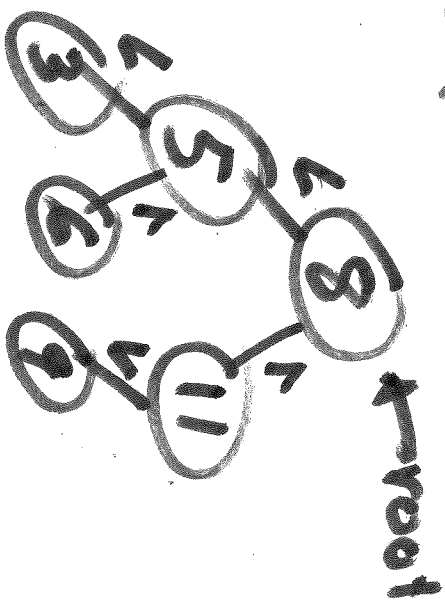


Binary Search Tree

(111)

- Each internal node stores an item (key, data)
- All the keys stored in the left subtree are less than "key"
- All the keys stored in the right subtree are greater than key



- Searching for a key will take time proportional to the distance from the root to the key.
- If we have a tree where height $\leq O(\log n)$ then we say that tree is balanced

and find(key) takes $O(\log n)$. (112)

treeSearch(6, root)

+ Operations

~~Insert~~

val = treeSearch(key, node)

if node is null // key not found
return null // ~~node~~ not found

$O(\log n)$
if (key == node.key) // key is in current node return data
return node → data

if (key < node.key) {

return treeSearch(key, node → left)

else // key is in right subtree

return treeSearch(key, node → right)

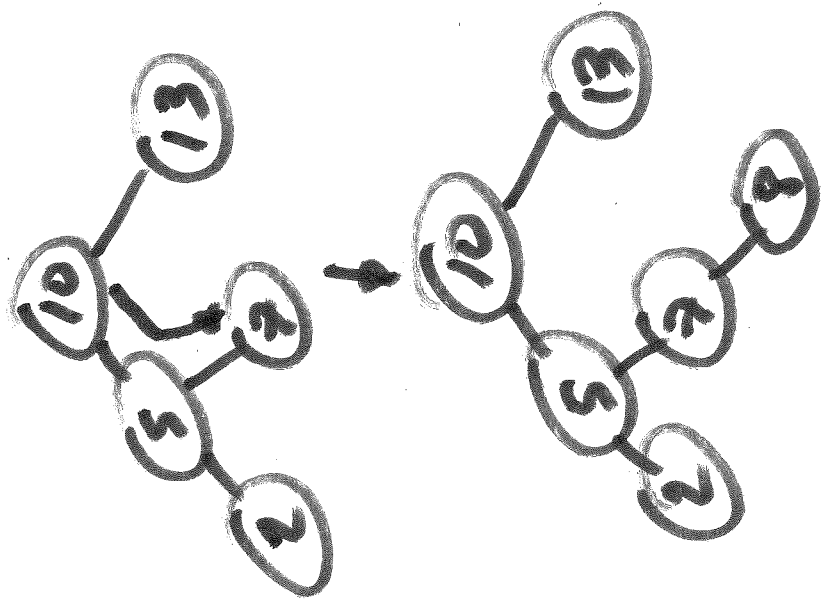
not balanced

$O(n)$
 $O(n)$

insert (key, data)

113

- Search key in tree.
- If found then replace data
- If not found then add a new node at the bottom where the search ended



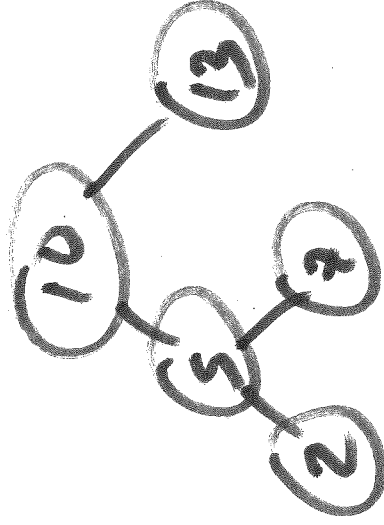
insert (8)

114

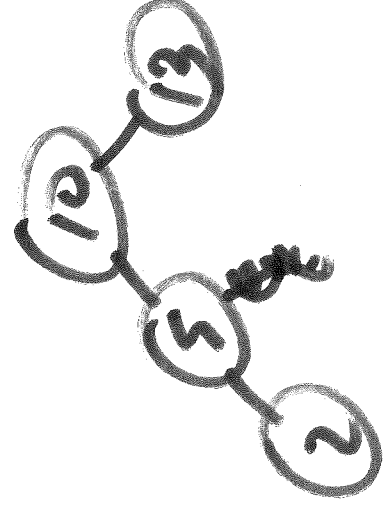
remove(key)

→ Search key in tree

→ If node is at the bottom just remove it.

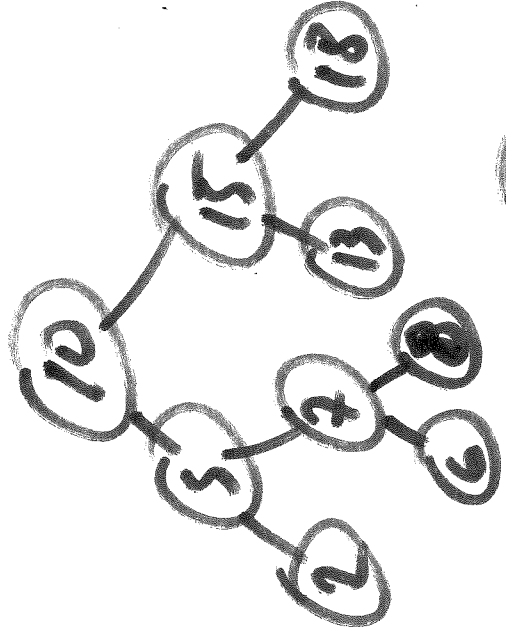


remove(7)



Quiz 12
Given the tree

115.1

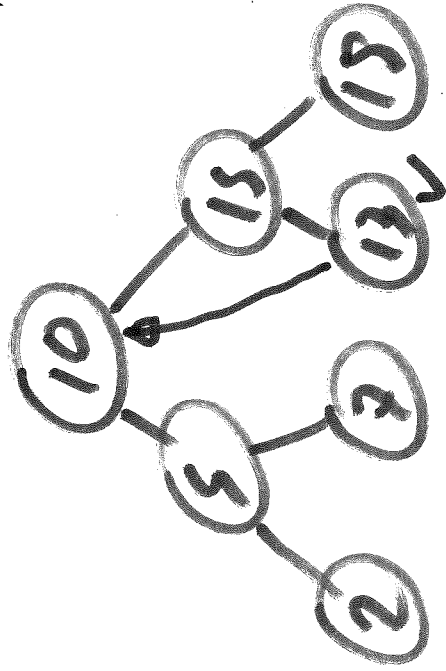


If we remove 5 what node will take its place

- a) 2
- b) 7
- c) 6
- d) 8

→ If node to remove is internal then

11.5.2



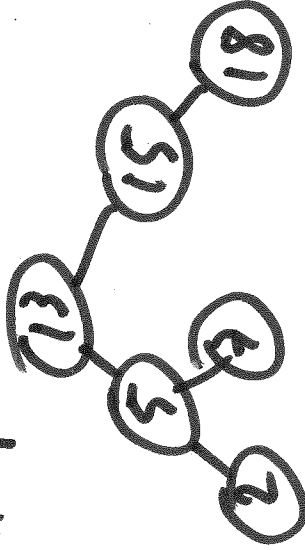
remove(10)

2, 5, 7, 10, 13, 15, 18
 ↳ successo

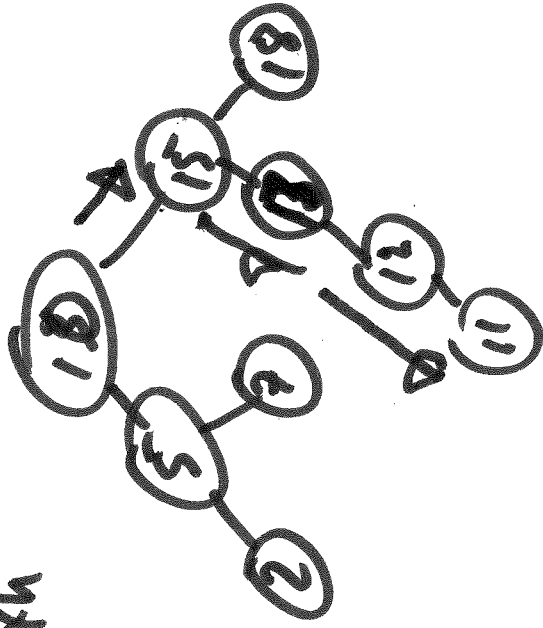
— Find the successor, that is the node that immediately follows after the node removed.

The successor of 10 is 13. In this tree, 13 is the successor of 10.

