

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
-



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 ☺
[Hash map]
- $O(\log n)$: log-time (good)
[Binary search]
- $O(n^p)$: $p = \text{constant}$: poly time
- $O(\sqrt{n})$



Asymptotic: "Big-o"

- Ignore constants
- Focus on large n
- $T(n) = 34n^2 + 15n \log n + 143$
- $T(n) = O(n^2)$



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?



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

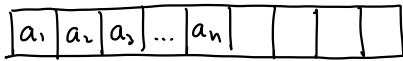
insert(i, x) - inserts x prior to i^{th} (moving others back)

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

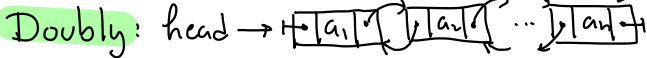
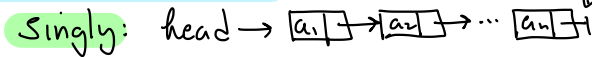
length() - returns num. of items

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

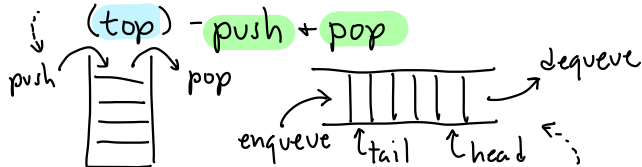
- allocate new array size $2n$
- copy old to new
- remove old array

Dynamic Lists + Sequential Allocation

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

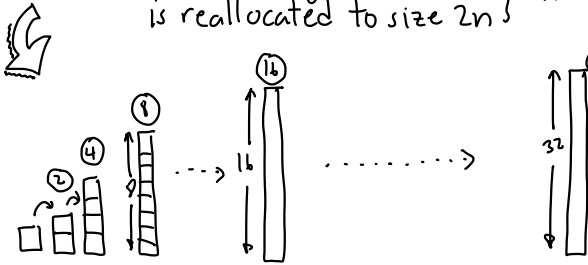


Queue: FIFO list: **enqueue** inserts at **tail** and **dequeue** deletes from **head**

Cost model (Actual cost)

Cheap: No reallocation \rightarrow 1 unit

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



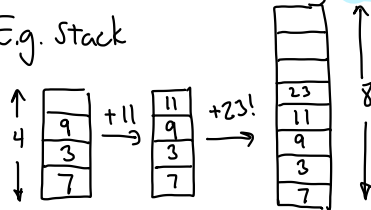
1 2 3 4 5 6 7 8 9 10 11 12 13 14 15 16 17
 $+1 +1 +1 +1 +1 +1 +1 +1 +1 +1 +1 +1 +1 +1 +1 +1 +1$

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.

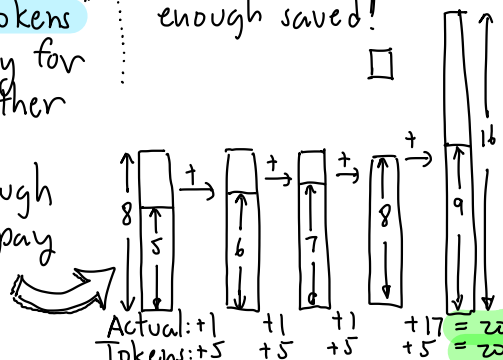
[i.e. any seq. of m ops has cost $\leq 5 \cdot m$]

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



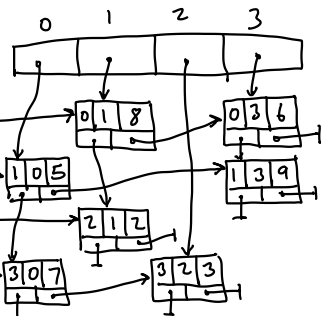
Fixed Increment: Increase by a fixed constant
 $n \rightarrow n + 100$

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

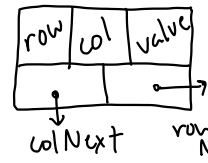
Dynamic Stack:
 - Showed doubling \Rightarrow Amortized $O(1)$
 - Other strategies?

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:

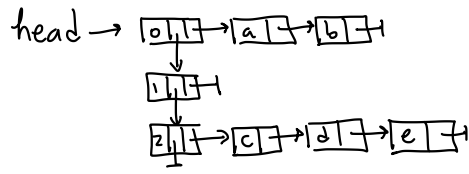


Idea: Store only non-zero entries linked by row and column

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

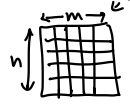
Leave as exercise ☹️
 (spoiler alert!)
 Fixed increment \rightarrow no
 Fixed factor \rightarrow yes
 Squaring \rightarrow yes

Multilists: Lists of lists



Sparse Matrices:

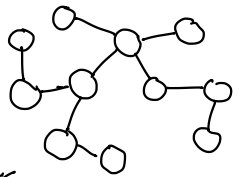
An $n \times m$ matrix has $n \cdot m$ entries and takes (naively) $O(n \cdot m)$ space



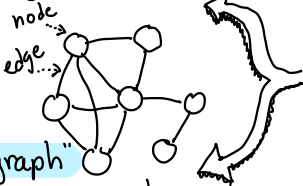
Sparse matrix: Most entries are zero

Tree (or "Free Tree")

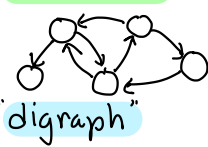
- undirected
- connected
- acyclic graph



Undirected



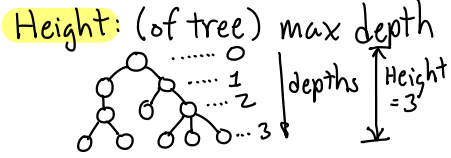
Directed



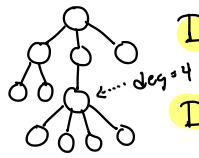
Graph: $G=(V,E)$
 V = finite set of vertices (nodes)
 E = set of edges (pairs of vertices)

Trees: Basic Concepts and Definitions

Depth: path length from root

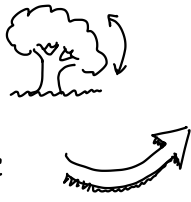
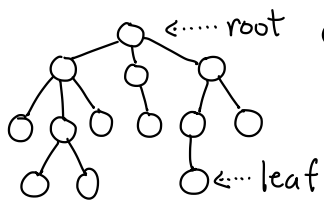


Degree (of node): number of children



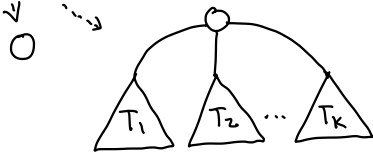
Degree (of tree): max. degree of any node

Rooted tree: A free tree with root node

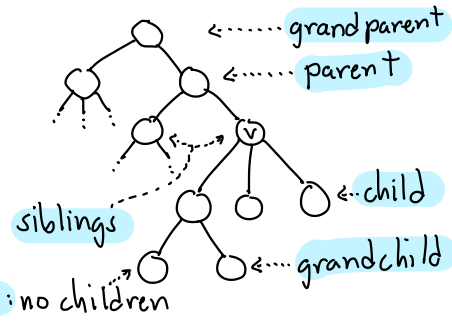


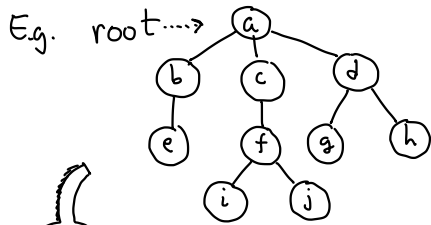
Formal definition:

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



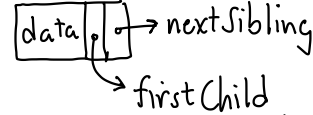
"Family" Relations





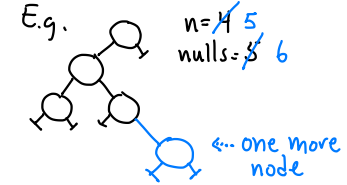
Representing rooted trees:
Each node stores a (linked) list of its children

Node structure:



Wasted space?

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

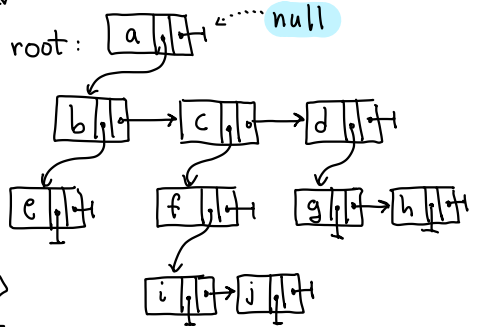


In Java:

```

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

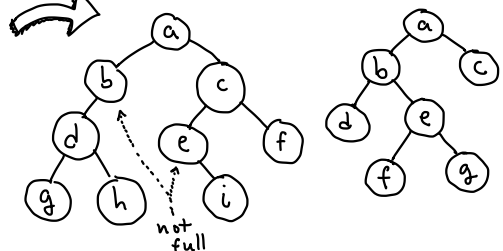
```



Trees Representation + Binary Trees (I)

(Not full)

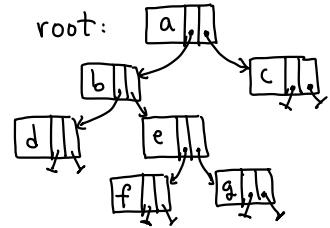
Full:



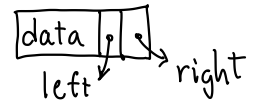
Full: Every non-leaf node has 2 children

called the **Binary representation**

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

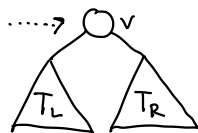
```



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

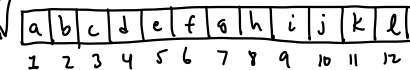
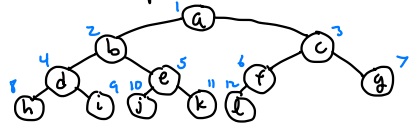
Binary Tree Traversals (can be generalized)

root → v



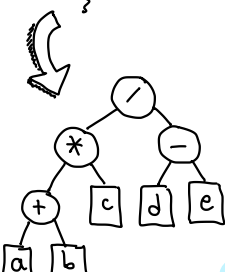
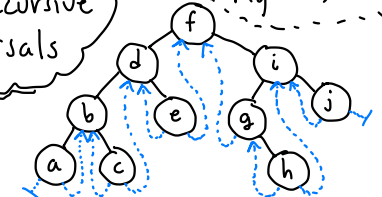
- process/visit v
 - traverse T_L
 - traverse T_R
- } recursive

Complete Binary Tree: All levels full (except last)



$parent(i) = \lfloor i/2 \rfloor$
 $left(i) = 2 \cdot i$
 $right(i) = 2 \cdot i + 1$

Challenge: Nonrecursive traversals



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

Binary Trees:
Traversals, Extension,
and More

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

Another way to save space...

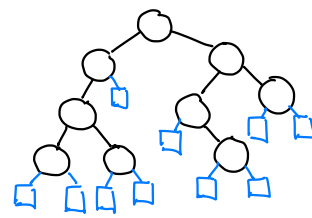
Threaded binary tree: Store (useful) links in the null links. (Use a mark bit to distinguish link types.)

Eg. **Inorder Threads:**
 Null left → inorder predecessor
 Null right → " successor



Those wasteful null links...

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



Observation: Every extended binary tree is full

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

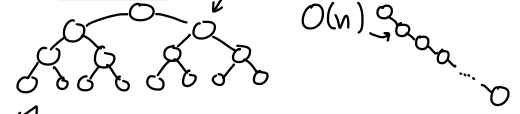
each associated with key x ;

- and value v_i
- store for quick access & updates

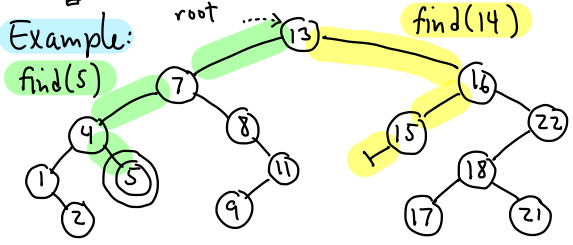
- Ordered: Assume that keys are totally ordered: $<, >, =$

Efficiency: Depends on tree's height

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

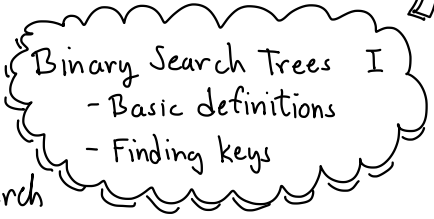


Example:



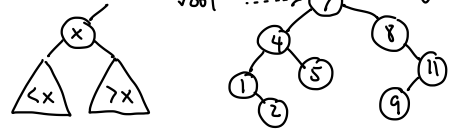
Sequential Allocation?

- 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**

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



Find: How to find a key in the tree?

- Start at root $p = \text{root}$
- 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!

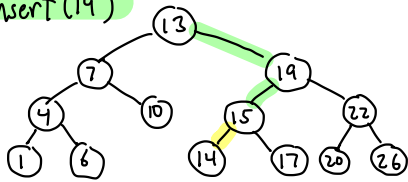


