

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 operations 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

Common:

$O(1)$: constant time ☺
[HashMap]

$O(\log n)$: log-time (good)
[Binary search]

$O(n^p)$: p=constant: poly time
 $O(\sqrt{n})$

Asymptotic: "Big-O"

- Ignore constants
- Focus on large n

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

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

Our approach:

- **Theoretical:** Algorithms + Asymptotic Analysis

- **Practical:** Implementation + practical efficiency

Asymptotic Analysis:

- Run time as function of n : no. of items

- Worst-case, average case, randomized,...

- **Amortized** - average over series of ops.

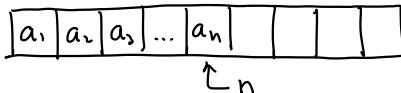
Linear List ADT:

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

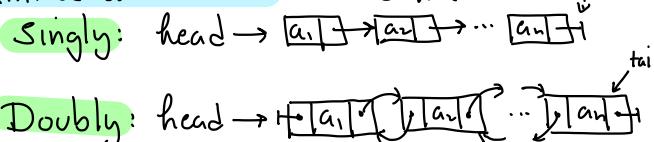
- `init()` - create an empty list
- `get(i)` - returns a_i
- `set(i, x)` - sets i^{th} element to x
- `insert(i, x)` - inserts x prior to i^{th} (moving others back)
- `delete(i)` - deletes i^{th} item (moving others up)
- `length()` - returns num. of items

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

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

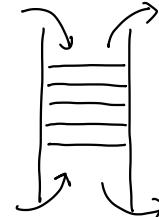
Deque ("deck"): Can insert or delete from either end

Stack: All access from one side

\downarrow (top) - push + pop



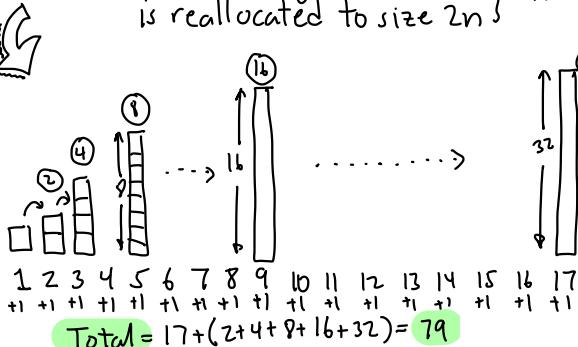
Queue: FIFO list: enqueue inserts at tail and dequeue deletes from head



Cost model (Actual cost)

Cheap: No reallocation \rightarrow 1 unit

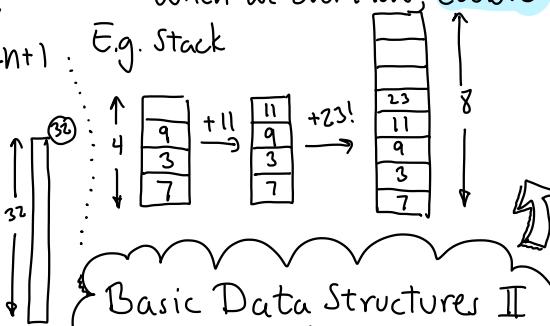
Expensive: Array of size $n \Rightarrow 2n+1$
is reallocated to size $2n$



Dynamic (Sequential) Allocation

- When we overflow, double

E.g. 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 $\leq 5 \cdot m$]

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

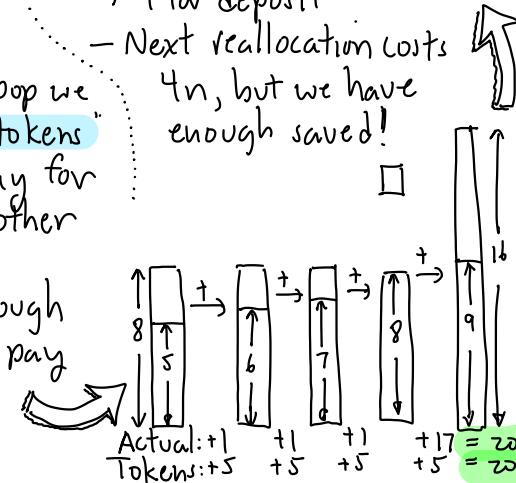
- During this time we collect at least $5n$ tokens
 $\rightarrow 1$ for each op
 $\rightarrow 4$ for deposit

- Next reallocation costs $4n$, but we have enough saved!

□

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.



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 $n \rightarrow \lceil n^{1.5} \rceil$

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

Leave as exercise 
 (Spoiler alert!)

Fixed increment → no

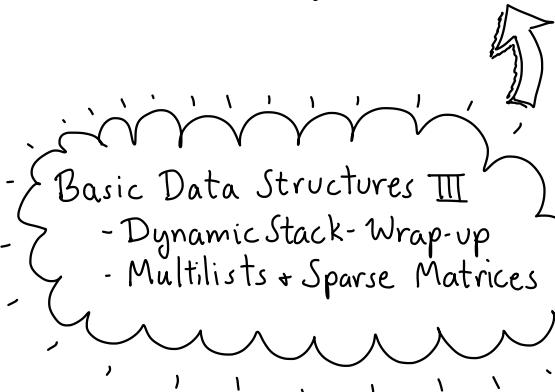
Fixed factor → yes

Squaring → yes

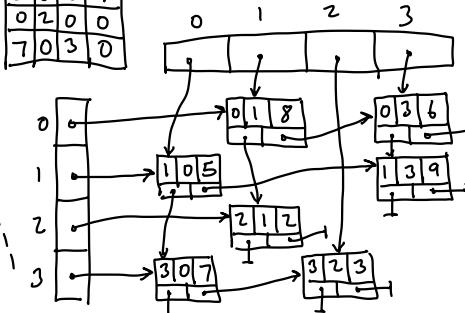
Dynamic Stack:

- Showed doubling \Rightarrow Amortized $O(1)$

- Other strategies?



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



Node:

row	col	value

row Next col Next

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

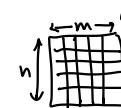
Multilists: Lists of lists

head → [0] → [a] → [b] →



Sparse Matrices:

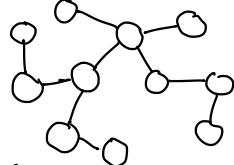
An $n \times m$ matrix has $n \cdot m$ entries and takes (naively) $O(n \cdot m)$ space



Sparse matrix: Most entries are zero

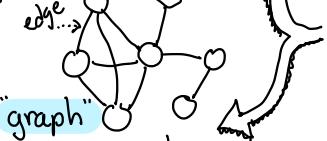
Tree (or "Free Tree")

- undirected
- connected
- acyclic graph

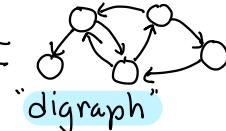


Undirected

node
edge



Directed



"digraph"

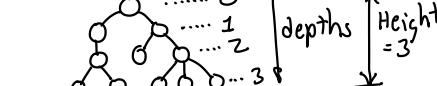
Graph: $G = (V, E)$

V = finite set of vertices
(nodes)

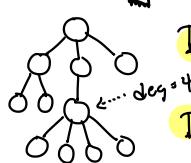
E = set of edges
(pairs of vertices)

Depth: path length from root

Height: (of tree) max depth

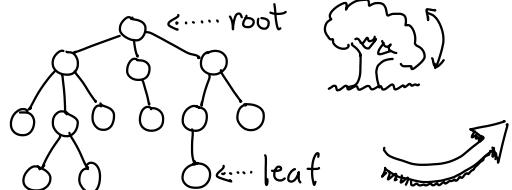


Degree (of node): number of children



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

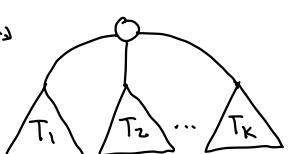
Rooted tree: A free tree with root node



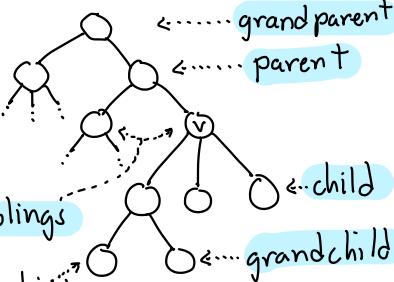
Formal definition:

Rooted tree: is either

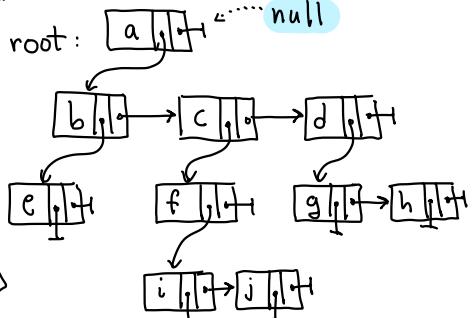
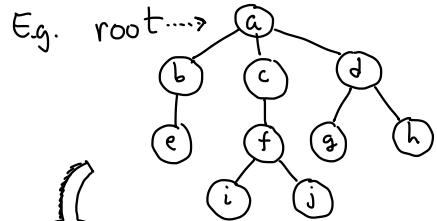
- single node (root)
- set of one or more rooted trees ("subtrees") joined to a common root



"Family" Relations



leaf: no children

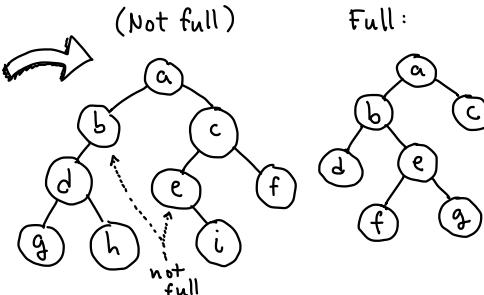
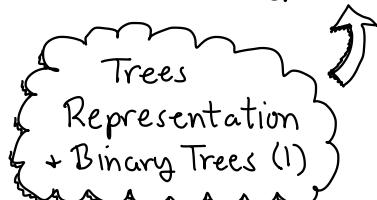
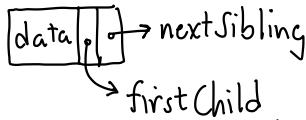


called the **Binary representation**

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

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

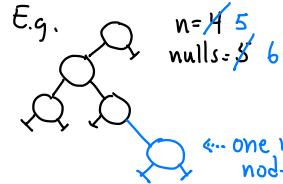
Node structure:



Full: Every non-leaf node has 2 children

Wasted space?

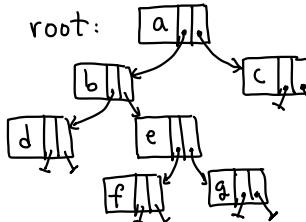
Theorem: A binary tree with n nodes has $n+1$ null links



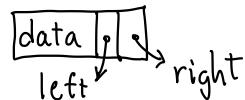
In Java: class BTNode<E> {

```

E data;
BTNode<E> left;
BTNode<E> right;
...
}
```



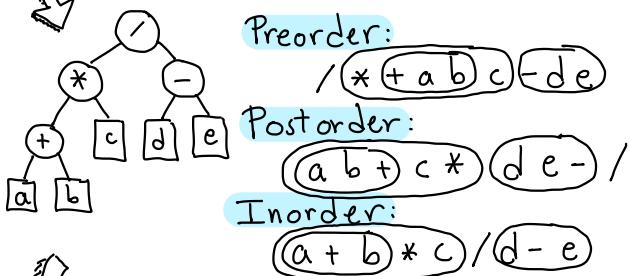
Node structure:



```

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
}

```

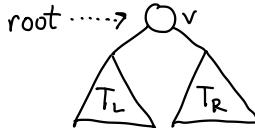


Those wasteful null links...

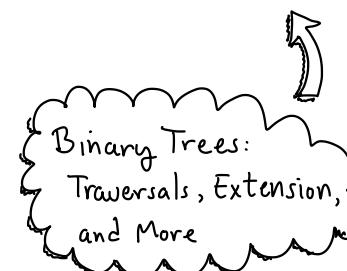
Extended binary tree: Replace each null link with a special leaf node: external node

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

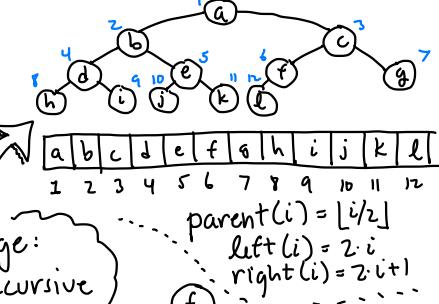
Binary Tree Traversals (can be generalized)



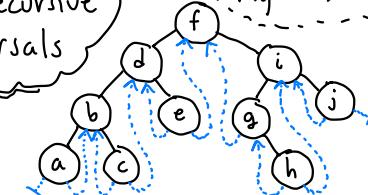
- process/visit v
- traverse T_L } recursive
- traverse T_R



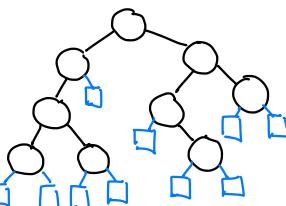
Complete Binary Tree: All levels full (except last)



Challenge:
Nonrecursive traversals



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



Observation: Every extended binary tree is full

Another way to save space...