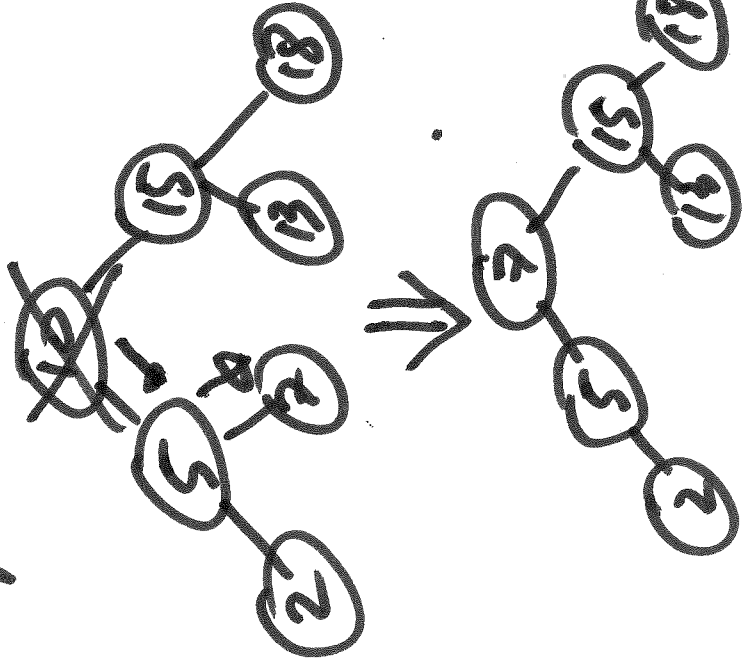
We move 13 in place of 10



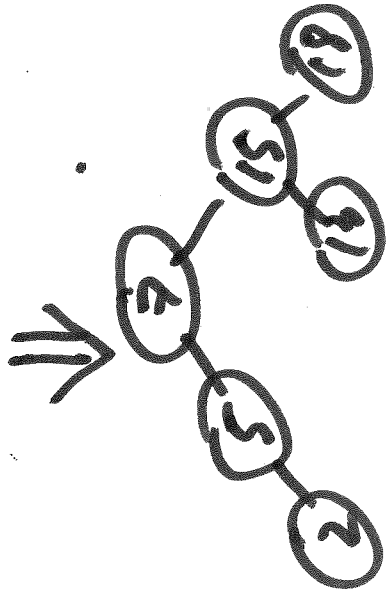
To obtain successor go to the right child of the node to remove and follow the left path



Also the predecessor could work

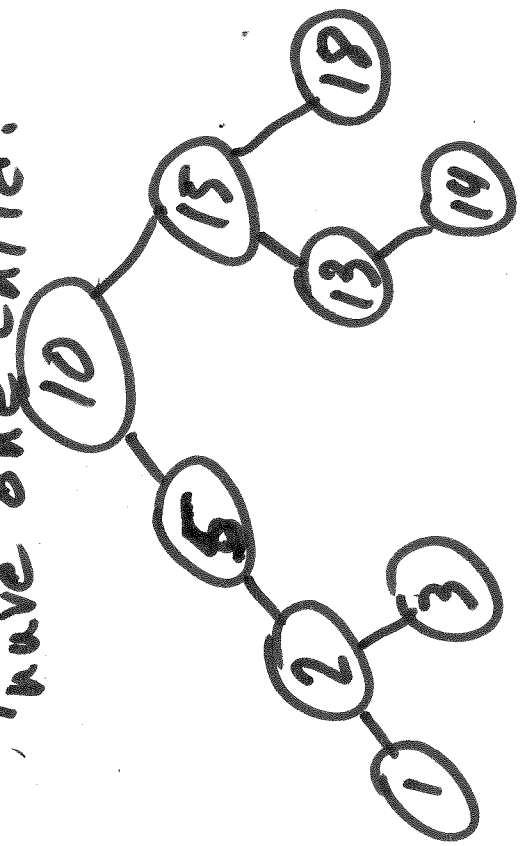


predecessor
8, 5, 7, 10, 12, 15, 18



(17)

Irregular case,
→ predecessor child's successor
have one child.

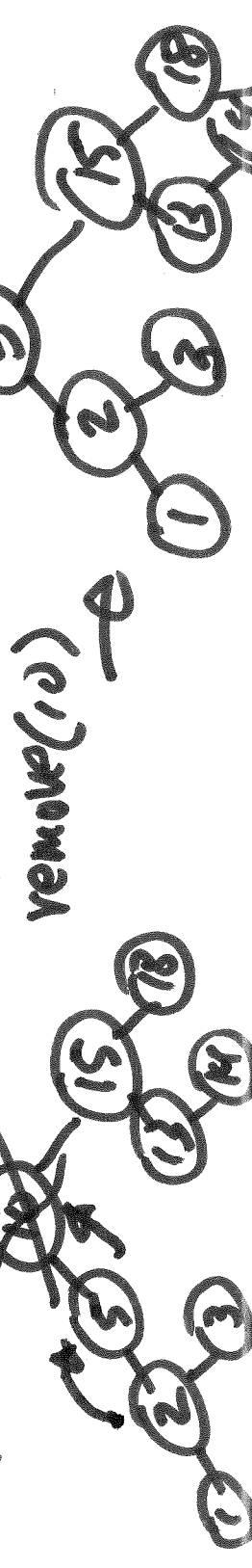


remove (10)

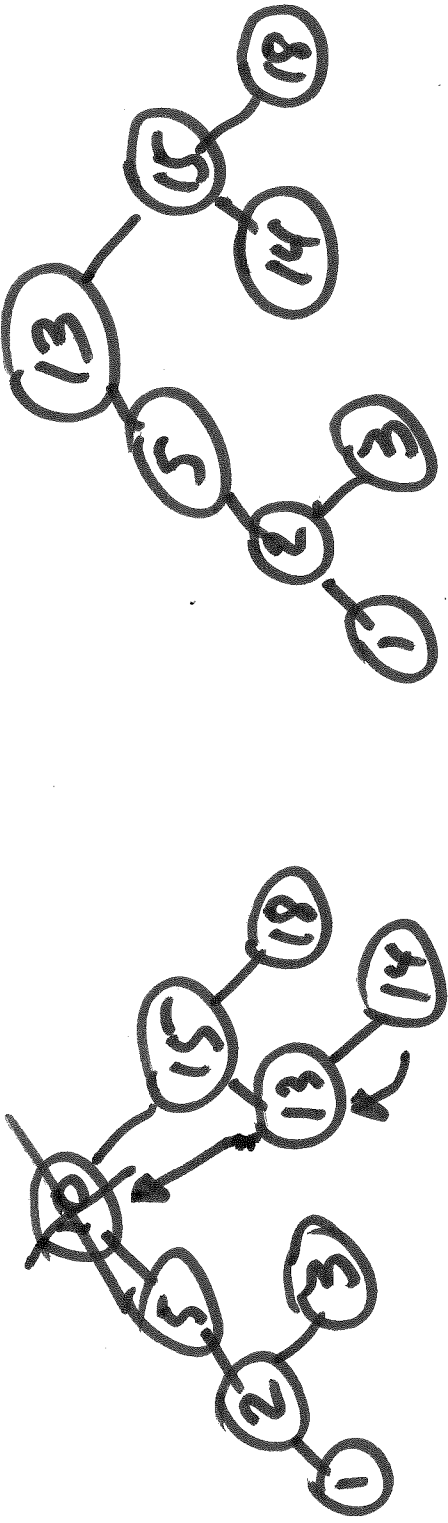
We still remove 10 and put the predecessor

5 or the successor 13 instead of 10
but ~~we~~ we promote the child
(should be one) instead of the
predecessor or successor to take place.

using predecessor



Using successor



The cost of removal is $O(h)$ where h is the maximum height of the tree.

If the tree is balanced then $h \approx O(\log n)$ and therefore removal is $O(\log n)$.

There is no guarantee that after a removal or insert the tree will be balanced.

AVL trees

(119)

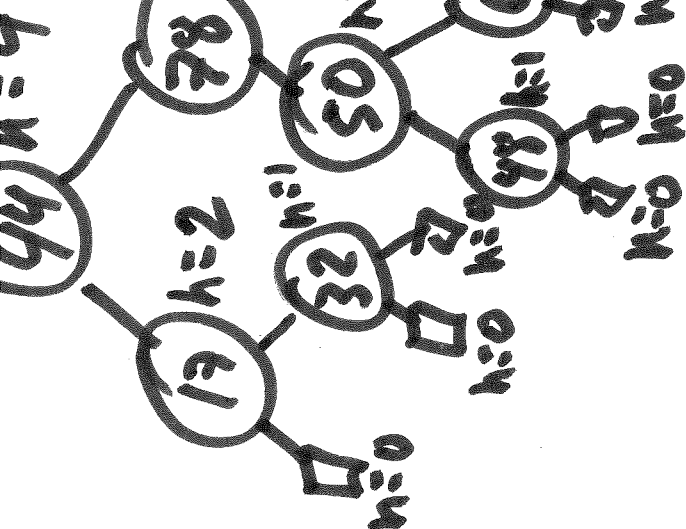
- They are binary trees that are balanced. That is, the height of the tree is always $O(\log n)$.
- To keep the tree balanced, we rebalance (readjust) the tree after every insert or remove operation.