```

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
}

```

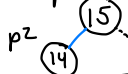
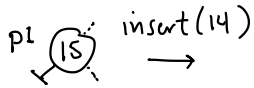
Binary Search Trees II

- insertion
- deletion



Why did we do:

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



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

$p2 = new BSTNode$
return $p2$

Be sure you understand this!

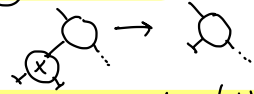
Delete (Key x)

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

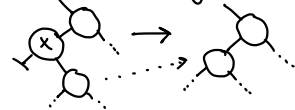


3 cases:

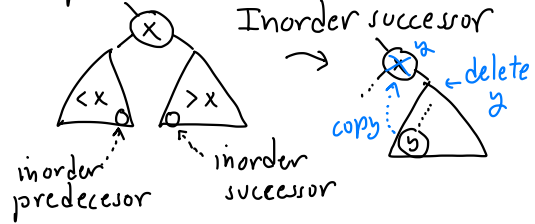
(1) x is a leaf



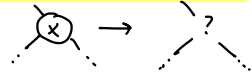
(2) x has single child



Replacement Node?



3. x has two children



Find replacement node

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

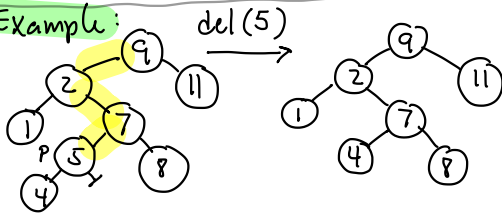


BSTNode delete (Key x, BSTNode p)

```

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

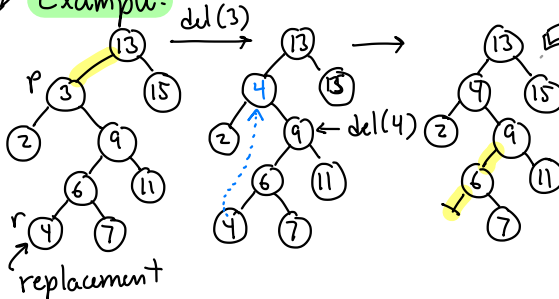
```

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

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

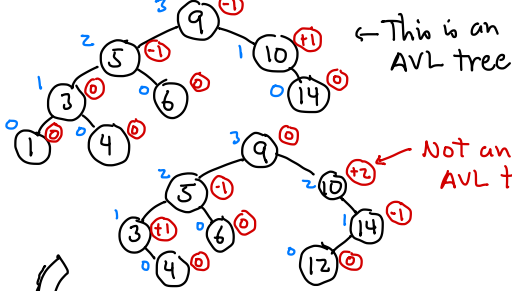
Thm: If n keys are inserted in random order, expected height is $O(\log n)$.

Analysis:

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

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



```
BSTNode rotateRight(BSTNode p) {
```

```
    BSTNode q = p.left
```

```
    p.left = q.right
```

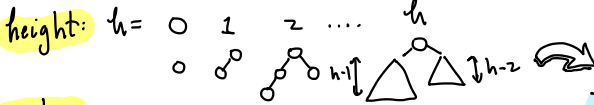
```
    q.right = p
```

```
    return q
```

```
}
```

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

Worst cases:



nodes: $n = 1, 2, 4, 7, 12, 20, \dots$
 $n+1 = 2, 3, 5, 8, 13, 21, \dots$

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)

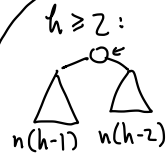
$$h=0: n(h) = 1 = F_3 - 1$$

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

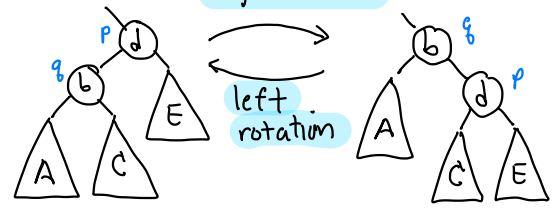
$$h \geq 2: n(h) = 1 + n(h-1) + n(h-2)$$

$$= 1 + (F_{h+2} - 1) + (F_{h+1} - 1)$$

$$= (F_{h+2} + F_{h+1}) - 1 = F_{h+3} - 1 \quad \square$$



How to maintain the AVL property?

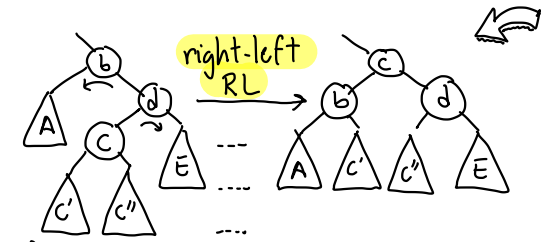


$A < b < C < d < E$

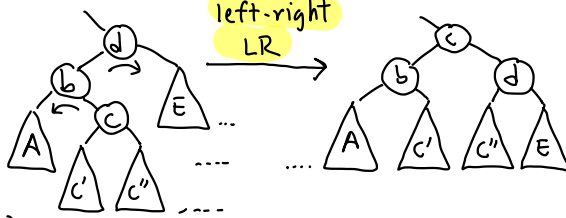
$A < b < C < d < E$

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



Double rotations:



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 ... (symmetrical)
updateHeight(p); return p

```

```

BSTNode rotateLeftRight (BSTNode p)
{
  p.left = rotateLeft(p.left)
  return rotateRight(p)
}

```

AVL Tree:

AVLNode: Same as BSTNode but + member **int height**

Utilities:

```

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

```

```

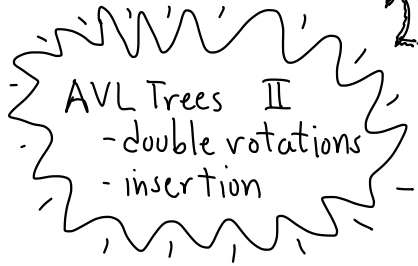
void updateheight (AVLNode p)
{
  p.height = 1 + max (height(p.left),
                    height(p.right))
}

```

```

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

```

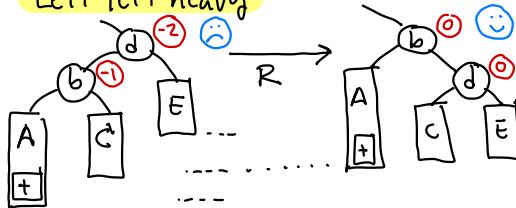


Find: Same as B.S.T.

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

How to rebalance? Bal = -2

Left-left heavy



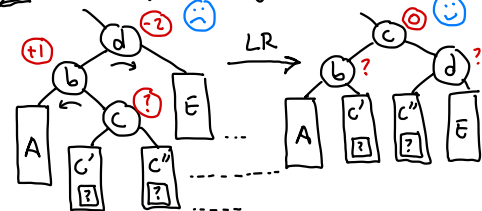
AVLNode insert (Key x, Value v, AVLNode p)

```

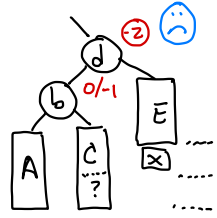
if (p == null) p = new AVLNode(x, v)
else if (x < p.key)
  p.left = insert(x, v, p.left)
else if (x > p.key)
  p.right = insert(x, v, p.right)
else throw -Error -Duplicate!
return rebalance(p)

```

Left-right heavy:

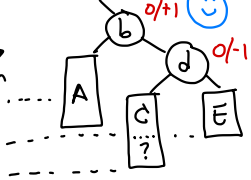


Cases: Balance factor -2

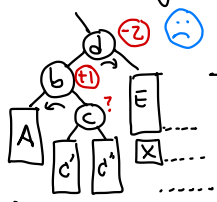


Left-left heavy

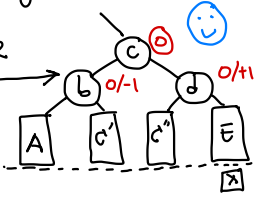
Right rotation



Left-right heavy



LR

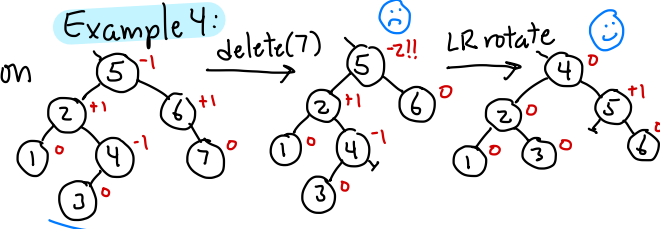


Deletion: Basic plan

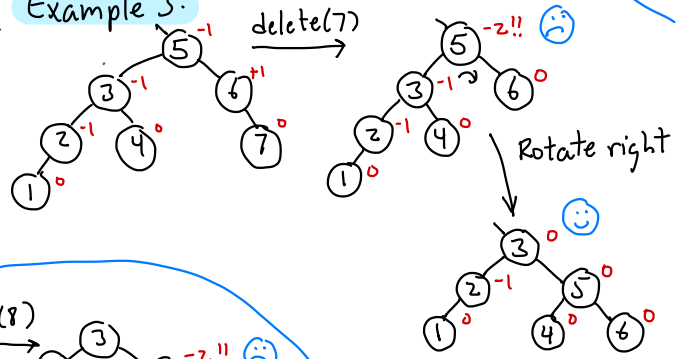
- Apply standard BST deletion
- find key to delete
- find replacement node
- copy contents
- delete replacement
- rebalance

AVL Trees III
- Deletion
- Examples

Example 4:



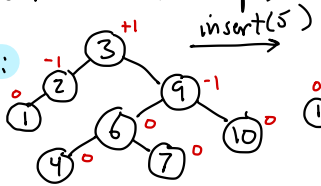
Example 3:



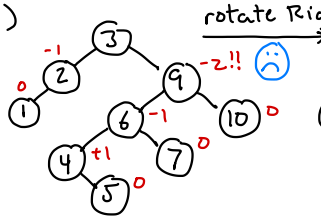
AVLNode delete (Key x, AVLNode p)

same as BST delete
return rebalance(p)

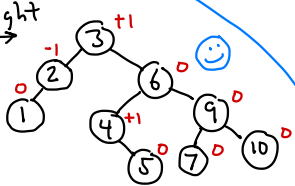
Examples:



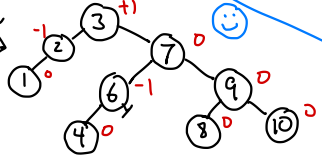
insert(5)



rotate Right



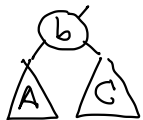
rotate LR



Node types:

2-Node

1 key
2 children



3-Node

2 keys
3 children



↑ Identical heights



Recap:

AVL: Height balanced

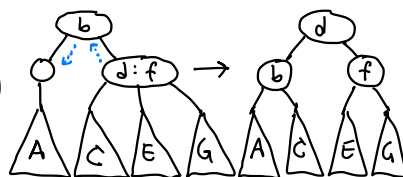
Binary

2-3 tree: Height exact
Variable width

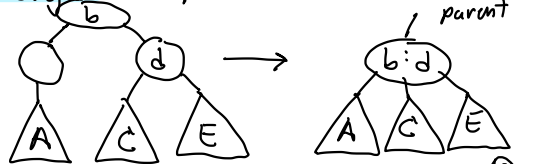


Adoption (Key-Rotation)

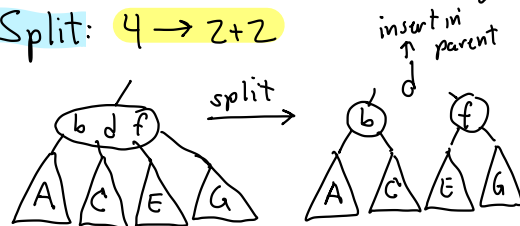
$$1+3 = 2+2$$



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 $O(\log n)$

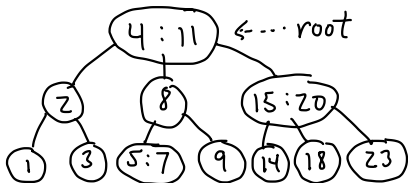
Roughly: $\log_3 n \leq h \leq \log_2 n$

How to maintain balance?