Threaded binary tree:

Store (useful) links in the null links. (Use a mark bit to distinguish link types.)

E.g. Inorder Threads:

Null left → inorder predecessor
 Null right → " successor

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: $<$, $>$, $=$



Sequential Allocation?

- Store in array sorted by key

→ **Find**: $O(\log n)$ by binary search

→ **Insert/Delete**: $O(n)$ time



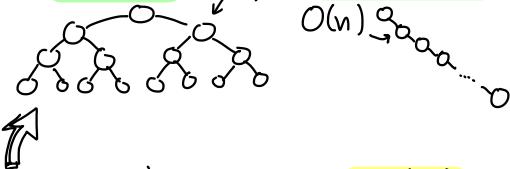
Can we achieve $O(\log n)$ time for all ops? **Binary Search Trees**

Binary Search Trees I

- Basic definitions
- Finding keys

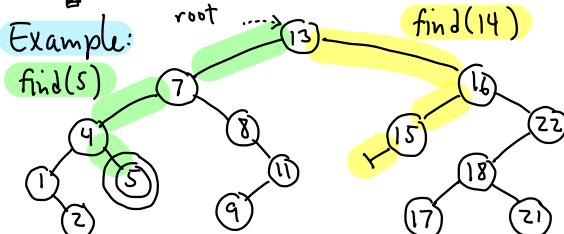
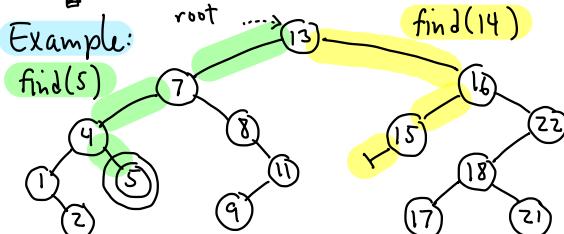
Efficiency: Depends on tree's height

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



Example:

find(5)

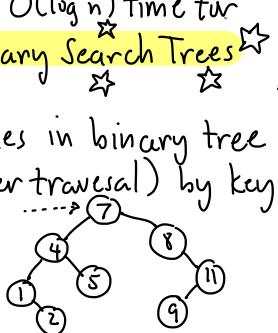


Find: How to find a key in the tree?

- 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!

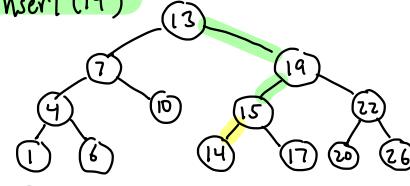


Idea: Store entries in binary tree sorted (inorder traversal) by key



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

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){}

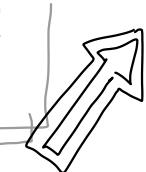
```

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



Delete (Key x)

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

How?

Why did we do:

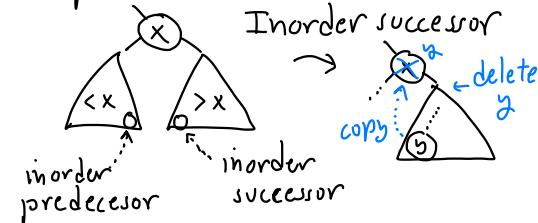
$p.left = \text{insert}(x, v, p.left)$?

p_1 $\text{insert}(14)$ \rightarrow p_2 $\text{new BSTNode}(14)$ \rightarrow p_2 $\text{return } p_2$

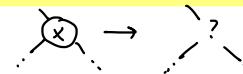
$p_1.\text{left} = \text{insert}(14, v, p_1.\text{left})$

Be sure you understand this!

Replacement Node?



3. \otimes has two children



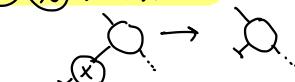
Find replacement node

y , copy to \otimes , and then delete y

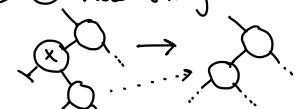


3 cases:

① \otimes is a leaf



② \otimes has single child

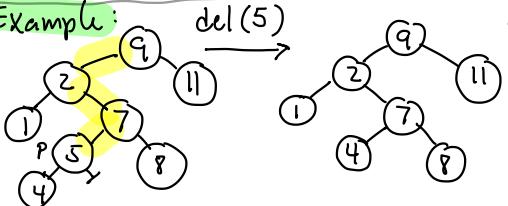


```

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

```

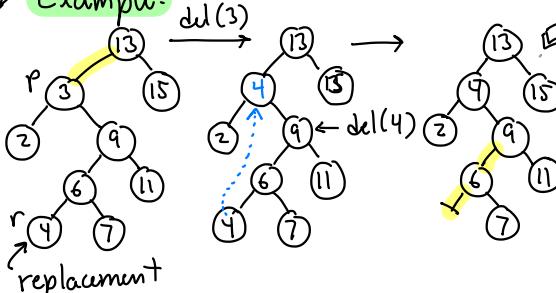
BSTNode findReplacement(BSTNode p) {
    BSTNode r = p.right
    while (r.left != null)
        r = r.left
    return r
}

```

Binary Search Trees III

- deletion
- analysis
- Java

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;
```

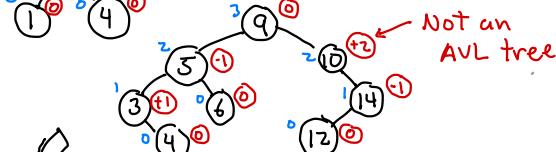
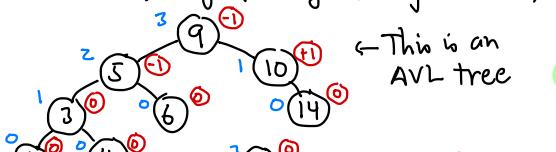
Data (private)

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

Public
members
(invoke
helpers)

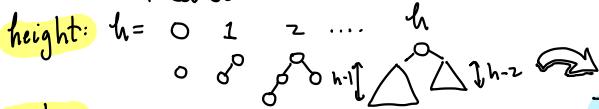
Balance factor:

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



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

Worst cases:



height: h	0	1	2	\dots	h
nodes:	0	1	2	\dots	n

$$n = 1, 2, 4, 7, 12, 20, \dots$$

$$n+1 = 2, 3, 5, 8, 13, 21, \dots$$

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

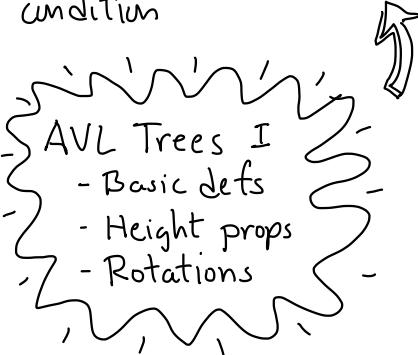
Conjecture: Min no. of nodes in AVL

tree of height h is $F_{h+3}-1$

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



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

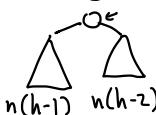
Proof: (Induct. on h)

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

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

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

$h \geq 2$:



BSTNode rotateRight(BSTNode p){

BSTNode q = p.left

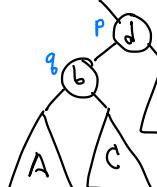
p.left = q.right

q.right = p

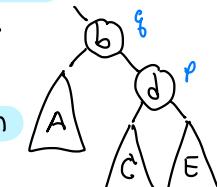
return q

}

How to maintain the AVL property?



$$A < b < C < j < E$$



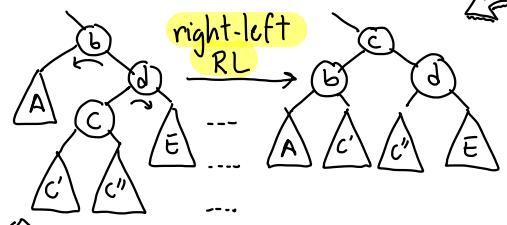
$$A < b < C < j < E$$

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

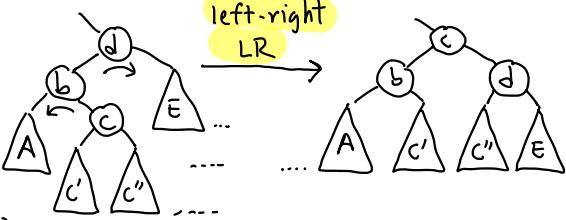
Proof: Fact: $F_h \approx \varphi^h / \sqrt{5}$ where

$$\varphi = (1 + \sqrt{5})/2 \quad \text{"Golden ratio"}$$

$$\begin{aligned} n &\geq \varphi^{h+3} = c \cdot \varphi^h \Rightarrow h \leq \log_{\varphi} n + c' \\ &\Rightarrow h \leq \log_2 n / \log_2 \varphi \\ &= O(\log n) \quad \square \end{aligned}$$



Double rotations:



`BSTNode rotateLeftRight(BSTNode p)`
 $p.left = \text{rotateLeft}(p.left)$
 return `rotateRight(p)`

AVL Tree:

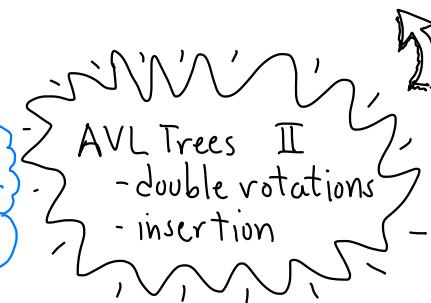
AVL Node: Same as BSTNode (from Lect 4) but add: int height

Utilities:

`int height(AVLNode p)`
 $\begin{cases} p == \text{null} \rightarrow -1 \\ \text{o.w. } \rightarrow p.height \end{cases}$

`void updateheight(AVLNode p)`
 $p.height = 1 + \max(\text{height}(p.left), \text{height}(p.right))$

`int balanceFactor(AVLNode p)`
 $\text{return height}(p.right) - \text{height}(p.left)$



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
```

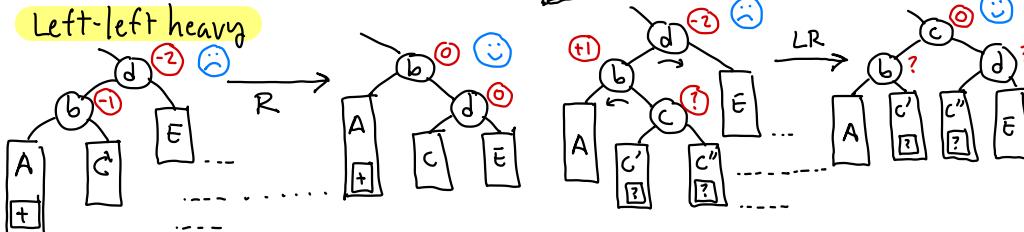
AVLNode insert(Key x, Value v, AVLNode p){
 $\begin{cases} p == \text{null} \rightarrow p = \text{new AVLNode}(x, v) \\ \text{else if } (x < p.key) \rightarrow p.left = \text{insert}(x, v, p.left) \\ \text{else if } (x > p.key) \rightarrow p.right = \text{insert}(x, v, p.right) \\ \text{else throw - Error - Duplicate!} \end{cases}$
 return rebalance(p)}

Find: Same as BST.