AVL Tree Property

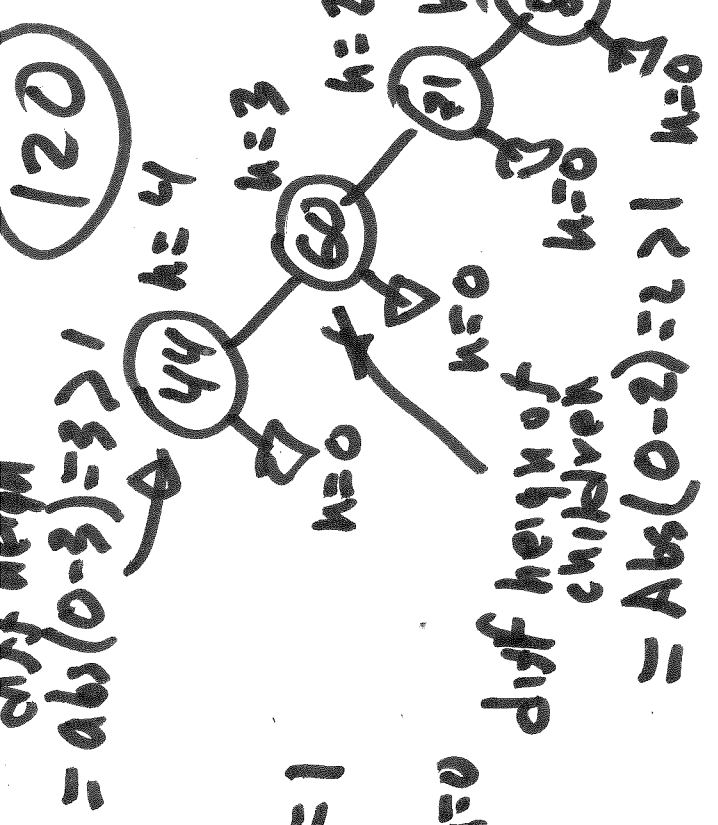
For every internal node v of T , the heights of the children can differ at most by 1.

~~height(v) =~~

$$\text{height}(v) = \begin{cases} 0 & v \text{ is a leaf} \\ 1 + \max(\text{height}(\text{left}), \text{height}(\text{right})) & \end{cases}$$



AVL tree



Not an AVL tree

$\text{diff height of children} = \text{Abs}(0-2) = 2 > 1$
 $\text{diff height of children} = \text{Abs}(0-2) = 2 > 1$

$\text{diff height of children} = \text{Abs}(0-3) = 3 > 1$
 (120)

Claim

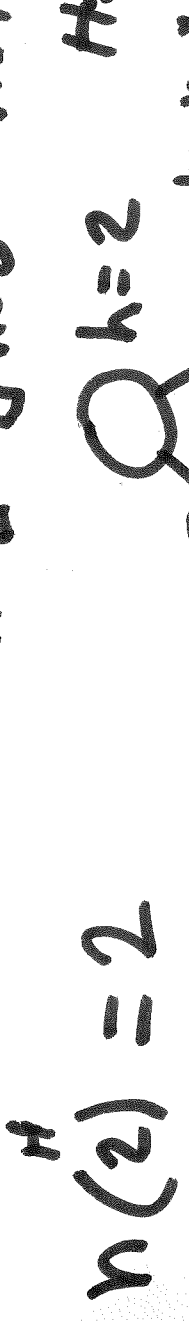
121

If we follow the AVL property then the tree's max height will be $O(\log n)$. i.e. the AVL Tree will be balanced $n = \#$ of nodes.

Proof

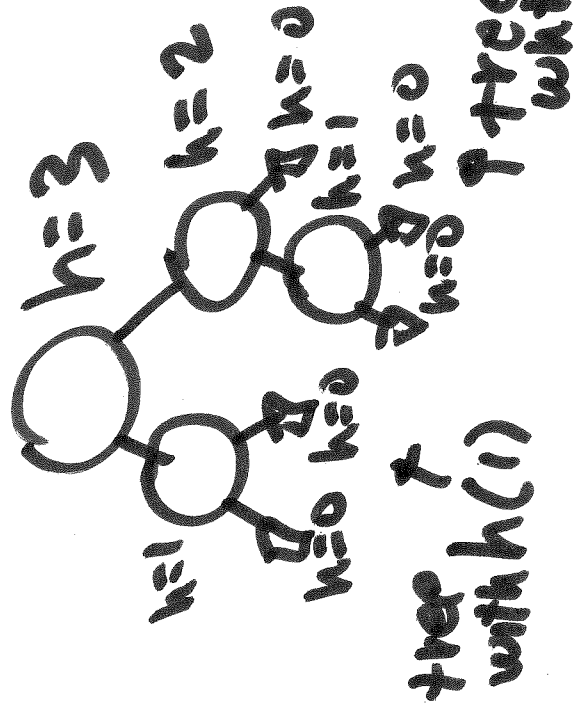
Find the minimum number of nodes of an AVL tree of height H . ($H=h$)

$n(H) =$ minimum number of internal nodes of an AVL tree of height H .



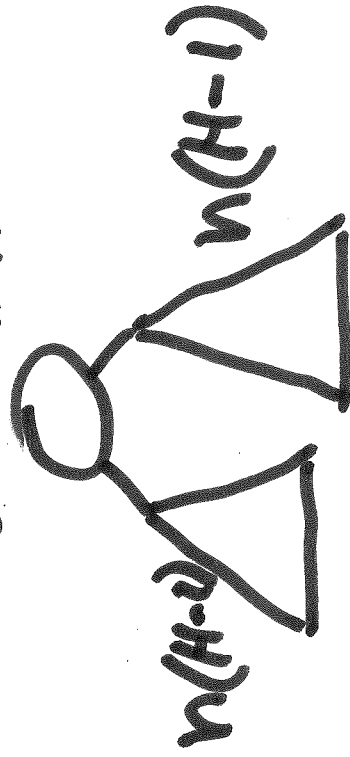
Why not minimum? Not minimum.

For n is the AVL tree
 with height H and minimum
 number of nodes n contains
 a subtree of height $H-1$ and
 a subtree of height $H-2$



In general

$$n(H)$$



(121) $(1-H)u + (2-H)u + 1 = (H)u$
 $n(H) > n(H-2) + n(H-1)$

since $n(H-1) > n(H-2)$

$$n(H) > 2n(H-2)$$

$$n(H) > 4n(H-4)$$

$$n(H) > 8n(H-6) \dots 2^i n(H-2i)$$

$$n(H) > 2^i n(1)$$

$$H - 2i = 1$$

$$H = 1 + 2i$$

$$i = \frac{H-1}{2}$$

$$n > 2^{\frac{H-1}{2}} (1)$$

$$\log_2 n > \log_2 2^{\frac{H-1}{2}} \rightarrow$$

$$\log_2 n > \frac{H-1}{2}$$

$$\boxed{\frac{n \log_2 n}{H} > 1} \rightarrow \boxed{\frac{n \log_2 n}{H} > 1} \rightarrow \boxed{2 \log_2 n + 1 > H}$$

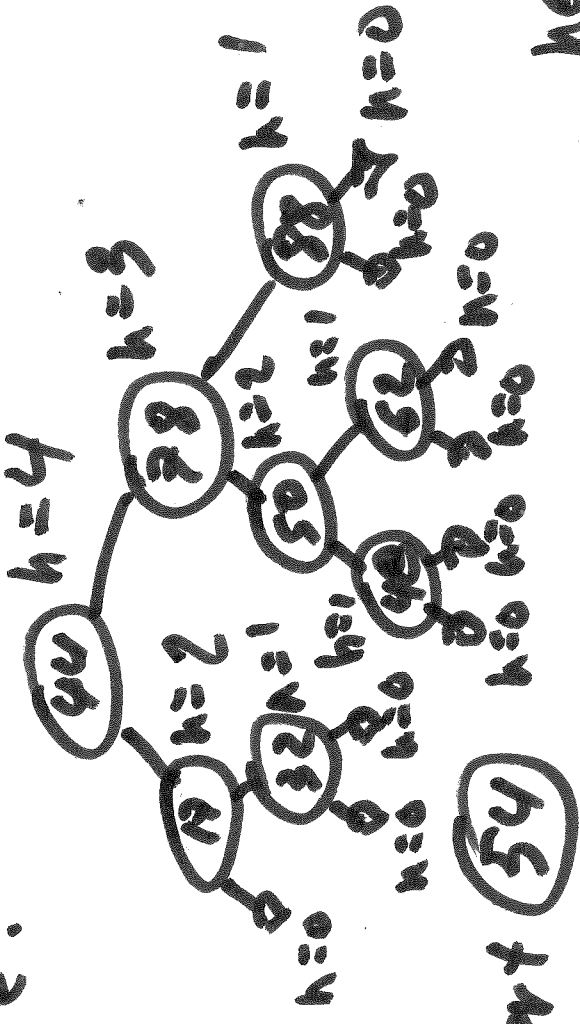
$$\boxed{\frac{n \log_2 n}{H} > 1} \rightarrow \boxed{H < 2 \log_2 n + 1} \rightarrow \boxed{H = O(\log n)}$$

(124)
If we enforce the property of the AVL tree (difference in height of ~~the~~ children of a node ≤ 1) then the height is $O(\log n)$ and therefore insert, remove, find are also $O(\log n)$.