- Split
- Merge
- Adoption (Key rotation)

Example:

2-3 tree of height 2

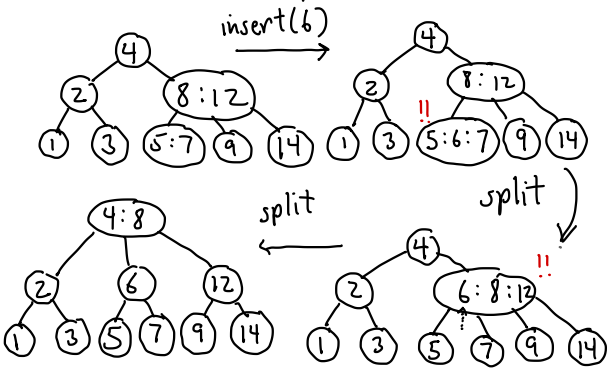


Conceptual tool:

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



Insertion example:



Dictionary operations:

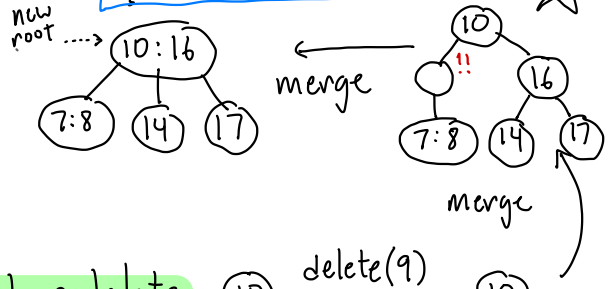
- Find** - straight forward
- Insert** - find leaf node where key "belongs" + add it (may split)
- Delete** - find/replacement/merge or adopt



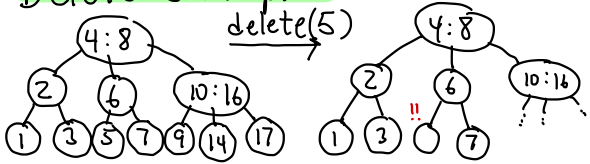
2-3 Trees II

Implementation?

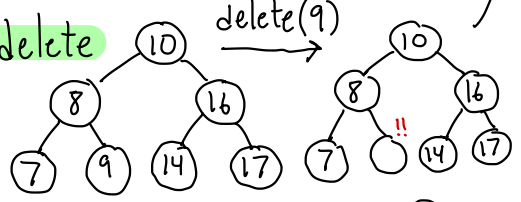
```
class TwoThreeNode {
    int nChildren
    TwoThreeNode children[3]
    Key key[2]
}
```



Delete Example:



Another delete example:

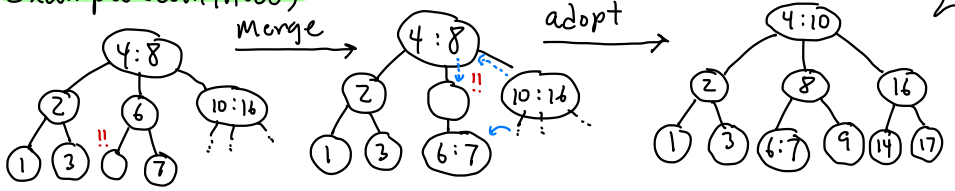


Deletion remedy:

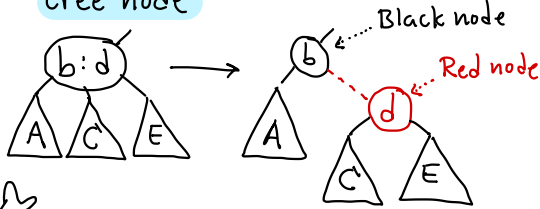
- Have a 3-node neighboring sibling → adopt
- o.w.: Merge with either siblings + steal key from parent



Example (continued)

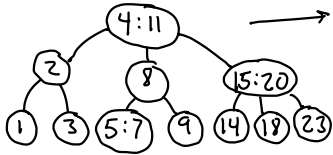


Encoding 3-node as binary tree node

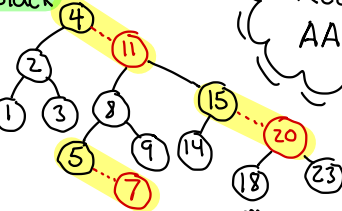


Example:

2-3 Tree:



Red-Black:



Rules:

- Every node labeled red/black
- Root is black
- Nulls treated as if black
- If node is red, both children are black
- 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!

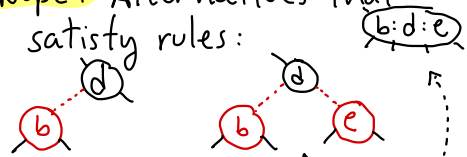
AA-tree:

Arne Anderson 1993

New rule:

- 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

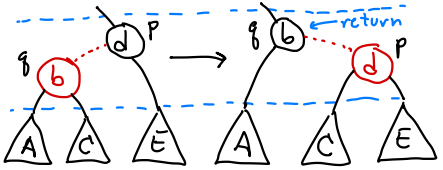
Lemma: A red-black tree with n keys has height $O(\log n)$

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?

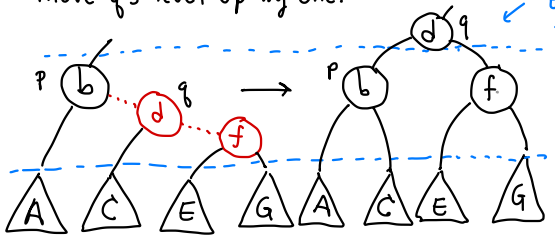
Restructuring Ops:

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 on its right child q , and move q 's level up by one.



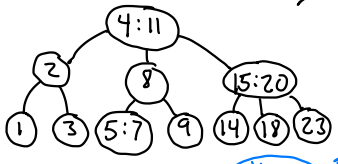
How to test?

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

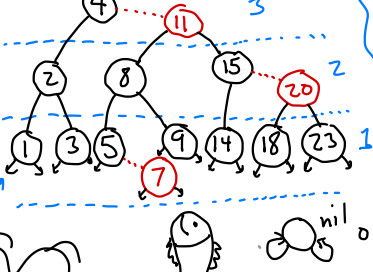


Example:

2-3 Tree:



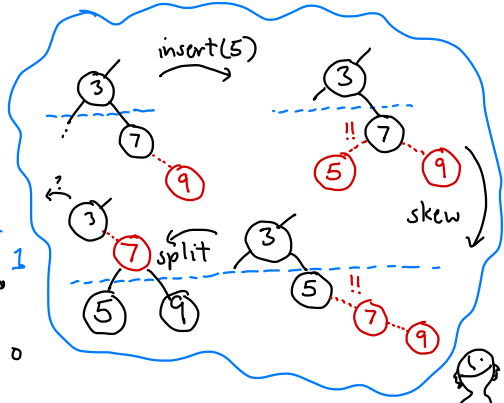
AA tree:



all to nil

Red-Black + AA Trees II

What 2-3 op does this remind you of?



AA Insertion:

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



```

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
}
    
```



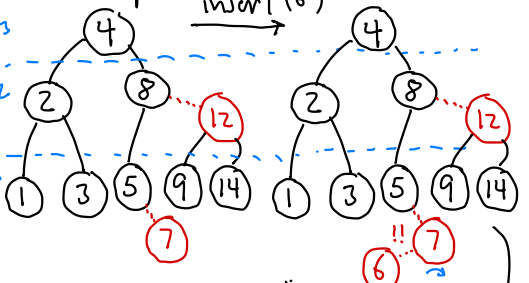
AA Node split (AANode p)

```

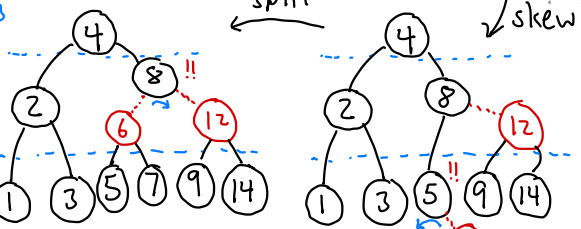
if (p == nil) return p
if (p.right.right.level == p.level) {
    AANode q = p.right
    p.right = q.left
    q.left = p
    q.level += 1
    return q
} else return p
    
```

Example:

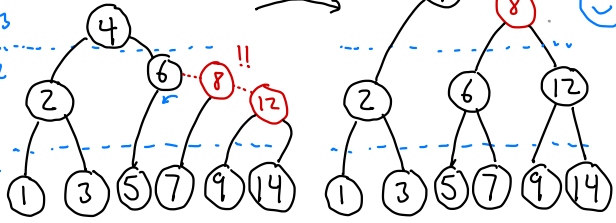
insert(6)



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:

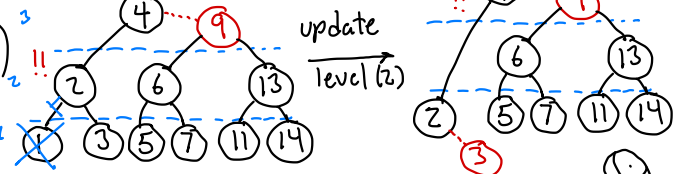
Two more helpers:

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

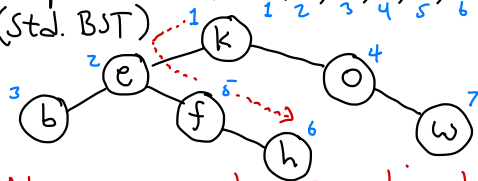
[Explosion of randomized algorithms]

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

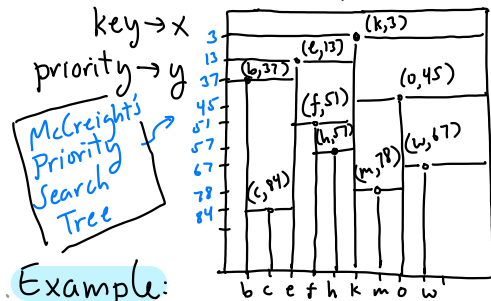
Randomized Data Structures

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



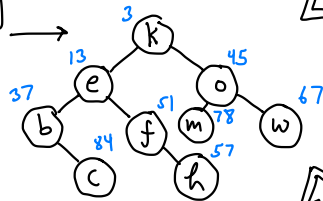
Obs: In a standard BST, keys are by inorder + insert times are in heap order (parent < child)

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?

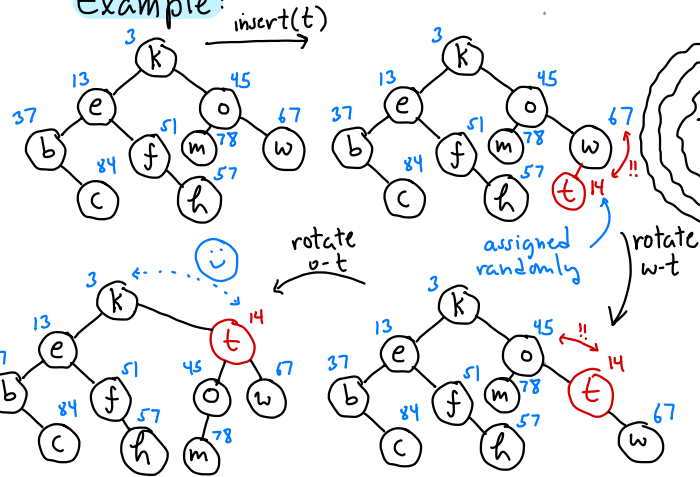
Yes: Just consider the corresponding BST

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

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



Example:



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

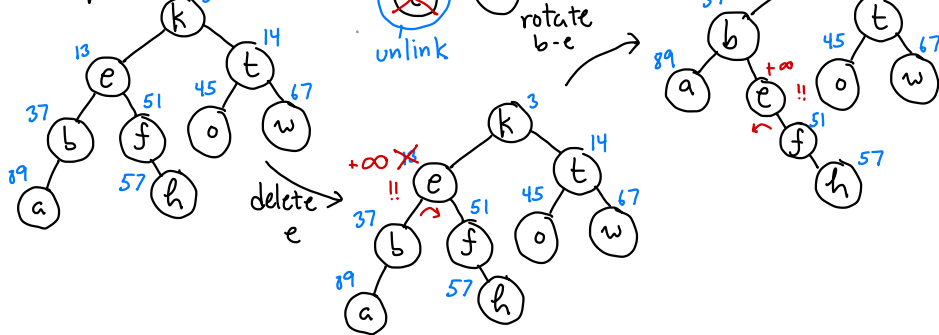


Theorem: A treap containing n entries has height $O(\log n)$ in expectation (averaged over all assignments of random priorities)

Proof: Follows directly from BST analysis



Example:



Implementation: (See pdf notes)

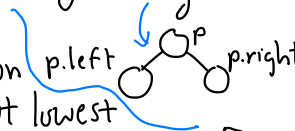
Node: Stores priority + usual...

