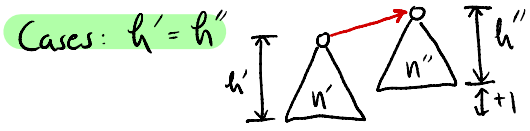
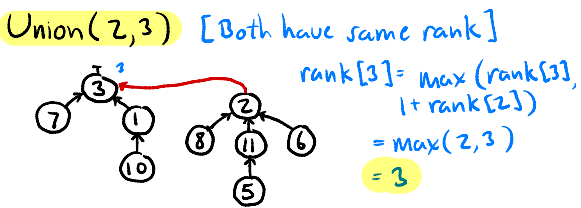
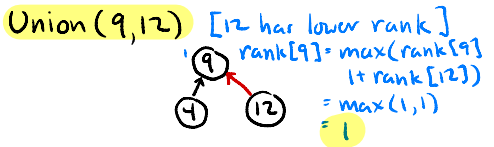
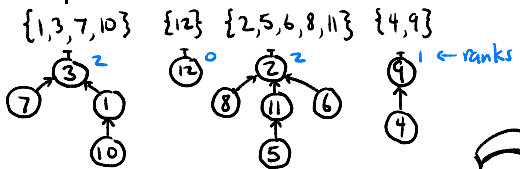
- Init: $O(n)$ - set a parents to null + ranks to 0
- Union: $O(1)$ - constant time
- Find: $O(\text{tree height})$

We'll show tree height = $O(\log n)$
 \Rightarrow Find takes $O(\log n)$ time

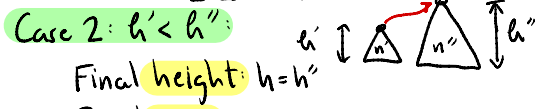
What is worst case?

Init: All ranks $\leftarrow 0$

Example:



Final tree height: $h = h' + 1 = h'' + 1$
 Final size: $n = n' + n'' \geq 2^{h'} + 2^{h''} = 2^{h-1} + 2^{h-1} = 2 \cdot 2^{h-1} = 2^h$



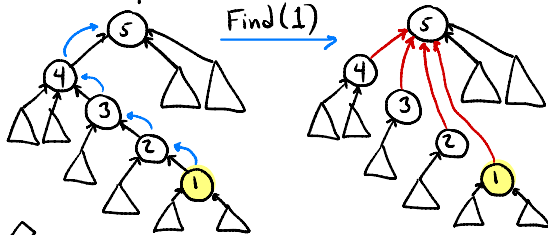
Final height: $h = h''$
 Final size: $n = n' + n'' \geq 2^{h'} + 2^{h''} \geq 2^{h''} = 2^h$

Case 3: $h' > h''$ (symmetrical) \square

Path Compression:

- Whenever we perform find, "short-cut" the links so they point directly to root
- This does not increase running by more than constant, but can speed up later finds

Example:



Does this little trick improve running times?

- Worst case - No. Find may take $O(\log n)$ time
- Amortized - Yes! Huge improvement! (But hard to prove)



Simple Union-Find performs a sequence of m Unions + Finds on set of size n in $O(m \log m)$ time.
 ⇒ Amortized time (average per op) is $O(\log m)$
 - Not bad - But can we do better?

Disjoint Set Union - Find III

Digression: Ackerman's Function (1926)

for $i, j \geq 0$

$$A(i, j) = \begin{cases} j+1 & \text{if } i=0 \\ A(i-1, 1) & \text{if } i>0, j=0 \\ A(i-1, A(i, j-1)) & \text{o.w.} \end{cases}$$

Looks innocent, but it's a monster!

Theorem: (Tarjan 1975) After init. any seq of m Union-Finds (with path compression) takes total time $O(m \cdot \alpha(m, n))$.
 ⇒ Amortized time = $O(\alpha(m, n))$

[For all practical purposes, this is constant time.]

From super big to super small
 Inverse of Ackerman:

$$\alpha(m, n) = \min \{ i \geq 1 \mid A(i, \lfloor L^{m/n} \rfloor) > \log m \}$$

Obs: $\alpha(m, n) \leq 4$ for any imaginable values of m, n ($m \geq n$)

i	$j=0$	$j=1$	$j=2$	$j=3$	$j=4$	$j=5$...
0	1	2	3	4	5	6	$A(0, j) = j+1$
1	2	3	4	5	6	7	$A(1, j) = j+2$
2	3	5	7	9	11	13	$A(2, j) = 2j+3$
3	5	13	29	...			$A(3, j) = 2^{j+3} - 3$
4	13	814	REAL BIG!	...			$A(4, j) = 2^{2^{j+3}} - 3$
...							

More than atoms in universe

Naive Solution:

- Store items in linear list
- Order?

Insert order -

fast insert / slow extract

Priority order -

fast extract / slow insert

Heap: Tree-based structure

(min) heap order: for all nodes, parent's key \leq node's key

[Reverse: max-heap order]

Many variants:

Binary, leftist, binomial, Fibonacci, pairing, quake, skew... heaps

Binary Heap:

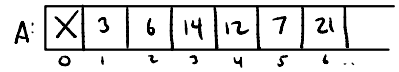
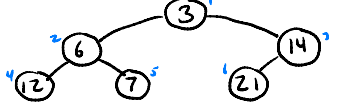
- simple, elegant, efficient
- old (1964)
- basic: insert/extract - $O(\log n)$
build - $O(n)$

Priority Queue:

- Stores key-value pairs
- Key \equiv priority
- Ops: insert(x, v) - insert value v with key x
- extract-min - remove/return pair with min key value



Pointerless tree



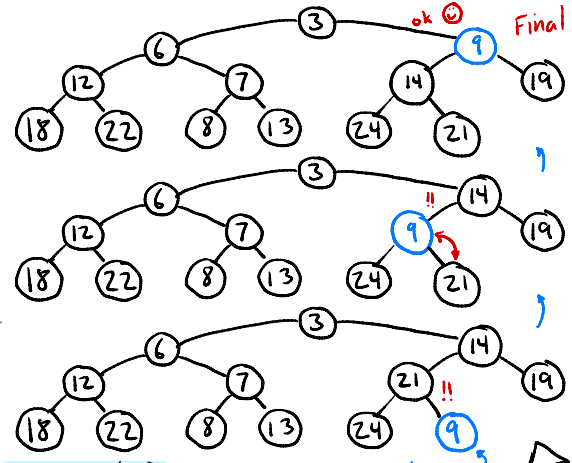
- left(i): if $(2 \cdot i \leq n)$ $2i$ else null
- right(i): if $(2 \cdot i + 1 \leq n)$ $2i + 1$ else null
- par(i): if $(i \geq 2)$ $\lfloor i/2 \rfloor$ else null

void insert(Key x) (ignore value)

```
n++; i ← sift-up(n, x)
A[i] ← x
```

int sift-up(int i, Key x)

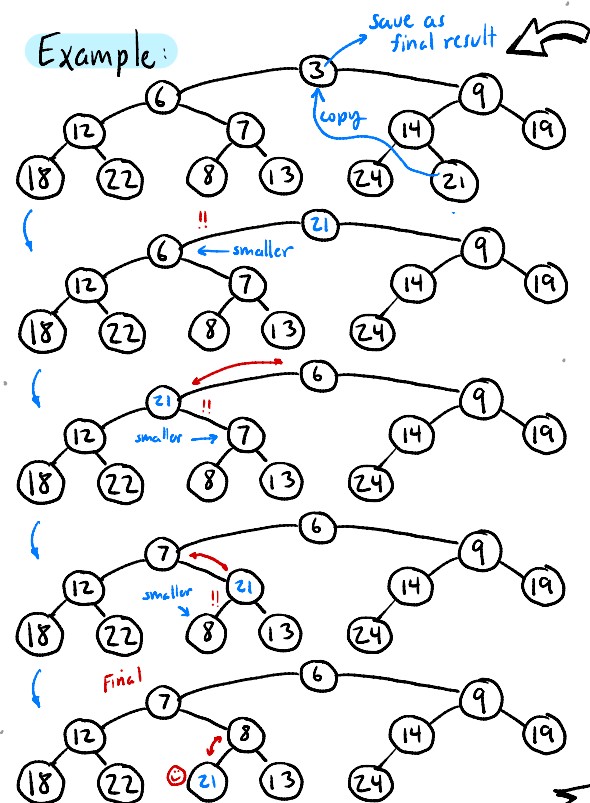
```
while (i > 1 && x < A[par(i)])
    A[i] ← A[par(i)]
    i ← par(i)
return i
```



Insert(x):

- Append x to end of array
- Sift x up until its parent's key is smaller (or reaching root)

Example:



Binary Heap - Extract Min

- Min key at root → save it
- Copy $A[n]$ to root ($A[1]$) + decrement n
- Sift the root key down
 - find smaller of two children
 - if larger, swap with this child
- Return saved root key

Priority Queues + Heaps II

Key extract-min()

```

if (n == 0) Error - Empty heap
result ← A[1]
z ← A[n--] // get replacement
i ← sift-down(i, z)
A[i] ← z
return result
    
```

int sift-down(int i, Key z)

```

while (left(i) ≠ null)
    u ← left(i); v ← right(i)
    if (v ≤ n && A[v] < A[u])
        u ← v // A[u] is smaller child
    if (A[u] < z)
        A[i] ← A[u]; i ← u
    else break
return i
    
```

Leftist Property: Null path length

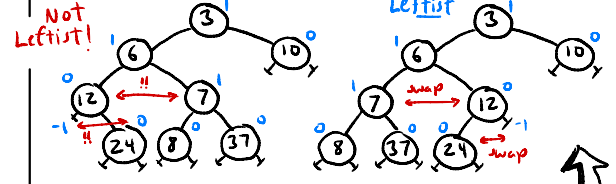
$npl(v)$ = length of shortest path to null

$$npl(v) = \begin{cases} -1 & \text{if } v = \text{null} \\ 1 + \min(npl(v.\text{left}), npl(v.\text{right})) & \text{o.w.} \end{cases}$$

Def: Leftist Heap is binary tree where:

- Keys are heap ordered
- \forall nodes v , $npl(v.\text{left}) \geq npl(v.\text{right})$

Examples:



Leftist Heaps: Meldable heaps

- Can merge two heaps into single heap
- Eg. One processor breaks. Awaiting jobs must be merged with another processor.

Analysis: Both insert + extract-min take time proportional to tree height
Tree is complete $\Rightarrow O(\log n)$ time

Class structure:

```

LeftistHeap<Key> {
    private class LNode {
        Key x
        LNode left, right
        int npl
    }
}
    
```

inner class - used only by LeftistHeap

```

private LNode root
public LeftistHeap() { root = null }
" void insert(Key x)
" Key extractMin()
" void mergeWith(LeftistHeap H2)
... (other private/protected utilities)
    
```

references root node

constructor public functions

```

public mergeWith(LeftistHeap H2) {
    root = merge(this.root, H2.root)
    H2.root = null
}
    
```

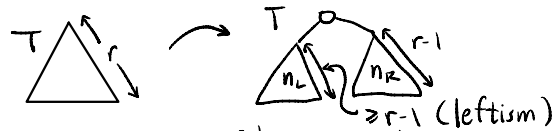
helper function merger destroys H2

Merge helper: 2 phases

- 1 Merge right paths by order of keys + update npl's
- 2 Check leftist property + swap

Lemma: A leftist tree with $r \geq 1$ nodes along its rightmost path has $n \geq 2^r - 1$ nodes

Proof: (Sketch - see latex notes)



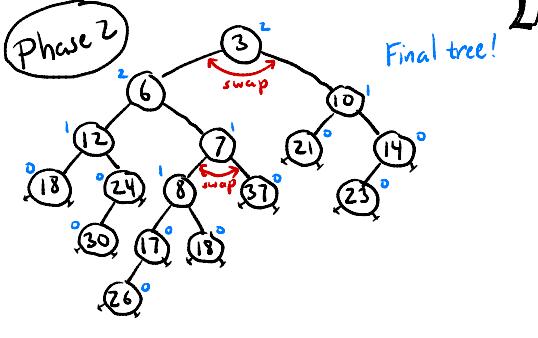
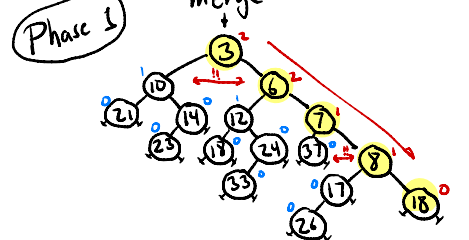
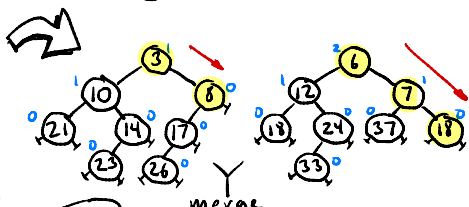
By induction: $n_L \geq 2^{r-1} - 1$, $n_R \geq 2^{r-1}$
 $n = 1 + n_L + n_R \geq (2 \cdot 2^{r-1} - 1) + 1 = 2^r - 1$ □

Priority Queues + Heaps III

Analysis: Time \sim Rightmost path = $O(\log n)$
 Insert + Extract-min? Exercises

```

LNode merge(LNode u, LNode v) {
    if (u == null) return v
    if (v == null) return u
    if (u.key > v.key) // swap so u is smaller
        swap u ↔ v
    if (u.left == null) u.left = v
    else
        u.right = merge(u.right, v)
        if (u.left.npl < u.right.npl)
            swap u.left ↔ u.right
        u.npl = u.right.npl + 1
    return u
}
    
```



Dictionary:

insert (Key x , Value v)

- insert (x, v) in dict. (No duplicates)

delete (Key x)

- delete x from dict. (Error if x not there)

find (Key x)

- returns a reference to associated value v , or null if not there.

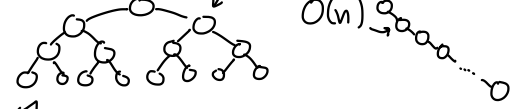


Search: Given a set of n entries each associated with key x and value v_i

- store for quick access & updates
- Ordered: Assume that keys are totally ordered: $<, >, =$

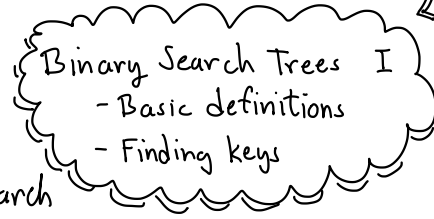
Efficiency: Depends on tree's height

Balanced: $O(\log n)$ Unbalanced: $O(n)$

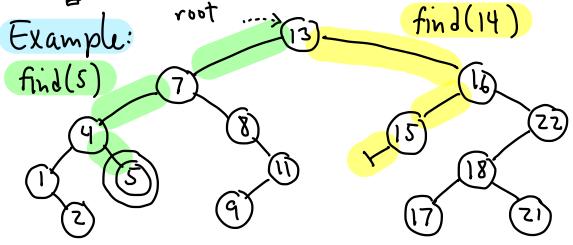


Sequential Allocation?

- Store in array sorted by key
- Find: $O(\log n)$ by binary search
- Insert/Delete: $O(n)$ time

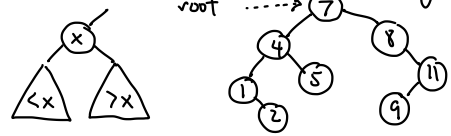


Example:



Can we achieve $O(\log n)$ time for all ops? **Binary Search Trees**

Idea: Store entries in binary tree sorted (inorder traversal) by key



Find: How to find a key in the tree?

- Start at root $p \leftarrow \text{root}$
- if $(x < p.\text{key})$ search left
- if $(x > p.\text{key})$ search right
- if $(x == p.\text{key})$ found it!
- if $(p == \text{null})$ not there!

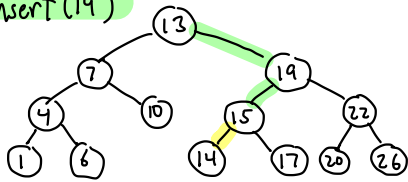


```

Value find (Key x, BSTNode p) {
  if (p == null) return null
  else if (x < p.key)
    return find(x, p.left)
  else if (x > p.key)
    return find(x, p.right)
  else return p.value
}

```

insert (14)



Insert (Key x , Value v)

- find x in tree
- if found \Rightarrow error! duplicate key
- else: create new node where we "fell out"



```

BSTNode insert (Key x, Value v, BSTNode p) {
    if (p == null)
        p = new BSTNode(x, v)
    else if (x < p.key)
        p.left = insert(x, v, p.left)
    else if (x > p.key)
        p.right = insert(x, v, p.right)
    else throw exception  $\rightarrow$  Duplicate!
    return p
}
  
```

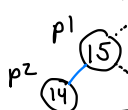
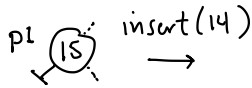
Binary Search Trees II

- insertion
- deletion



Why did we do:

$p.left = insert(x, v, p.left)$?



Be sure you understand this!