→ At every insert, remove we rebalance tree to make sure property is followed. rebalancing also has to be at most $O(\log n)$ to keep operations $O(\log n)$.

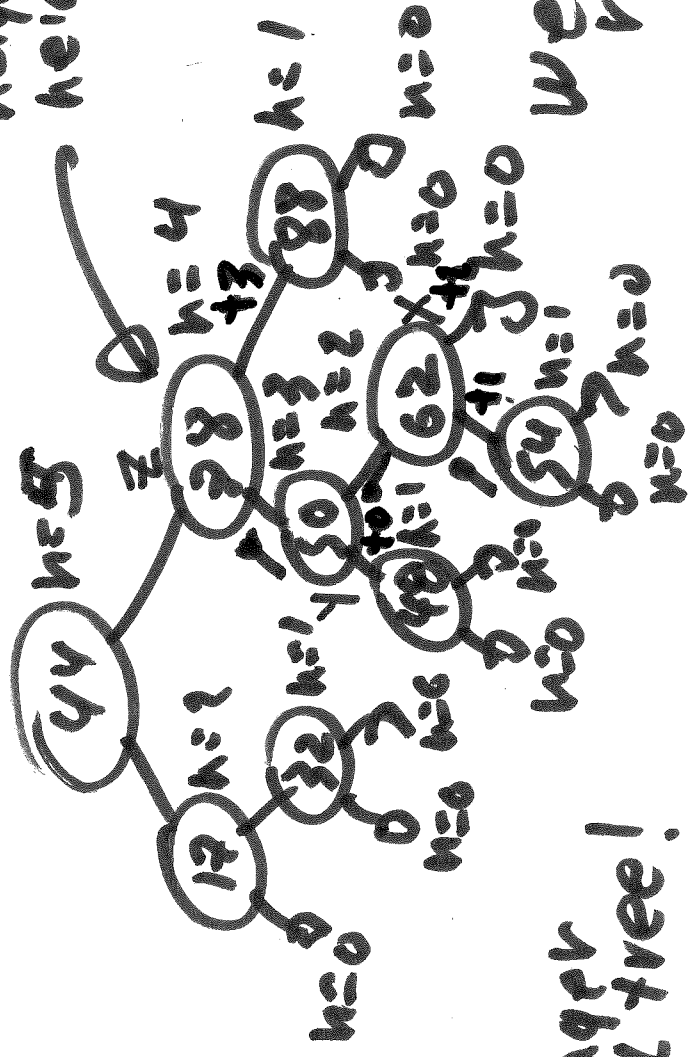
Insertion in an AVL tree.

Insert item like we usually insert in a binary tree:



Insert 54

height(left) = 3
height(right) = 1
| 3 - 1 | = 2

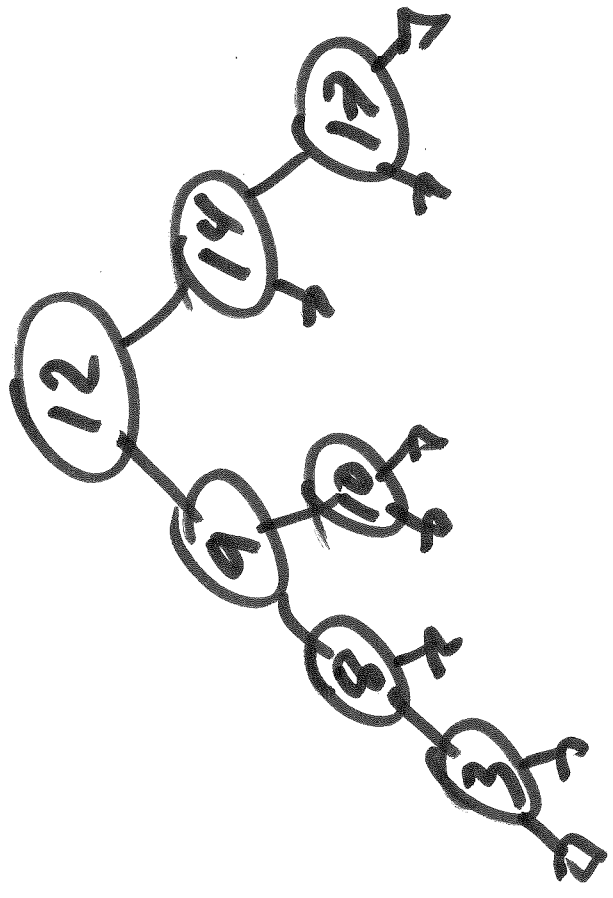


→ No longer an AVL tree!

We need to rebalance.

926

Quiz 13 AVL
What is the height
of the root node in
this ~~tree~~ tree (



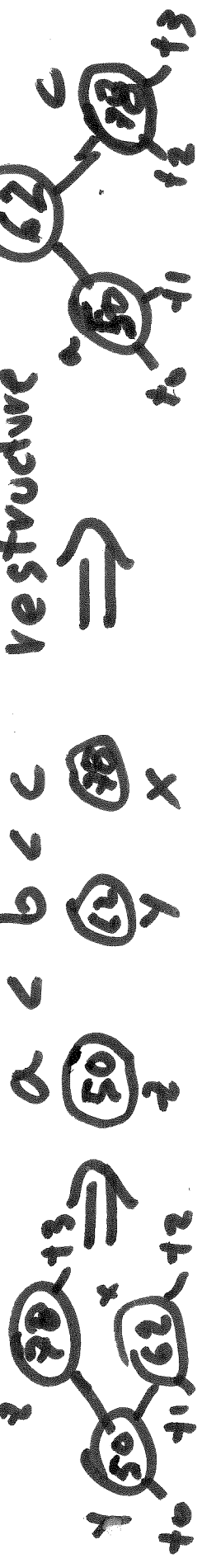
- a) 5
 - b) 4
 - c) 3
 - d) 2
- ~~tree~~

(122)

→ To rebalance tree, from inserted node go up until we find an ancestor that is root of the unbalanced tree, then we do the following labeling

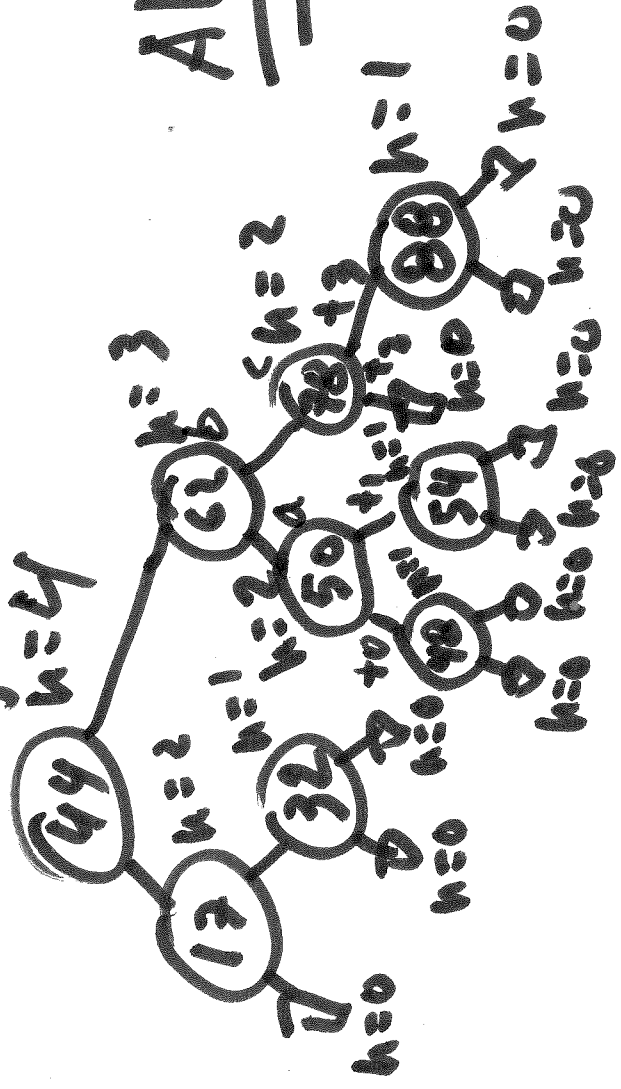
- z → Node that is the root of an unbalanced subtree. i.e. the difference in the height of the left subtree and right subtree is larger than 1
- y → the child of z with largest height.
- x → the child of y with largest height.