Helpers:

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

restructure:

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

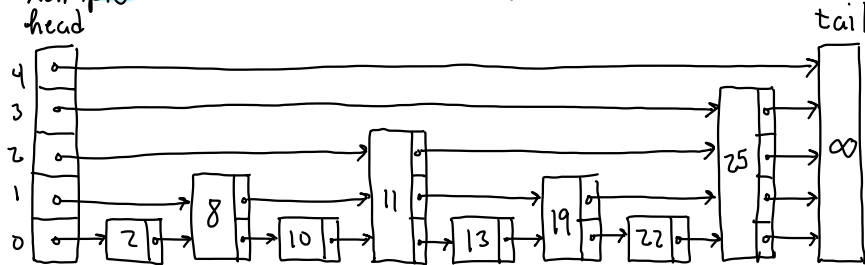


Ideal Skip list:

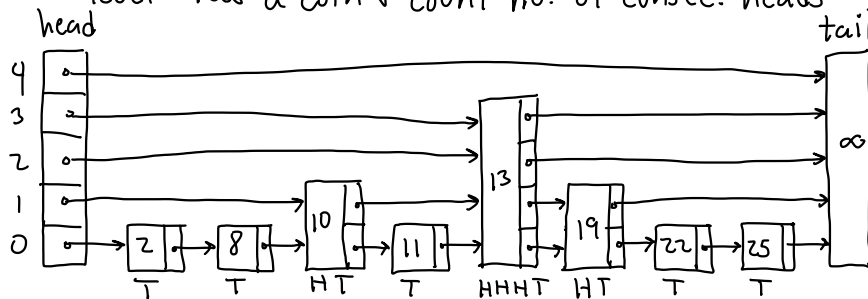
- Organize list in levels
- Level 0: Everything
 - 1: Every other $\circ \rightarrow \circ \rightarrow \circ \rightarrow \circ$
 - 2: Every fourth $\circ \rightarrow \circ \rightarrow \circ \rightarrow \circ$
 - \vdots
 - i : Every 2^i $\circ \rightarrow \circ \rightarrow \circ \rightarrow \circ \rightarrow \circ \rightarrow \circ \rightarrow \circ$



Example:



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



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

```

current node
until we hit base level
advance horizontal
drop down a level
we are at base level

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 $c \cdot \lg n$ is:

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

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

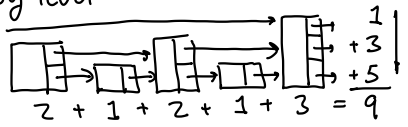
$$= n/n^c = 1/n^{c-1}$$

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 $i = n/2^i$
- ⇒ $E(n_i) = n/2^i$

Total space (expected) is:

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

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

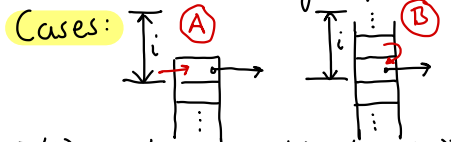
Thm: Expected search time is $O(\log n)$

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.



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

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

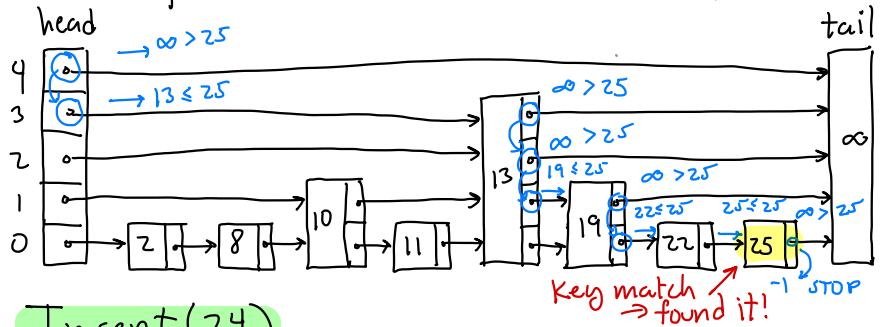
Delete:

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

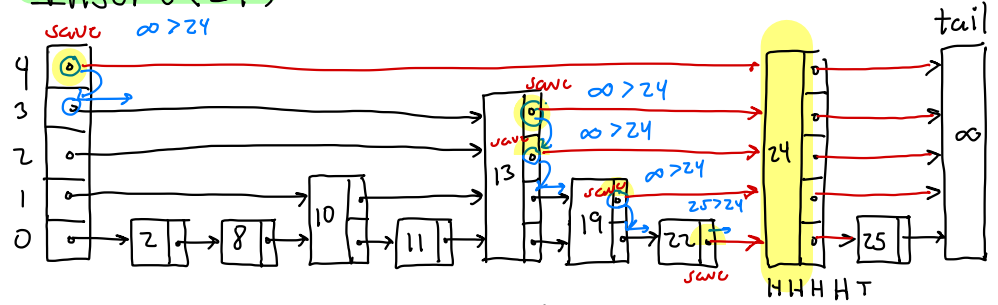
Insert: (Similar to linked lists)

- 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

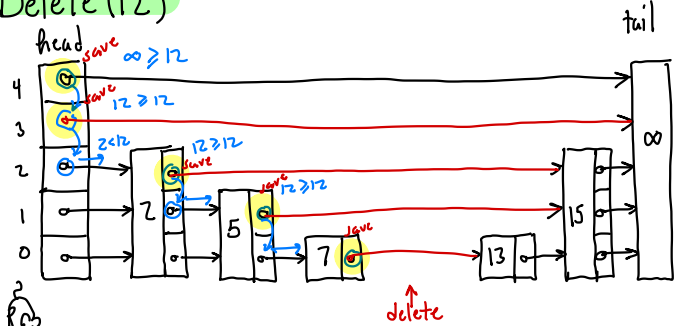
Example: find(25)



Insert(24)



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

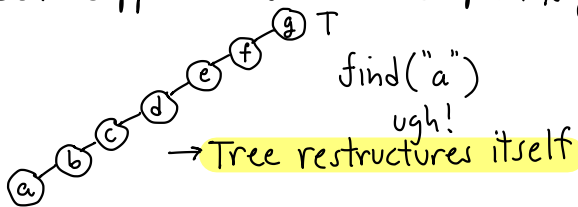
Self-adjusting: Tree adapts as popularity changes

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



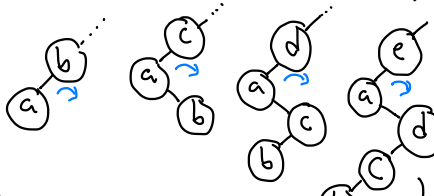
Recap: Lots of search trees

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

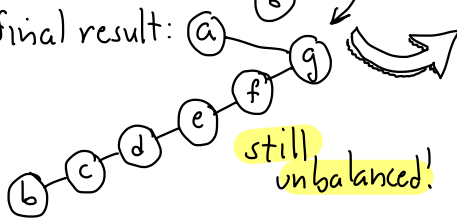
→ **Focus:** Worst-case or randomized expected case

SPLAY TREES I

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

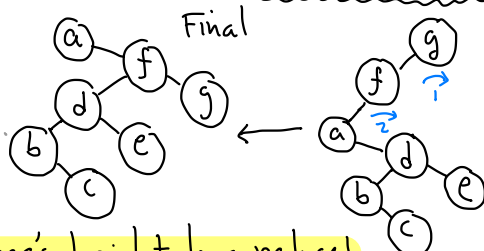


... final result:

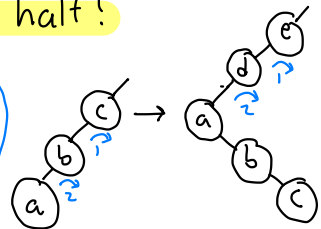


Lesson: Different combinations of rotations can:

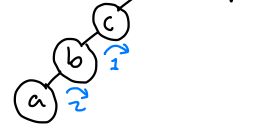
- bring given node to root
- significantly change (improve) tree structure.



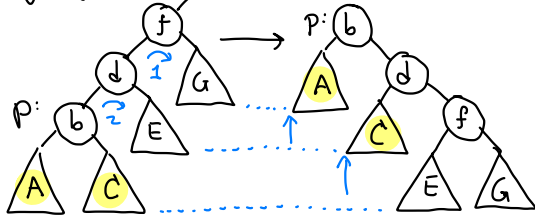
Tree's height has reduced by ~ half!



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

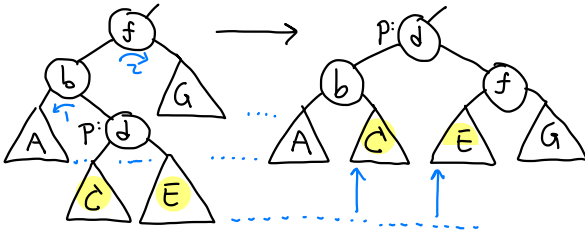


ZigZig(p): [LL case]



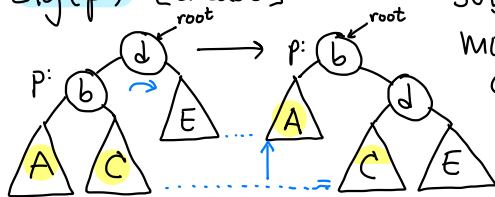
Subtrees A, C move up ↑

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) [nearest node]

```

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):

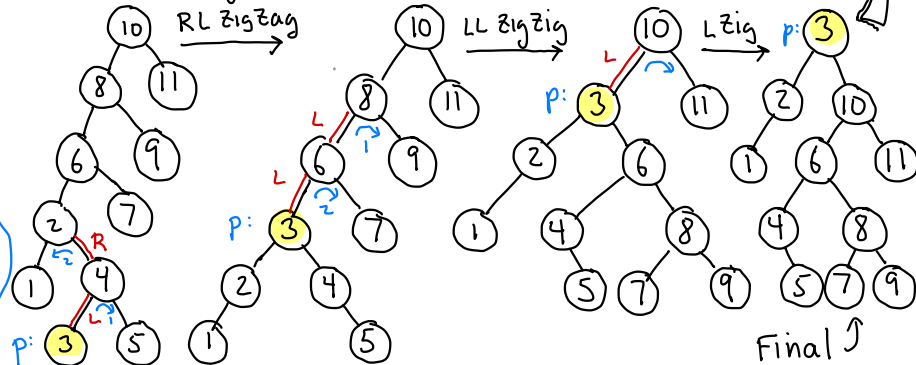
```

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

root.key ≠ x or error!

Splay Trees II

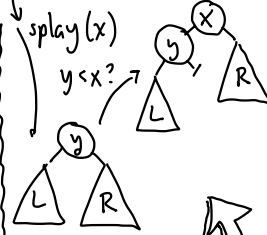
Example: splay(3)

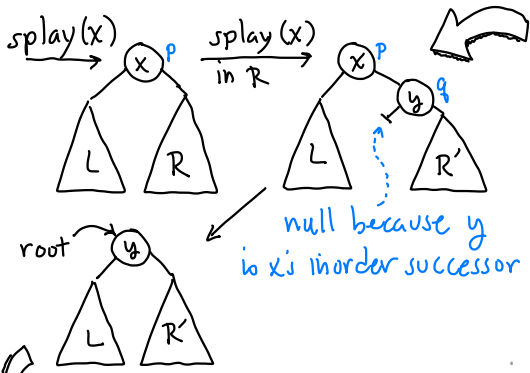


find(x):

```

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





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

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

SPLAY TREES III

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 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.

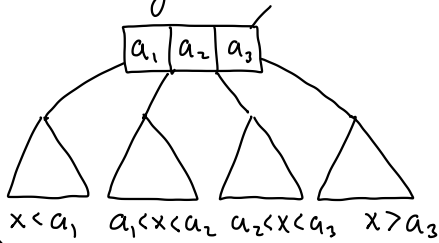
Static Optimality:

- 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 $O(m \sum p_i \lg 1/p_i)$

Multiway Search Trees:



B-Tree:

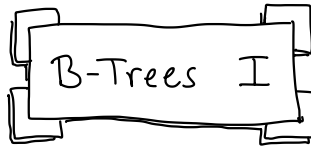
- Perhaps the most widely used search tree
- 1970 - Bayer & McCreight
- Databases
- Numerous variants

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

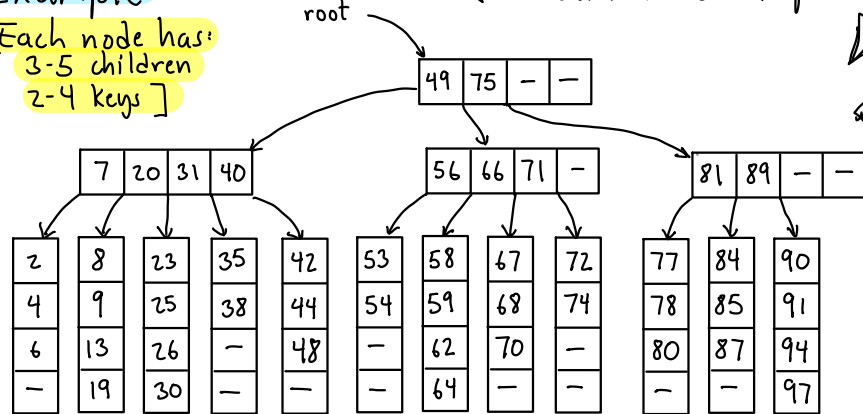
Secondary Memory:

- 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 $M=...$

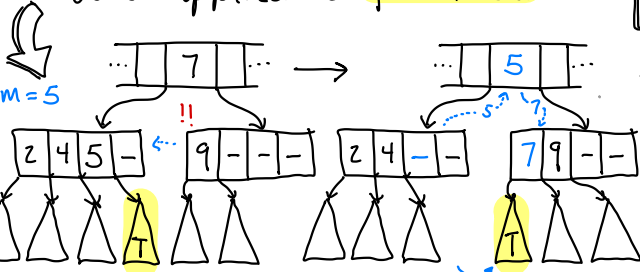
```
class BTreeNode {
    int nChild // no. of children
    BTreeNode child[M] // children
    Key key[M-1] // keys
    Value value[M-1] // values
}
```

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)

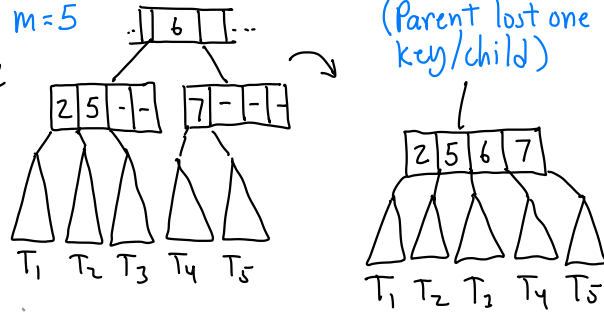
Key Rotation (Adoption)

- 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**.