$p1.left = insert(14, v, p1.left)$

$p2 = new BSTNode$
return $p2$

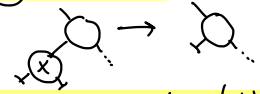
Delete (Key x)

- find x
 - if not found \rightarrow error
 - else: remove this node + restore BST structure
- How?

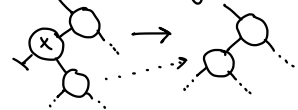


3 cases:

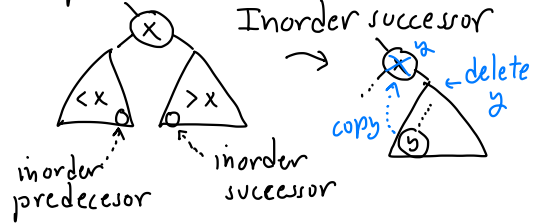
(1) x is a leaf



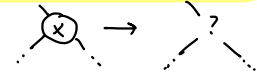
(2) x has single child



Replacement Node?



3. x has two children



Find replacement node

(y), copy to (x), and then delete (y)

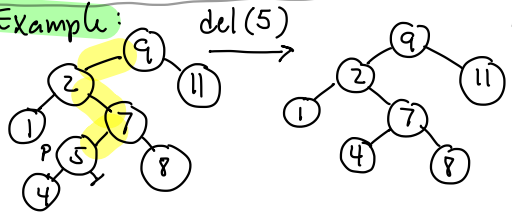


BSTNode delete (Key x, BSTNode p)

```

if (p == null) error! Key not found
else
  if (x < p.key)
    p.left = delete(x, p.left)
  else if (x > p.key)
    p.right = delete(x, p.right)
  else if (either p.left or p.right null)
    if (p.left == null)
      return p.right
    if (p.right == null)
      return p.left
  else
    r = findReplacement(p)
    copy r's contents to p
    p.right = delete(r.key, p.right)
  return p
  
```

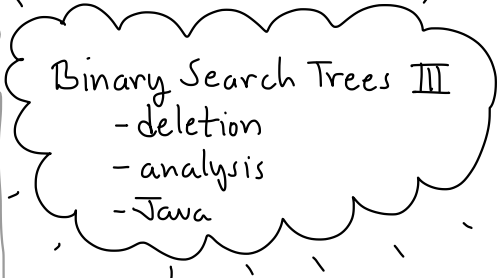
Example:



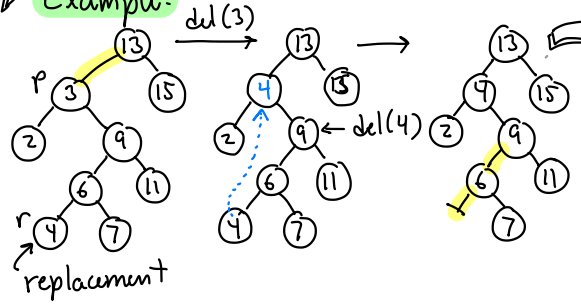
Find Replacement Node

```

BSTNode findReplacement (BSTNode p)
  BSTNode r = p.right
  while (r.left != null)
    r = r.left
  return r
  
```



Example:



Java Implementation:

- Parameterize Key + Value types: `extends Comparable`
- class `BinSearchTree<K,V>`
- BSTNode - inner class
- Private data: `BSTNode root`
- `insert, delete, find`: local
- provide public fns `insert, delete, find`

But height can vary from $O(\log n)$ to $O(n)$..

Expected case is good

Thm: If n keys are inserted in random order, expected height is $O(\log n)$.

Analysis:

All operations (find, insert, delete) run in $O(h)$ time, where h = tree's height

Java implementation (see notes for details)

```
public class BSTree <Key extends Comparable, Value> {
```

```
class Node {  
    Key key  
    Value value  
    Node left, right  
  
    ... constructor, toString...  
}
```

Inner class
for node
(protected)

Local helpers
(private or protected)

```
Value find (Key x, Node p) {...}  
Node insert (Key x, Value v, Node p) {...}  
Node delete (Key x, Node p) {...}
```

```
private Node root;
```

Data (private)

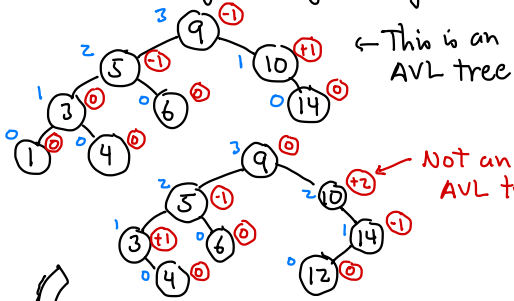
```
public Value find (Key x) {...}  
public void insert (Key x, Value v) {...}  
public void delete (Key x) {...}
```

Public
members
(invoke
helpers)

```
}
```

Balance factor:

$$\text{bal}(v) = \text{hgt}(v.\text{right}) - \text{hgt}(v.\text{left})$$



AVL Height Balance

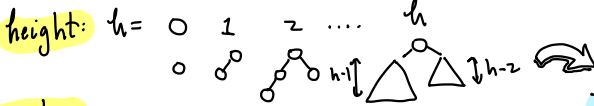
- for each node v , the heights of its subtrees differ by ≤ 1 .

AVL tree: A binary search tree that satisfies this condition



Does this imply $O(\log n)$ height?

Worst cases:



nodes:

$n =$	1	2	4	7	12	20	...
$n+1 =$	2	3	5	8	13	21	...

Recall: $F_0 = 0, F_1 = 1, F_h = F_{h-1} + F_{h-2}$

Conjecture: Min no. of nodes in AVL tree of height h is $F_{h+3} - 1$

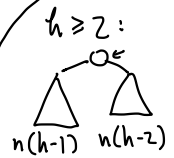
Theorem: An AVL tree of height h has at least $F_{h+3} - 1$ nodes.

Proof: (Induct. on h)

$$h = 0: n(h) = 1 = F_3 - 1$$

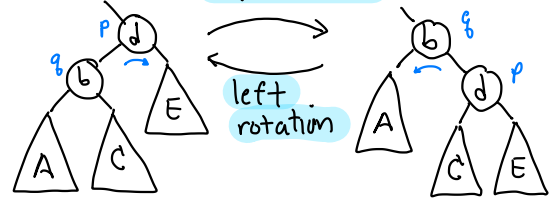
$$h = 1: n(h) = 2 = F_4 - 1$$

$$\begin{aligned} h \geq 2: n(h) &= 1 + n(h-1) + n(h-2) \\ &= 1 + (F_{h+2} - 1) + (F_{h+1} - 1) \\ &= (F_{h+2} + F_{h+1}) - 1 = F_{h+3} - 1 \quad \square \end{aligned}$$



```
BSTNode rotateRight(BSTNode p) {
    BSTNode q = p.left
    p.left = q.right
    q.right = p
    return q
}
```

How to maintain the AVL property?

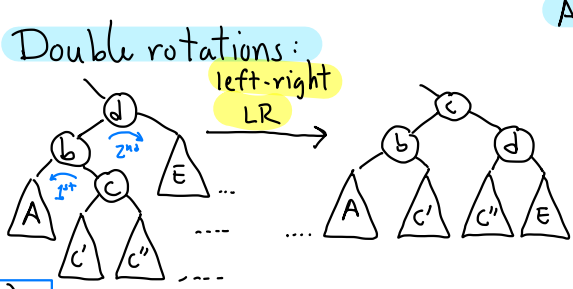
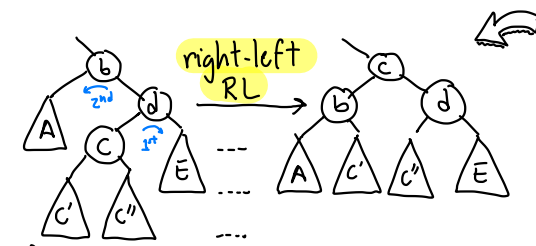


$$A < b < c < d < E$$

$$A < b < c < d < E$$

Corollary: An AVL tree with n nodes has height $O(\log n)$

Proof: Fact: $F_h \approx \varphi^h / \sqrt{5}$ where $\varphi = (1 + \sqrt{5})/2$ "Golden ratio"
 $n \geq \varphi^{h+3} = c \cdot \varphi^h \Rightarrow h \leq \log_{\varphi} n + c$
 $\Rightarrow h \leq \log_2 n / \log_2 \varphi = O(\log n) \quad \square$



```

AVLNode rebalance (AVLNode p)
if (p == null) return p
if (balanceFactor(p) < -1)
  if (ht(p.left.left) >= ht(p.left.right))
    p = rotateRight(p)
  else p = rotateLeftRight(p)
else if (balanceFactor(p) > +1)
  ... (symmetrical)
updateHeight(p); return p
  
```

```

BSTNode rotateLeftRight (BSTNode p)
p.left = rotateLeft(p.left)
return rotateRight(p)
  
```

AVL Tree:

AVL Node: Same as BSTNode (from Lect 4) but add: **int height**

Utilities:

```

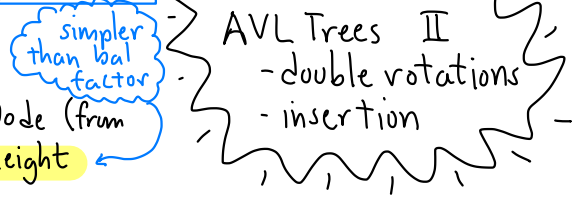
int height (AVLNode p)
return { p == null -> -1
        { ow. -> p.height
  
```

```

void updateheight (AVLNode p)
p.height = 1 + max (height (p.left),
                  height (p.right))
  
```

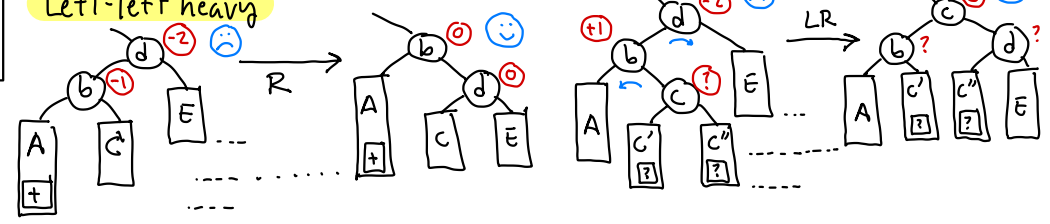
```

int balanceFactor (AVLNode p)
return height (p.right) -
       height (p.left)
  
```



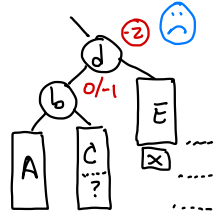
Find: Same as B.S.T.
Insert: Same as BST but as we "back out" rebalance

How to rebalance? Bal = -2
Left-left heavy



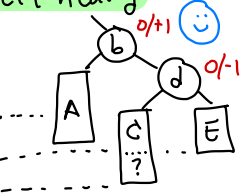
Left-right heavy:

Cases: Balance factor -2

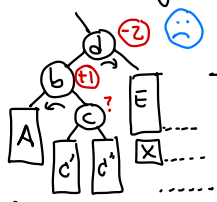


Left-left heavy

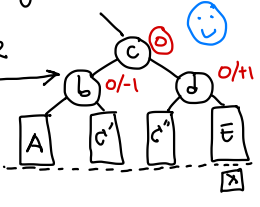
Right rotation



Left-right heavy



LR

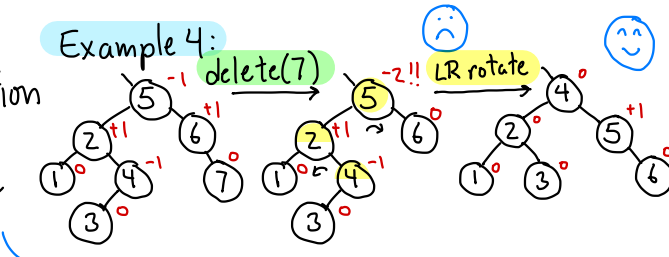


Deletion: Basic plan

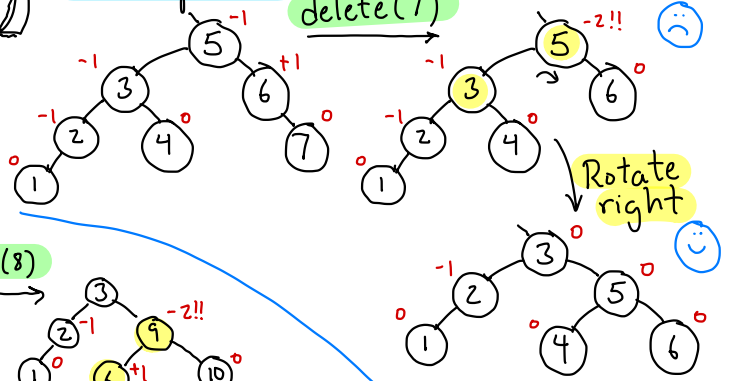
- Apply standard BST deletion
- find key to delete
- find replacement node
- copy contents
- delete replacement
- rebalance

AVL Trees III
- Deletion
- Examples

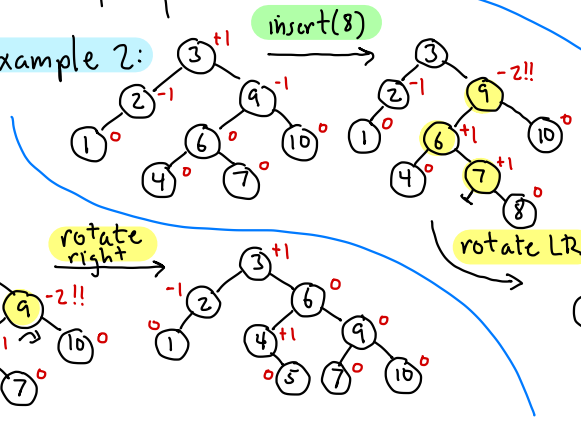
Example 4:



Example 3:

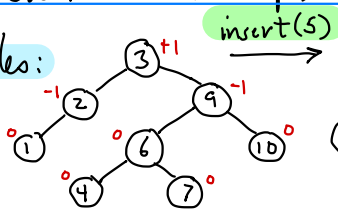


Example 2:

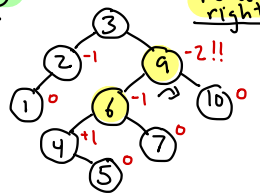


AVLNode delete (Key x, AVLNode p)
same as BST delete
return rebalance(p)

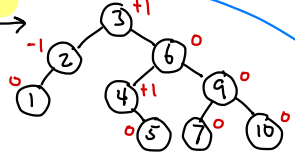
Examples:



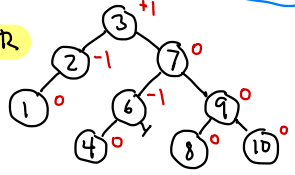
insert(5)



rotate right



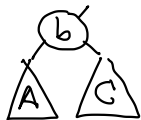
rotate LR



Node types:

2-Node

1 key
2 children



3-Node

2 keys
3 children



↑ Identical heights

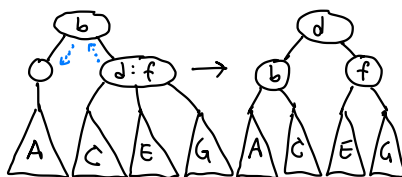


Recap:

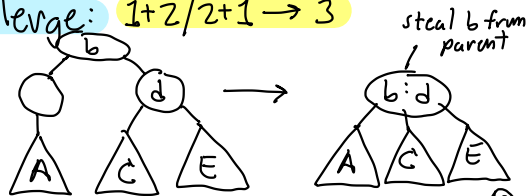
AVL: Height balanced
Binary

2-3 tree: Height exact
Variable width

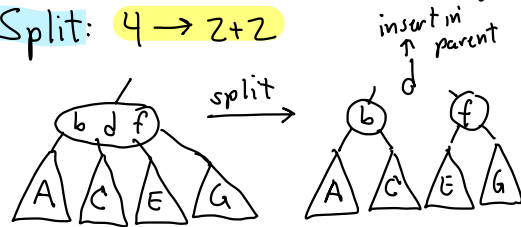
Adoption (Key-Rotation)
 $1+3 = 2+2$



Merge: $1+2/2+1 \rightarrow 3$



Split: $4 \rightarrow 2+2$



Def: A 2-3 tree of height h is either:

- Empty ($h = -1$)
- A 2-Node root and two subtrees, each 2-3 tree of height $h-1$
- A 3-Node root and three subtrees... height $h-1$.



Thm: A 2-3 tree of n nodes has height $O(\log n)$

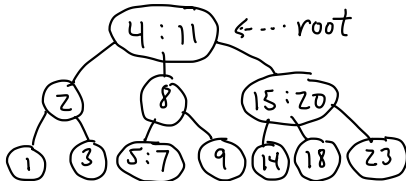


Roughly: $\log_3 n \leq h \leq \log_2 n$



Example:

2-3 tree of height 2



How to maintain balance?