→ Rename z, y, x as a, b, c based on the order of the nodes in an in-order traversal, that is $a < b < c$



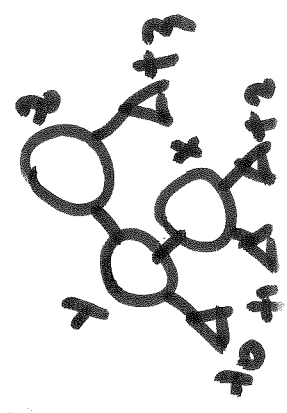
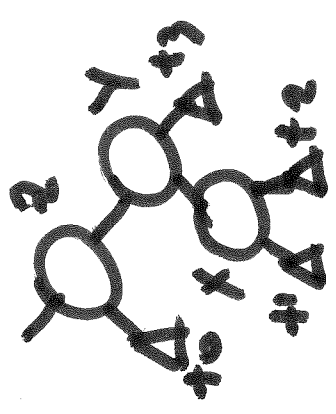
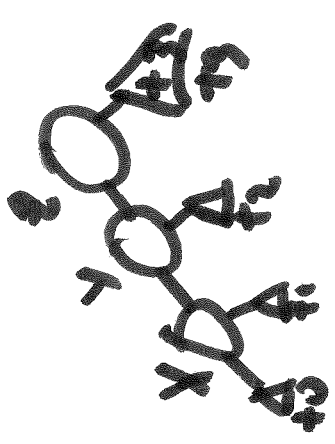
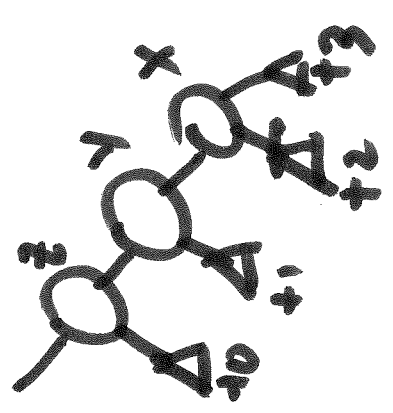
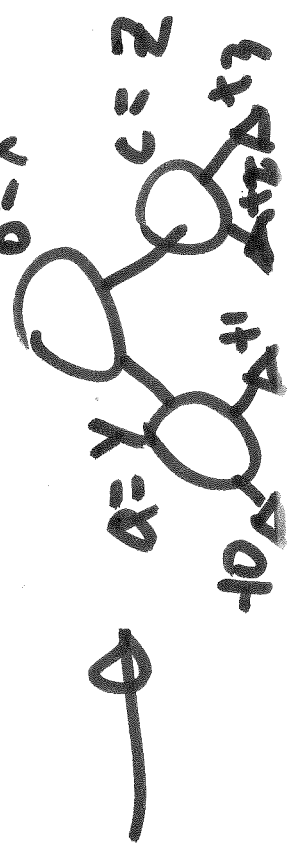
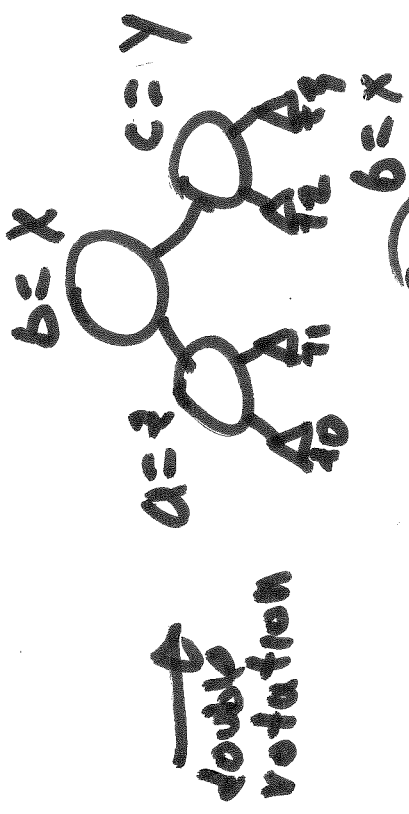
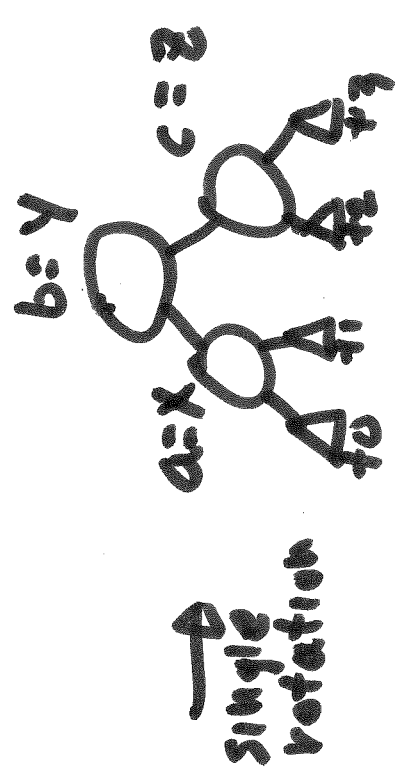
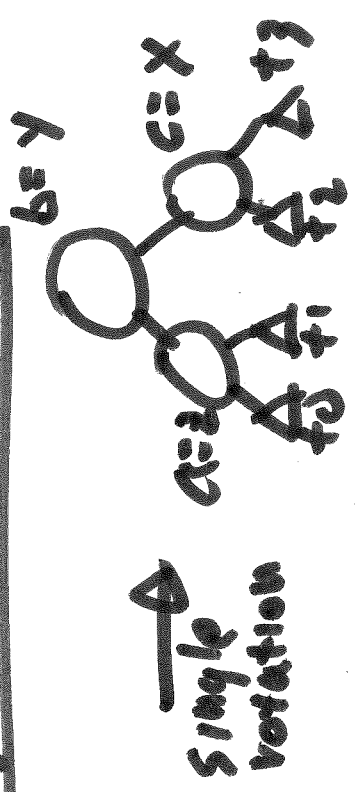
128

Rearranging a, b, c in the tree we get:



AVL tree!