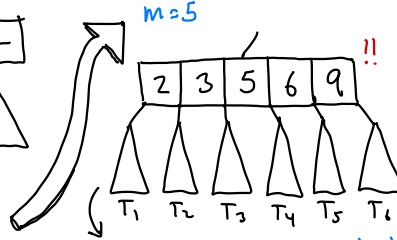


B-Tree restructuring:

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



B-Trees II



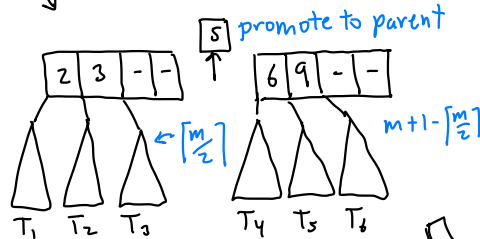
Lemma: For all $m \geq 2$,
 $\lceil m/2 \rceil \leq 2\lceil m/2 \rceil - 1 \leq m$
 \Rightarrow Resulting node is valid

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



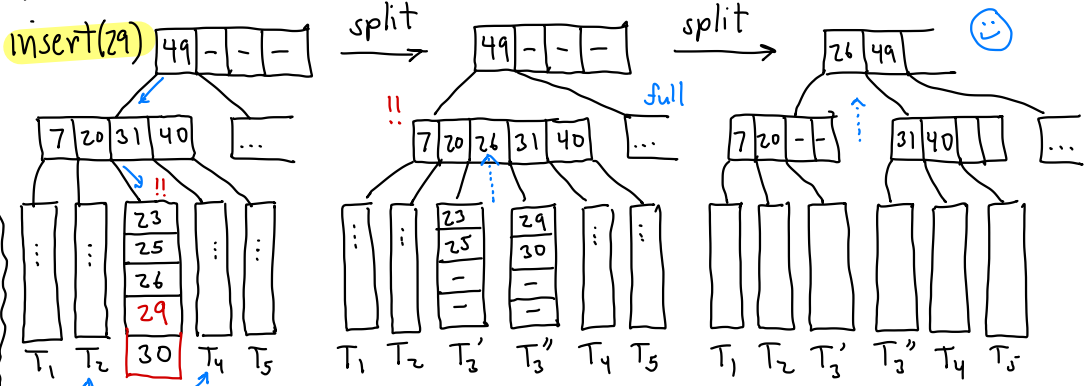
Node Merging:

- 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: $m=5$

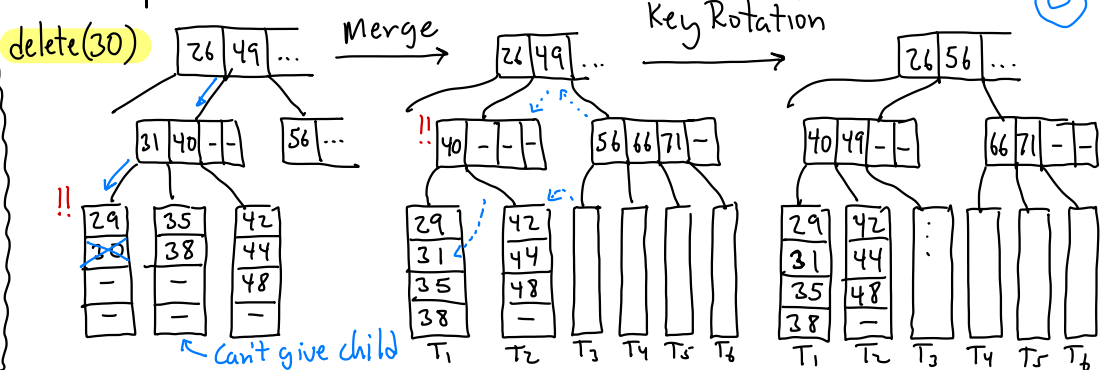


B-Trees III

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: $m=5$



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

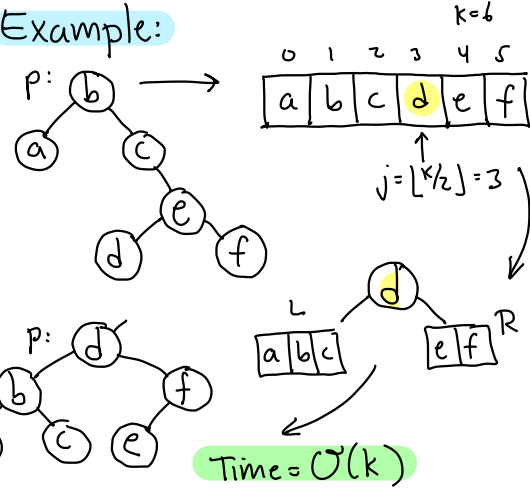


Recap:

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



Example:



Overview:

Insert:

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

Find: Same as std BST

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



Delete:

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

How? Maintain $n, m \leftarrow 0$

Insert: $n++$, $m++$

Delete: $n--$... If $m > 2n$ rebuild

How to rebuild?

rebuild(p):

- inorder traverse p 's subtree \rightarrow 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
- $L \leftarrow$ buildSubtree($A[0..j-1]$)
- $R \leftarrow$ buildSubtree($A[j+1..k-1]$)
- return Node(x, L, R)



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, $size(p)$ = no. of nodes in p 's subtree
 - if $\frac{size(p.child)}{size(p)} > \frac{2}{3}$
- $p \leftarrow rebuild(p)$
- break



Details of Operations:

Init: $n \leftarrow m \leftarrow 0$ $root \leftarrow null$

Delete:

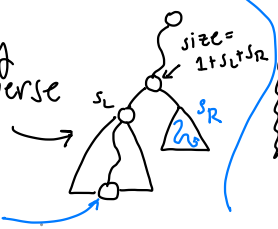
- Same as std BST
 - $n--$
 - if $m > 2n$,
- $rebuild(root)$

Time: $O(n)$



How to compute $size(p)$?

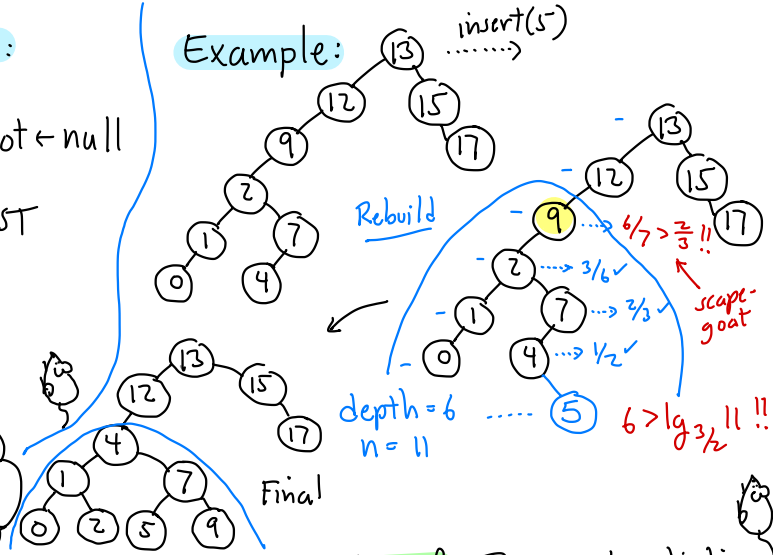
- Can compute it on the fly
- While backing out, traverse "other sibling"
- Too slow? No!
- \rightarrow Change to rebuild.



Must there be a scapegoat? yes!

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

Example:



Proof: By contradiction

- Suppose p 's depth $> \log_{3/2} n$ but \forall ancestors u , $size(u.child) \leq \frac{2}{3} \cdot size(u)$

depth 0: $size = n$

depth 1: $size \leq \frac{2}{3}n$

depth 2: $size \leq \frac{4}{9}n$

...

depth $d > \log_{3/2} n$: $size \leq (\frac{2}{3})^d n$

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

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

Weight Balance:

- Given a set of **keys**

$$X = \{x_0, \dots, x_{n-1}\}$$

- and **values**

$$V = \{v_0, \dots, v_{n-1}\}$$

- and **weights**

$$W = \{w_0, \dots, w_{n-1}\}$$

- Assume:

$$x_0 < x_1 < \dots < x_{n-1} \text{ sorted}$$

$$w_i > 0 \text{ positivity}$$

Pseudo-Probability:

- Let: $\bar{W} = \sum_{i=0}^{n-1} w_i$ **total weight**

- Let: $p_i = w_i / \bar{W}$ **pseudo-prob**

- Obs: $0 < p_i \leq 1$ } discrete prob. distribution
- $\sum_i p_i = 1$ }

Shannon's Theorem: If p_i is the prob. of accessing x_i , any BST has expected search at least $\sum_i p_i \lg(1/p_i)$ ← (called the **entropy** of distribo)



Overview:

- Splay trees - **Static**

Optimality

- More frequently accessed keys closer to root


⇒ **Weight-balanced trees**

Weight-Balanced Trees I

Implementation: (as extended BST)

Internal node:

Stores:
 Key key → splitter
 float wt → total weight of entries in subtree
 left, right



External Node:

Key key, ← x_i
 Value value
 float wt ← w_i



How to (Nearly) Achieve Shannon's bound

→ Weight-balanced tree

→ For each node p :

$wt(p)$ = total weight of keys in p 's subtree

$$\text{balance}(p) = \frac{\max(wt(p.\text{left}), wt(p.\text{right}))}{wt(p)}$$



Given $\frac{1}{2} \leq \alpha \leq 1$, a BST is **α -balanced** if for all internal nodes p , $\text{balance}(p) \leq \alpha$

$\alpha = \frac{1}{2}$: Perfectly balanced
 $= 1$: Arbitrarily bad

$\alpha = \frac{2}{3}$: A reasonable compromise

Balance by Rebuilding:

Given an array $A[0..k-1]$ of external nodes:

$A[i].key$, $A[i].value$,
 $A[i].wt$

- Assume keys are sorted
- Assume weights > 0



Weight-based median:

- Select splitter to minimize left-right weight difference

- Let $\bar{w} = \sum_{i=0}^{k-1} A[i].wt \leftarrow$ Total wt

- Let $\bar{w}_{i,j} = \sum_{m=i}^{j-1} A[m].wt \leftarrow$ Total wt of $A[i..j-1]$

- Let $\Delta_j = |\bar{w}_{0,j} - \bar{w}_{j,k}| \leftarrow$ Absolute diff if we split $A[0..j-1] \{A[j..k-1]$

- Goal: Split at $0 \leq j < k$ that minimizes:

$$\Delta_{min} = \min_{0 \leq j < k} \Delta_j \leftarrow \text{Most balanced split}$$

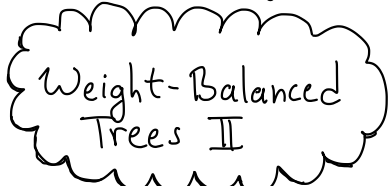


How to maintain balance?

Options:

- Rotations: Similar to AVL trees (single + double)
... \rightarrow BB[x] trees

- Rebuild subtrees: Similar to scapegoat



Example:

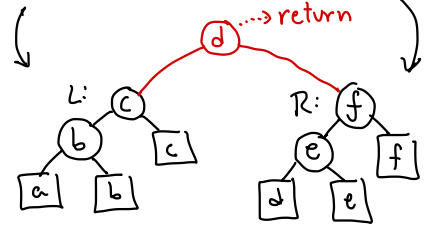
	0	1	2	3	4	5	$k=6$
key:	a	b	c	d	e	f	
wt:	$\leftarrow 3$	$\leftarrow 2$	$\leftarrow 4$	$\rightarrow 1$	$\rightarrow 2$	$\rightarrow 4$	

$$\Delta_{min} = |(1+2+4) - (3+2+4)| = 2$$

$\bar{w} = 16$

$L = \text{buildTree}(A[0..2])$

$R = \text{buildTree}(A[3..5])$



buildTree(A[0..k-1])

if ($k == 1$) return $A[0]$ /* base case */

$\bar{w} = \sum_{i=0}^{k-1} A[i].wt$ /* total weight */

Init: $b = 0$; $Lwt = 0$; $Rwt = \bar{w}$; $\Delta_{min} = \bar{w}$

for ($i = 0 \dots k-1$)

$Lwt += A[i].wt$; $Rwt -= A[i].wt$

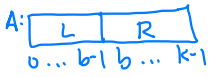
$\Delta = |Rwt - Lwt|$ /* weight difference */

if ($\Delta < \Delta_{min}$) { $b = i+1$; $\Delta_{min} = \Delta$ }

$L = \text{buildTree}(A[0..b-1])$

$R = \text{buildTree}(A[b..k-1])$

return new IntNode($A[b].key, L, R$)



But it is pretty close! 😊 ↩

Theorem: (Mehlhorn '77)
The above balanced split algorithm produces a tree whose exp. search time is

$$\leq H + 3$$

where H = entropy bound.