- Split
- Merge
- Adoption (Key rotation)



Conceptual tool:

We'll allow 1-nodes + 4-nodes temporary



Insertion example:

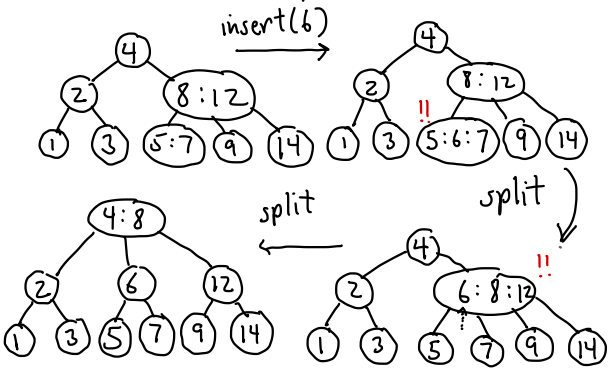


Dictionary operations:

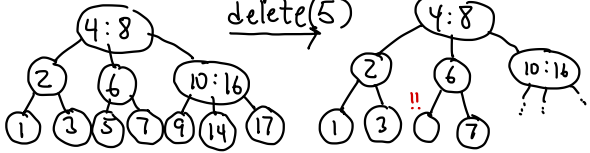
- Find** - straight forward
- Insert** - find leaf node where key "belongs" + add it (may split)
- Delete** - find/replacement/merge or adopt

Implementation?

```
class TwoThreeNode {
    int nChildren
    TwoThreeNode children[3]
    Key key[2]
}
```

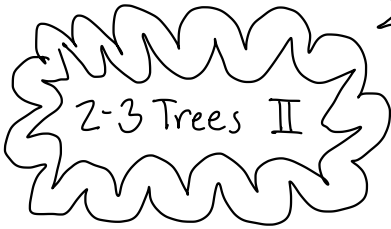


Delete Example:

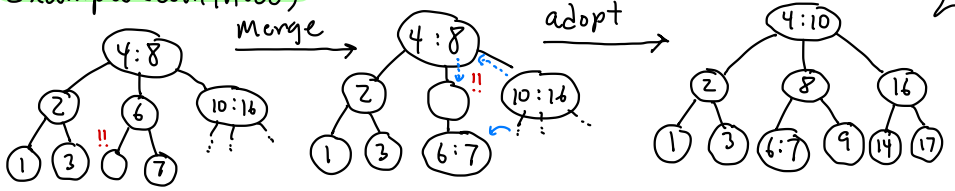


Deletion remedy:

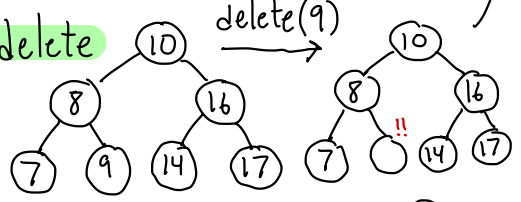
- Have a 3-node neighboring sibling → adopt
- o.w.: Merge with either siblings + steal key from parent



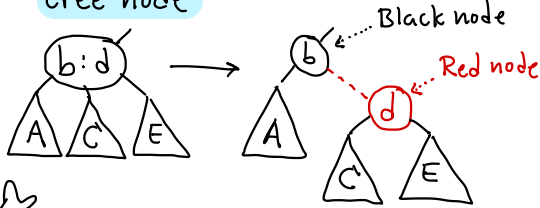
Example (continued)



Another delete example:

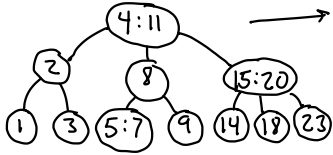


Encoding 3-node as binary tree node

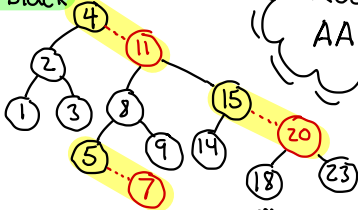


Example:

2-3 Tree:



Red-Black:



Rules:

- ① Every node labeled red/black
- ② Root is black
- ③ Nulls treated as if black
- ④ If node is red, both children are black
- ⑤ Every path from root to null has same no. of black

Some history:

2-3 Trees: Bayer 1972

Red-black Trees: Guibas & Sedgwick 1978 (a binary variant of 2-3)

Rumor - Guibas had two pens - red & black to draw with

Red-Black and AA-Trees I

AA-Trees: Simpler to code

- No null pointers: Create a sentinel node, nil, and all nulls point to it → nil
- No colors: Each node stores level number. Red child is at same level as parent. q is red \Leftrightarrow q.level == p.level

What we need are stricter rules!

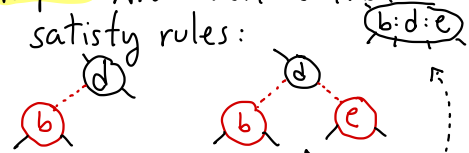
AA-tree:

Arne Anderson 1993

New rule:

- ⑥ Each red node can arise only as right child (of a black node)

Nope! Alternatives that satisfy rules:



A "left-skewed" encoding

Corresponds to 2-3-4 trees

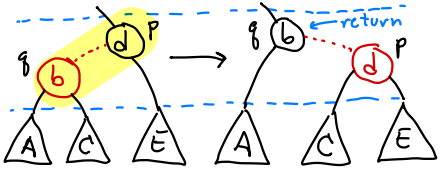
Lemma: A red-black tree with n keys has height $O(\log n)$

Proof: It's at most twice that of a 2-3 tree.

Q: Is every Red-Black Tree the encoding of some 2-3 tree?

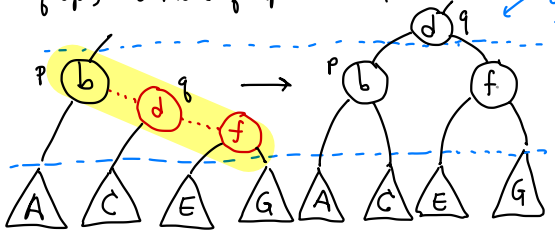
Restructuring Ops:

Skew: Restore right skew
 → If black node has red left child, rotate



How to test? $p.\text{left.level} == p.\text{level}$

Split: If a black node has a right-right red chain, do a left rotation at p (bringing its right child q up) and move q up one level.

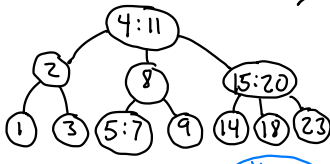


How to test?
 $p.\text{level} == p.\text{right.level} == p.\text{right.right.level}$
 not needed (levels are monotone)

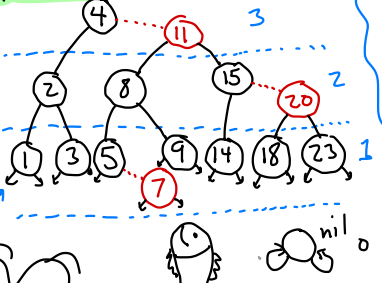


Example:

2-3 Tree:



AA tree:

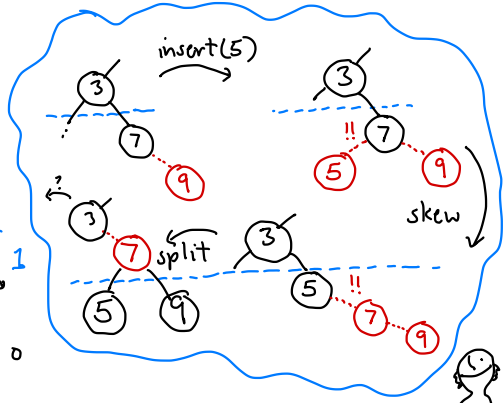


Red-Black + AA Trees II



```

AANode skew(AANode p) {
    if (p == nil) return p
    if (p.left.level == p.level) {
        AANode q = p.left
        p.left = q.right; q.right = p
        return q
    } else return p
}
    
```



AA Insertion:

- Find the leaf (as usual)
- Create new red node
- Back out applying skew + split

AA Node split (AANode p)

```

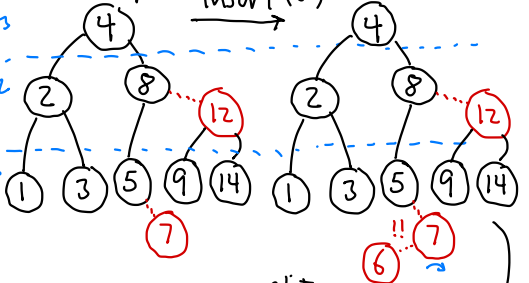
if (p == nil) return p
if (p.right.right.level == p.level) {
    AANode q = p.right
    p.right = q.left
    q.left = p
    q.level += 1
    return q
} else return p
    
```

What 2-3 op does this remind you of?

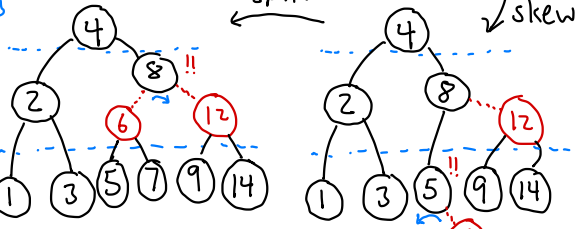


Example:

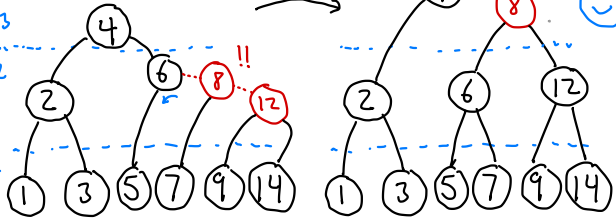
insert(6)



split



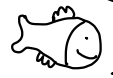
split



```

AANode insert(Key x, Value v, AANode p)
{
    if (p == nil)
        p = new AANode(x, v, 1, nil, nil)
    else if (x < p.key) ... insert on left
    else if (x > p.key) ... insert on right
    else Duplicate Key!
    return split(skew(p))
}
    
```

Red-Black and AA Trees III



Deletion:

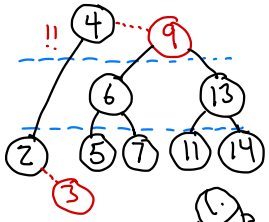
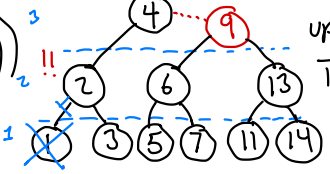
Two more helpers:

updateLevel: If p's level exceeds $l = 1 + \min(p.\text{left.level}, p.\text{right.level})$ then set p's level to l + also p's right child

Example:

delete(1)

update level(2)



fix After Delete(p):

- update p's level
- skew(p), skew(p.right)
- split(p), split(p.right)

deletion: Same as AVL deletion, but end with: **return fix After Delete(p)**



History:

1989: Seidel + Aragon

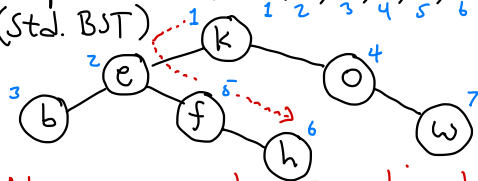
[Explosion of randomized algorithms]

Later discovered this was already known: Priority Search Trees from different context (geometry)
McCreight 1980

Intuition:

- Random insertion into BSTs $\Rightarrow O(\log n)$ expected height
- Worst case can be very bad $O(n)$ height
- Treap: A tree that behaves as if keys are inserted in random order

Example: Insert: k, e, b, o, f, h, w (std. BST)



Along any path - Insertion times increase

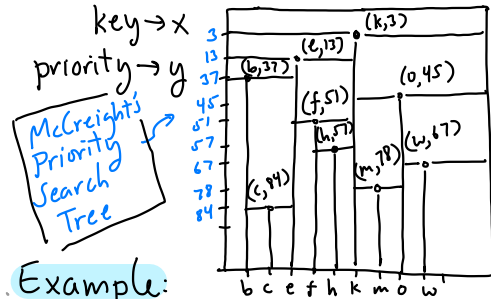
Randomized Data Structures

- Use a random number generator
- Running in expectation over all random choices
- Often simpler than deterministic



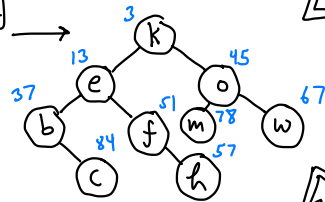
Obs: In a standard BST, keys are by inorder + insert times are in heap order (parent < child)

Geometric Interpretation:



Example:

Key	Priority
b	37
c	84
e	13
f	51
h	57
k	3
m	78
o	45
w	67



Treap: Each node stores a key + a random priority. Keys are in inorder. Priorities are in heap order

? Is it always possible to do both?

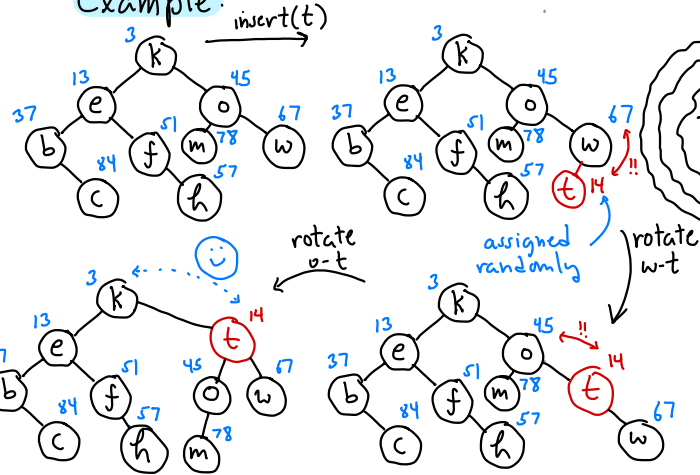
Yes: Just consider the corresponding BST

Insertion: As usual, find the leaf + create a new leaf node.

- Assign random priority
- On backing out - check heap order + rotate to fix.



Example:



rotate o-t

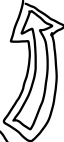
assigned randomly

rotate w-t

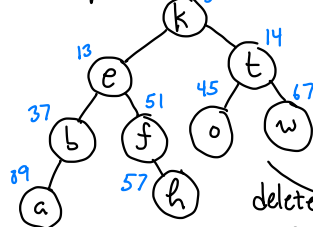


Theorem: A treap containing n entries has height $O(\log n)$ in expectation (averaged over all assignments of random priorities)

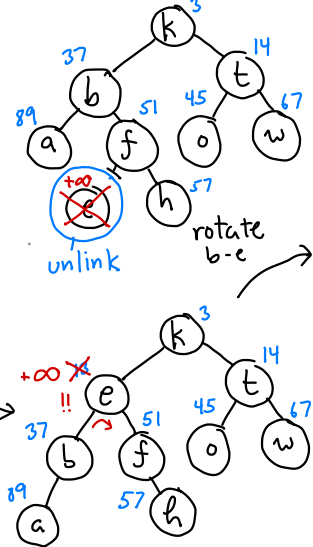
Proof: Follows directly from BST analysis



Example:



delete e



unlink

rotate b-e

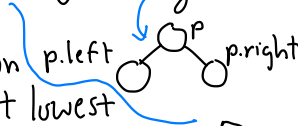
Implementation: (See pdf notes)

Node: Stores priority + usual...

Helpers:

lowest priority (p)

returns node of lowest priority among:



restructure:

performs rotation (if needed) to put lowest priority node at p.



Deletion: (cute solution) Find node to delete. Set its priority to $+\infty$. Rotate it down to leaf level + unlink.

Ideal Skip list:

- Organize list in levels

- Level 0: Everything

1: Every other $\circ \rightarrow \circ \rightarrow \circ \rightarrow \circ$

2: Every fourth

i : Every 2^i $\circ \rightarrow \circ \rightarrow \circ \rightarrow \circ \rightarrow \circ \rightarrow \circ \rightarrow \circ$



Sorted linked lists:

- Easy to code

- Easy to insert/delete

- Slow to search... $O(n)$



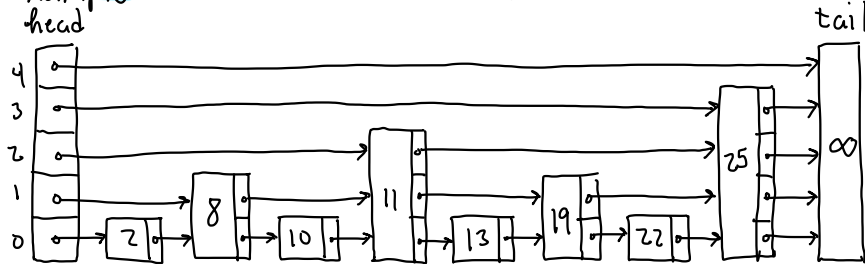
Skip Lists I

Idea: Add extra links to skip

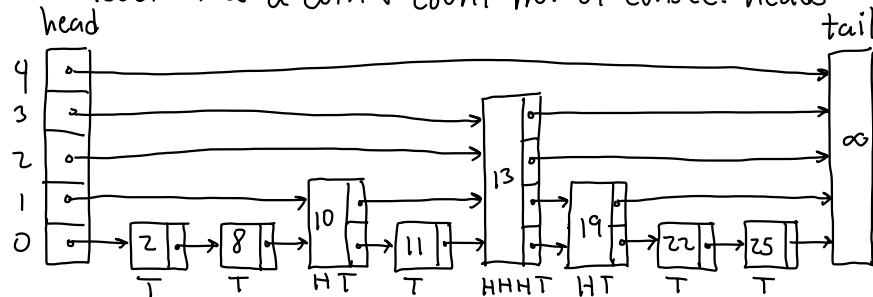


How to generalize?