Restructuring has 4 cases $a < b < c$

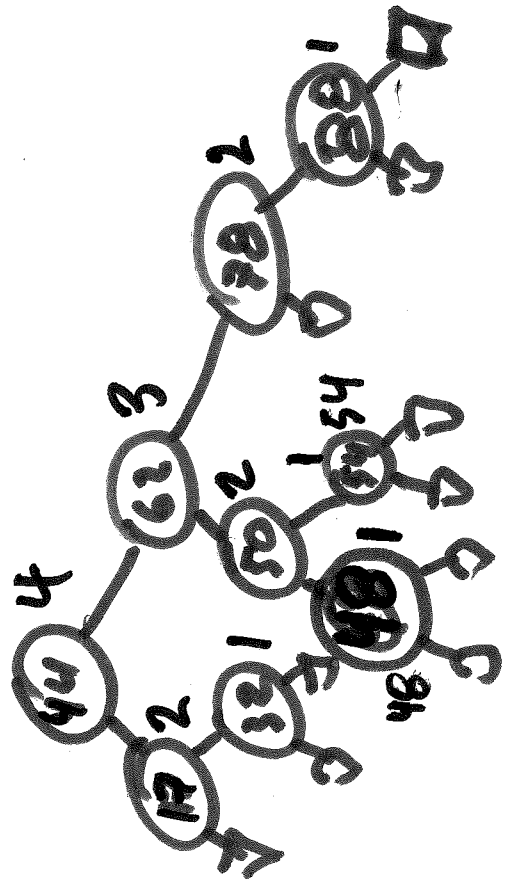


Removal

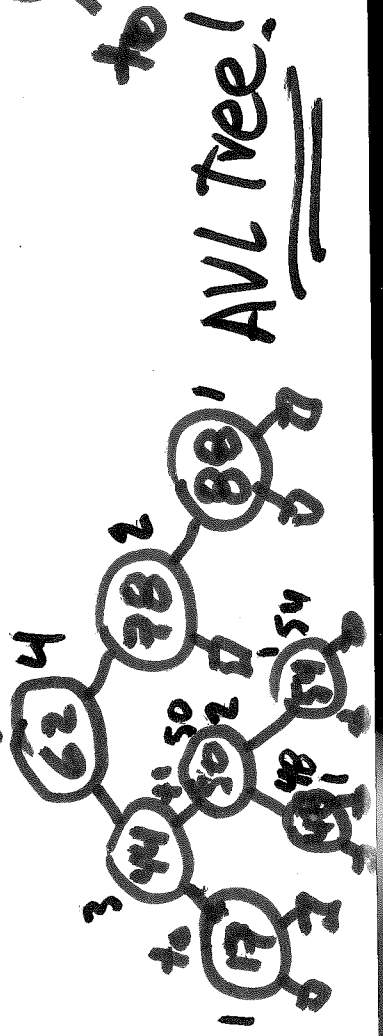
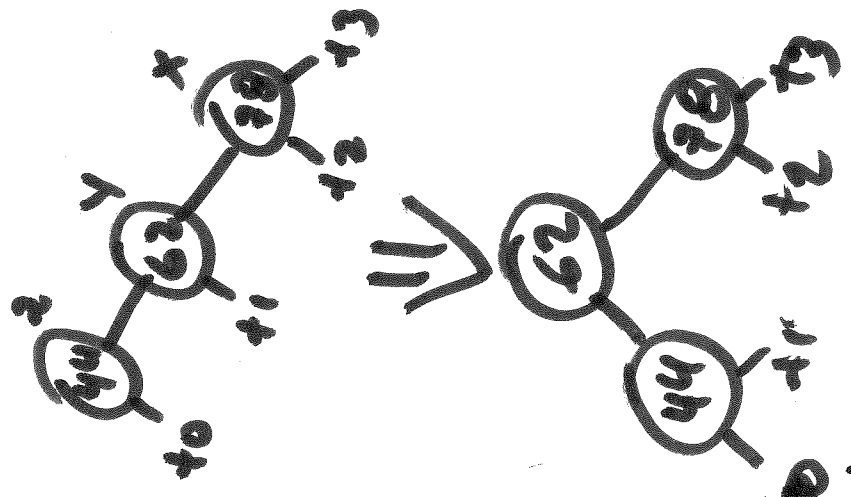
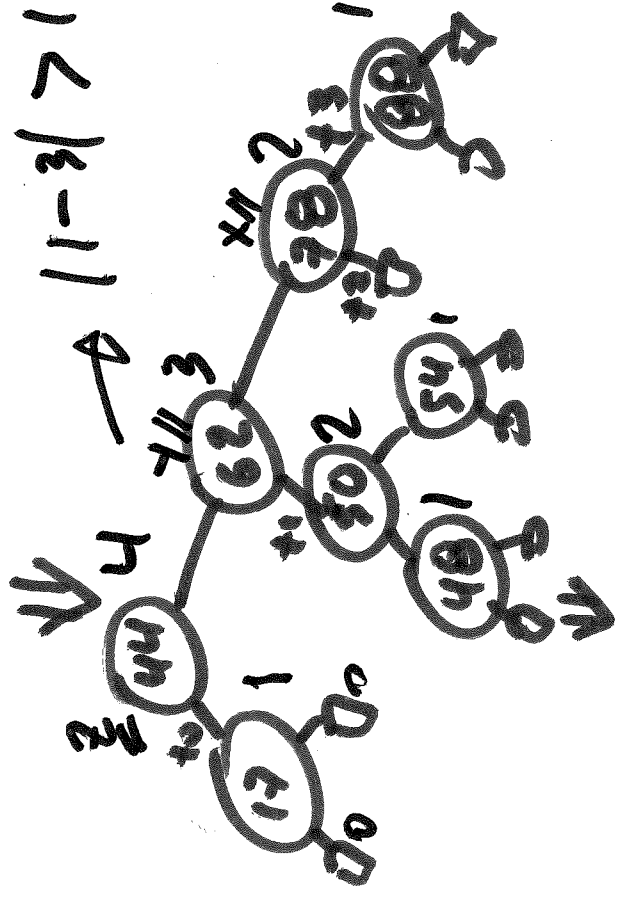
(130)

- Removing a node can cause a tree to become unbalanced
- Assuming that the node removed is w , we find the first unbalanced subtree by going up in the tree.
- Similar to insertion, we find Z, Y, X
 - Z - Parent of unbalanced subtree
 - Y - child of Z with largest height
 - X - child of Y with largest height (if both children have same height, the use either one as X).
- Apply same steps of restructuring used in insertion.
- As this restructuring may upset the balance of another node higher in the tree, we continue trying to restructure until it is not needed or root is reached.

(131)



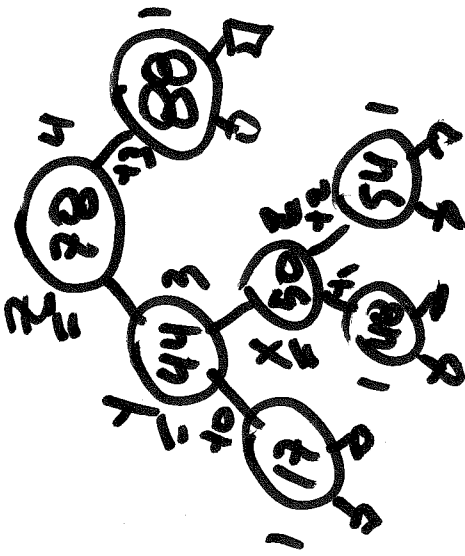
remove(32)



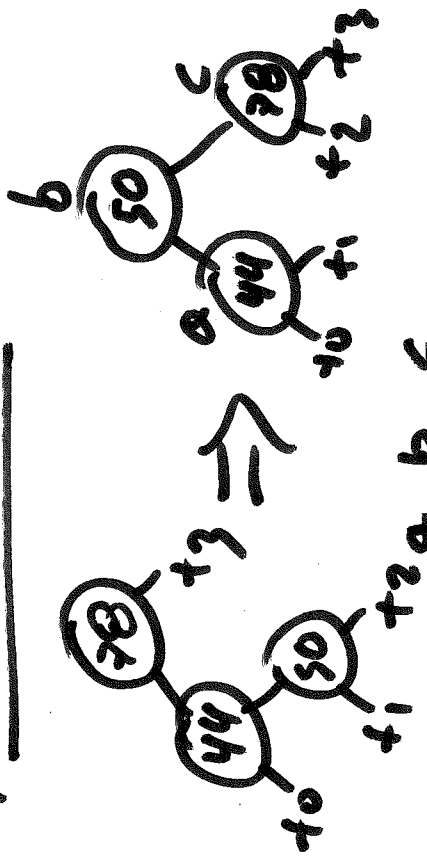
remove(62)

132

Remove (62) and it leaves a hole.
 we put in its place either the predecessor (54)
 or the successor (78). We will choose (78)

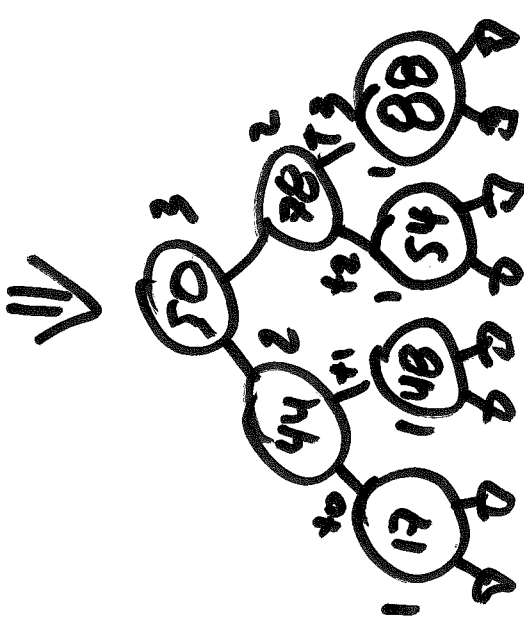


Not balanced



88, 50, 78
 44, 50, 78

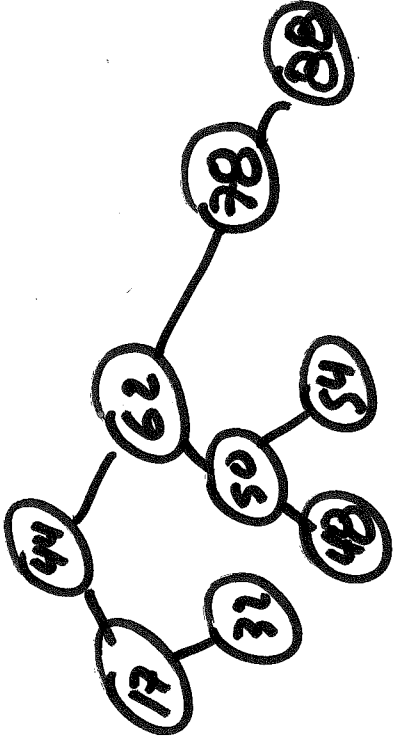
balanced!



Quiz 14

133

Starting at tree



If we remove (17) what are the values of z, y, x

- a) $z=62, y=50, x=54$
- b) $z=62, y=78, x=88$
- c) $z=44, y=62, x=78$
- d) $z=44, y=62, x=54$

AVL Tree Implementation

AVLTree.h

~~int~~ typedes int KeyType; ~~int~~ typedes const char * DataType;

```
struct AVLNode {
    int height;
```

```
    KeyType key key;
```

```
    DataType data data;
```

```
    AVLNode *left;
```

```
    AVLNode *right;
```

```
    AVLNode *parent;
```

// Type of key is int.

// The ~~key~~ can be other types like
// const char * etc.