Dictionary Operations:

→ Balance by destroying + rebuilding - **Jackhammer Trees**

Find: Same as usual. Tree height $\leq \log_{3/2} n$, so $O(\log n)$ time-guaranteed.

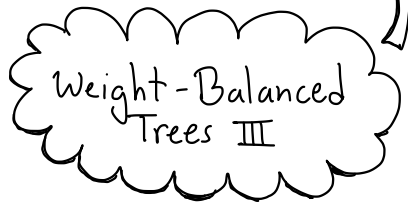
Insert/Delete: Start same as standard BST
→ After operation completes **check + rebuild**

Analysis:

Does this algorithm produce the **optimal tree** (w.r.t. expected case search time)?

- No. 😞 The optimal BST can be computed by **dynamic programming**

CMSC 451



Check + Rebuild:

- When returning from recursive calls, update each node's weight
 $p.wt \leftarrow p.left.wt + p.right.wt$
- Starting at root, walk down search path. Stop at first node p s.t.

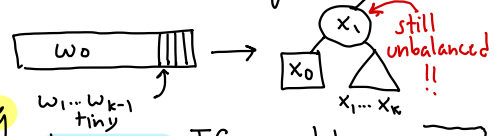
Recall def earlier

$$balance(p) > \alpha$$

Given by designer e.g. $\alpha = 2/3$

Bad weight distributions?

- If a weight is very large relative to neighbors, rebalance may be ineffective



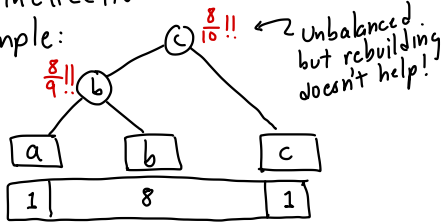
Lemma: If weights are "nice" (not too much variation), insert + delete run in $O(\log n)$ amortized time.

→ If no such p found - Great! Tree is balanced

Else: **Jackhammer!**
- Traverse p 's subtree inorder, store extern nodes in array $A[0..k-1]$
- Replace p 's subtree with $buildTree(A)$

Very heavy entries:

- If an entry's weight is too high, rebuilding is ineffective
- Example:



- This tree is best possible!

- Exemption: Don't rebuild if a key's weight is very high

For node p : $\max(p) =$
max weight in p 's subtree

$$\text{max-ratio}(p) = \frac{\max(p)}{\text{weight}(p)}$$

Given parameter $0 < \beta < 1$,
a node is β -exempt if
 $\text{max-ratio}(p) > \beta$

Dictionary Operations:

- **find**: as usual
- **insert**: insert as usual but rebuild if needed
- **delete**: delete as usual but rebuild if needed

Weight-Balanced
Trees IV

(α, β) -balance: Every
internal node p is either
 α -balanced or
 β -exempt

Lemma: For any set of
weighted entries, \exists an
 (α, β) -balanced BSTree if
 $\frac{1}{2} < \alpha < 1$ and $\beta < 2\alpha - 1$

When to rebuild?

- When "backing out" from insert/delete, update node weights
- Walk down search path from root [opposite from scapegoat!]

- If any node p is out of balance:

$$\text{balance}(p) > \alpha$$

-and-

$$\text{max-ratio}(p) \leq \beta$$

then:

- **Rebuild p** :

- Traverse p 's subtree in order
- Collect external nodes in array $A[0..k-1]$
- replace p with $\text{buildTree}(A)$

Eg.
 $\alpha = 2/3$
 $\beta = 1/4$

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 \downarrow simple/slow
 - Quadratic probing \downarrow complex/fast
 - Double hashing



Collision Resolution:

If there were **no collisions** hashing would be trivial!

- insert**(x, v) \rightarrow table[$h(x)$] = v
- find**(x) \rightarrow return table[$h(x)$]
- delete**(x) \rightarrow 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

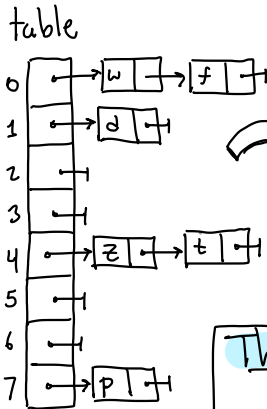
Separate Chaining:

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

$m=8$



Hashing II

Token-based - See latex notes!

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

- 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

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

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 $\lambda > \lambda_{max}$...

Open Addressing:

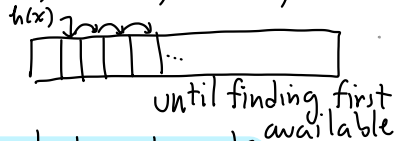
- Special entry ("empty") means this slot is unoccupied
- Assume $\lambda \leq 1$
- To insert key: check: $h(x)$ if not empty try
 - $h(x) + i_1$
 - $h(x) + i_2$
 - \vdots

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

- What's the best probe sequence?

Linear Probing:

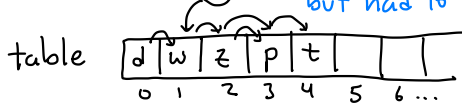
$h(x), h(x)+1, h(x)+2, \dots$



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 "

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 $m=4, j^2 \bmod 4 = 0 \text{ or } 1 \text{ but not } 2 \text{ or } 3$)
- 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.

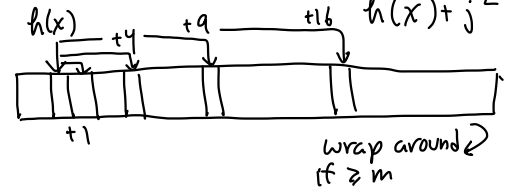
Pf: See latex notes.

Clustering

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

Quadratic Probing:

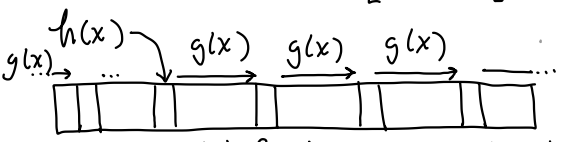
$h(x), h(x)+1, h(x)+4, h(x)+9, \dots, h(x)+j^2$



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

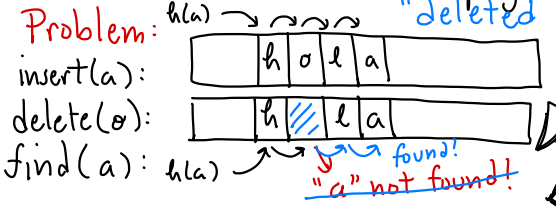
\rightarrow Proof is nontrivial (skip)

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

very efficient!

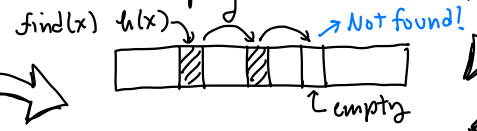
Delete(x): Apply find(x)
 \rightarrow Not found \Rightarrow error
 \rightarrow Found \Rightarrow set to "empty"

Is this right??



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!)

Geometric Search:

- Nearest neighbors \rightarrow 
- Range searching \rightarrow 

- Point Location




- Intersection Search

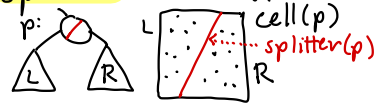


Sofar: 1-dimensional keys

- 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.

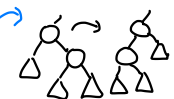
Point: A d -vector in \mathbb{R}^d
 $p = (p_1, \dots, p_d)$ $p_i \in \mathbb{R}$

Multi-Dim vs. 1-dim Search?

Similarities:

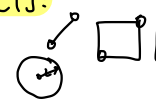
- Tree structure
- Balance $\mathcal{O}(\log n)$
- Internal nodes - split
- External nodes - data

Differences:

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

Quadtrees & kd-Trees I

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 

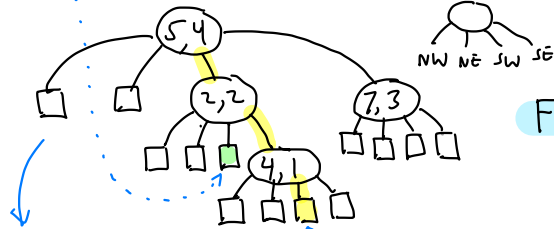
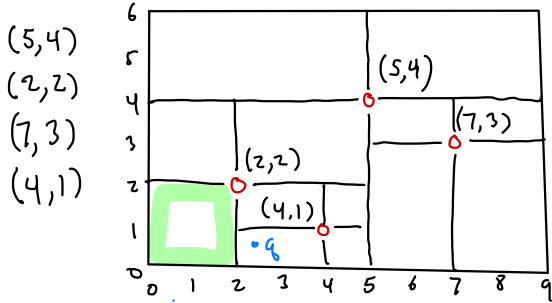
class Point {

```

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?

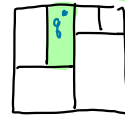
→ Follow path from root down (generalizing BST find)

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



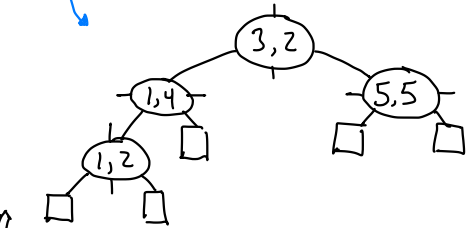
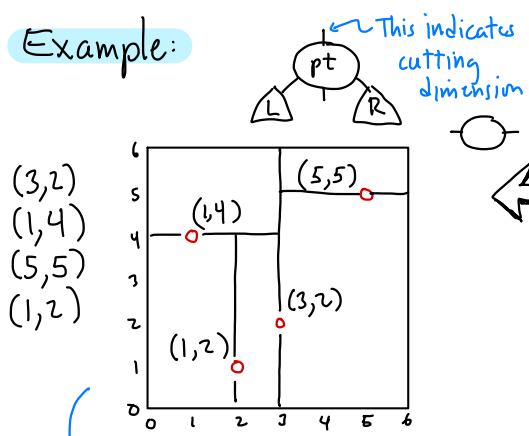
left: left/below
right: right/above

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:

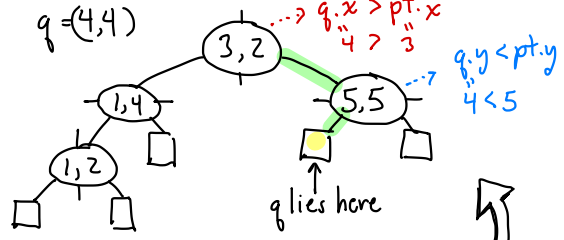


Kd-Tree Node:

```
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) → 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) → all coords match?
        return p.value
    else if (p.onLeft(q))
        return find(q, p.left)
    else
        return find(q, p.right)
}
```

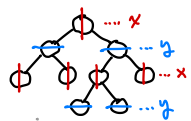
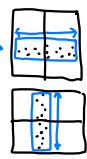
Find:

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

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

How do we choose cutting dim?

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



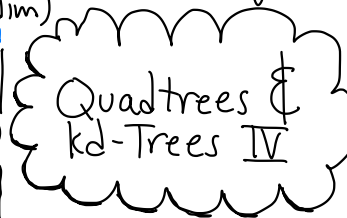
KDNode insert (Point x, Value v, KDNode p, int cd)