Example:



Too rigid \rightarrow **Randomize!** To determine level - toss a coin + count no. of consec. heads:



Node Structure: (Variable sized)

```
class SkipNode {
    Key key
    Value value
    SkipNode[] next
}
```

In constructor, set size (height)

```
Value find(Key x) {
    i = topmost level
    SkipNode p = head
    while (i >= 0) {
        if (p.next[i].key <= x) p = p.next[i]
        else i--
    }
    if (p.key == x) return p.value
    else return null
}
```

Annotations:
 - current node (points to p)
 - until we hit base level (points to while loop)
 - advance horizontal (points to p = p.next[i])
 - drop down a level (points to i--)
 - we are at base level (points to })

Thm: A skip list with n nodes has $O(\log n)$ levels in expectation

Proof: Will show that probability of exceeding $c \cdot \lg n$ is $\leq 1/n^{c-1}$

→ Prob that any given node's level exceeds l is $1/2^l$

[l consecutive heads]

→ Prob that any of n nodes' level exceeds l is $\leq n/2^l$

[n trials with prob $1/2^l$]

→ Let $l = c \cdot \lg n$ ($\lg \equiv \log_2$)

Prob that max level exceeds

$$c \cdot \lg n \text{ is:}$$

$$\leq n/2^l = n/2^{(c \cdot \lg n)}$$

$$= n/(2^{\lg n})^c$$

$$= n/n^c = 1/n^{c-1}$$

□

Obs: Prob. level exceeds

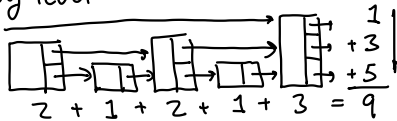
$$3 \cdot \lg n \text{ is } \leq 1/n^2$$

(If $n \geq 1,000$, chances are less than 1 in million!)

Skip Lists II

Thm: Total space for n -node skip list is $O(n)$ expected.

Proof: Rather than count node by node, we count level by level:



- Let $n_i =$ no. of nodes that contrib. to level i .

- Prob that node at level $\geq i$ is $1/2^i$

- Expected no. of nodes that contrib. to level $i = n/2^i$

$$\Rightarrow E(n_i) = n/2^i$$

Total space (expected) is:

$$E\left(\sum_{i=0}^{\infty} n_i\right) = \sum_{i=0}^{\infty} E(n_i) = \sum_{i=0}^{\infty} n/2^i$$

$$= n \sum_{i=0}^{\infty} 1/2^i = 2n$$

□

Thm: Expected search time is $O(\log n)$

Proof:

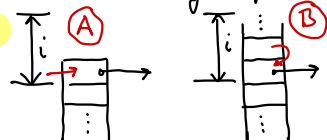
- We have seen no. levels is $O(\log n)$

- Will show that we visit 2 nodes per level on average

Obs: Whenever search arrives first time to a node, it's at top level. (Can you see why?)

Def: $E(i) =$ Expect. num. nodes visited among top i levels.

Cases:



$$E(i) = 1 + (\text{Prob(A)})E(i) + (\text{Prob(B)})E(i-1)$$

current node
same level
from prior level

$$= 1 + 1/2 E(i) + 1/2 E(i-1)$$

$$\Rightarrow E(i)(1 - 1/2) = 1 + 1/2 E(i-1)$$

$$\Rightarrow E(i) = [1 + 1/2 E(i-1)] \cdot 2 = 2 + E(i-1)$$

Basis: $E(0) = 0 \Rightarrow E(i) = 2 \cdot i$

Let $l =$ max level. Total visited = $E(l)$

\Rightarrow We visit 2 nodes per level on average. □

Skip Lists III

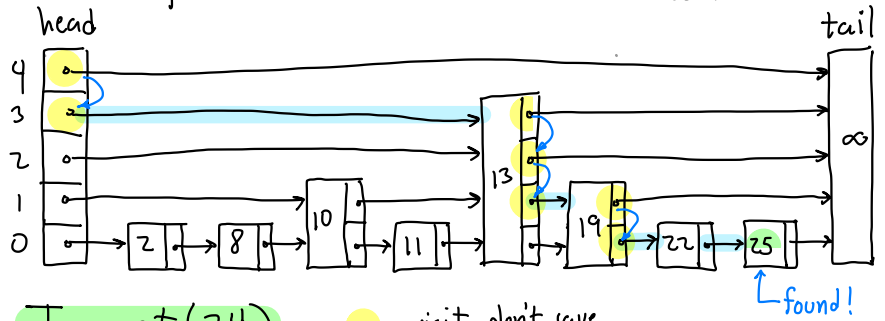
Delete:

- Start at top
- Search each level saving last node < key
- On reaching node at level 0, remove it and unlink from saved pointers

Insert: (Similar to linked lists)

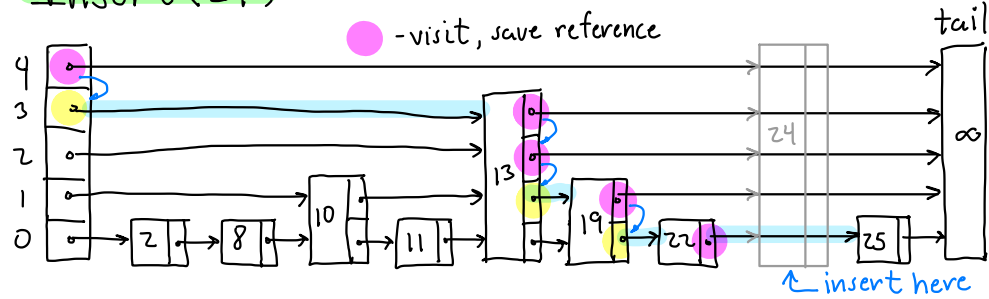
- Start at top level
- At each level:
 - Advance to last node \leq key
 - Save node + drop level
- At level 0:
 - Create new node (flip coins to determine height)
 - Link into each saved node

Example: find(25)

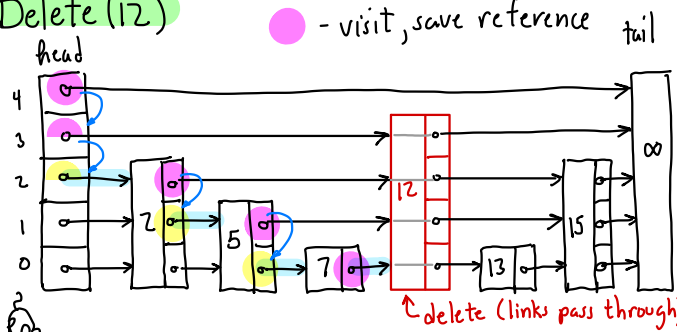


Insert(24)

- visit, don't save (yellow dot)
- visit, save reference (pink dot)



Delete(12)



- visit, don't save (yellow dot)
- visit, save reference (pink dot)

Analysis: All operations run in time \sim find $\Rightarrow O(\log n)$ expected
Note: Variation in running times due to randomness only - not sequence
 \Rightarrow User cannot force poor performance.

Other/Better Criteria?

Expected case: Some keys more popular than others

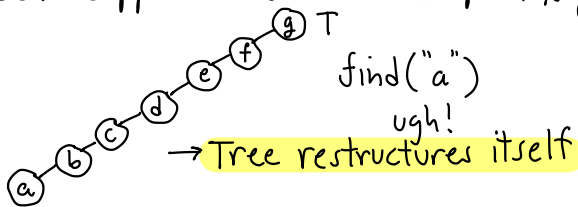
Self-adjusting: Tree adapts as popularity changes

How to design/analyze?

Splay Tree: A self-adjusting binary search tree

- **No rules!** (yay anarchy!)
 - No balance factors
 - No limits on tree height
 - No colors/levels/priorities
- **Amortized efficiency:**
 - Any single op - slow
 - Long series - efficient on avg.

Intuition: Let T be an unbalanced BST + suppose we access its deepest key



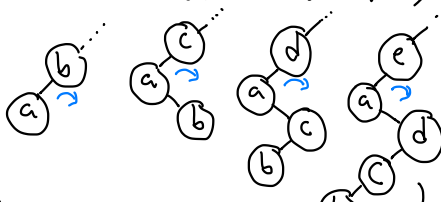
Recap: Lots of search trees

- Unbalanced BSTs
- AVL Trees
- 2-3, Red-black, AA Trees
- Treaps + Skip lists

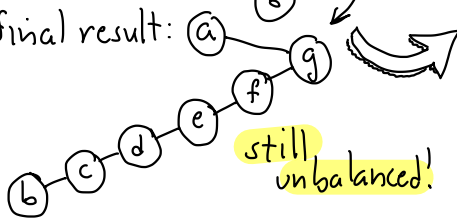
→ **Focus:** Worst-case or randomized expected case

SPLAY TREES I

Idea I: Rotate "a" to top (Future accesses to "a" fast)

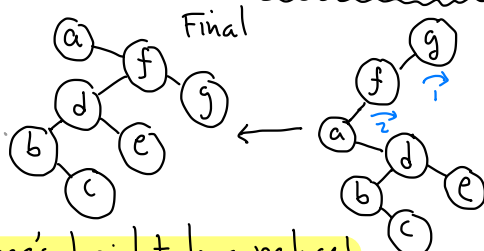


... final result:

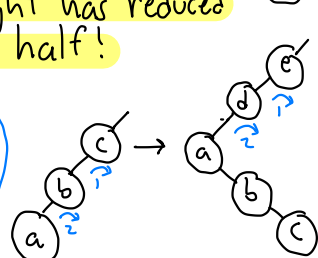


Lesson: Different combinations of rotations can:

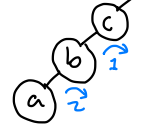
- bring given node to root
- significantly change (improve) tree structure.



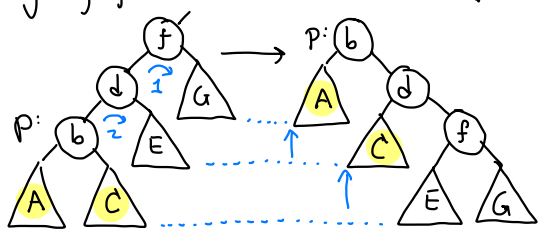
Tree's height has reduced by ~ half!



Idea II: Rotate 2 at a time - upper + lower

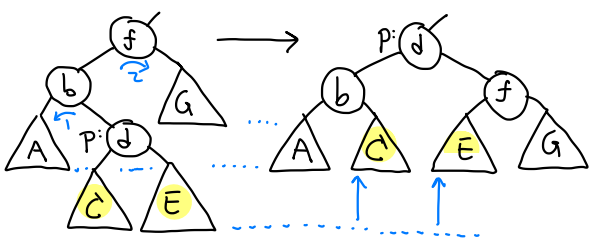


ZigZig(p): [LL case]



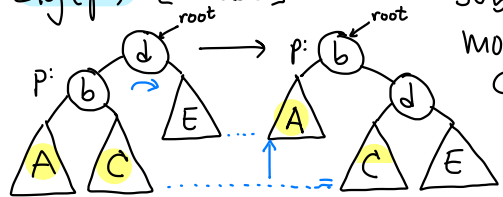
Subtrees A, C move up ↑

ZIGZAG(P): [LR case]



Subtrees C, E of p move up ↑

Zig(p): [L case]



Subtree A moves up ↑
C unchanged

Splay(Key x):

```

Node p ← find x by standard BST search
while (p ≠ root) {
  if (p == child of root) zig(p)
  else /* p has grand parent */
    if (p is LL or RR grand child) zigzig(p)
    else /* p is LR or RL gr. child */ zigzag(p)
}
  
```

insert(x):

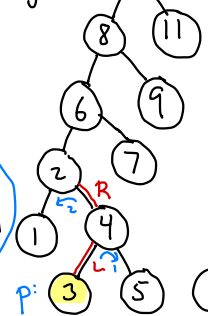
```

Node p ← splay(x)
if (p.key == x) Error!!
q ← new Node(x)
if (p.key < x)
  q.left ← p
  q.right ← p.right
  p.right ← null
else ... (symmetrical)...
root ← q
  
```

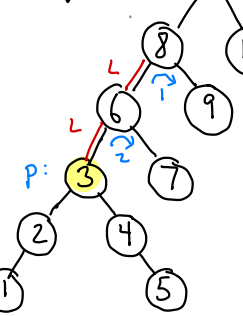
Splay Trees II

Example:

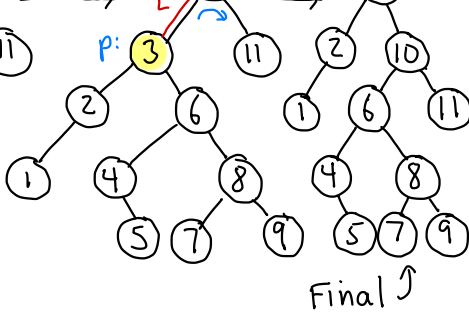
splay(3) RL zigzag



LL zigzig



L zig p: 3



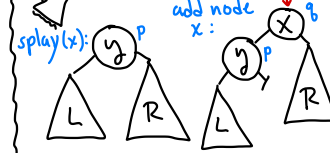
Final ↑

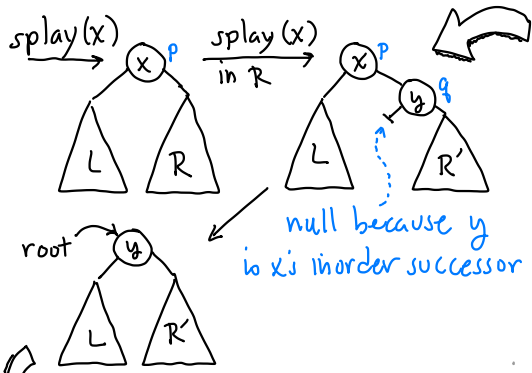
find(x):

```

root ← splay(x)
if (root.key == x)
  return root.value
else return null
  
```

insert(x):





delete(x):
 splay(x) [x now at root]
 p = root
 if (p.key ≠ x) **error!**
 splay(x) in p's right subtree
 q = p.right [q's key is xi successor]
 q.left = p.left [q.left == null]
 root = q

Dynamic Finger Theorem:
 Keys: $x_1 < \dots < x_n$. We perform accesses $x_{i_1}, x_{i_2}, \dots, x_{i_m}$
 Let $\Delta_j = i_j - i_{j-1}$: distance between consecutive items

 • Δ_1 Δ_2 Δ_3 Δ_4

Thm: Total access time is $O(m + n \log n + \sum_{j=1}^m (1 + \lg \Delta_j))$

Analysis:

- Amortized analysis
- Any one op might take $O(n)$
- Over a long sequence, average time is $O(\log n)$ each
- Amortized analysis is based on a sophisticated **potential argument**
- Potential: A function of the tree's structure
- **Balanced** \Rightarrow Low potential.
- **Unbalanced** \Rightarrow High potential.
- Every operation tends to reduce the potential

SPLAY TREES III

Splay Trees are **Amazingly Adaptive!**

Balance Theorem: Starting with an empty dictionary, any sequence of m accesses takes total time $O(m \log n + n \log n)$ where $n = \max.$ entries at any time.


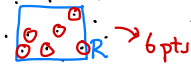
Static Optimality:

- Suppose key x_i is accessed with prob p_i ($\sum p_i = 1$)
- **Information Theory:** Best possible binary search tree answers queries in expected time $O(H)$ where $H = \sum p_i \lg 1/p_i$ **Entropy**

Static Optimality Theorem:

Given a seq. of m ops. on splay tree with keys x_1, \dots, x_n , where x_i is accessed g_i times. Let $p_i = g_i/m$. Then total time is $O(m \sum p_i \lg 1/p_i)$

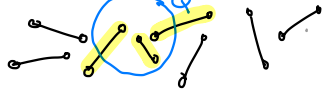
Geometric Search:

- Nearest neighbors \rightarrow 
- Range searching \rightarrow 

- Point Location



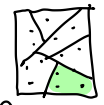
- Intersection Search

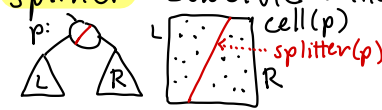


So far: 1-dimensional keys

- Multi-dimensional data
- Applications:
 - Spatial databases + maps
 - Robotics + Auton. Systems
 - Vision/Graphics/Games
 - Machine Learning

Partition Trees:



- Tree structure based on hierarchical space partition
 - Each node is associated w. a region - **cell**
 - Each internal node stores a **splitter** - subdivides the cell
- 



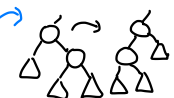
Quadtrees & kd-Trees I

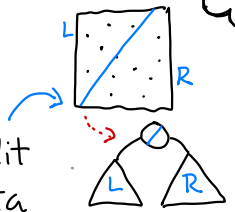
Multi-Dim vs. 1-dim Search?