// Also data can be of other type
// different than int.

```
}
class AVLTree {
```

```
    AVLNode *root;
```

```
public:
    AVLTree();
```

```
    bool insert (KeyType key, DataType data);
```

```
    bool remove (KeyType key);
```

```
    bool find (KeyType key, DataType &data);
```

```
}
```

AVL.cpp

#include "AVL.h"

AVLTree::AVLTree() {

 root = NULL;

}

bool

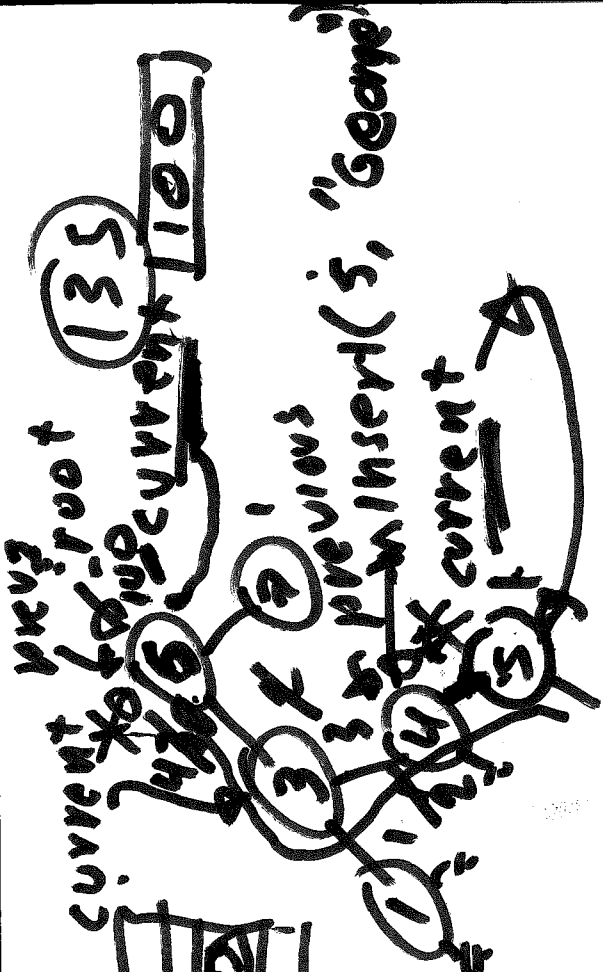
AVLTree::insert(keyType key, DataType data) {

 // Find node to insert into

 AVLNode *current = root;

 AVLNode *previous = NULL;

 while (current != NULL) {
 previous = current;
 if (key < current->key) {
 current = current->left;
 }
 else if (key > current->key) {
 current = current->right;
 }
 else { // key found
 return true;
 }
 }



100



prev
current
root
insert(5, "George")
get to the bottom
or find key.

136

```

// change data
current->data = data;
return false; // true if key found.
} // else key found

```

```

} // while
// key was not found. Create new node.
ALNode *n = new AVINode();

```

```

n->key = key;
n->data = data;
n->left = NULL;
n->right = NULL;
n->parent = NULL;
n->height = 1;

```

```

// Test if tree is empty
if (previous == NULL) { // -root is NULL
  // insert n as root node
  -root = n;
  return false;
}

```

```
// insert at the left or right of previous.
if (key < previous -> key) {
  // insert at the left
  previous -> left = n;
}
```

```
else {
  // insert at the right
  previous -> right = n;
}
```

```
n -> parent = previous;
// Adjust height of nodes in path
// from inserted node to root.
```

```
AVLNode * m = n -> parent;
while (m != NULL) { // iterate until we reach root.
```

```
  int maxHeight = 0;
  if (m -> left != NULL) { // check's left child
    maxHeight = m -> left -> height;
  }
  if (m -> right != NULL) { // check's right child
    maxHeight = m -> right -> height;
  }
}
```

O(log n)

138

```
// compute height of m
m->height = max(height + 1,
m = m->parent; // go to parent
} // while
restructure(n); // restructure tree starting
return false; // at inserted node n.
}
```

Quiz 15

1851

~~What~~

In the following code that finds

the node `current` if `previous = null`
`while (current != null) {`

`if (key < current->key)`

`current = current->left;`

`}`

`else if (key > current->key)`
`current = current->right;`

`}`

`else {`
`current->data = data;`

`}`

`// what is the missing line?`

a) `current = previous;`

b) `current = previous->next;`

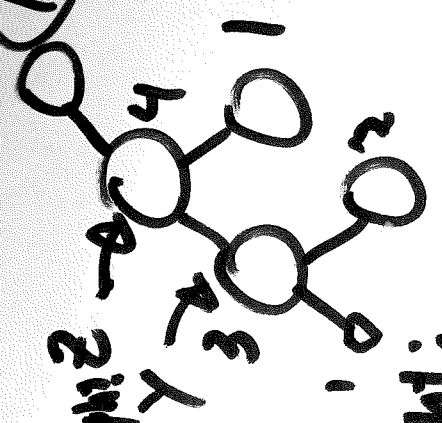
c) `previous = current;` `current->left;`

Q142

// z is unbalanced

// Find y. child z with max height
 AVLNode *y = NULL;
 // find y

int maxh = 0;
 if (z->left != NULL) {
 y = z->left;
 maxh = z->left->height;
 }
 if (z->right != NULL) {
 y = z->right;
 maxh = z->right->height;
 }
 return y;



if (z->right != NULL && maxh < z->right->height) {
 y = z->right;
}

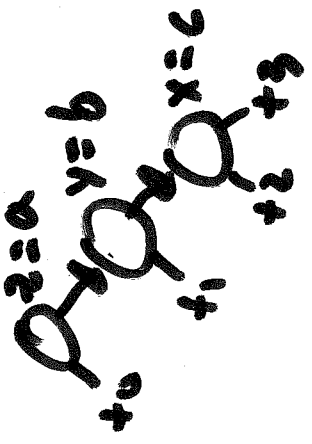
assert(y != NULL); // should never happen
 // but just in case
 // we have a bug.

// Find x
 AVLNode *x = NULL;

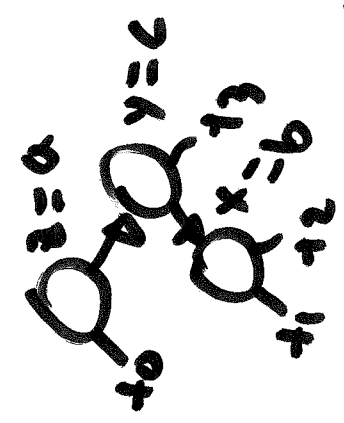
maxh = 0;
 if (y->left != NULL) {
 x = y->left;
 maxh = x->height;
 }
 if (y->right != NULL && y->right->height > maxh) {
 x = y->right;
 }
 assert(x != NULL);

4 cases:

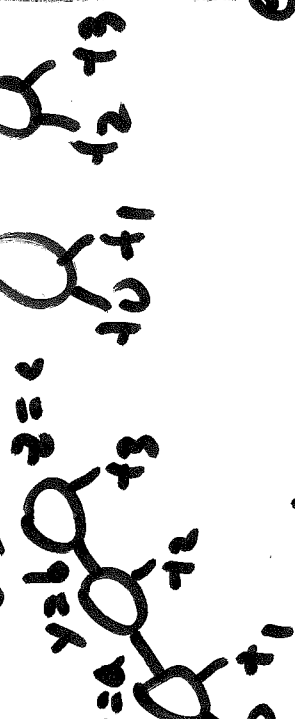
case 1:



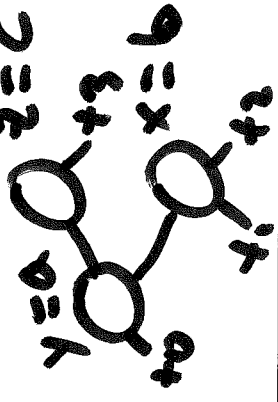
Case 2



Case 3



Case 4



// Identify a, b, c

AVLNode *a, *b, *c, *t0, *t1, *t2, *t3;

if (z->right == y) {

// case 1 or 2

if (y->right == x) {

// case 1

a = z; b = y; c = x;

t0 = z->left; t1 = y->left;

t2 = x->left; t3 = x->right;

else {

// case 2

a = z; c = y; b = x;

t0 = z->left; t1 = x->left;

t2 = x->right; t3 = y->right;

else {

// case 3 or 4

if (y->left == x) {

// case 3

a = x; b = y; c = z;

t0 = x->left; t1 = x->right;

(141)

```

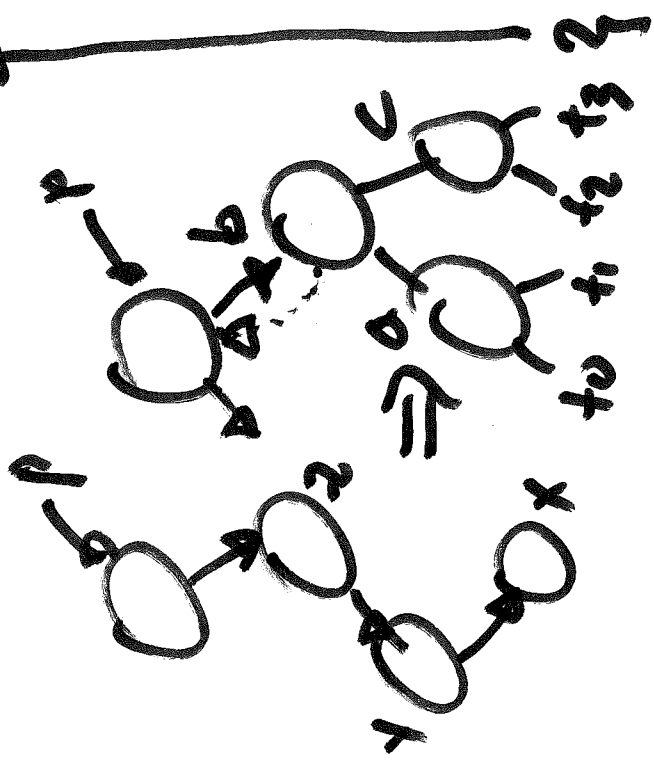
1 t2 = y -> right; t3 = z -> right;
2 // case 3
3 else {