Insert: Same as BST but as we "back out" rebalance

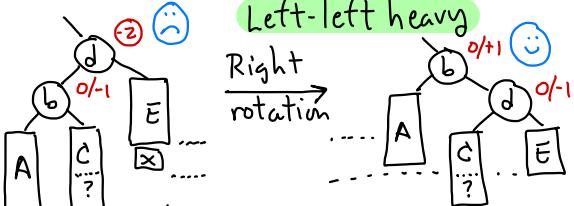
How to rebalance? Bal = -2

Left-left heavy

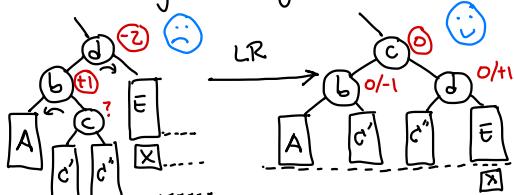


Left-right heavy:

Cases: Balance factor -2



Left-right heavy



Deletion: Basic plan

- Apply standard BST deletion

- find key to delete

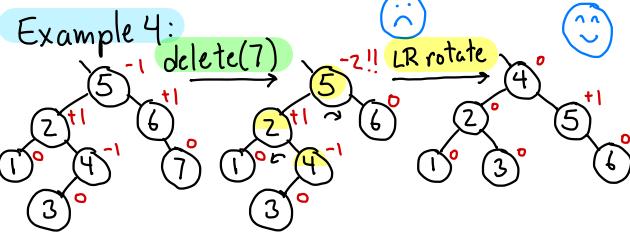
- find replacement node

- copy contents

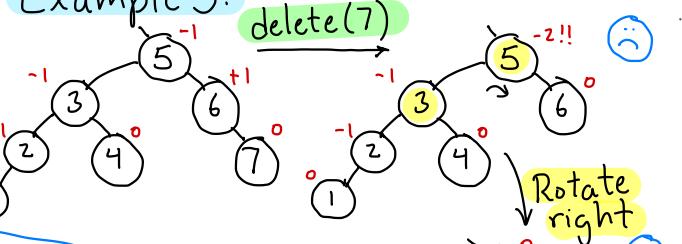
- delete replacement

- rebalance

Example 4:



Example 3:

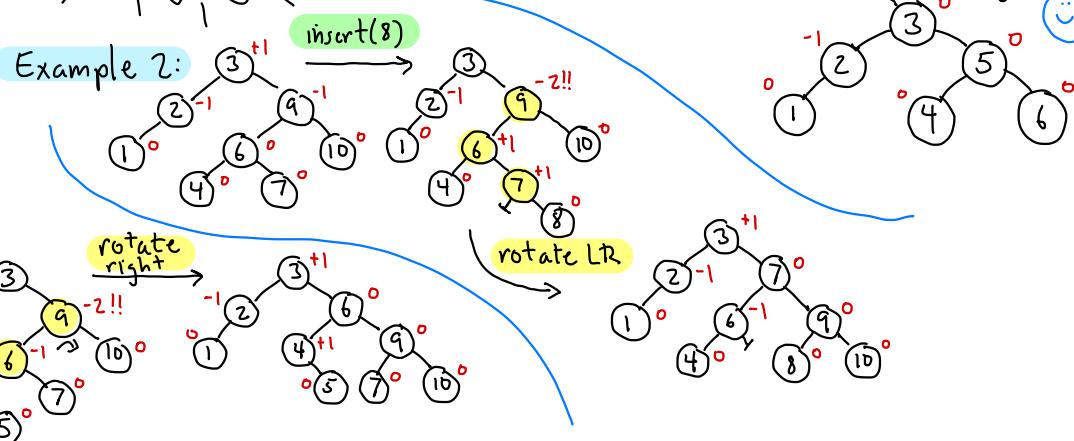


AVLNode delete (Key x, AVLNode p)

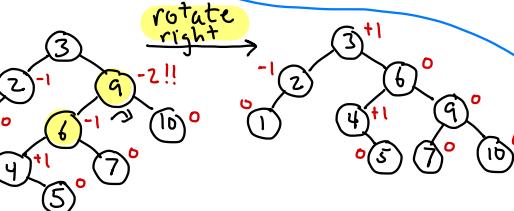
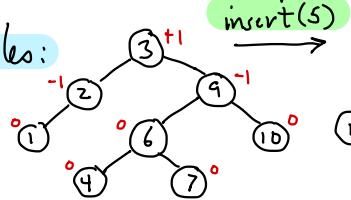
: same as BST delete

: return rebalance(p)

Example 2:



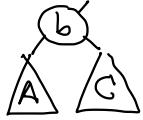
Examples:



Node types:

2-Node

1 key
2 children



3-Node

2 keys
3 children



Recap:

AVL: Height balanced
Binary

2-3 tree: Height exact
Variable width



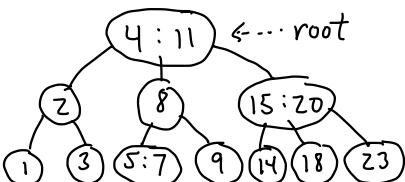
Def: A 2-3 tree of height h is either:

- 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 $h-1$.



Example:

2-3 tree of height 2



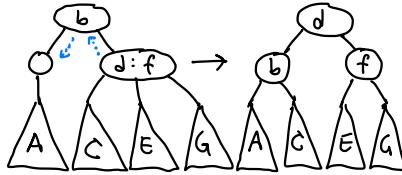
Recap:

AVL: Height balanced
Binary

2-3 tree: Height exact
Variable width

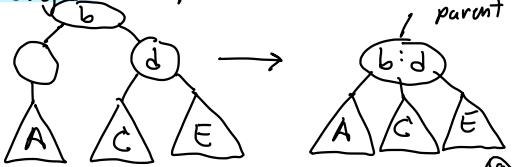
Adoption
(Key-Rotation)

$$1+3 = 2+2$$



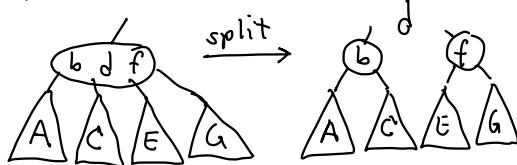
Merge:

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



steal b from parent

Split: $4 \rightarrow 2+2$



insert in parent
↑
d

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

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

How to maintain balance?

- Split
- Merge
- Adoption (Key rotation)

Conceptual tool:

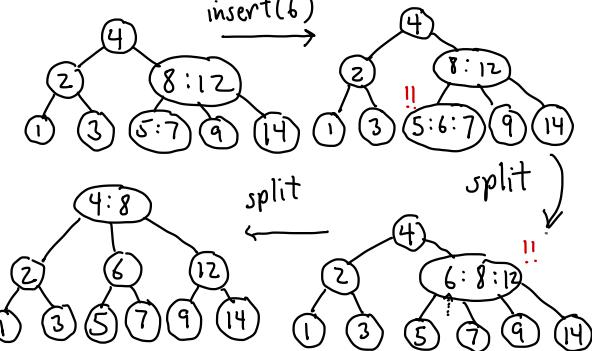
We'll allow 1-nodes + 4-nodes temporary

1-node



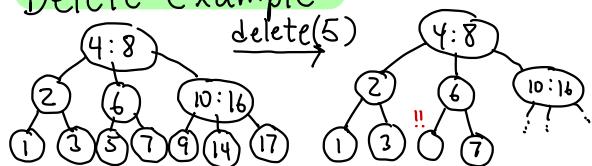
Insertion example:

insert(6)



Delete Example:

delete(5)



Deletion remedy:

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

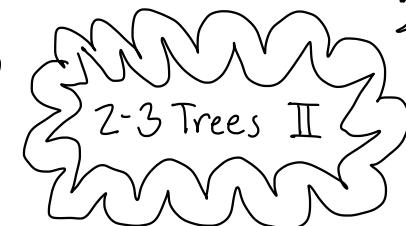
Dictionary operations:

Find - straightforward

Insert - find leaf node

where key "belongs"
+ add it (may split)

Delete - find / replacement/
merge or adopt



Implementation?

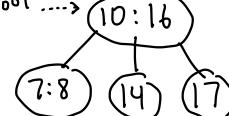
```
class TwoThreeNode {
```

```
int nChildren
```

```
TwoThreeNode children[3]
```

```
Key key[2]
```

new
root



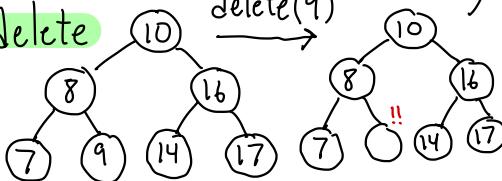
merge



merge

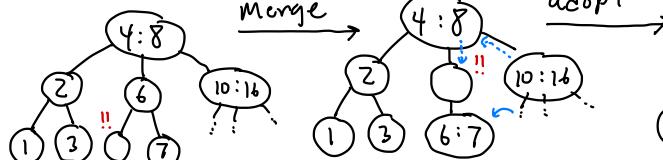
Another delete example:

delete(9)

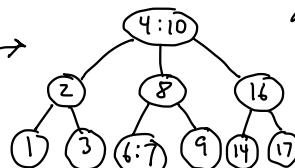


Example (continued)

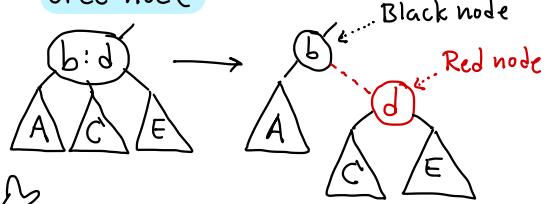
merge



adopt

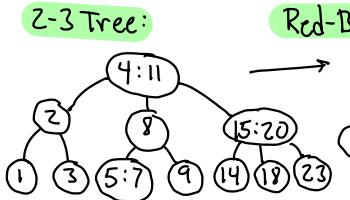


Encoding 3-node as binary tree node

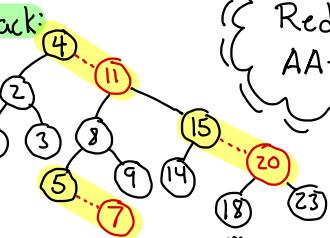


Example:

2-3 Tree:



Red-Black:



Some history:

2-3 Trees: Bayer 1972

Red-black Trees: Guibas + Sedgewick 1978 (a binary variant of 2-3)

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

Red-Black and AA-Trees I

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

AA-Trees: Simpler to code

- No null pointers: Create a sentinel node, nil, and all nulls point to it \rightarrow 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:

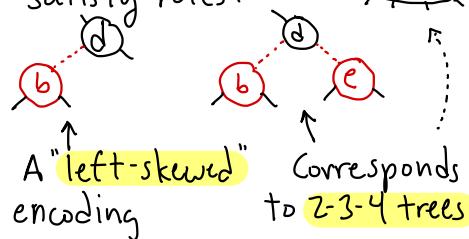
- ⑥ Each red node can arise only as right child (of a black node)

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?

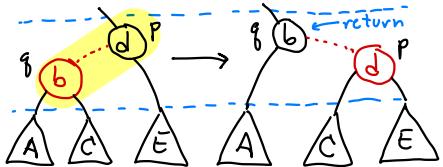
Nope! Alternatives that satisfy rules:



Restructuring Ops:

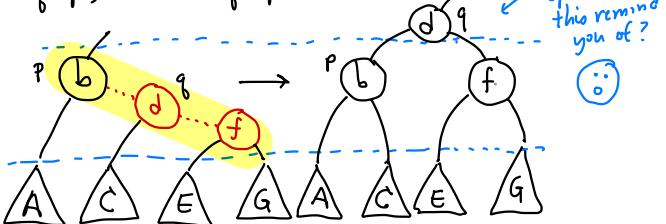
Skew: Restore right skew

→ If black node has red left child, rotate



How to test? $p.left.level == p.level$

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.

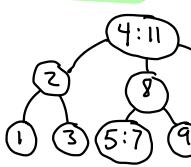


How to test?

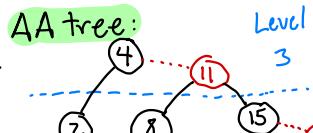
$p.level == p.right.level == p.right.right.level$
not needed (levels are monotone)

Example:

2-3 Tree:



AA tree:



Level
3
2
1
0

all to nil

nil

Red-Black + AA Trees II

insert(5)

?

!

?

!

?

!

?

!

?

!

?

!

?

!

?

!

?

!

?

!

?

!

?

!

?

!

?

!

?

!