Similarities:

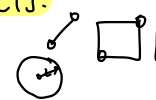
- Tree structure
- Balance $O(\log n)$
- Internal nodes - split
- External nodes - data

Differences:

- No (natural) total order
- Need other ways to discriminate + separate
- Tree rotation may not be meaningful 



Representations:

- **Scalars:** Real numbers for coordinates, etc. float
- **Points:** $p = (p_1, \dots, p_d)$ in real d-dim space \mathbb{R}^d
- **Other geom objects:** Built from these 

- External nodes store pts.

Point: A d-vector in \mathbb{R}^d
 $p = (p_1, \dots, p_d)$ $p_i \in \mathbb{R}$

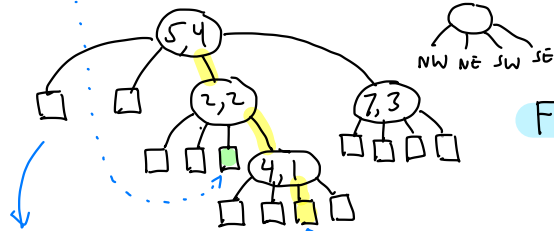
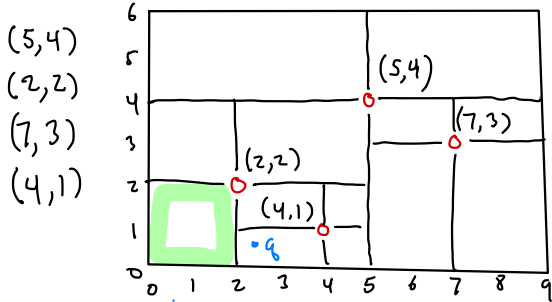
class Point {

```

float[] coord // coords
Point(int d)
    ...  $\rightarrow$  coord = new float[d]
int getDim()  $\rightarrow$  coord.length
float get(int i)  $\rightarrow$  coord[i]
... others: equality, distance
toString...
  
```

Point Quadtree:

- Each internal node stores a point
- Cell is split by horiz. + vertic. lines through point



Each external node corresponds to cell of final subdivision



Quadtrees: (abstractly)

- Partition trees
- Cell: Axis-parallel rectangle [AABB - Axis-aligned bounding box]
- Splitter: Subdivides cell into four (genly 2^d) subcells



Quadtrees & kd-Trees II



Find/Pt Location:

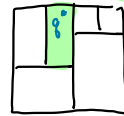
Given a query point q , is it in tree, and if not which leaf cell contains it?
 → Follow path from root down (generalizing BST find)

History: Bentley 1975

- called it 2-d tree (\mathbb{R}^2)
- 3-d tree (\mathbb{R}^3)
- In short kd-tree (any dim)
- Where/which direction to split? → next

kd-Tree: Binary variant of quadtree

- splitter: Horiz. or vertic. line in 2-d (orthogonal plane ow.)
- cell: Still AABB



left: left/below
right: right/above

Quadtrees - Analysis

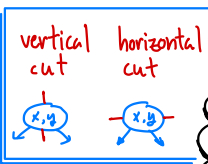
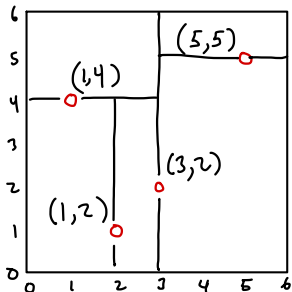
- Numerous variants! PR, PMR, QR, QX, ... see Samet's book
- Popular in 2-d apps (in 3-d, **outtrees**)
- Don't scale to high dim
 - out degree = 2^d
- What to do for higher dims?



Example:



- (3,2)
- (1,4)
- (5,5)
- (1,2)

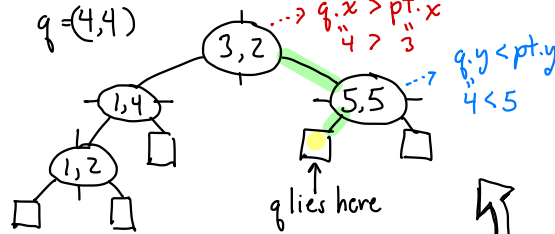


Kd-Tree Node:

```
class KDNode {
    Point pt // splitting point
    int cutDim // cutting coordinate
    KDNode left // low side
    KDNode right // high side
}
```

Quadtrees & kd-Trees III

Example: find(q) ^{calls} find(q, root)



Analysis: Find runs in time $O(h)$, where h is height of tree.

Theorem: If pts are inserted in random order, expected height is $O(\log n)$

```
Value find(Point q, KDNode p) {
    if (p == null) return null;
    else if (q == p.pt) return p.value; // all coords match?
    else if (p.onLeft(q)) return find(q, p.left);
    else return find(q, p.right);
}
```

Find point q in subtree rooted at p with cutDim cd :

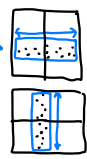
- if $q == p.point \Rightarrow$ found!
- if $q[cd] < p.point[cd] \Rightarrow$ left
- if $q[cd] \geq p.point[cd] \Rightarrow$ right

Helper:

```
class KDNode {
    boolean onLeft(Point q) {
        return q[cutDim] < pt[cutDim];
    }
}
```

How do we choose cutting dim?

- Standard kd-tree: cycle through them (eg. $d=3: 1,2,3,1,2,3,\dots$) based on tree depth
- Optimized kd-tree: (Bentley) Based on widest dimension of pts in cell.



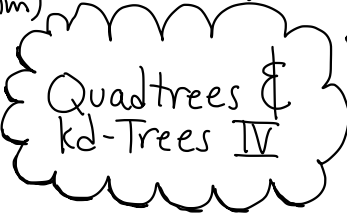
KDNode insert (Point pt, KNode p, int cd)

```

if (p == null) // fell out?
    p = new KNode(pt, cd) // new leaf node
else if (p.point == pt)
    Error! Duplicate key
else if (p.onLeft(pt))
    p.left = insert(pt, p.left, (cd+1)%dim)
else
    p.right = insert(pt, p.right, (cd+1)%dim)
return p
    
```

Kd-Tree Insertion:

- (Similar to std. BSTs)
- Descend tree until
 - find pt → Error - duplicate
 - falling out ← (Although we draw extended trees, lets assume standard trees)
 - create new node
 - set cutting dim



Deletion:

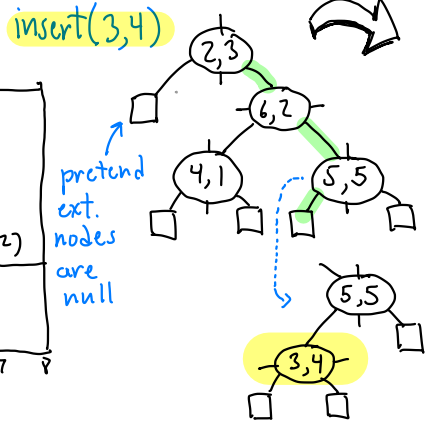
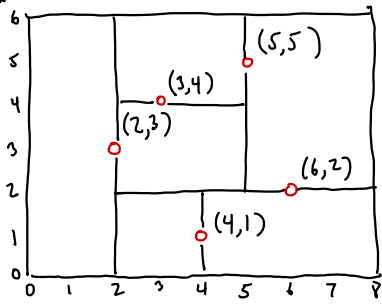
- Descend path to leaf
- If found:
 - leaf node → just remove
 - internal node
 - find replacement
 - copy here
 - recur. delete replacement

This is the hardest part. See Latex notes.

Rebalance by Rebuilding:

- Rebuild subtrees as with scapegoat trees
- $O(\log n)$ amortized
- Find: $O(\log n)$ guaranteed.

Example:

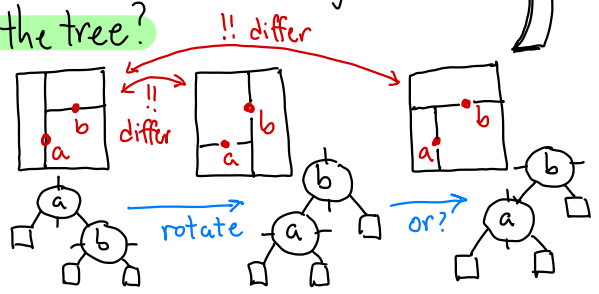


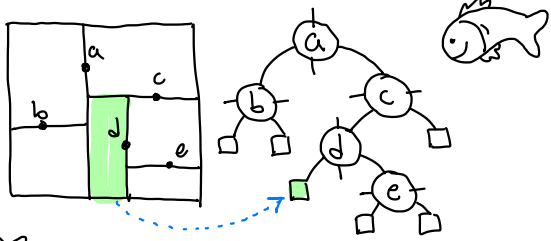
Analysis:

Run time: $O(h)$

Can we balance the tree?

- Rotation does not make sense !!





Kd-Trees:

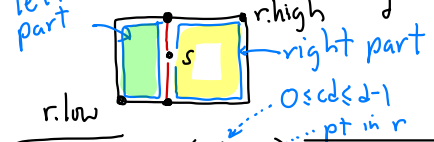
- Partition trees → vert

L	R
---	---
- Orthogonal split → horz

R	L
---	---
- Alternate cutting dimension x, y, x, y, \dots
- Cells are axis-aligned rectangles (AABB)

Rectangle methods for kd-cells:

- Split a cell r by a split pt $s \in r$, along cut dim cd



$r.\text{leftPart}(cd, s)$
 → returns rect with $low = r.low$
 + $high = r.high$ but
 $high[cd] \leftarrow s[cd]$

$r.\text{rightPart}(cd, s)$
 → $high = r.high$ + $low = r.low$ but
 $low[cd] \leftarrow s[cd]$

Queries?

- **Orthogonal range queries**
 - Given query rect. (AABB) count/report pts in this rect.
- Other range queries?
 - Circular disks
 - Halfplane
- **Nearest neighbor queries**
 - Given query pt, return closest pt in the set
 - Find k^{th} closest point
 - Find farthest point from q

Kd-Tree Queries I

Axis-Aligned Rect in \mathbb{R}^d

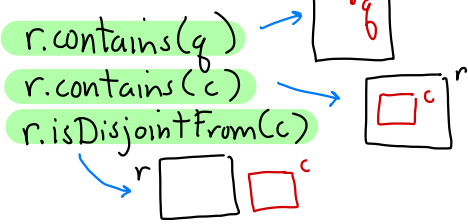
- Defined by two pts: $low, high$



- Contains pt $q \in \mathbb{R}^d$ iff $low_i \leq q_i \leq high_i$

Useful methods:

Let r, c - Rectangle
 q - Point



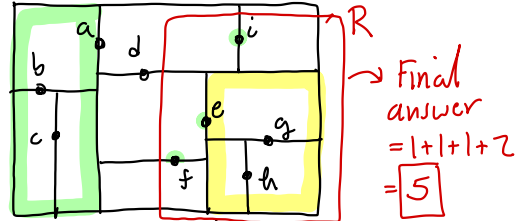
This Lecture: $O(\sqrt{n})$ time alg for orthog. range counting queries in \mathbb{R}^2
 → General \mathbb{R}^d : $O(n^{1-1/d})$

Orthog. Range Query

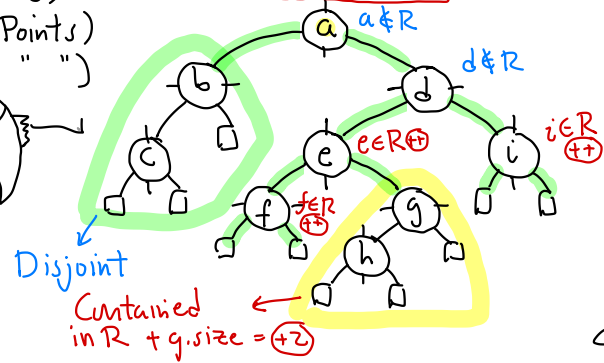
- Assume: Each node p stores:
 - $p.pt$: splitting point
 - $p.cutDim$: cutting dim
 - $p.size$: no. of pts in p 's subtree
- Tree stores ptr. to **root** and **bounding box** for all pts.
- Recursive helper stores current. **node p + p 's cell.**

```

class Rectangle {
private Point low, high
public Rect (Point l, Point h)
    " boolean contains(Point q)
    " boolean contains(Rect c)
    " Rect leftPart (int cd, Points)
    " Rect rightPart (" " " " " ")
}
    
```



Kd-Tree Queries II



Cases:

- $p == null \rightarrow$ fell out of tree $\rightarrow 0$
- Query rect is **disjoint** from p 's cell
 - \rightarrow return 0
 - \rightarrow no point of p contributes to answer
- Query rect **contains** p 's cell
 - \rightarrow return $p.size$
 - \rightarrow every point of p 's subtree contributes to answer.
- Otherwise: Rect. + cell **overlap** - Recurse on both children

```

int rangeCount(Rect R, KDNode p, Rect cell)
{
if (p == null) return 0 // fell out of tree
else if (R.isDisjointFrom(cell)) return 0 // no overlap
else if (R.contains(cell)) return p.size // take all
else { int ct = 0
    if (R.contains(p.pt)) ct++ // p's pt in range
    ct += rangeCount(R, p.left, cell.leftPart(p.cutDim, p.pt))
    ct += rangeCount(R, p.right, cell.rightPart... )
}
}
    
```

Theorem: Given a balanced kd-tree storing n pts in \mathbb{R}^2 (using alternating cut dim), orthog. range queries can be answered in $O(\sqrt{n})$ time.

→ Slower than $\log n$. Faster than n

Analysis: How efficient is our algorithm?

- Tricky to analyze
- At some nodes we recurse on both children $\Rightarrow O(n)$ time?
- At some we don't recurse at all!

Solving the Recurrence:

- Macho: Expand it
- Wimpy: Master Thm (CLRS)

Master Thm:

$$T(n) = aT\left(\frac{n}{b}\right) + n^d + d \log_b a$$

$$\Rightarrow T(n) = n^{\log_b a}$$

For us: $a=2$
 $b=4$
 $d=0$

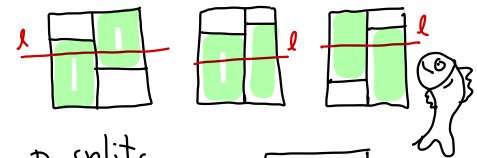
$$\Rightarrow T(n) = n^{\log_4 2} = n^{1/2} = \sqrt{n}$$

Since tree is **balanced** a child has half the pts + grandchild has quarter.

Recurrence: $T(n) = 2 + 2T(n/4)$

2 cells stabbed
 Recurse on 2 grandchildren
 Each has $n/4$ pts

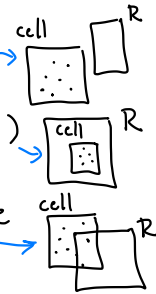
If we consider 2 consecutive levels of kd-tree, l stabs at most 2 of 4 cells:



p splits horizontally
 l stabs only one

Stabbing: 3 cases

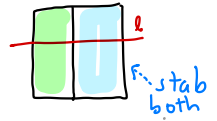
- cell is disjoint (easy)
- cell is contained (easy)
- cell partially overlaps or is stabbed by the query range (hard!)



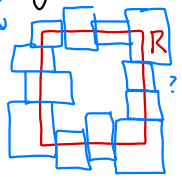
Kd-Tree Queries III

Lemma: Given a kd-tree (as in Thm above) and horiz. or vert. line l , at most $O(\sqrt{n})$ cells can be stabbed by l

Proof: w.l.o.g. l is horiz.
Cases: p splits vertically



How many cells are stabbed by R ? (worst case)



Simpler: Extend R 's sides to 4 lines + analyze each one.



Range Tree Applications:

- Range trees can be applied to a variety of query problems

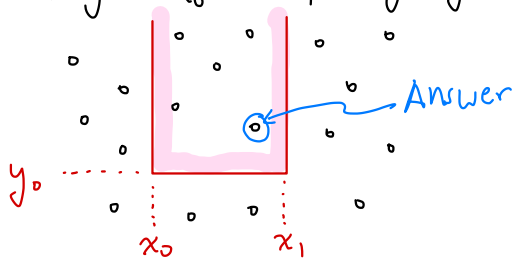
- Methods:

- Minimization/Maximization
- Transform coordinates
- Adding new coordinates

Minimization/Maximization -

3-Sided Min Query

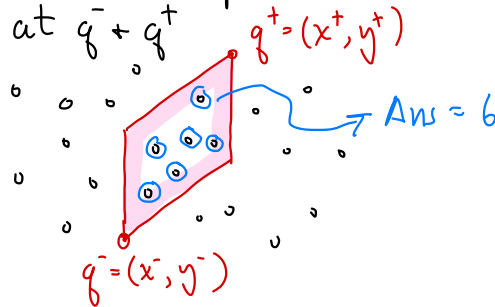
Given a set P of n pts in \mathbb{R}^2 , a query consists of x -interval $[x_0, x_1]$ and y value y_0 . Return the lowest pt in 3-sided region $x_0 \leq x \leq x_1$, $y \geq y_0$



Transforming coordinates:

Skewed rectangle query:

Given a set P of n pts in \mathbb{R}^2 , a skewed rectangle is given by 2 pts $q^- = (x^-, y^-)$ and $q^+ = (x^+, y^+)$ and consists of pts in parallelogram with two vertical sides and two with slope $+1$ corners at $q^- + q^+$

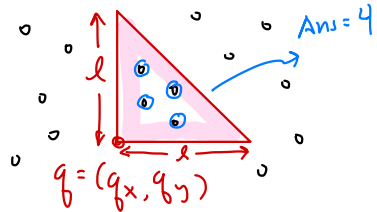


Return a count of the number of pts of P inside the skewed rectangle.