```

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
    
```

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

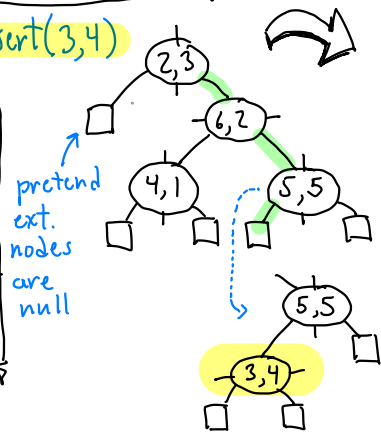
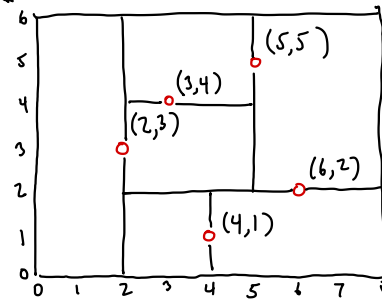
This is the hardest part. See Latex notes.

Rebalance by Rebuilding:

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

Example:

insert(3,4)



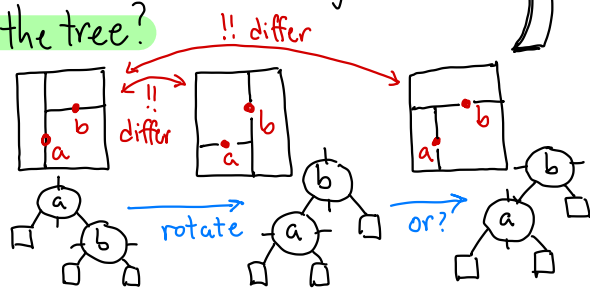
Analysis:

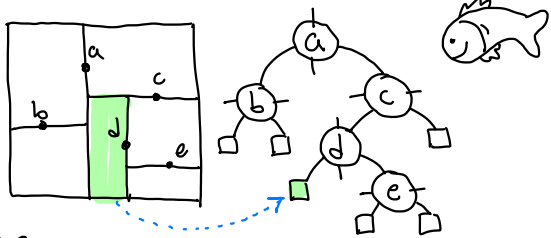
Run time: $O(h)$

Tree height

Can we balance the tree?

- Rotation does not make sense





Kd-Trees:

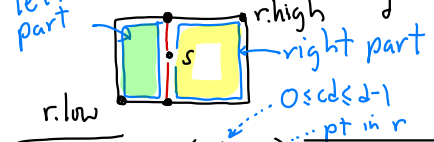
- Partition trees \rightarrow vert

L	R
---	---
- Orthogonal split \rightarrow 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



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

$r.\text{rightPart}(cd, s)$
 $\rightarrow high = r.high + low = r.low$ 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

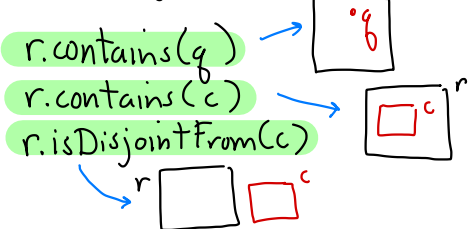
- Defined by two pts: $low, high$



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

Useful methods:

Let r, c - Rectangle
 g - Point



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

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
 ⇒ $O(n)$ time?
 → At some we don't recurse at all!

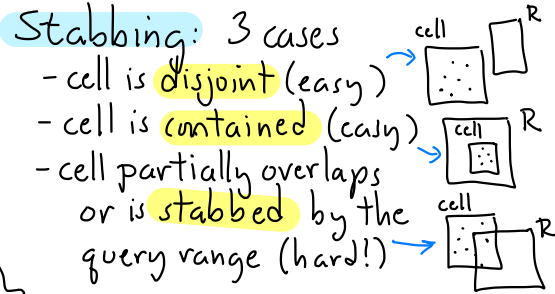
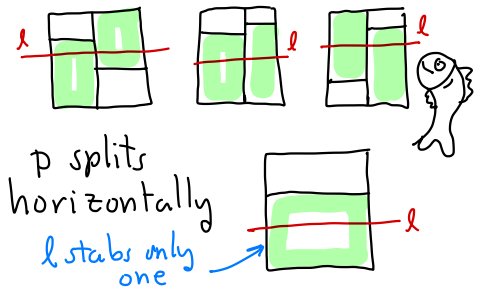
Solving the Recurrence:
 - Macho: Expand it
 - Wimpy: Master Thm (CLRS)

Master Thm:
 $T(n) = aT(\frac{n}{b}) + n^d + d \log_b a$
 $\Rightarrow T(n) = n^{\log_b a}$
 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) → (Each has $n/4$ pts)
 Recurse on 2 grandchildren

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

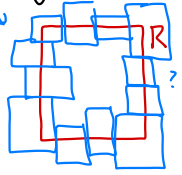


Kd-Tree Queries III

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.

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

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

Recap:

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

- Orthogonal range query: Count/report pts in axis-aligned rect.  Ans=4

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

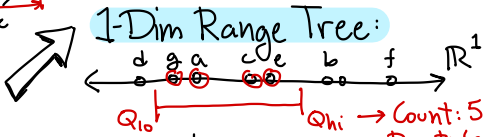
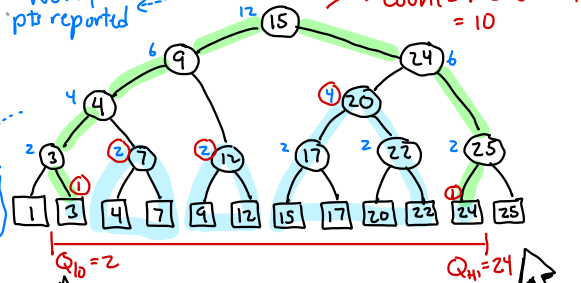
No. of pts reported



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 $O(k + \log n)$)

No. of pts reported ←  Count = 1+2+2+4+1 = 10



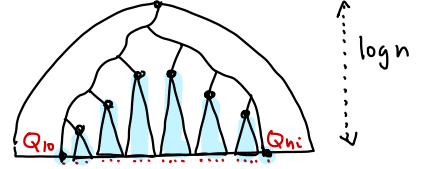
1-Dim Range Tree:

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 $Q_{lo} + Q_{hi}$ + 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:

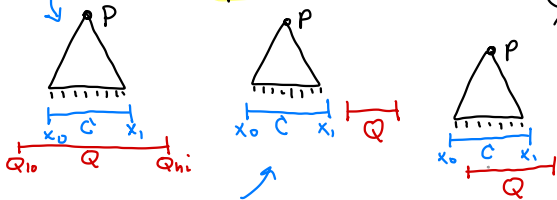
p is external:

- 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

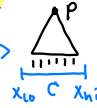
More details:

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
```