AA Insertion:

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

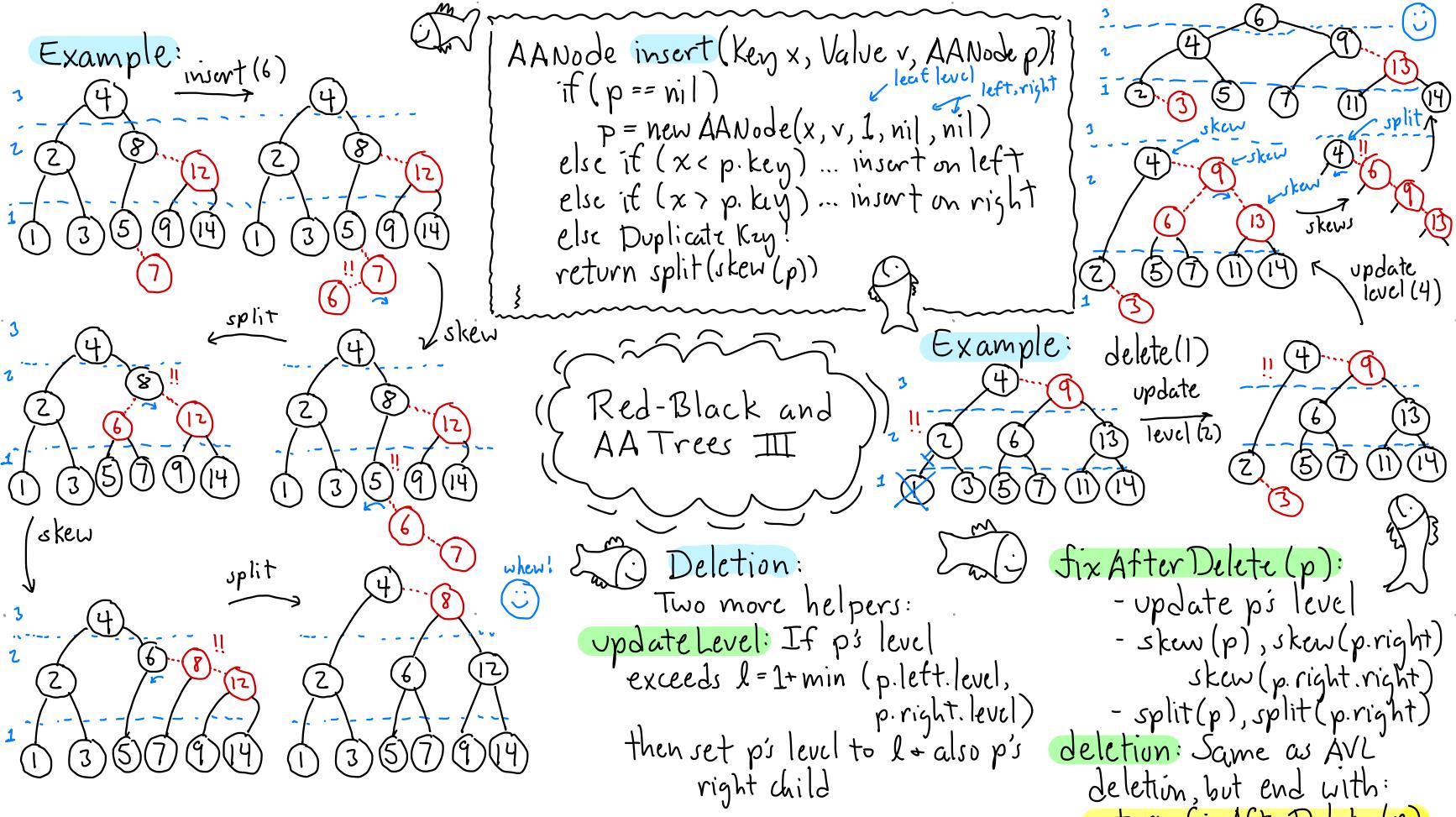
AA Node split (AANode p)

```
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
```

← all okay

```
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
```

← everything's fine



History:

1989: Seidel + Aragon

[Explosion of randomized algorithms]

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

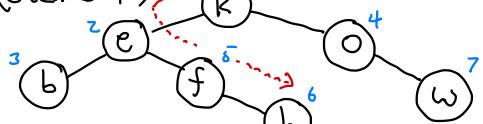
McCreight 1980

Intuition:

- Random insertion into BSTs
⇒ $O(\log n)$ 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

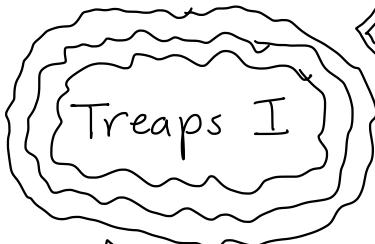
(Std. BST)



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



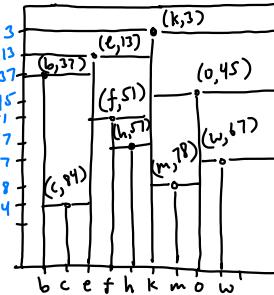
Obs: In a standard BST, keys are by inorder + insert times are in heap order (parent < child)



Geometric Interpretation:

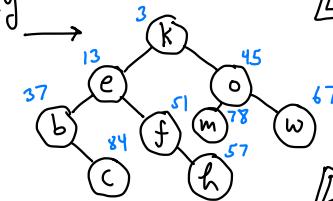
key → x
priority → y

McCreight's Priority Search Tree



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:

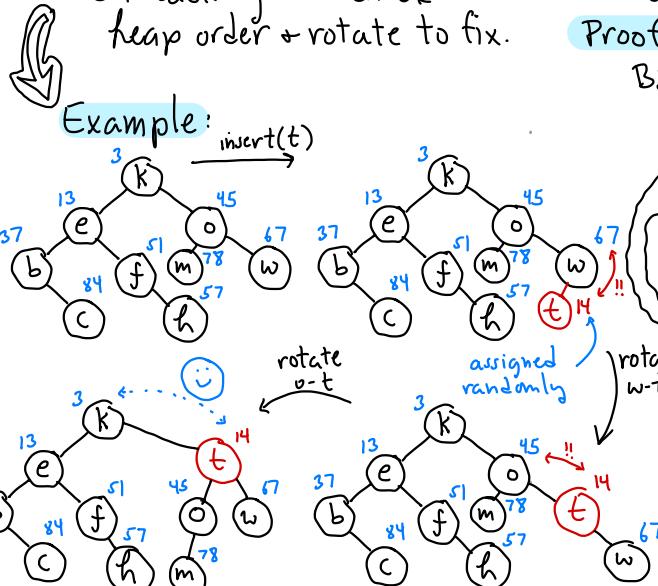
Priorities are in heap order

? Is it always possible to do both?

Yes: Just consider the corresponding BST

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

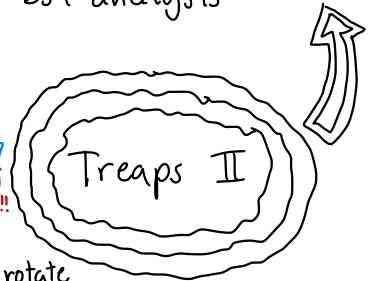
- Assign random priority
- On backtracking - check heap order + rotate to fix.



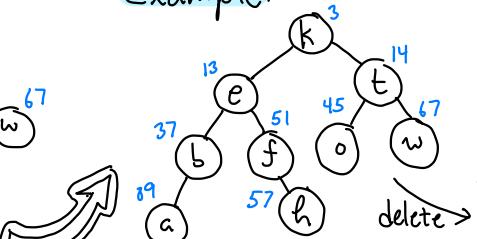
Deletion: (Cute solution) Find node to delete. Set its priority to $+\infty$. Rotate it down to leaf level + unlink.

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:



Implementation: (See pdf notes)

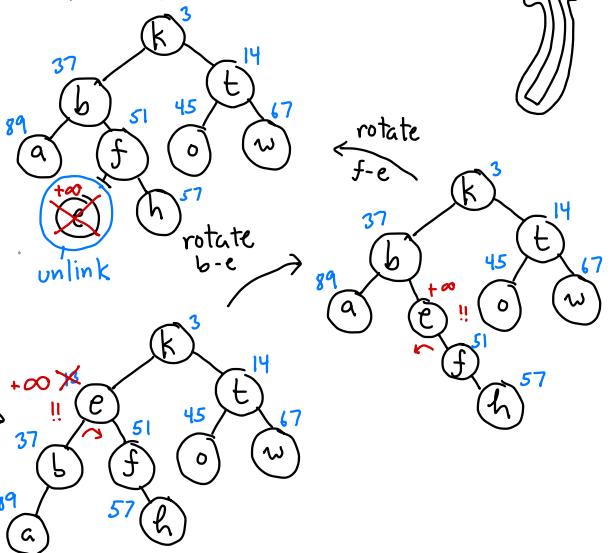
Node: Stores priority + usual...

Helpers:

lowest priority (p)
returns node of lowest priority among:

restructure:

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



Ideal Skip List:

- Organize list in levels

- Level 0: Everything

1: Every other

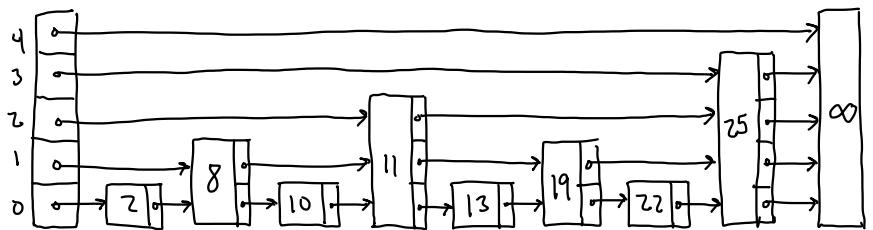
2: Every fourth

i : Every 2^i



Example:

head



Sorted linked lists:

- Easy to code

- Easy to insert/delete

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

Idea: Add extra links to skip



How to generalize?

Skip Lists I

Node Structure: (Variable sized)

class SkipNode{

Key key

Value value

SkipNode[] next

In constructor,
set level and
size

Value find(Key x){

i = topmost Level

SkipNode p = head

while(i >= 0) {

if(p.next[i].key <= x) p = p.next[i]

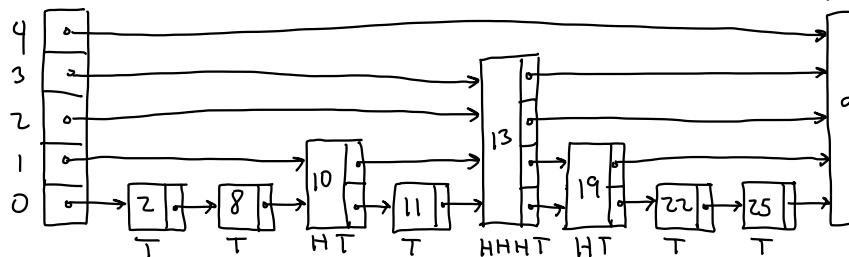
else i-- drop down a level

} we are at base level

if(p.key == x) return p.value
else return null

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

head



Thm: A skip list with n nodes has $O(\lg n)$ 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 node's 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

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

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

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

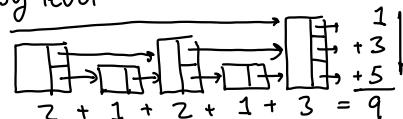
Obs: Prob. level exceeds $3 \cdot \lg n$ is $\leq 1/n^2$.

(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 $O(\lg n)$

Proof:

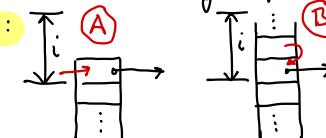
- We have seen no. levels is $O(\lg n)$

- 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

$$= 1 + \frac{1}{2}E(i) + \frac{1}{2}E(i-1)$$

$$\Rightarrow E(i)(1 - \frac{1}{2}) = 1 + \frac{1}{2}E(i-1)$$

$$\Rightarrow E(i) = [1 + \frac{1}{2}E(i-1)]2 = 2 + E(i-1)$$

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

Let $l = \max \text{ level}$. Total visited = $E(l)$

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

Delete:

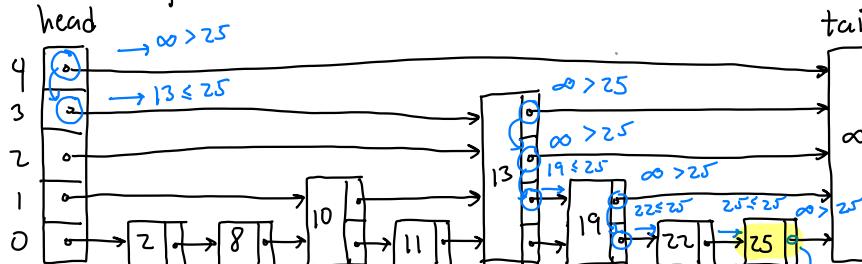
- Start at top
- Search each level saving last node $<$ key
- On reaching node at level 0, remove it and unlink from saved pointers

Insert: (Similar to linked lists)

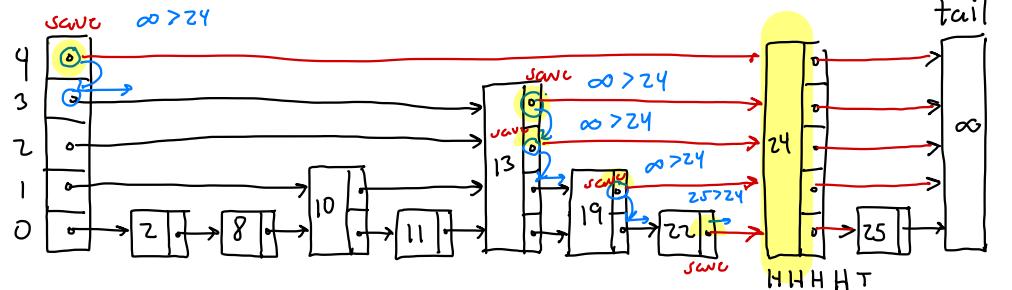
- Start at top level
- At each level:
 - Advance to last node \leq key
 - Save node + drop level
- At level 0:
 - Create new node (flip coins to determine height)
 - Link into each saved node

Skip Lists III

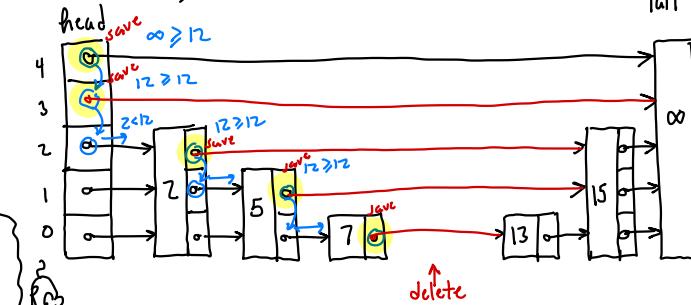
Example: find(25)



Insert(24)



Delete(12)



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

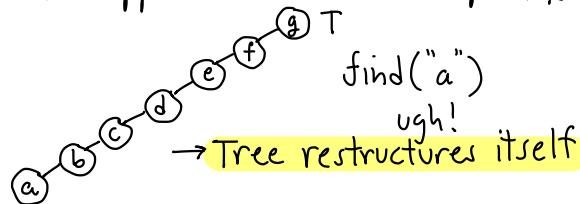
How to design/analyze?