Adding New Coordinates:

NE Right Triangle Query

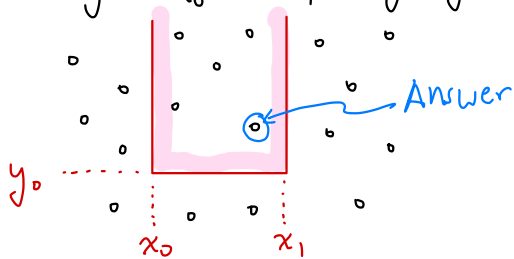
Given a set P of n pts in \mathbb{R}^2 and scalar $l > 0$, a NE triangle is a 45-45 right triangle with lower left corner at q and side length l .



Return a count of the number of pts of P lying within the triangle.

3-Sided Min Query

Return lowest in region
 region $x_0 \leq x \leq x_1, + y \geq y_0$



Data structure:

- Build a range tree for x
- Aux. trees are range trees for y that support find larger

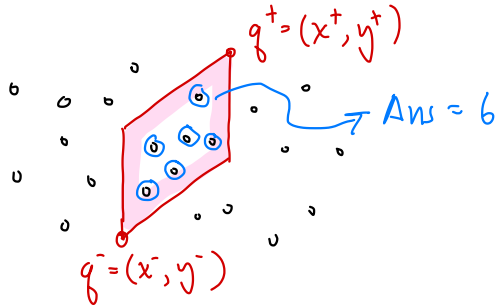
Query Processing:

- Do 1D range search in main tree for interval $[x_0, x_1]$
- For each maximal subtree in range, do find larger (y_0)
- Return smallest of these.

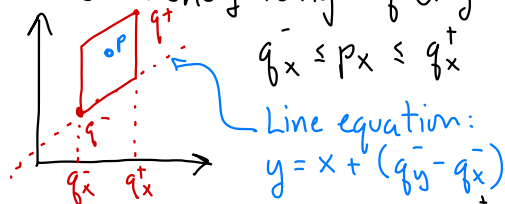
Analysis:

- Same as 2D range tree
- Space: $O(n \log n)$ Time: $O(\log^2 n)$

Skewed rectangle query:



Transform coordinates to
 make orthog range query



$$q'_x \leq p_x \leq q'_y$$

Line equation:
 $y = x + (q'_y - q'_x)$

$$p_x + (q'_y - q'_x) \leq p_y \leq p_x + (q'_y - q'_x)$$

$$\Leftrightarrow q'_y - q'_x \leq p_y - p_x \leq q'_y - q'_x$$

$\underbrace{\hspace{2cm}}_{p'_y = p_y - p_x}$

Map each $p = (p_x, p_y) \in P$
 to $p' = (p'_x, p'_y) \triangleq (p_x, p_y - p_x)$

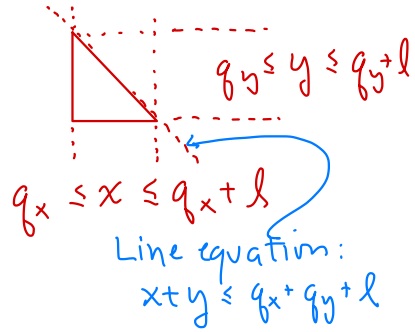
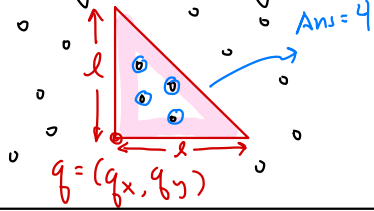
Let P' be resulting set.

Build std. range tree for P' . Return ans. to query

$$q'_x \leq x \leq q'_x$$

$$q'_y - q'_x \leq y \leq q'_y - q'_x$$

NE Right Triangle Query



- Add new coord:

$$z = x + y$$

- Map pts:

$$p = (p_x, p_y) \rightarrow p' = (p_x, p_y, p_x + p_y)$$

- Let P' be resulting set

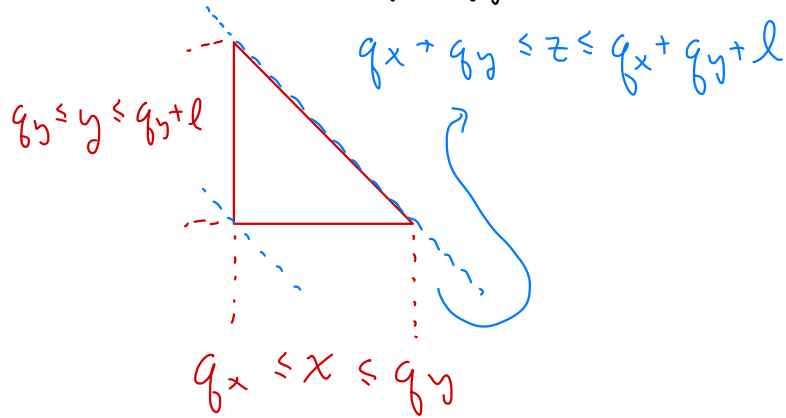
Build a 3D range tree on P'

NE triangle query becomes:

$$q_x \leq x \leq q_x + l$$

$$q_y \leq y \leq q_y + l$$

$$q_x + q_y \leq z \leq q_x + q_y + l$$



Space:

$$O(n \log^2 n)$$

Query time:

$$O(\log^3 n)$$

Can we do better?



Recap:

- **kd-Tree**: General-purpose data structure for pts in \mathbb{R}^d
- **Orthogonal range query**: Count/report pts in axis-aligned rect. \rightarrow Ans = 4
- **kd-Tree**: **Counting**: $O(\sqrt{n})$ time
Report: $O(k + \sqrt{n})$ time

Call this a **1-Dim Range Tree**:

Claim: A 1-Dim range tree with n pts has space $O(n)$ and answers 1-D range count/rept queries in time $O(\log n)$ (or $O(k + \log n)$)

- Space is $O(n \log^{d-1} n)$
- Query time: **Counting**: $O(\log^d n)$
Reporting: $O(k + \log^d n)$
- \rightarrow In \mathbb{R}^2 : $\log^2 n$ much better than \sqrt{n} for large n
- \rightarrow Range trees are more limited

Layering: Combining search structures

- Suppose you want to answer a **composite query** w. multiple criteria:

- Medical data: Count subjects
 - Age range**: $a_{l_0} \leq \text{age} \leq a_{h_1}$
 - Weight range**: $w_{l_0} \leq \text{weight} \leq w_{h_1}$

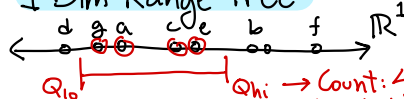
- Design a data structure for each criterion **individually**
- **Layer** these structures together to answer full query

\rightarrow **Multi-Layer Data Structures**

Range Trees I

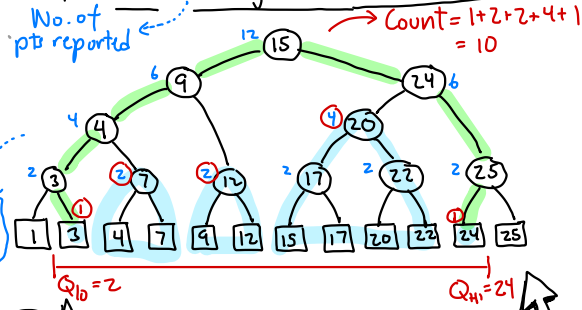


1-Dim Range Tree:



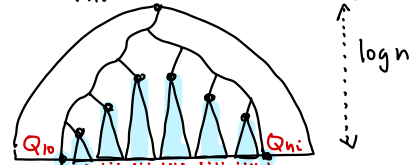
Approach:

- Balanced BST (eg. AVL, RB, ...)
- Assume **extended tree**
- Each node p stores no. of entries in subtree: $p.size$



Canonical Subsets:

- **Goal**: Express answer as disjoint union of subsets
- **Method**: Search for $Q_{l_0} + Q_{h_1} +$ take maximal subtrees



Recursive helper:

```
int range1Dx(Node p,
    Intv Q=[Qlo, Qhi], Intv C=[xo, xi])
```

initial call: range1Dx(root, Q, C_o)

Cases:

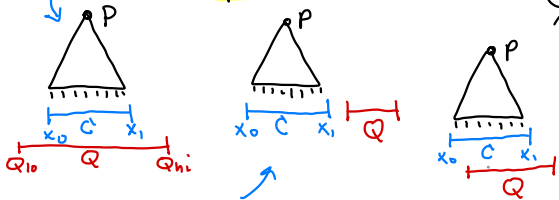
p is external:

- if p.pt.x ∈ Q → 1 else → 0

p is internal:

- C ⊆ Q ⇒ all of p's pts lie within query

→ return p.size



- C is disjoint from Q ⇒ none of p's pts lie in Q

→ return 0

- Else partial overlap

→ Recurse on p's children + trim the cell

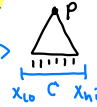
More details:

Given a 1-D range tree T:

- Let Q=[Q_{lo}, Q_{hi}] be query interval

- For each node p, define interval cell C=[x_o, x_i] s.t. all pts of p's subtree lie in C

- Root cell: C_o=[-∞, +∞]



Range Trees II

```
int range1Dx(Node p,
```

```
Intv Q, Intv C=[xo, xi]) {
```

```
if (p is external) return 1
```

```
else if (C ⊆ Q) return p.size
```

```
else if (Q ∩ C disjoint) return 0
```

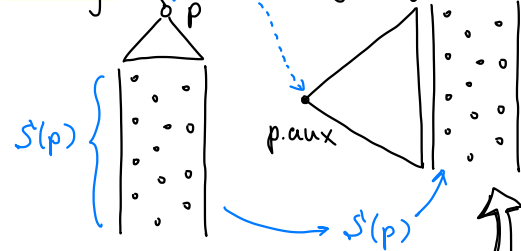
```
else return:
```

```
    range1Dx(p.left, Q, [xo, p.x])
```

```
    + range1Dx(p.right, Q, [p.x, xi])
```

x-range:

y-range:



2-D Range Searching:

- "layer" a range tree for x with range tree for y

- For each node p ∈ 1D-x tree, let S(p) = set of pts in p's subtree

- Def: p.aux: A 1D-y tree for S'(p)

Analysis:

Lemma: Given a 1-D range tree with n pts, given any interval Q, can compute O(log n) subtrees whose union is answer to query.

Thm: Given 1-D range tree...

can answer range queries in time O(log n) ... → (+k to report)

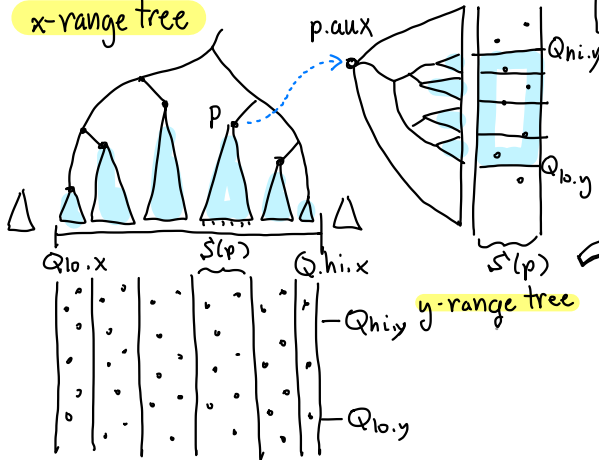
Answering Queries?

Given query range

$$Q = [Q_{lo.x}, Q_{hi.x}] \times [Q_{lo.y}, Q_{hi.y}]$$

- Run range1D_x to find all subtrees that contribute
- For each such node p, run range1D_y on p.aux
- Return sum of all result

x-range tree



Intuition: The x-layer finds subtrees p contained in x-range + each aux tree filters based on y .

2D Range Tree:

- Construct 1D range tree based on x coord for all pts
- For each node p :
 - Let $S(p)$ be pts of p 's tree
 - Build 1D range tree for $S(p)$ based on $y \rightarrow p.aux$
- Final structure is union of x -tree + $(n-1)$ y -trees

Range trees III

```
int range2D(Node p, Rect Q, Intv C=[x0, x1]) {
    if (p is external) return p.pt ∈ Q? 1 : 0
    else if (Q.x contains C) { // C ⊆ Q's x-projection
        [y0, y1] = [-∞, +∞] // init y-cell
        return range1Dy(p.aux, Q, [y0, y1])
    } else if (Q.x is disjoint of C) return 0
    else // partial x-overlap
        return range2D(p.left, Q, [x0, p.x])
        + range2D(p.right, Q, [p.x, x1])
}
```

Analysis:

Invoked $O(\log n)$ times - once per maximal subtree

Invoked $O(\log n)$ times - once for each ancestor of max subtree

Higher Dimensions?

- In d -dim space, we create d -layers
- Each recurses one dim lower until we reach 1-d search
- Time is the product:

$$\log n \cdot \log n \cdot \dots \log n = O(\log^d n)$$

Analysis: The 1D x search takes of $O(\log n)$ time + generates $O(\log n)$ calls to 1D y search

$$\Rightarrow \text{Total: } O(\log n \cdot \log n) = O(\log^2 n)$$

Hashing: (Unordered)

dictionary

- stores key-value pairs in **array table** $[0..m-1]$
- supports basic dict. ops. (insert, delete, find) in **$O(1)$ expected time**
- does not support ordered ops (getMin, findUp, ...)
- simple, practical, widely used

Overview:

- To store n keys, our table should (ideally) be a bit larger (e.g., $m \geq c \cdot n$, $c = 1.25$)
- **Load factor:**
 $\lambda = n/m$
- Running times increase as $\lambda \rightarrow 1$
- **Hash function:**
 $h: \text{Keys} \rightarrow [0..m-1]$
→ Should **scatter** keys random.
→ Need to handle **collisions**

Recap: So far, **ordered dicts.**

- insert, delete, find
 - **Comparison-based:** $<, =, >$
 - getMin, getMax, getK, findUp...
 - Query/Update time: $O(\log n)$
→ Worst-case, amortized, random.
- Can we do better? $O(1)$?

Hashing I

Good Hash Function:

- Efficient to compute
- Produce few collisions
- Use every bit in key
- Break up natural clusters

Eg. Java variable names: temp1, temp2, temp3

table:



$x \neq y$
but
 $h(x) = h(y)$

Universal Hashing:

Even better → randomize!

- Let H be a **family** of hash fns
 - Select $h \in H$ randomly
 - If $x \neq y$ then $\text{Prob}(h(x) = h(y)) = \frac{1}{m}$
- Eg. Let p - large prime, $a \in [1..p-1]$
 $b \in [0..p-1]$ **all random**
- $h_{a,b}(x) = ((ax + b) \bmod p) \bmod m$

Why "mod p mod m"?

- modding by a large prime scatters keys
- m may not be prime (e.g. power of 2)

Common Examples:

- **Division hash:**
 $h(x) = x \bmod m$
- **Multiplicative hash:**
 $h(x) = (ax \bmod p) \bmod m$
 a, p - large prime numbers
- **Linear hash:**
 $h(x) = ((ax + b) \bmod p) \bmod m$
 a, b, p - large primes

Assume keys can be interpreted as ints

Overview:

- Separate Chaining
- Open Addressing:
 - Linear probing
 - Quadratic probing
 - Double hashing

simpl./slow
↓
complex/fast

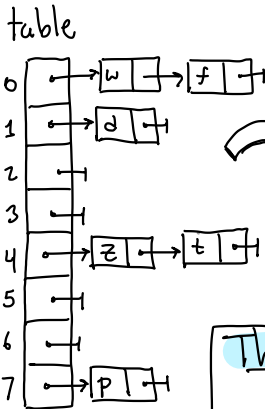
Separate Chaining:

table[i] is head of linked list of keys that hash to i.

Example:

Keys(x)	h(x)
d	1
z	4
p	7
w	0
t	4
f	0

m=8



Collision Resolution:

If there were no collisions hashing would be trivial!

- insert(x, v) → table[h(x)] = v
- find(x) → return table[h(x)]
- delete(x) → table[h(x)] = null

If $\lambda < \lambda_{min}$ or $\lambda > \lambda_{max}$? Rehash!

- Alloc. new table size = n/λ_0
- Compute new hash fn h
- Copy each x, v from old to new using h
- Delete old table

Hashing II

Token-based - See latex notes!

Thm: Amortized time for rehashing is $1 + (2\lambda_{max} / (\lambda_{max} - \lambda_{min}))$

How to control λ ?