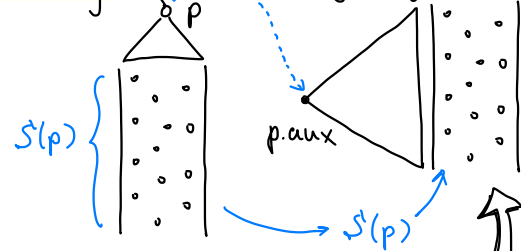
```
else return:
```

```
    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)

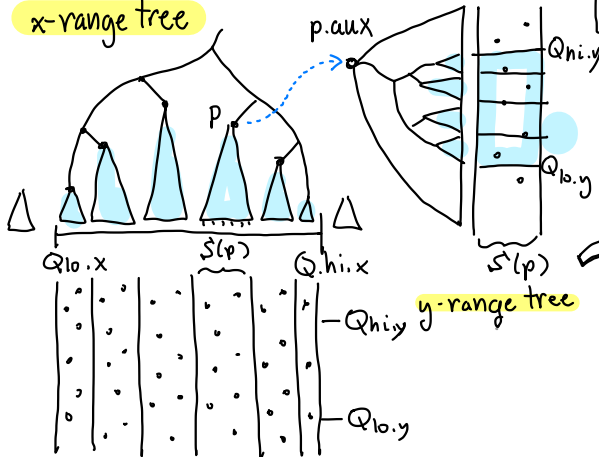
Answering Queries?

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

⇒ Total: $O(\log n \cdot \log n) = O(\log^2 n)$

History:

- Driscoll, Sarnak, Sleator, Tarjan (1986) - First serious theoretical analysis
- Applied to **geometric search** (time is coordinate)



Splay trees

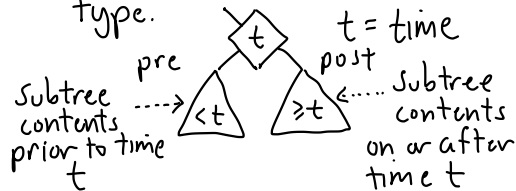
Approaches:

Copy-on-write: Whenever you modify, make a copy
→ very inefficient!

Change-log: Make a list of recent updates
→ slow to process

Fat nodes / Node copying: When changes occur, save just modified portions

→ **Temporal node:** New node type.



Persistent Data Structures:

- Preserves **prior states** of data structure
- Allows **searches in history**
→ "Was Jane Smith enrolled in UMD in Spring 2016?"
- **Full/Partial Persistence:**
- change can be made to any/current version



Case study: PBJ Tree Persistent Weight-Balanced Jackhammer trees

- A **partially persistent** BJ tree
- Uses **rebuilding** to balance subtrees
- Allows **weighted entries**

Example: Rebuild subtree T_1 at time t . Let T_2 be new subtree:

BJ Tree:



PBJ Tree:



Approach: Whenever a modification made - save old contents + use temporal node to distinguish

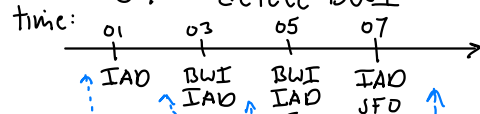
Change at t :



Example:

Time: **Command:**

01	insert IAD
03	insert BWI
05	insert SFO
07	delete BWI



find IAD at 00: not found
find IAD at 02: found
get-min at 04: BWI
get-min at 08: IAD

PBJ Tree (Private) Data:

- final float ALPHA, BETA
↳ same as BJT Tree
- Node root - root → init: null
- int firstInsertTime } init: -1
- int lastInsertTime }

Insertion (without rebalancing).

Helper: Node insert(x, v, w, t)

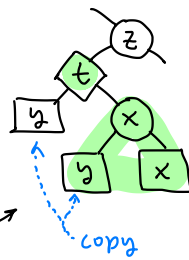
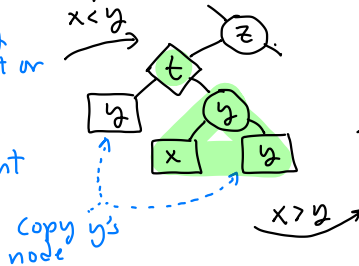
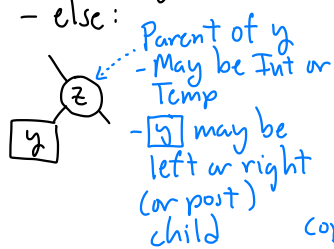
→ insert (x, v) of wt w at t

- If root == null: (first insertion)
 - firstInsertTime ← t
 - root = new ExtNode (x, v, w)

- else

- search for x in tree
- reach ExtNode y
 - if $x == y \Rightarrow$ Error: Duplicate key
 - else:

Int: left ← left.insert(...)
right ← right.insert(...)
Temp: post ← post.insert(...)



→ Update node weight + maxWt (later)

Weight, maxWt

Persistent Search Trees II

PBJ Tree: Tech Specs

Node structure:

- Node: weight + maxWt (float)
 - ↳ TempNode: time (int)
 - pre, post (Node)
 - Int Node: key, left, right
 - Ext Node: key, value

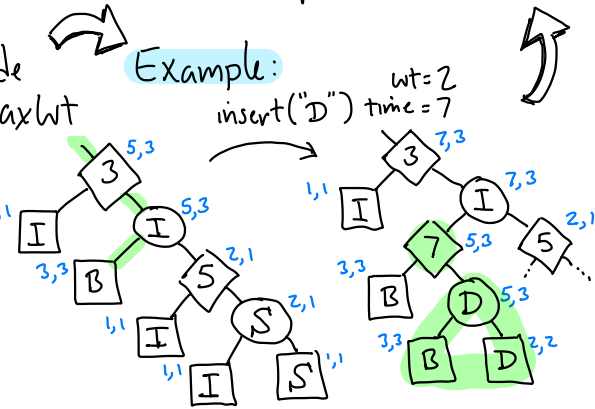
Few More Things:

- Before insert, check that $t >$ lastInsertTime → else Error!
- Set lastInsertTime ← t
- Note: Partial persistence
- Rebalance → Next

Updating Weights?

- Int Node: $w_t \leftarrow \text{left}.w_t + \text{right}.w_t$
 $\text{maxWt} \leftarrow \max(\text{left}.w_t, \text{right}.w_t)$
- Temp Node: $w_t \leftarrow \text{post}.w_t$ (only post!)
 $\text{maxWt} \leftarrow \text{post}.w_t$

Example:



Internal Node Rebalance:

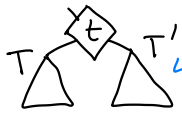
- Test α, β condition (same as B+ tree)
- If unbalanced, compile list A of extern nodes for **current tree**

For temporal recurse only on post side

The new tree has no temporal nodes

- T = original tree
- T' = buildTree(A)

- return:



Temporal Node Rebalance:

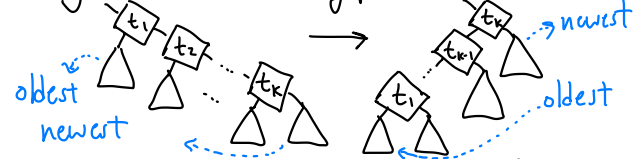
- Recurse on post side: **post = post.rebalance(x, t)**

- On return:

if (post child is temporal)...

- perform left rotation
 - update weights
- return root of subtree

Why? Don't like long post chains



Rebalancing after insertion:

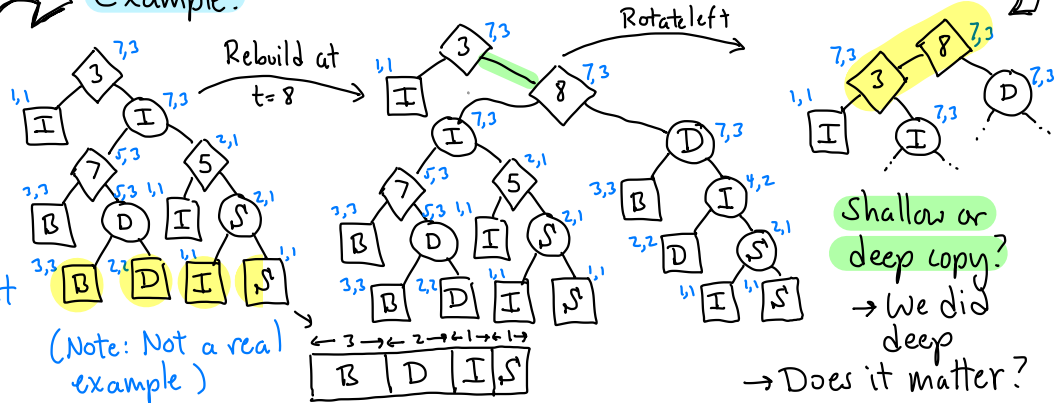
- Starting at root, retrace search path
- Recursive helper:

Node rebalance(x, t)

- **Internal**: Apply to left/right based on $x \leftrightarrow \text{key}$
- **Temporal**: Apply to post only.



Example:



find(x, t), findUp(x, t), getMin(t)...

- same as before but for temporal nodes visit pre/post based on t
- check whether $t < \text{firstUpdateTime}$
→ if so → null

getPreorderList(t)/getFullPreorderList()

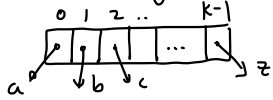
- First gets preorder list at time t (no temporal nodes!)
- Second gets full preorder list (all nodes)

Delete/Clear: Not implemented/Not required! (A bit messy)

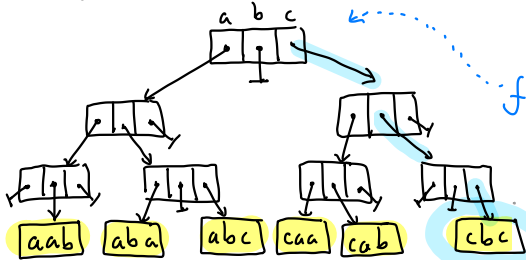
Tries: History

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

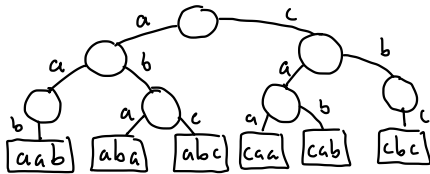
Node: Multiway of order k



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:

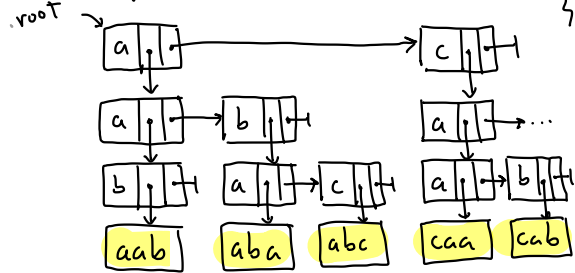
- No. of nodes \sim total no. of chars in all strings
- Space $\sim k \cdot (\text{no. of nodes})$

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

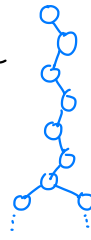
Patricia Tries:

- 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.



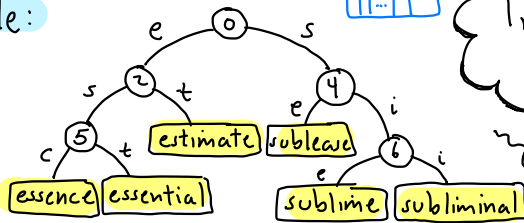
Example:

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

Path compression!

Example:

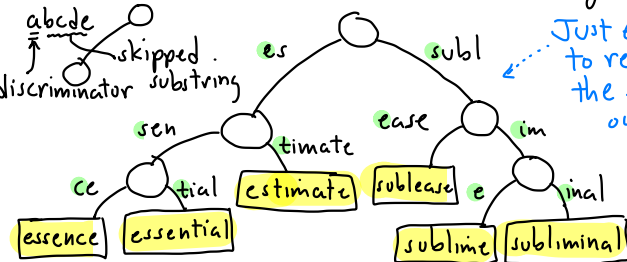
- essence
- essential
- estimate
- sublease
- sublime
- subliminal



Branch based on i^{th} char of string

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

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

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

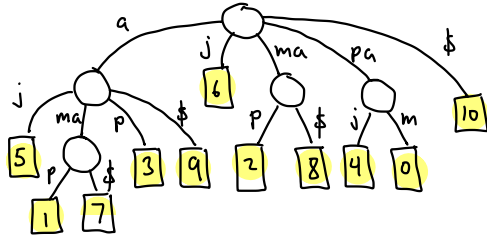
Def: **Substring identifier** for

- 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:** $S = a_0 a_1 a_2 \dots a_{n-1} \#$ (special terminal)
- **Suffix:** $S_i = a_i a_{i+1} \dots a_{n-1} \#$
- **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. $ID(S, a) = \text{amap}$ $ID(S, a) = \text{ama\$}$

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

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

Example:

Search("ama") → End at intern node
 Report: 2 occs. ← 1, 7

Search("amapaj") → End at extern node
 Go to S_i + verify ← 1

Tries and Digital Search Trees III

Analysis:

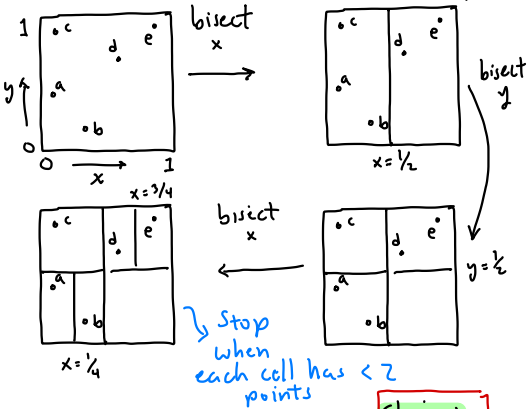
- Space: $O(n)$ nodes
 $O(n \cdot k)$ total space
 $(k = |Σ| = O(1))$
- Search time: \sim to length of target string
- Construction time: $\sim O(n \cdot k)$ [nontrivial]

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

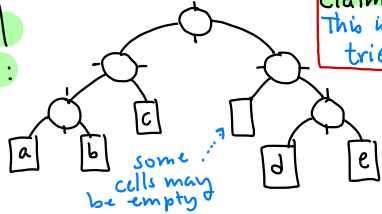
Geometric Applications:

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

Assume points lie in unit square



Final tree:



Claim: This is a trie!

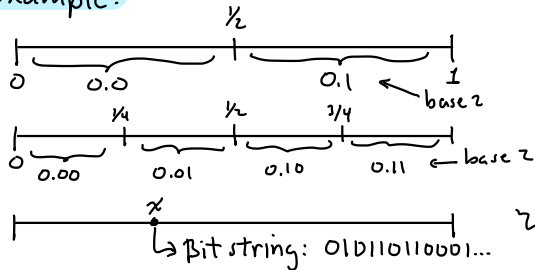
Binary Encoding:

- 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**:

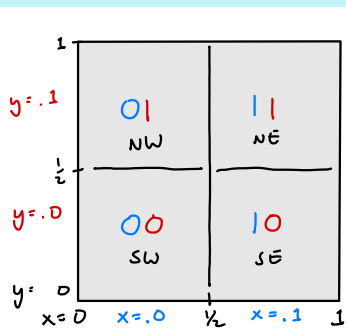
$$x = 0.a_1a_2a_3\dots \quad a_i \in \{0,1\}$$

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

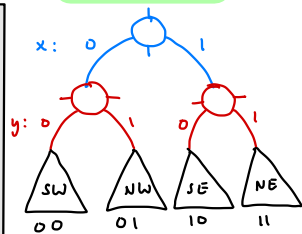
Example:



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**
 $x = 0.a_1a_2\dots$
 $y = 0.b_1b_2\dots$
 $z = 0.c_1c_2\dots$

Lemma: Given a pt set $P \subseteq \mathbb{R}^2$

(in unit square $[0,1]^2$) let

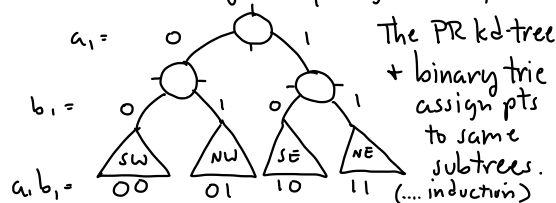
$P = \{p_1, \dots, p_n\}$ where $p_i = (x_i, y_i)$

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.

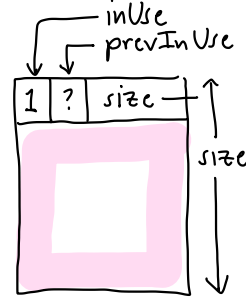
- **Stack** - local vars, recursion
- **Heap** - for "new" objects

Don't confuse with heap data structure/heap sort

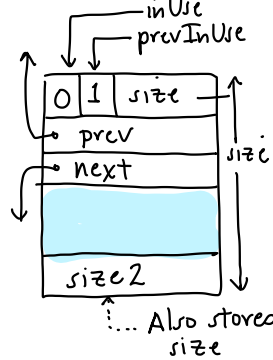


Block Structure:

Allocated:

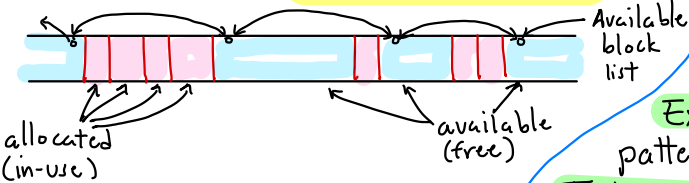


Available:



Explicit Allocation/Deallocation

- Heap memory is split into **blocks** whenever requests made
- **Available blocks**:
 - merged when contiguous
 - stored in **available block list**



Memory Management I

Fragmentation:

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



External: Caused by pattern of alloc/dealloc

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

Guide:

prevInUse: 1 if prev. contig. block is allocated

prev/next: links in avail. list

size/size2: total block size (includes headers)



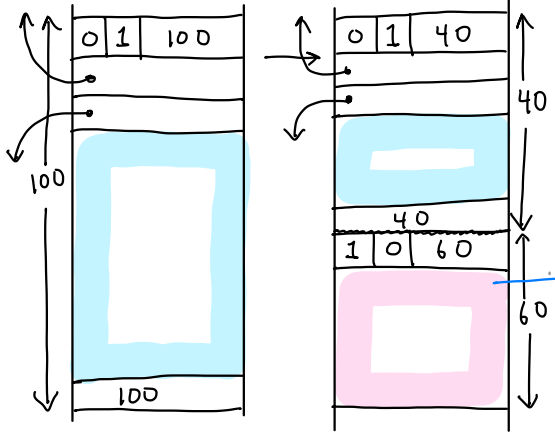
How to select from available blocks?

- **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: $\text{malloc}(b)$

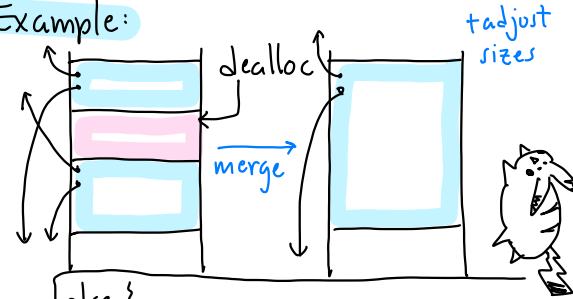
- Search avail. list for block of size $b' \geq b+1$
- If b' close to b : alloc entire block (unlink from avail list)
- Else: split block

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

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

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)

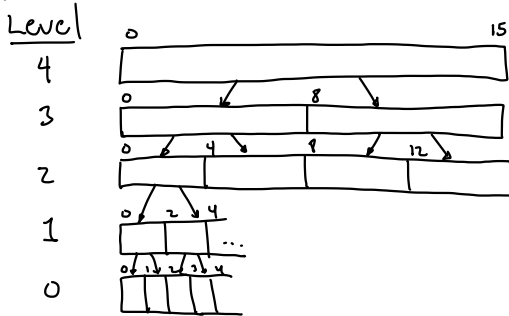
```

else {
    p.size -= b // remove allocation
    *(p+p.size-1) = p.size // size 2
    q = p+p.size // start of new block
    q.size = b
    q.prevInUse = 0 // new block header
}
q.inUse = 1
(q+q.size).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.

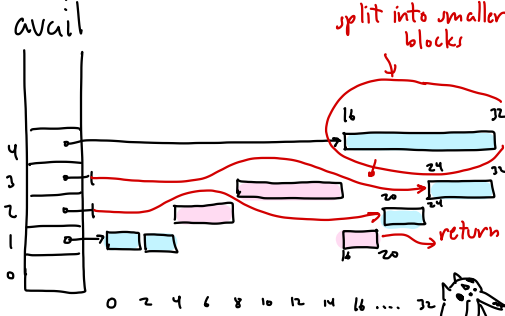
Memory Management III

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