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

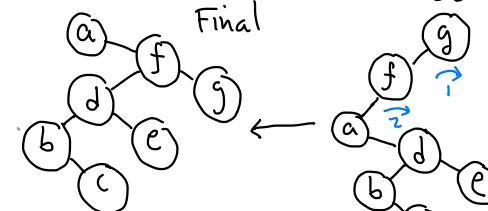
- Unbalanced BSTs
- AVL Trees
- 2-3, Red-black, AA Trees
- Treaps + Skip lists

→ Focus: Worst-case or randomized expected case

SPLAY TREES I

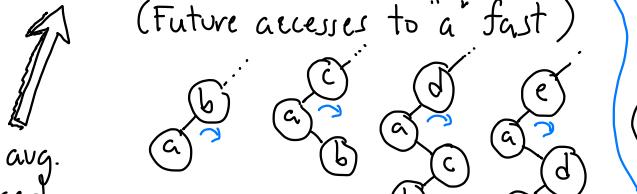
Lesson: Different combinations of rotations can:

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

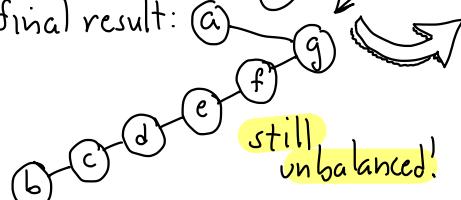


Tree's height has reduced by ~ half!

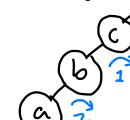
Idea I: Rotate "a" to top
(Future accesses to "a" fast)



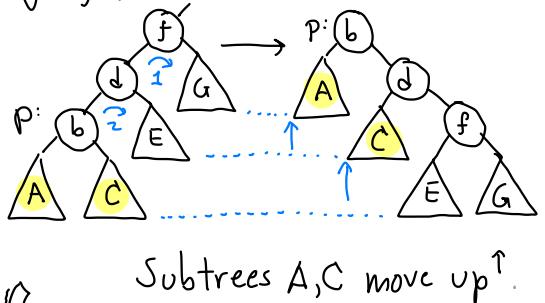
....final result:



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



ZigZig(p): [LL case]



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) zigzag(p)
        else /* p is LR or RL gr. child */ zigzag(p)
    }
  
```

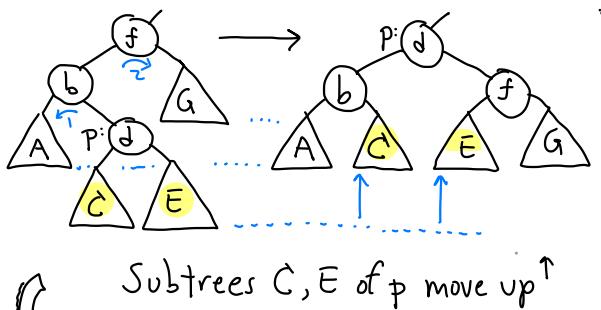
insert(x):

```

{splay(x)
q = new Node(x)
if (root.key < x)
    x.left = root
    x.right = root.right
    root.right = null
else ... symmetrical...
}
  
```

Subtrees A,C move up ↑.

ZIGZAG(p): [LR case]

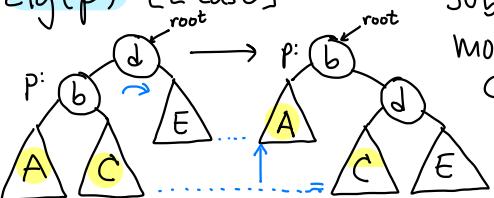


Splay Trees II

Example: splay(3)

Subtree A moves up ↑
C unchanged

Zig(p): [L case]

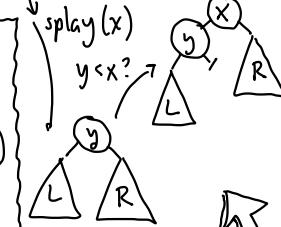


find(x):

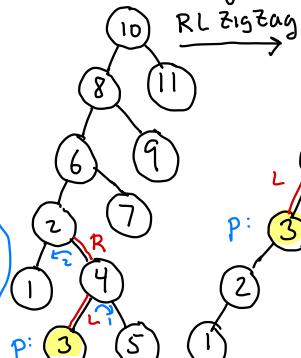
```

{splay(x)
if (root.key == x)
    found!
else not found
}
  
```

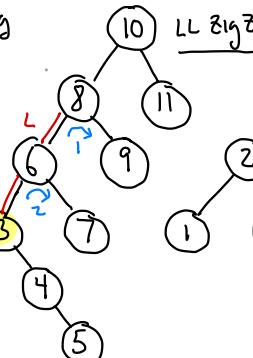
splay(x)



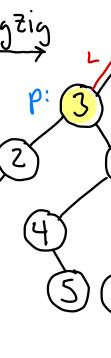
RL zigzag



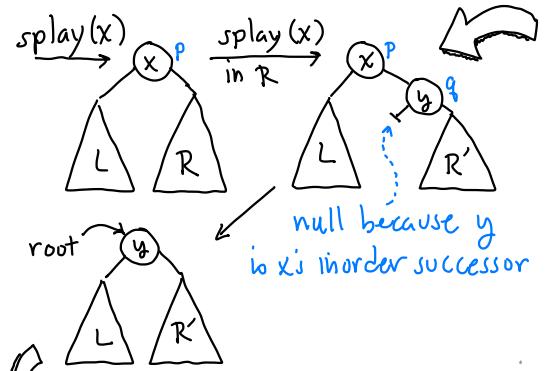
LL zigzag



L zig



Final ↴



Analysis:

- Amortized analysis
- Any one op might take $\Theta(n)$
- Over a long sequence, average time is $\mathcal{O}(\log n)$ 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

delete(x):

- splay(x) [x now at root]
- $p = \text{root}$
- if ($p.\text{key} \neq x$) error!
- splay(x) in p's right subtree
- $q = p.\text{right}$ [q's key is x's successor]
- $q.\text{left} = p.\text{left}$ [$q.\text{left} == \text{null}$]
- 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

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

SPLAY TREES III

Splay Trees are Amazingly Adaptive!

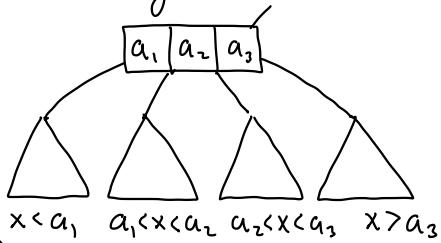
Balance Theorem: Starting with an empty dictionary, any sequence of m accesses takes total time $\mathcal{O}(m \log n + n \log n)$ where $n = \max.$ entries at any time.

Static Optimality:

- Suppose key x_i is accessed with prob p_i : $(\sum_{i=1}^n p_i = 1)$
- **Information Theory:** Best possible binary search tree answers queries in expected time $\mathcal{O}(H)$ where $H = \sum p_i \lg \frac{1}{p_i}$ \leftarrow 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 q_i times. Let $p_i = q_i/m$. Then total time is $\mathcal{O}(m \sum p_i \lg \frac{1}{p_i})$

Multiway Search Trees:



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

Node Structure: constant int M = ...

class BTNode {

```

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

B-Tree:

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

B-Tree: of order m (≥ 3)

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

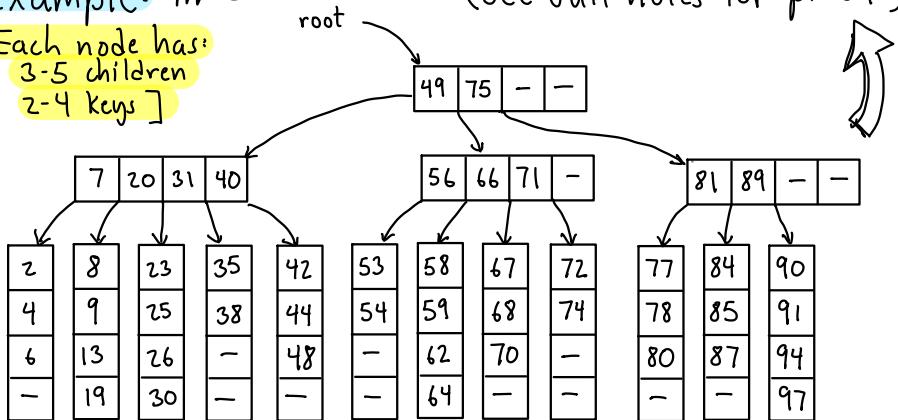
B-Trees I

Example: m=5

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

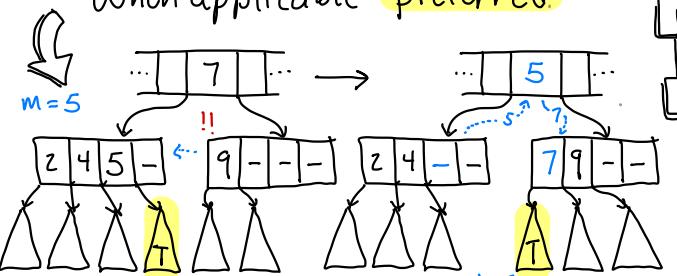
Theorem: A B-tree of order m with n keys has height at most $(\lg n)/\gamma$, where $\gamma = \lg(m/2)$

(See full notes for proof)



Key Rotation (Adoption)

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



Node Splitting:

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

Lemma: For all $m \geq 2$,

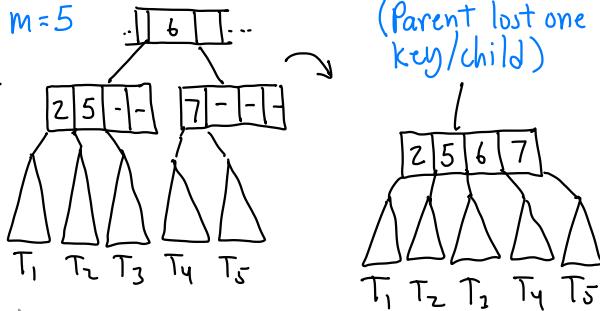
$$\lceil \frac{m}{2} \rceil \leq m+1 - \lceil \frac{m}{2} \rceil \leq m$$

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

B-Tree restructuring:

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

$m=5$



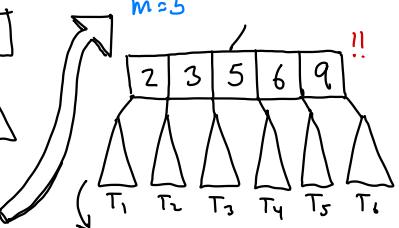
(Parent lost one key/child)

Lemma: For all $m \geq 2$,

$$\lceil \frac{m}{2} \rceil \leq 2\lceil \frac{m}{2} \rceil - 1 \leq m$$

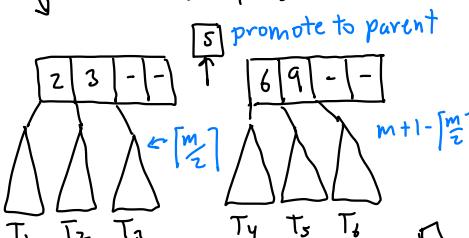
\Rightarrow Resulting node is valid

B-Trees II



Node Merging:

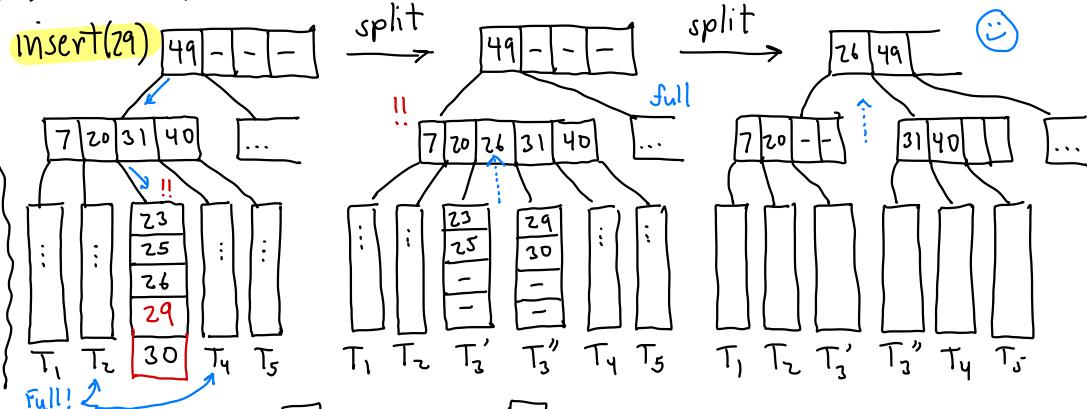
- A node has **too few** children $\lceil \frac{m}{2} \rceil - 1$
- Neither sibling has extra (both $\lceil \frac{m}{2} \rceil$)
- Merge with either sibling to produce node with $(\lceil \frac{m}{2} \rceil - 1) + \lceil \frac{m}{2} \rceil$ child



Insertion:

- Find insertion point (leaf level)
- Add key/value here
- If node **overfull** (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: $m=5$

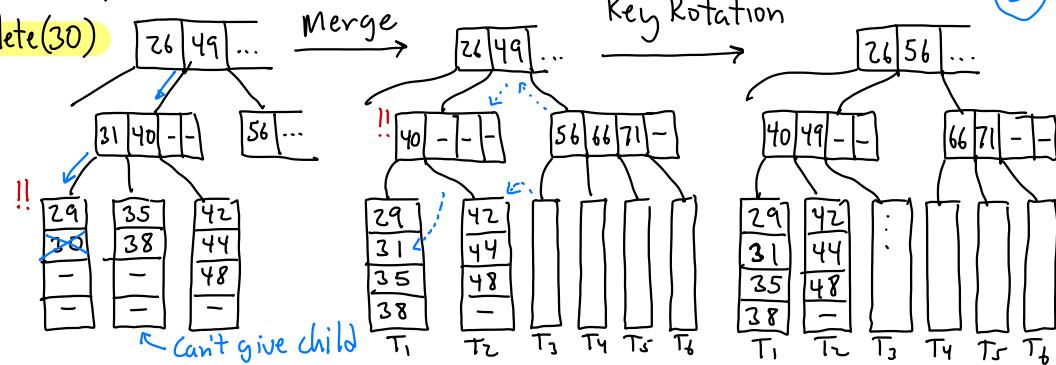


B-Trees III

Deletion:

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

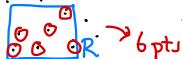
Example: $m=5$



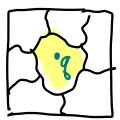
Geometric Search:

- Nearest neighbors

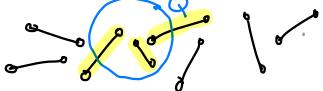
- Range searching



- 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



- External nodes store pts.

Multi-Dim vs. 1-dim Search?

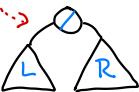
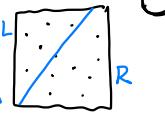
Similarities:

- Tree structure

- Balance $O(\log n)$

- Internal nodes - split

- External nodes - data



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

Differences:

- No (natural) total order

- Need other ways to discriminate + separate

- Tree rotation may not be meaningful



Quadtrees & kd-Trees I

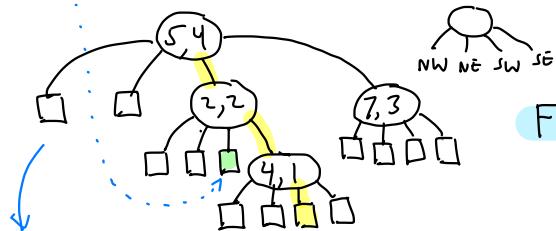
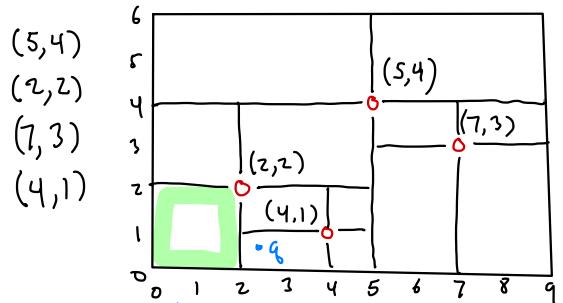
Point: A d -vector in \mathbb{R}^d

$$p = (p_1, \dots, p_d) \quad p_i \in \mathbb{R}$$

```
class Point {
    float[] coord // coords
    Point(int d)
        ... > coord = new float[d]
    int getDim() ... > coord.length
    float get(int i) ... > coord[i]
    ... others: equality, distance
    to_string...
}
```

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 (gently 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 cut.)

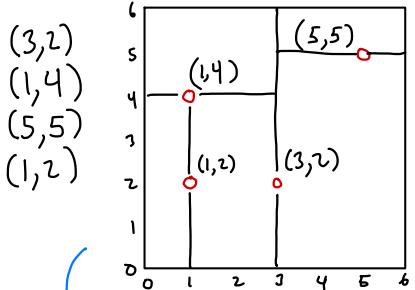
- cell: Still AABB
left: left/below
right: right/above



Quadtrees- Analysis

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

Example:



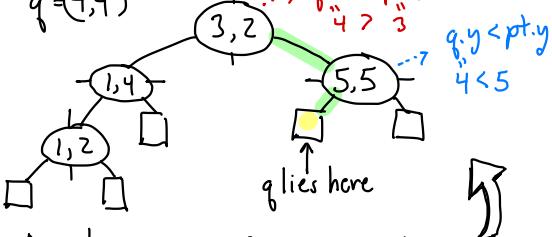
Kd-Tree Node:

class KDNode {

Point pt // splitting point
int cutDim // cutting coordinate
KDNode left // low side
KDNode right // high side

Example: $\text{find}(q) \xrightarrow{\text{calls}} \text{find}(q, \text{root})$

$q = (4,4)$



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

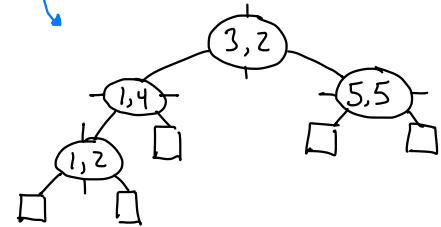
Theorem: If pts are inserted in random order, expected height is $O(\log n)$

How do we choose cutting dim?

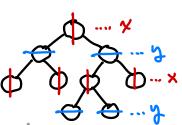
- Standard kd-tree: cycle through them (e.g. d=3: 1,2,3,1,2,3,...) based on tree depth

- Optimized kd-tree: (Bentley)

- Based on widest dimension of pts in cell.



Quadtrees &
kd-Trees III

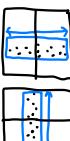


Find:

- Descend the tree
- Compare query pt with node pt along cutDim

class KDNode {

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



Value $\text{find}(\text{Point } q, \text{KDNode } p)$

```
if (p == null) return null;
else if (q == p.pt) ...all coords match?
    return p.value
else if (p.onLeft(q))
    return find(q, p.left)
else
    return find(q, p.right)
```

KDNode insert (Point x, Value v, KDNode p, int cd) {

```

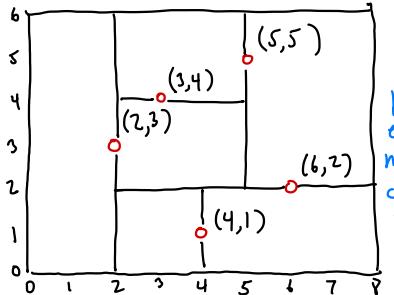
if(p == null) // fell out?
  [p = new KDNode(x, v, cd)] // new leaf node
else if(p.pt == x)
  Error! Duplicate key
else if(p.onLeft(x))
  p.left = insert(x, v, p.left, (cd+1)%dim)
else
  p.right = insert(x, v, p.right, (cd+1)%dim)
return p
}
  
```

(Similar to std. BSTs)

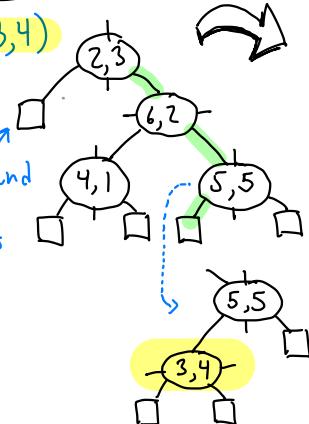
- Descend tree until
→ find pt → Error - duplicate
→ falling out (Although we draw extended trees, let's assume standard trees)
→ create new node
→ set cutting dim

Quadtrees & kd-Trees IV

Example:



insert(3,4)



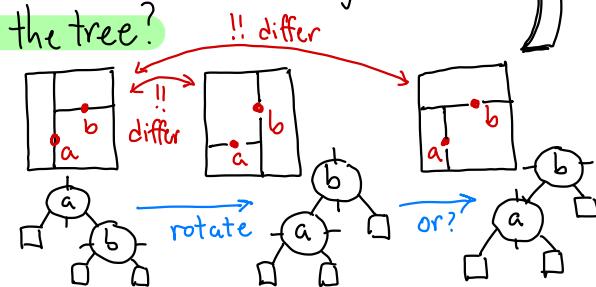
pretend ext. nodes are null

Analysis:

Runtime: $O(h)$

(Can we balance the tree?)

- Rotation does not make sense

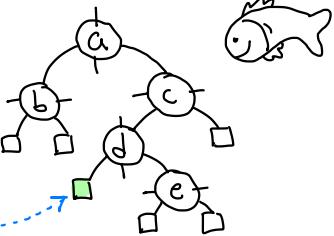
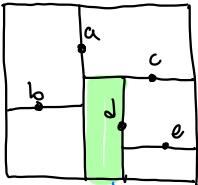


Deletion:

- 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.



kd-Trees:

- Partition trees
- Orthogonal split → vert [L|R]
- Alternate cutting → horz [R|L]
- Cells are axis-aligned rectangles (AABB)

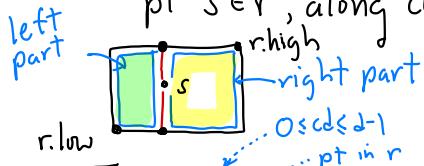
- Queries?
- Orthogonal range queries
 - Given query rect. (AABB) count/report pts in this rect.
 - Other range queries?
 - Circular disks
 - Halfplane
 - Nearest neighbor queries
 - Given query pt, return closest pt in the set
 - Find k^{th} closest point
 - Find farthest point from q

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

Kd-Tree Queries I

Rectangle methods for kd-cells:

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



r.leftPart(cd, s)

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

r.rightPart(cd, s)

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

Useful methods:

Let $r, c - \text{Rectangle}$
 $q - \text{Point}$



r.contains(q)



r.contains(c)



r.isDisjointFrom(c)

1 ≤ i ≤ 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.

Cases:

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

class Rectangle {

private Point low, high

public Rect (Point l, Point h)

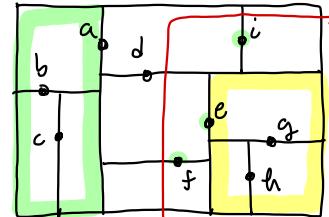
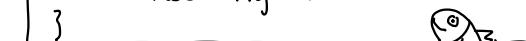
- " boolean contains (Point q)

- " boolean contains (Rect c)

- " Rect leftPart (int cd, Point s)

- " Rect rightPart (" " " ")

}



Final answer
= 1+1+1+2
= 5

Kd-Tree Queries II

II

Disjoint

Contained
in R + g.size = +2



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



Stabbing: 3 cases

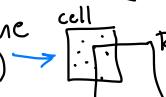
- cell is disjoint (easy)



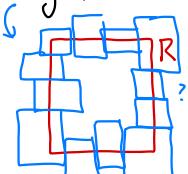
- cell is contained (easy)



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

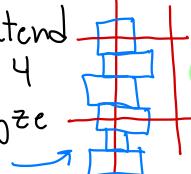


How many cells are stabbed by R ? (worst case)



Simpler: Extend R 's sides to 4

lines + analyze each one.



Analysis: How efficient is our algorithm?

→ Tricky to analyze

→ At some nodes we recurse on both children
⇒ $O(n)$ time?

→ At some we don't recurse at all!