- Rehashing: If table is too dense / too sparse, realloc. to new table of ideal size

Designer: $\lambda_{min}, \lambda_{max}$ - allowed λ values
 $\lambda_0 = \frac{\lambda_{min} + \lambda_{max}}{2}$ "ideal"

If $\lambda < \lambda_{min}$ or $\lambda > \lambda_{max}$...

Analysis: Recall load factor $\lambda = n/m$
 $n = \#$ of keys
 $m =$ table size

S_{sc} = Expected search time if x found (successful)
 U_{sc} = Expect. search time if x not found (unsuccessful)

Thm: $S_{sc} = 1 + \lambda/2$ $U_{sc} = 1 + \lambda$
 Proof: On avg. each list has $n/m = \lambda$
 success: 1 for head + half the list
 unsuccess: 1 " " + all the list

Open Addressing:

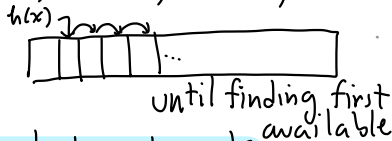
- Special entry ("empty") means this slot is unoccupied
- Assume $\lambda \leq 1$
- To insert key: check: $h(x)$ if not empty try
 - $h(x) + i_1$
 - $h(x) + i_2$
 - \vdots

$\langle i_1, i_2, i_3, \dots \rangle$ - Probe sequence

- What's the best probe sequence?

Linear Probing:

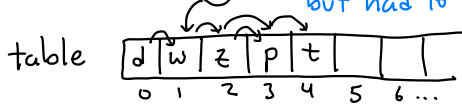
$h(x), h(x)+1, h(x)+2, \dots$



Simple, but is it good?

$x: d, z, p, w, t$

$h(x): 0, 2, 2, 0, 1$



Collision Resolution: (cont.)

- Separate chaining is efficient, but uses extra space (nodes, pointers, ...)
- Can we just use the table itself?

→ Open Addressing

Hashing III

Analysis:

Let S_{LP} = expected time for successful search

U_{LP} = " " unsuccessful "

$$\text{Thm: } S_{LP} = \frac{1}{2} \left(1 + \frac{1}{1-\lambda} \right)$$

$$U_{LP} = \frac{1}{2} \left(1 + \frac{1}{1-\lambda} \right)^2$$

Obs: As $\lambda \rightarrow 1$ times increase rapidly

Analysis: Improves secondary clustering

- Many fail to find empty entry (Try $m=4, j^2 \bmod 4 = 0 \text{ or } 1$ but not $2 \text{ or } 3$)
- How bad is it? It will succeed if $\lambda < 1/2$.

Thm: If quad. probing used + m is prime, the the first $\lfloor m/2 \rfloor$ probe locations are distinct.

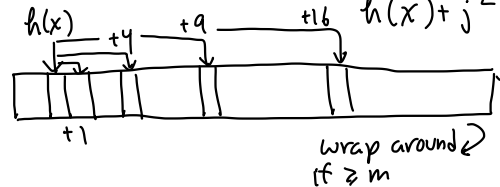
Pf: See latex notes.

Clustering

- Clusters form when keys are hashed to nearby locations
- Spread them out!

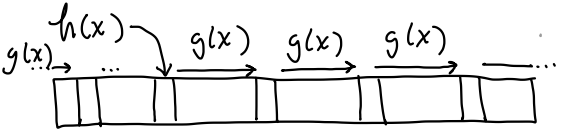
Quadratic Probing:

$h(x), h(x)+1, h(x)+4, h(x)+9, \dots, h(x)+j^2$



Double Hashing:
(Best of the open-addressing methods)

- Probe sequence det'd by second hash fn. - $g(x)$
 $h(x) + \{0, g(x), 2 \cdot g(x), 3 \cdot g(x) \dots\}$
 $[\text{mod } m]$



(until finding an empty slot)

Why does bust up clusters?
 Even if $h(x) = h(y)$ [collision] it is **very unlikely** that $g(x) = g(y)$
 \Rightarrow Probe sequences are entirely different!

Analysis: Defs:
 S_{DH}^v = Expected search time of doub. hash. if successful
 U_{DH} = Exp. if unsuccessful
 Recall: **Load factor** $\lambda = n/m$

Recap:
Separate Chaining:
 Fastest but uses extra space (linked list)
Open Addressing:
 Linear probing: } clustering
 Quadratic probing: }
 probing: }

Hashing IV

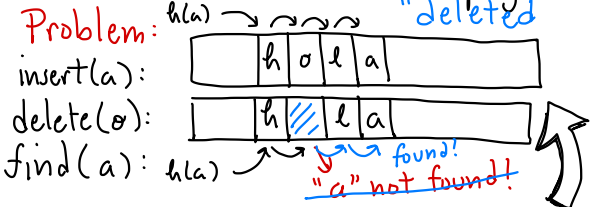
Thm: $S_{DH} = \frac{1}{\lambda} \ln(\frac{1}{1-\lambda})$
 $U_{DH} = 1/(1-\lambda)$

\rightarrow Proof is nontrivial (skip)

λ :	0.5	.075	0.95	0.99
U_{DH} :	2	4	20	100
S_{DH}^v :	1.39	1.89	3.15	4.65

very efficient!

Delete(x): Apply find(x)
 \rightarrow Not found \Rightarrow error
 \rightarrow Found \Rightarrow set to "empty"
 "deleted"



Find(x): Visit entries on probe sequence until:
 - found $x \Rightarrow$ return v
 - hit empty \Rightarrow return null

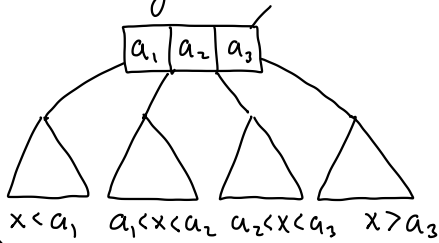
find(x) $h(x)$ \rightarrow Not found!
 empty

Dictionary Operations:

Insert(x,v): Apply probe sequence until finding first empty slot.
 - Insert (x,v) here.
 (If x found along the way \Rightarrow duplicate key error!)

Is this right??

Multiway Search Trees:



B-Tree:

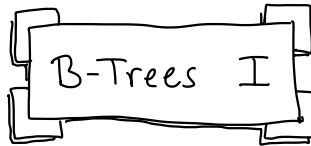
- Perhaps the most widely used search tree
- 1970 - Bayer & McCreight
- Databases
- Numerous variants

B-Tree: of order $m (\geq 3)$

- Root is leaf or has ≥ 2 children
- Non-root nodes have $\lceil m/2 \rceil$ to m children [null for leaves]
- k children $\Rightarrow k-1$ key-values
- All leaves at same level

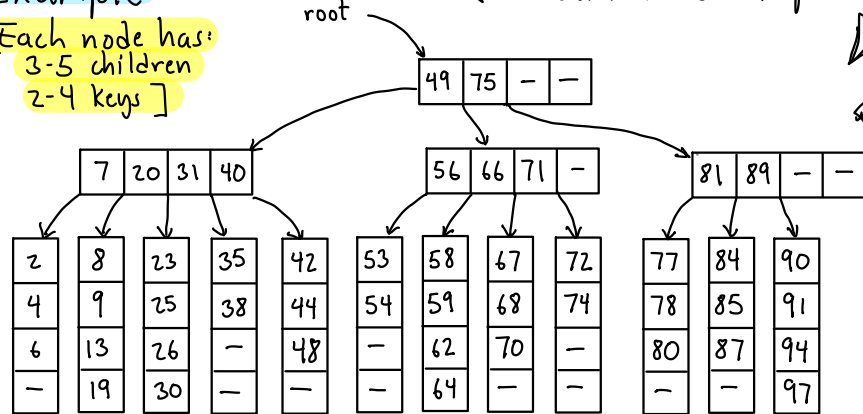
Secondary Memory:

- Most large data structures reside on disk storage
- Organized in blocks/pages
- Latency: High start-up time
- Want to minimize no. of blocks accessed



Example: $m=5$

[Each node has:
3-5 children
2-4 keys]



Node Structure: constant int $M=...$

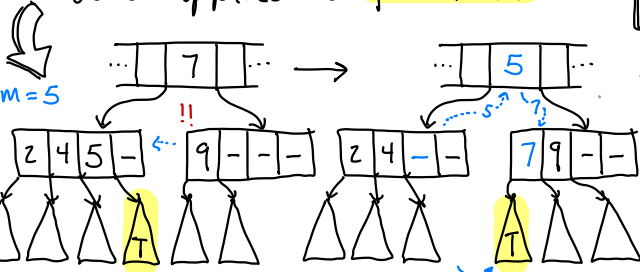
```
class BTreeNode {
    int nChild // no. of children
    BTreeNode child[M] // children
    Key key[M-1] // keys
    Value value[M-1] // values
}
```

Theorem: A B-tree of order m with n keys has height at most $(\lg n)/\gamma$, where $\gamma = \lg(m/2)$

(See full notes for proof)

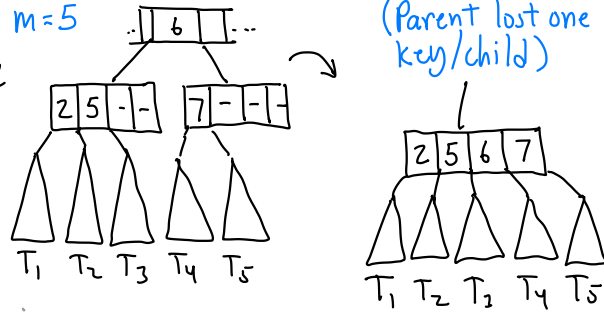
Key Rotation (Adoption)

- A node has **too few** children $\lceil m/2 \rceil - 1$
- Does either immediate sibling have **extra**? $\geq \lceil m/2 \rceil + 1$
- Adopt child from sibling & rotate keys
- When applicable - **preferred**.

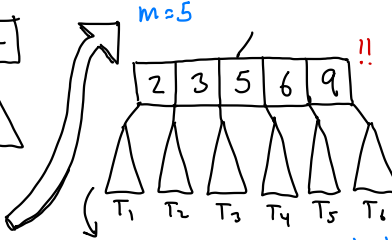


B-Tree restructuring:

- Generalizes 2-3 restructure
- Key rotation (Adoption)
- Splitting (insertion)
- Merging (deletion)



B-Trees II



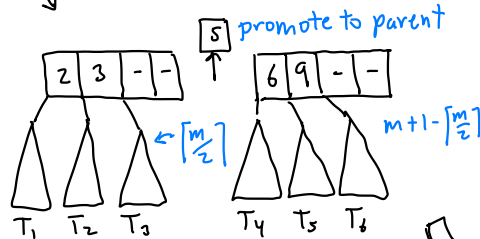
Lemma: For all $m \geq 2$,
 $\lceil m/2 \rceil \leq 2\lceil m/2 \rceil - 1 \leq m$
 \Rightarrow Resulting node is valid

Node Splitting:

- After insertion, a node has too many children... $m+1$
- We split into two nodes of sizes $m' = \lceil m/2 \rceil$ and $m'' = m+1 - \lceil m/2 \rceil$

Lemma: For all $m \geq 2$,
 $\lceil m/2 \rceil \leq m+1 - \lceil m/2 \rceil \leq m$

\Rightarrow **$m' + m''$ are valid node sizes**



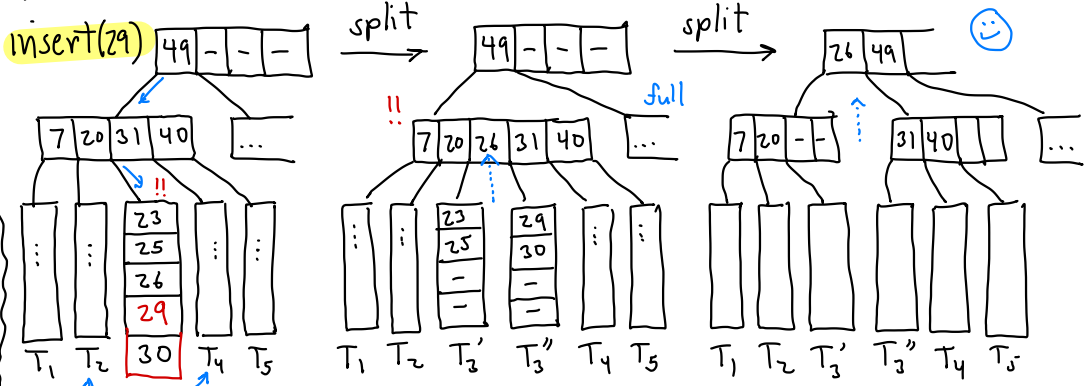
Node Merging:

- A node has too few children $\lceil m/2 \rceil - 1$
- Neither sibling has extra (both $\lceil m/2 \rceil$)
- Merge with either sibling to produce node with $(\lceil m/2 \rceil - 1) + \lceil m/2 \rceil$ child

Insertion:

- Find insertion point (leaf level)
- Add key/value here
- If node **overflow** (m keys, $m+1$ children)
 - Can either sibling take a child ($< m$)?
 - ⇒ **Key rotation** [done]
 - Else, **split**
 - Promotes key
 - If root splits add new root

Example: $m=5$

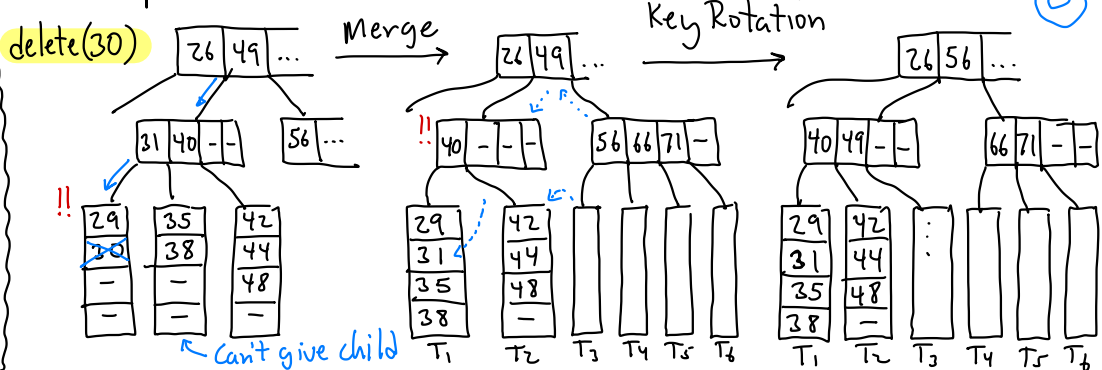


B-Trees III

Deletion:

- Find key to delete
- Find replacement/copy
- If **underfull** ($\lceil m/2 \rceil - 1$) child
 - If sibling can give child
 - **Key rotation**
 - Else (sibling has $\lceil m/2 \rceil$)
 - **Merge** with sibling
 - Propagates → If root has 1 child → collapse root

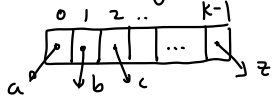
Example: $m=5$



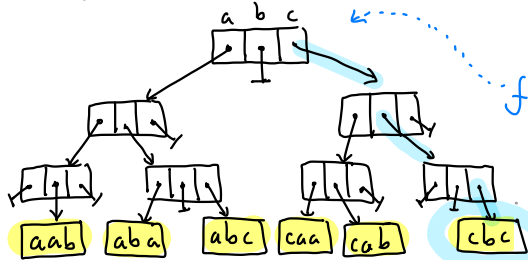
Tries: History

- de la Briandais (1959)
- Fredkin - "trie" from "retrieval"
- Pronounced like "try"

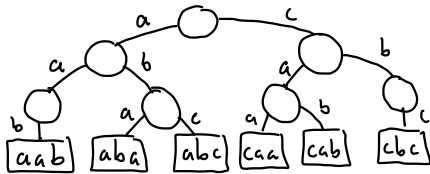
Node: Multiway of order k



Example: $\Sigma = \{a=0, b=1, c=2\}$
 Keys: $\{aab, aba, abc, caa, cab, cbc\}$



Same structure/Alt. Drawing



Digital Search:

- Keys are strings over some alphabet Σ
- E.g. $\Sigma = \{a, b, c, \dots\}$
 $\Sigma = \{0, 1\}$ Let $k = |\Sigma|$
- Assume chars coded as ints: $a=0, b=1, \dots, z=k-1$

Tries and Digital Search Trees I

Analysis:

Search: \sim length of query string $[O(1)$ time per node]

Space:

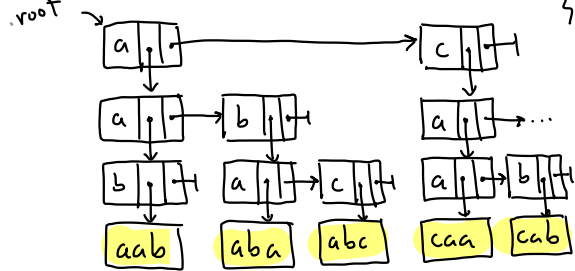
- No. of nodes \sim total no. of chars in all strings
- Space $\sim k \cdot (\text{no. of nodes})$

Large!