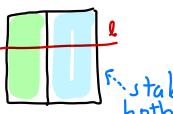
Kd-Tree Queries

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
 p splits horizontally
First stab both



Solving the Recurrence:

- Macho: Expand it

- Wimpy: Master Thm (CLRS)

Master Thm:

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

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

$$\text{For us: } a=2, b=4, d=0 \Rightarrow T(n) = n^{\log_4 2} = n^{1/2} = \sqrt{n}$$



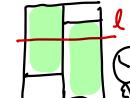
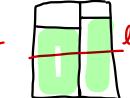
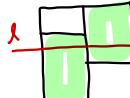
Since tree is balanced a child has half the pts + grandchild has quarter.

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

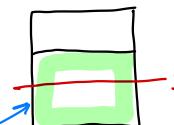
2 cells stabbed
Recursive on 2 grand children

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



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

Universal Hashing:

Even better → randomize!

- Let H be a family of hash fns
- Select $h \in H$ randomly
- If $x \neq y$ then $\text{Prob}(h(x) = h(y)) = 1/m$

E.g. 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)

Assume
keys can
be interpreted as
ints

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

E.g. Java variable names:



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

Overview:

- Separate Chaining
 - Open Addressing:
 - Linear probing
 - Quadratic probing
 - Double hashing
- simple/slow complex/fast

Separate Chaining:

$\text{table}[i]$ is head of linked list of keys that hash to i .

Example:

table	
Keys (x)	$h(x)$
d	1
z	4
p	7
w	0
t	4
f	0
m=8	

Collision Resolution:

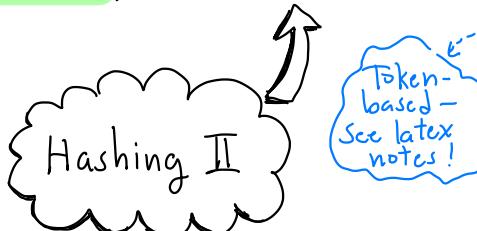
If there were no collisions hashing would be trivial!

$\text{insert}(x, v) \rightarrow \text{table}[h(x)] = v$
 $\text{find}(x) \rightarrow \text{return } \text{table}[h(x)]$
 $\text{delete}(x) \rightarrow \text{table}[h(x)] = \text{null}$

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

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



Token-based -
See latex notes!

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
 $\lambda = n/m$ $n = \# \text{ of keys}$
 $m = \text{table size}$

Proof: On avg. each list has $n/m = \lambda$
 success: 1 for head + half the list
 unsuccessful: 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
 $h(x) + i_1$
 $h(x) + i_2$
 $h(x) + i_3$

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

- What's the best probe sequence?

Linear Probing:

- $h(x), h(x)+1, h(x)+2, \dots$
- until finding first available

Simple, but is it good?

- $x: d, z, p, w, t$
 - $h(x): 0, 2, 2, 0, 1$
 - t did not collide directly but had to probe 3 times!
- | | | | | | | |
|---|---|---|---|---|---|------|
| d | w | z | p | t | | |
| 0 | 1 | 2 | 3 | 4 | 5 | 6... |

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: Improves secondary clustering

- May 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
 \Leftrightarrow if $\lambda < \frac{1}{2}$.

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

Pf: See latex notes.

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

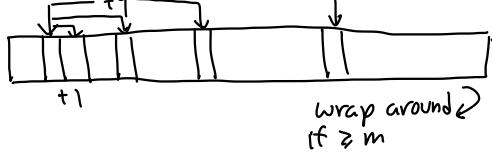
Clustering

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

Quadratic Probing:

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

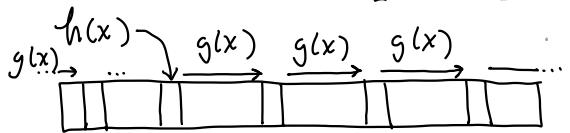
$$h(x) \xrightarrow{+4} h(x)+a \xrightarrow{+16} h(x)+j^2$$



Double Hashing:

(Best of the open-addressing methods)

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



(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} = 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:



$$\text{Thm: } S_{DH} = \frac{1}{\lambda} \ln \left(\frac{1}{1-\lambda} \right)$$

$$U_{DH} = 1/(1-\lambda)$$

\rightarrow Proof is nontrivial (skip)

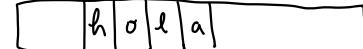
$\lambda:$	0.5	.075	0.95	0.99
$U_{DH}:$	2	4	20	100
$S_{DH}:$	1.39	1.89	3.15	4.65

\curvearrowleft Very efficient!

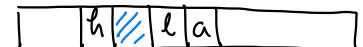
Delete(x): Apply find(x)

- \rightarrow Not found \Rightarrow error
- \rightarrow Found \Rightarrow set to "empty"

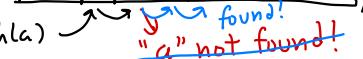
Problem:



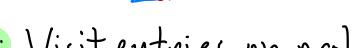
insert(a):



delete(a):

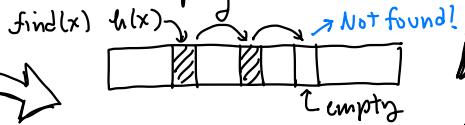


find(a):



Find(x): Visit entries on probe sequence until:

- found $x \Rightarrow$ return v
- hit empty \Rightarrow return null



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??

Scapegoat Trees:

- Arne Anderson (1989)
- Galperin + Rivest (1993)
rediscovered/extended
- Amortized analysis
 - $O(\log n)$ for dictionary ops amortized (guaranteed for find)
 - Just let things happen
 - If subtree unbalanced
 - rebuild it



Overview:

Insert:

- same as standard BST
- if depth too high
 - trace search path back
- find unbalanced node - **scapegoat**
- rebuild this subtree

Find:

- Tree height $\leq \log_{3/2} n \approx 1.71 \lg n$



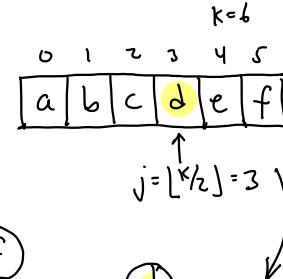
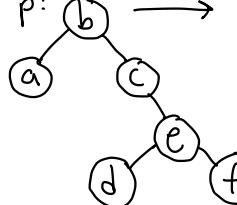
Recap:

- Seen many search trees
- Restructure via **rotation**
- Today: Restructure via **rebuilding**
- Sometimes rotation not possible
- Better mem. usage

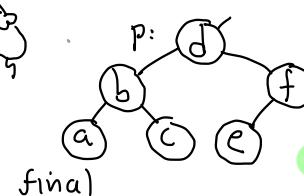


Example:

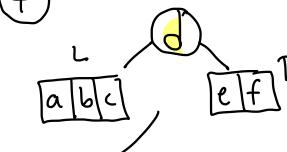
p: b →



$j = \lfloor x/2 \rfloor = 3$



final)



Time = $O(k)$

How to rebuild?

rebuild(p):

- inorder traverse p's subtree → array A[]
- buildSubtree(A)

buildSubtree(A[0..k-1]):

- if k=0 return null
- $j \leftarrow \lfloor k/2 \rfloor$; $x \leftarrow A[j]$ median
- $L \leftarrow \text{buildSubtree}(A[0..j-1])$
- $R \leftarrow \text{buildSubtree}(A[j+1..k-1])$
- return Node(x, L, R)

Delete:

- Same as std. BST
- If num. of deletions is large rel. to n - rebuild entire tree!

How? Maintain $n, m \leftarrow 0$

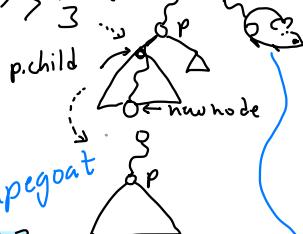
Insert: $n++, m++$

Delete: $n-- \dots \rightarrow$ If $m > 2n$ rebuild



Insert: _____

- $n++$; $m++$
- same as std BST but keep track of inserted node's depth $\rightarrow d$
- if $(d > \log_{3/2} m)$ {
 - /* rebuild event */
 - trace path back to root
 - for each node p visited, $\text{size}(p) = \text{no. of nodes in } p\text{'s subtree}$
 - if $\frac{\text{size}(p.\text{child})}{\text{size}(p)} > \frac{2}{3}$
 - $p \leftarrow \text{rebuild}(p)$
 - break



How to compute $\text{size}(p)$?

- Can compute it on the fly
- While backing out, traverse "other sibling"
- Too slow? No!
→ Charge to rebuild.

Details of Operations:

Init: $n \leftarrow m \leftarrow 0$ root $\leftarrow \text{null}$

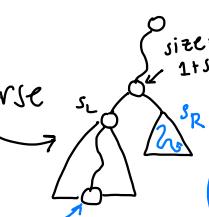
Delete:

- Same as std BST
- $n--$
- if $m > 2n$, rebuild(root)

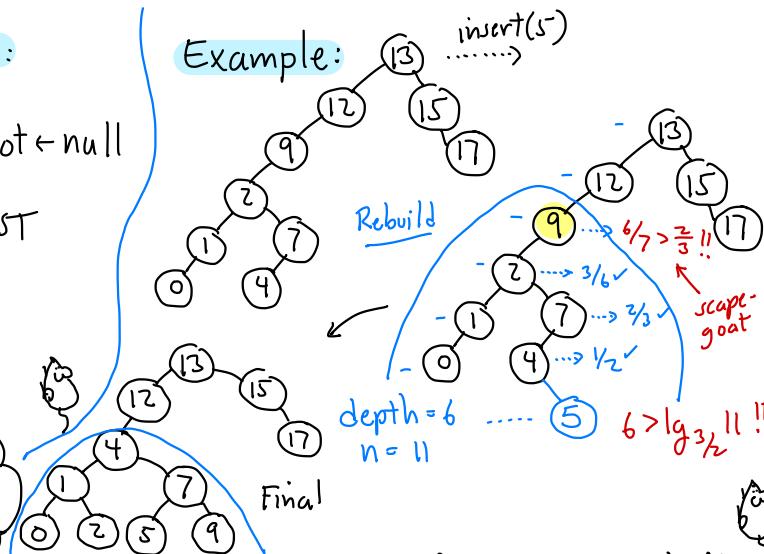
Time: $O(n)$

Scapegoat Trees II

Must there be a scapegoat? Yes!



Example: $\text{insert}(5)$



Proof: By contradiction

- Suppose p 's depth $> \log_{3/2} n$ but \forall ancestors

depth 0: $\text{size} \leq n$
1: $\text{size} \leq \frac{2}{3}n$
 \vdots
 $d > \log_{3/2} n$: $\text{size} \leq \frac{2}{3}^d n$

\Rightarrow Since p has 1 node:
 $1 \leq \text{size}(p) \leq (\frac{2}{3})^d n$

$\Rightarrow (\frac{3}{2})^d \leq n$
 $\Rightarrow d \leq \log_{3/2} n$ \square

Lemma: Given a binary tree with n nodes, if \exists node p of depth $> \log_{3/2} n$, then \exists ancestor of p that satisfies scapegoat condition

Scapegoat Trees

III

Theorem: Starting with an empty tree,
any sequence of m dictionary operations
on a scapegoat tree take time
 $O(m \log m)$ [Amortized: $O(\log m)$]

Proof: (Sketch)

Find: $O(\log n)$ guaranteed [Height = $O(\log n)$]

Delete: In order to induce a rebuild,
number of deletes \sim number of
nodes in tree

→ Amortize rebuild time against
delete ops

Insert: Based on potential argument

→ It takes $\sim k$ ops to cause a
subtree of size k to be unbalanced.

→ Charge rebuild time to these
operations

Can we do better?

Range Trees:

- Space is $O(n \log^{d-1} n)$

- Query time:

Counting: $O(\log^d n)$

Reporting: $O(k + \log^d n)$

→ In \mathbb{R}^2 : $\log^2 n$ much better than \sqrt{n} for large n

→ Range trees are more limited



Recap:

- **kd-Tree**: General-purpose data structure for pts in \mathbb{R}^d

- **Orthogonal range query**: Count/report pts in axis-aligned rect.