Analysis:

- **Space:** Smaller by factor k
- **Search Time:** Larger by factor of k

Example:



How to save space?

de la Briandais trees:

- Store 1 char. per node
- $\boxed{x} \rightarrow \neq x \Rightarrow$ try next char in Σ
 $= x \Rightarrow$ advance to next character of search string
- First-child/next-sibling

Patricia Tries:

- Improves trie by compressing degenerate paths
- PATRICIA = Practical Alg. to Retrieve Info. Coded in Alpha...
- Late 1960's: Morrison + Guchenberger
- Each node has **index field**, indicates which char to check next (Increase with depth)



Dealing with long Paths:

- To get both good spaces + query time efficiency, need to avoid long, degenerate paths.



Example:

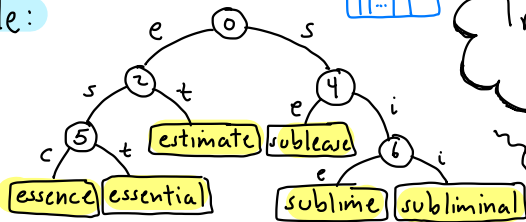
ID	String	Prefix	Identifier
S_5	ajam...	aj	aj
S_{10}	\$		
S_4	pajam...	paj	paj
S_9	a#	a#	a#
S_3	apaja...	ap	ap
S_8	ma#	ma#	ma#
S_2	mapaj...	map	map
S_7	ama#	ama#	ama#
S_1	amapaj...	amap	amap
S_6	jama#	j	j
S_0	pamapa...	pam	pam

(Increase with depth)



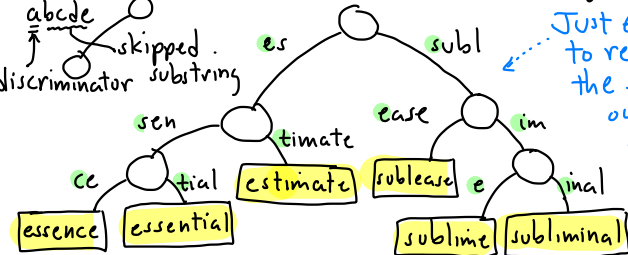
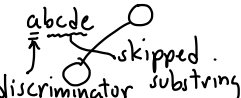
Example:

- essence
- essential
- estimate
- sublease
- sublime
- subliminal



Tries and Digital Search Trees II

Same data structure - Drawn differently



Just easier to read the strings out... same data struct.

Analysis:

- **Query time:** (Same as std trie) \sim search string length (may be less)

Space:

- **No. nodes:** \sim No. of strings (irres. of length)
- **Total space:** $K \cdot$ (No. of nodes) + (Storage for strings)

Example: $S = \text{pamapajama}\#$

- $S_{10} = \#$
- $S_9 = a\#$
- $S_8 = ma\#$
- $S_7 = ama\#$
- \vdots

Def: Substring identifier for

- S_i is shortest prefix of
- S_i unique to this string
- Eg. $ID(S_1) = \text{"amap"}$
- $ID(S_7) = \text{"ama\#"}$

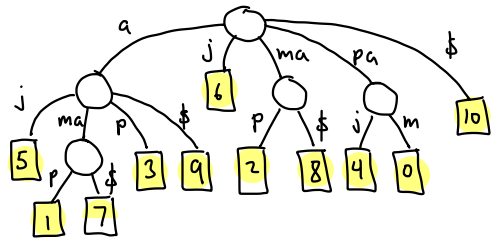
Suffix Trees:

- Given single large **text** S
- Substring queries: "How many occurrences of "tree" in CMSC 420 notes"

Notation: $S = a_0 a_1 a_2 \dots a_{n-1} \#$

- **Suffix:** $S_i = a_i a_{i+1} \dots a_{n-1} \#$ (special terminal)
- **Q:** What is minimum substring needed to identify suffix S_i ?

Example: $S = \overset{0}{p}\overset{1}{a}\overset{2}{m}\overset{3}{a}\overset{4}{p}\overset{5}{a}\overset{6}{j}\overset{7}{a}\overset{8}{m}\overset{9}{a}\overset{10}{\$}$



E.g. $ID(S, a) = \text{amap}$ $ID(S, a) = \text{ama\$}$

Substring Queries:

How many occurrences of t in text?

- Search for target string t in trie
- if we end in internal node (or midway on edge) - return no. of extern. nodes in this subtree
- else (fall out at extern. node)
 - compare target with string
 - if matches - found 1 occurrence
 - else - no occurrences

Example:

Search("ama") → End at intern node Report: Zocci.

Search("amapaj") → End at extern node Go to S_i + verify

Suffix Trees (cont.)

S - text string $|S| = n$

$S_i = i^{\text{th}}$ suffix

Substring ID = min substr. needed to identify S_i

A suffix tree is a Patricia trie of the $n+1$ substring identifiers

Tries and Digital Search Trees III

Analysis:

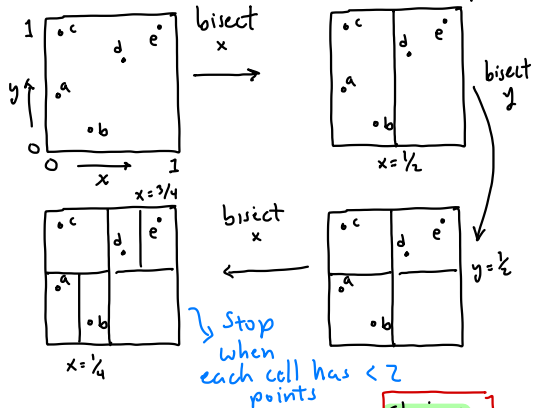
- Space: $O(n)$ nodes $O(n \cdot k)$ total space ($k = |S| = O(l)$)
- Search time: \sim total length of target string
- Construction time: $O(n \cdot k)$ [nontrivial]

PR k-d tree: Can be used for answering same queries as point kd-tree (orth. range, near. neigh)

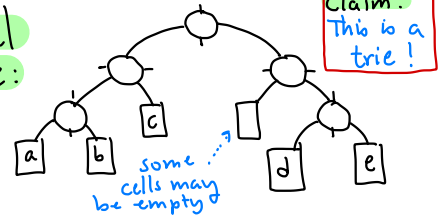
Geometric Applications:

PR kd-Tree: kd-tree based on midpoint subdivision

Assume points lie in unit square



Final tree:



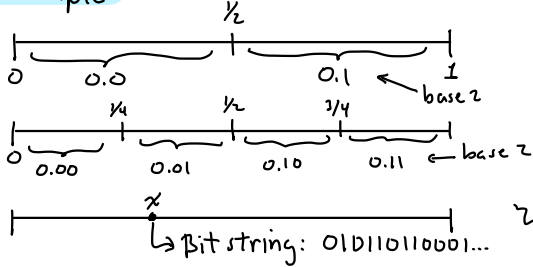
Binary Encoding:

- Assume our points are scaled to lie in **unit square**
 $0 \leq x, y < 1$ (can always be done)
- Represent each coordinate as **binary fraction**:

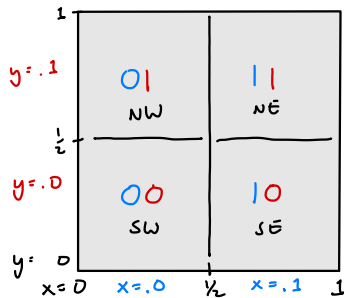
$$x = 0.a_1a_2a_3\dots \quad a_i \in \{0,1\}$$

$$x = \sum a_i \cdot \frac{1}{2^i}$$

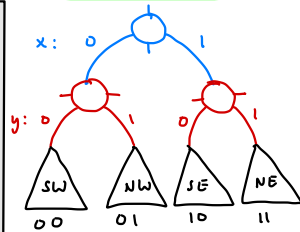
Example:



How do we extend to 2-D?



PR kd-tree



Bit Interleaving:

Given a point $p=(x,y)$
 $0 \leq x, y < 1$

let: $x = 0.a_1a_2\dots$ in binary
 $y = 0.b_1b_2\dots$

Define:

$$\phi(x,y) = a_1b_1a_2b_2a_3b_3\dots$$

Called **Morton Code** of p

PR kd-Tree \equiv Trie ??

- Approach: Show how to map any point in \mathbb{R}^2 to bit string
- Store bit strings in a trie (alphabet $\Sigma = \{0,1\}$)
- Prove that this trie has same structure as kd-tree

Tries and Digital Search Trees IV

Further Remarks:

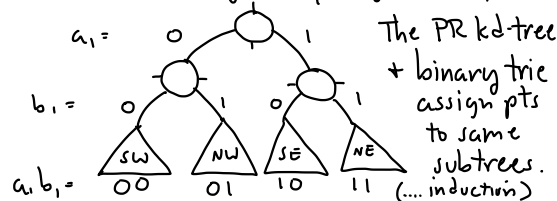
- Techniques for efficiently encoding, building, serializing, compressing... tries **apply immediately** to PR kd-tree
- Can generalize to **any dimension**
 $x = 0.a_1a_2\dots$
 $y = 0.b_1b_2\dots$
 $z = 0.c_1c_2\dots$

$\phi = a_1b_1c_1a_2b_2c_2\dots$

Lemma: Given a pt set $P \subseteq \mathbb{R}^2$ (in unit square $[0,1]^2$) let $P = \{p_1, \dots, p_n\}$ where $p_i = (x_i, y_i)$
 Let $\Phi(P) = \{\phi(p_1), \phi(p_2), \dots, \phi(p_n)\}$ (n binary strings)
 Then the PR kd-tree for P is equivalent to binary trie for $\Phi(P)$.

Proof: By induction on no. of bits

Let $x = 0.a_1a_2\dots$ $y = 0.b_1b_2\dots$
 and consider just $\phi(x,y) = a_1b_1\dots$



Deallocation Models:

Explicit: (C, C++)

- programmer deletes
- may result in **leaks** if not careful

Implicit: (Java, Python)

- runtime system deletes
- **Garbage collection**
- Slower runtime
- Better memory compaction



What happens when you do

- new (Java)
- malloc/free (C)
- new/delete (C++) ?

Runtime System Mem. Mgr.

- **Stack** - local vars, recursion
- **Heap** - for "new" objects

Don't confuse with heap data structure/heap sort



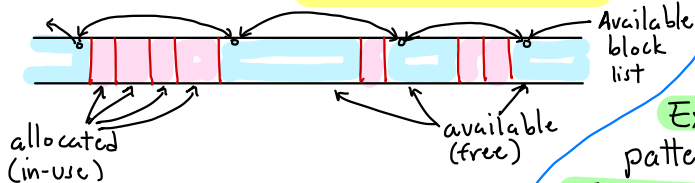
Memory Management I

Explicit Allocation/Deallocation

- Heap memory is split into **blocks** whenever requests made

Available blocks:

- merged when contiguous
- stored in **available block list**



Fragmentation:

- Results from repeated allocation + deallocation
- (**Swiss-cheese effect**)

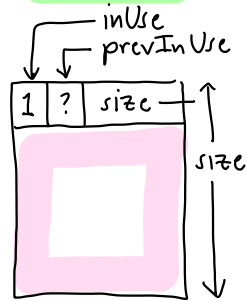


External: Caused by pattern of alloc/dealloc

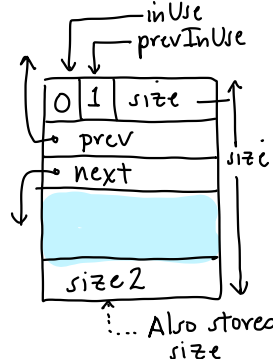
Internal: Induced by mem. manage. policies (not user)

Block Structure:

Allocated:



Available:



Guide:

prevInUse: 1 if prev. contig. block is allocated

prev/next: links in avail. list

size/size2: total block size (includes headers)



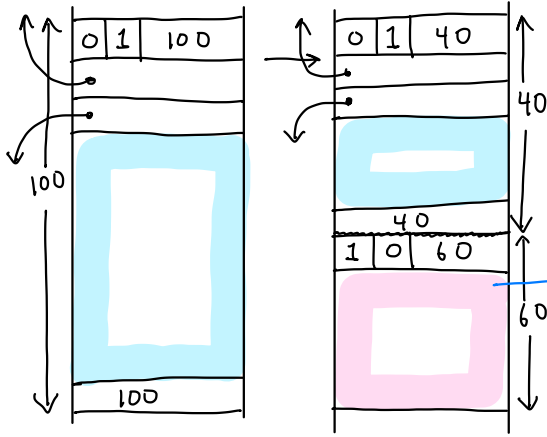
How to select from available blocks?

- **First-fit:** Take first block from avail. list that is large enough

- **Best fit:** Find closest fit from avail list

Surprise: First-fit is usually better - faster + avoids small fragments

Example: Alloc $b=59$



Allocation: $\text{malloc}(b)$

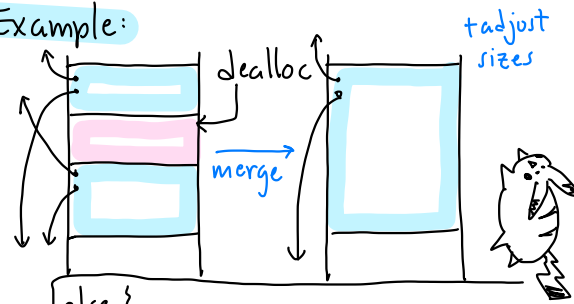
- Search avail. list for block of size $b' \geq b+1$
- If b' close to b : alloc entire block (unlink from avail list)
- Else: split block

Memory Management II

Deallocation:

- If prev + next contiguous blocks are allocated \rightarrow add this to avail
- Else - merge with either/both to make max. avail block

Example:



Some C-style pointer notation

void^* - pointer to generic word of memory

Let p be of type void^* :

$p+10$ - 10 words beyond p

$*(p+10)$ - contents of this

Let p point to head of block:

$p.\text{inUse}$, $p.\text{prevInUse}$, $p.\text{size}$

- We omit bit manipulation

$*(p+p.\text{size}-1)$ - references last word in this block



$(\text{void}^*) \text{alloc}(\text{int } b) \{$

$b+=1$ // add +1 for header

$p = \text{search avail list for block}$

$\text{size} \geq b$

if ($p == \text{null}$) Error- Out of mem!

if ($p.\text{size} - b \leq \text{TOO_SMALL}$)

 | unlink p from avail. list

 | $q = p$

else (continued)

else {

$p.\text{size} -= b$ // remove allocation

$*(p+p.\text{size}-1) = p.\text{size}$ // size 2

$q = p + p.\text{size}$ // start of new block

$q.\text{size} = b$

$q.\text{prevInUse} = 0$ // new block header

$q.\text{inUse} = 1$

$(q+q.\text{size}).\text{prevInUse} = 1$

// update prevInUse for next contig. block

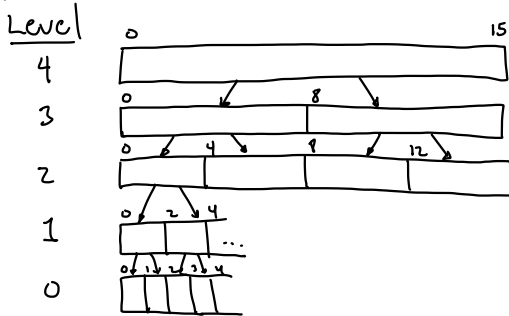
return $q+1$ // skip over header

}

Buddy System:

- Block sizes (including headers) are power of 2
- Requests are rounded up (internal fragmentation)
- Block size 2^k starts at address that is multiple of 2^k
- k = level of a block

Structure:



In practice: There is a minimum allowed block size

Buddy system only allows allocations aligning with these blocks



Coping with External Fragmentation

- Unstructured allocation can result in severe external fragmentation
- Can we compress? Problem of pointers
- By adding more structure we can reduce extern frag. at cost of internal frag.

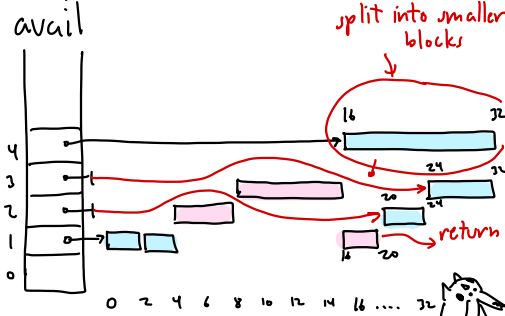
Memory Management III

Merging:

- When two adjacent blocks are available, we don't always merge them
- Must have same size: 2^k
- Must be buddies - siblings in this tree structure

Def: $buddy_k(x) = \begin{cases} x + 2^k & \text{if } 2^{k+1} \text{ divides } x \\ x - 2^k & \text{otherwise} \end{cases}$
 $\equiv buddy_k(x) = (1 \ll k) \oplus x$ [Bit manipulation]

Example: $alloc(2)$ ^{round up} $\rightarrow alloc(4)$



Allocation: $alloc(b)$

- $k = \lceil \lg(b+1) \rceil$ ^{add +1 for header}
- if $avail[k]$ non empty - return entry + delete
- else: find $avail[j] \neq \emptyset$ for $j > k$
- split this block

Big Picture:

- Avail list is organized by level: $avail[k]$
- Block header structure same as before except: $prevInUse$ } not needed size 2