- **kd-Tree**: Counting: $O(\sqrt{n})$ time
Report: $O(k + \sqrt{n})$ time

No. of pts reported



Layering:

 Combining search structures

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

- Medical data: Count subjects

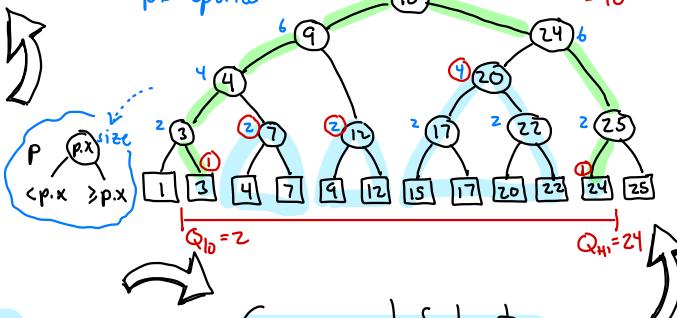
Age range: $a_{lo} \leq \text{age} \leq a_{hi}$

Weight range: $w_{lo} \leq \text{weight} \leq w_{hi}$

- Design a data structure for each criterion individually

- Layer these structures together to answer full query

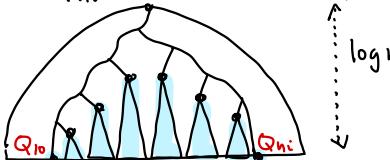
→ Multi-Layer Data Structures



Canonical Subsets:

- **Goal**: Express answer as disjoint union of subsets

- **Method**: Search for Q_{lo} + Q_{hi} + take maximal subtrees



Approach:

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

Recursive helper:

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

Intv Q = [Q_{lo}, Q_{hi}], Intv C = [x₀, x₁])

initial call: range1Dx(root, Q, C₀)

Cases:

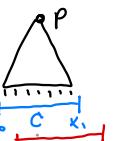
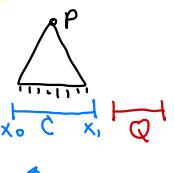
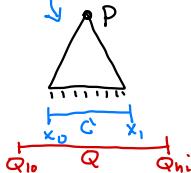
p is external:

- if p.pt.x ∈ Q → 1 else → 0

p is internal:

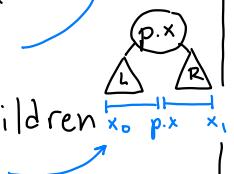
- C ⊆ Q ⇒ all of p's pts lie within query

→ 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:

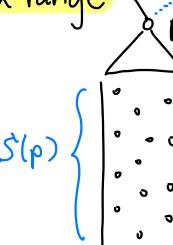
- Let Q = [Q_{lo}, Q_{hi}] be query interval

- For each node p, define interval cell C = [x₀, x₁] s.t. all pts of p's subtree lie in C

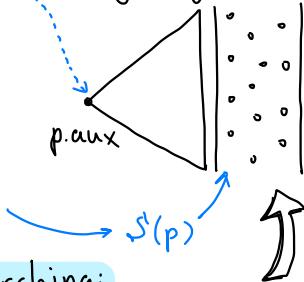
- Root cell: C₀ = [-∞, +∞]

Range Trees II

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:

```
int range1Dx(Node p,  
Intv Q, Intv C = [x0, x1]) {  
    if(p is external) → 1  
    return p.pt.x ∈ Q → 0  
    else if (C ⊆ Q) return p.size  
    else if (Q+C disjoint) return 0  
    else return:  
        range1Dx(p.left, Q, [x0, p.x])  
        + range1Dx(p.right, Q, [p.x, x1])
```

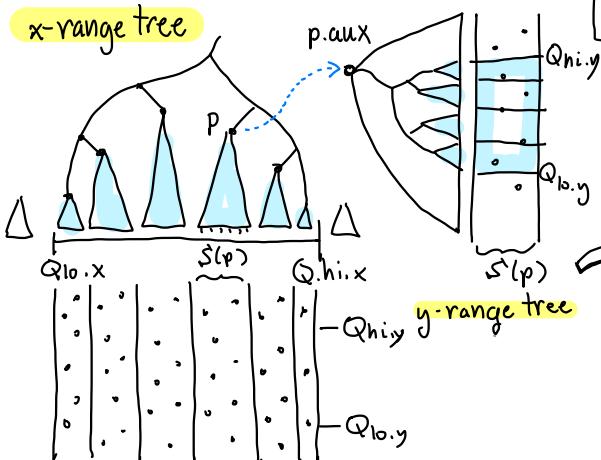
Lemma: Given a 1-D range tree with n pts, given any interval Q, can compute O(log n) subtrees whose union is answer to query.

Thm: Given 1-D range tree...
can answer range queries in time O(log n) → (+k to report)

Answering Queries?

Given query range $Q = [Q_{lo,x}, Q_{hi,x}] \times [Q_{lo,y}, Q_{hi,y}]$

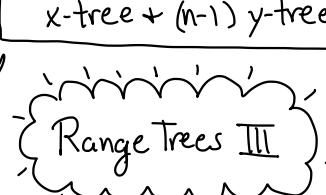
- Run range1D_x to find all subtrees that contribute
- For each such node p,
 - run range1D_y on p.aux
- Return sum of all result



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



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:

$$\underbrace{\log n \cdot \log n \cdots \log n}_{d} = O(\log^d n)$$

Analysis: The 1D x search takes of $O(\log n)$ time + generates $O(\log n)$ calls to 1D y search
 \Rightarrow Total: $O(\log n \cdot \log n) = O(\log^2 n)$

```
int range2D(Node p, Rect Q, Intrv C=[x0, x1]) {
```

```
    if (p is external) return p.pt ∈ Q? 1 : 0
    else if (Q.x contains C) { // C ⊆ Q; 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:

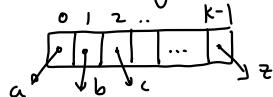
Invoked $O(\log n)$ times - once per maximal subtree

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

Tries: History

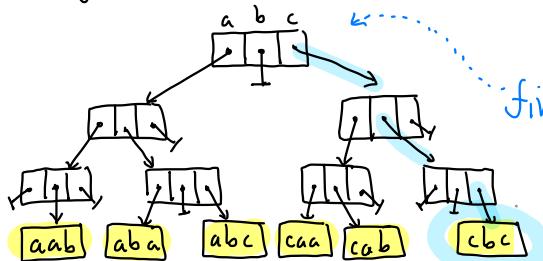
- de la Briandais (1959)
- Fredkin - "trie" from "retrieval"
- Pronounced like "try"

Node: Multiway of order k

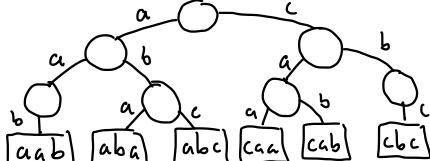


Example: $\Sigma = \{a=0, b=1, c=2\}$

Keys: {aab, aba, abc, caa, cab, cbc}



Same structure/Alt. Drawing



Digital Search:

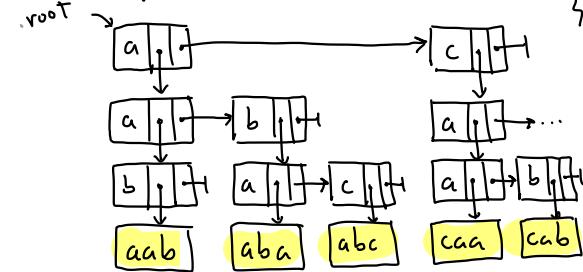
- Keys are strings over some alphabet Σ
- E.g. $\Sigma = \{a, b, c, \dots\}$
- $\Sigma = \{0, 1\}$ Let $k = |\Sigma|$
- Assume chars coded as ints: $a=0, b=1, \dots, z=k-1$

Tries and Digital Search Trees I

Analysis:

- 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
- $x \rightarrow \#x \Rightarrow$ try next char in Σ
- $= x \Rightarrow$ advance to next character of search string
- First-child/next-sibling

Space:

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

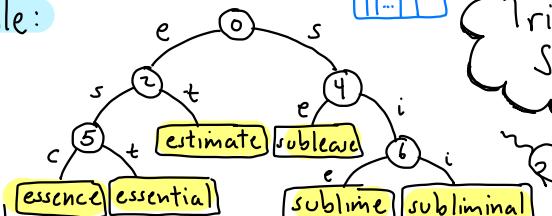
Large!

Patricia Tries:

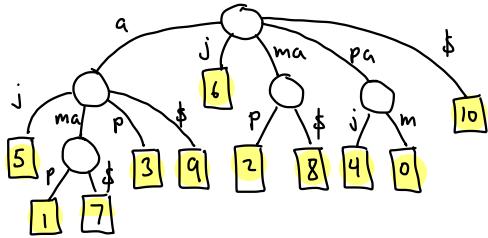
- Improves trie by compressing degenerate paths
- PATRICIA = Practical Alg. to Retrieve Info. Coded in Alpha...
- Late 1960's: Morrison & Gwachberger
- Each node has index field, indicates which char to check next (Increase with depth)

Example:

essence
essential
estimate
sublease
sublime
subliminal



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



E.g. $\text{ID}(S_1) = \text{amap}$ $\text{ID}(S_2) = \text{ama\$}$.

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 (full or on extern node)
 - compare target with string
 - if matches - found 1 occurrence
 - else - no occurrences

Example:

$\text{Search("ama")} \rightarrow$ End at intern node ama

Report: 2 occ's.

$\text{Search("amapaj")} \rightarrow$ End at extern node amap

Goto S_1 + 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
($k = |\Sigma| = O(1)$)

- **Search time:** \sim to length of target string

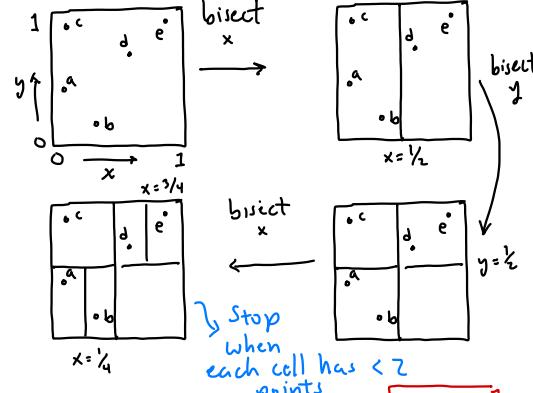
- **Construction time:**
 $-O(n \cdot k)$ [nontrivial]

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

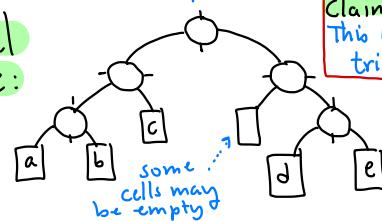
Assume points lie in unit square



Stop when each cell has < 2 points

Claim:
This is a trie!

Final tree:



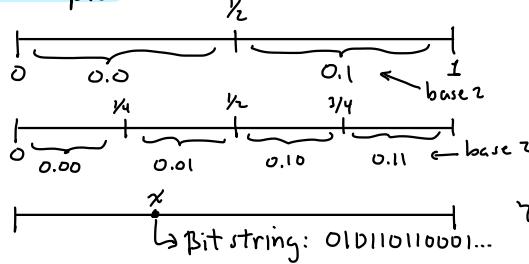
Binary Encoding:

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

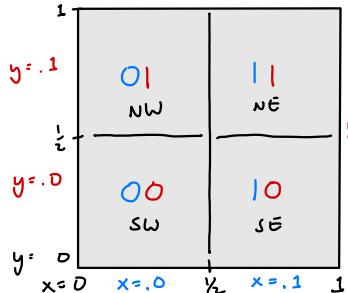
$$x = 0.a_1 a_2 a_3 \dots \quad a_i \in \{0, 1\}$$

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

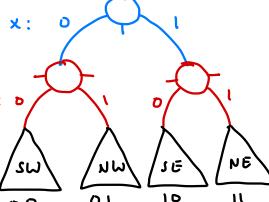
Example:



How do we extend to 2-D?



PR kd-tree



Bit Interleaving:

Given a point $p = (x, y)$

$$0 \leq x, y \leq 1$$

let: $x = 0.a_1 a_2 \dots$ in binary

$$y = 0.b_1 b_2 \dots$$

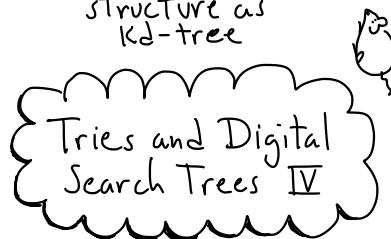
Define:

$$\phi(x, y) = a_1, b_1, a_2 b_2, a_3 b_3, \dots$$

Called Morton Code of p

PR kd-Tree \equiv Trie ??

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



Further Remarks:

- Techniques for efficiently encoding, building, serializing, compressing... tries apply immediately to PR kd-tree
- Can generalize to any dimension

$$\begin{aligned} x &= 0.a_1 a_2 \dots \\ y &= 0.b_1 b_2 \dots \\ z &= 0.c_1 c_2 \dots \end{aligned} \quad \left. \begin{array}{l} \phi = a_1 b_1 c_1 a_2 b_2 c_2 \dots \\ \vdots \end{array} \right\}$$

Lemma: Given a pt set $P \subseteq \mathbb{R}^2$ (in unit square $[0, 1]^2$) let

$$P = \{p_1, \dots, p_n\} \text{ where } p_i = (x_i, y_i)$$

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 $x = 0.a_1 a_2 \dots$ $y = 0.b_1 b_2 \dots$

and consider just $\phi(x, y) = a_1, b_1, \dots$