```

```

4 // case 4
5 a = y; b = x; c = z;
6 t0 = y -> left; t1 = x -> left;
7 t2 = x -> right; t3 = z -> right;
8 }

```



```

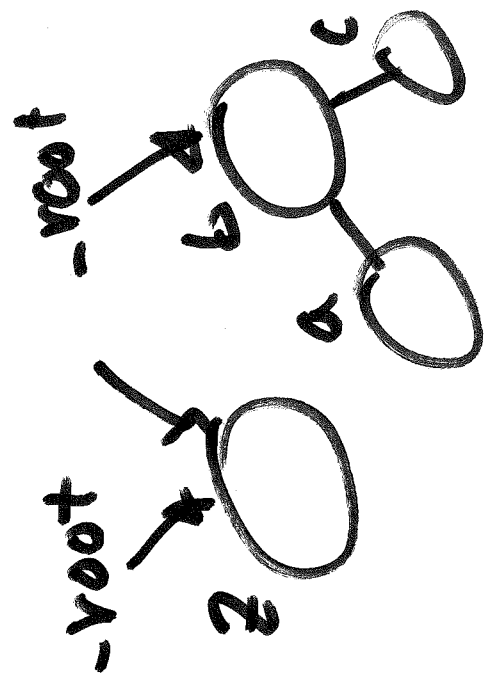
// Do rotation
AVLNode *p = z -> parent;

```

```

1 if (p != NULL) {
2     if (p -> left == z) {
3         attach b to left of p
4         p -> left = b;
5     }
6     else {
7         attach b to right of p
8         p -> right = b;
9     }
10    else { // z is root
11        root = b;
12    }
13 }

```



$b \rightarrow \text{-parent} = p;$
 $b \rightarrow \text{-left} = a;$
 $b \rightarrow \text{-right} = c;$
 $a \rightarrow \text{-parent} = b;$
 $a \rightarrow \text{-left} = t_0;$
 $a \rightarrow \text{-right} = t_1;$

$b \rightarrow \text{-parent} = b$
 $c \rightarrow \text{-left} = t_2$
 $c \rightarrow \text{-right} = t_3$

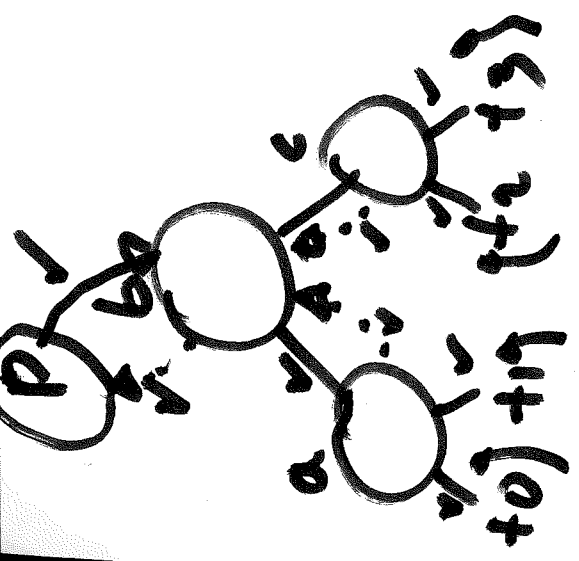
// connect parent of t_0, t_1, t_2, t_3

if ($t_0 != \text{null}$) {
 $t_0 \rightarrow \text{-parent} = a;$

}
 if ($t_1 != \text{null}$) {
 $t_1 \rightarrow \text{-parent} = a;$

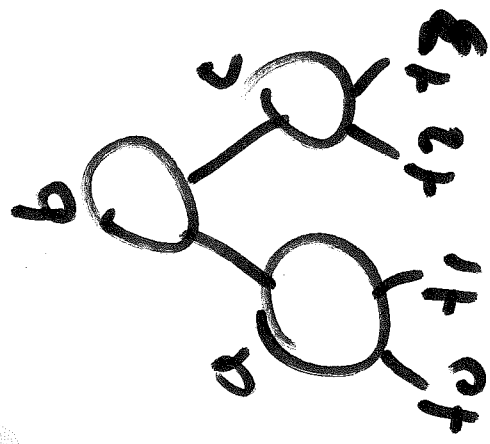
}
 if ($t_2 != \text{null}$) {
 $t_2 \rightarrow \text{-parent} = c;$

}
 if ($t_3 != \text{null}$) {
 $t_3 \rightarrow \text{-parent} = c;$



```
// Fix height of a
int maxHeight = 0
if (a->left != NULL) {
    maxHeight = a->left->height;
}
if (a->right != NULL) {
    maxHeight = a->right->height;
}
a->height = maxHeight + 1;
```

```
// Fix height of b
maxHeight = 0;
if (c->left != NULL) {
    maxHeight = c->left->height;
}
if (c->right != NULL) {
    maxHeight = c->right->height;
}
c->height = maxHeight + 1;
```



```

// Fix height of b
max height = 0
if (b -> left != null) {
  max height = 0 -> left -> height;
}
if (b -> right != null) {
  max height = 0 -> right -> height;
}
max height = b -> right -> height +
max height = max height + 1;
z = p; // go up.
} // while z != null.
} // restructure.

```

The AVL remove is similar to insert

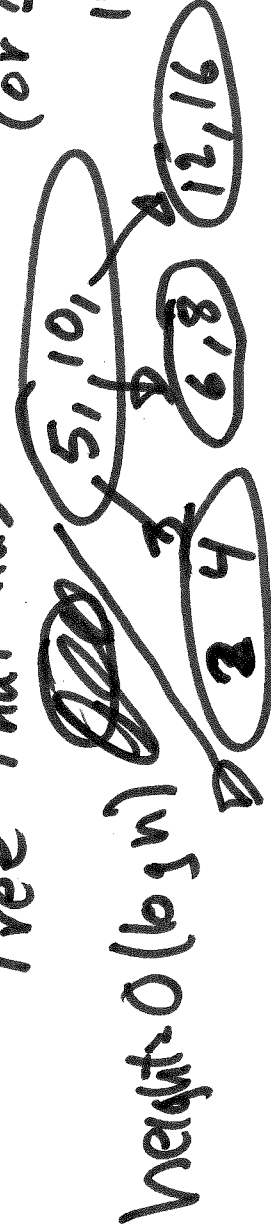
- find node
- remove node
- restructure

Other types of balanced trees

→ AVL Trees

→ 2-3 tree

Tree that has 2 or 3 children (or 1 or 2 keys in each node)



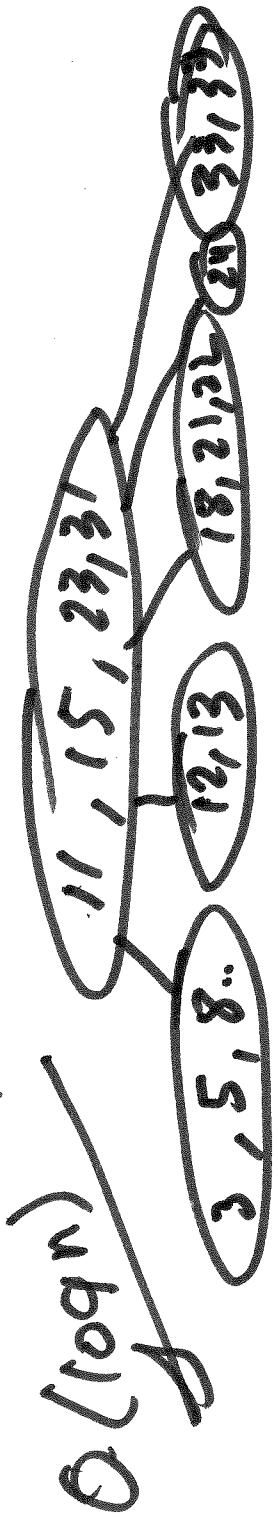
→ Red-black Tree. Each edge has a color red or black



b-trees

148

→ They may have many keys in each node



- The more keys are in each node, the height will be smaller.
- B-trees are important in databases.
- Database are stored in disks.
- In a disk the read and write operations are in entire blocks, where each block is 1KB or 4KB, 8KB, 16KB.
- Even if you only want to read 1 byte, a whole block (E.g 8KB) has to be read.
- A database access will take as many seconds as the number of block reads in disk.

In AVL tree shown, access to bottom of the tree will take 4 disk accesses.

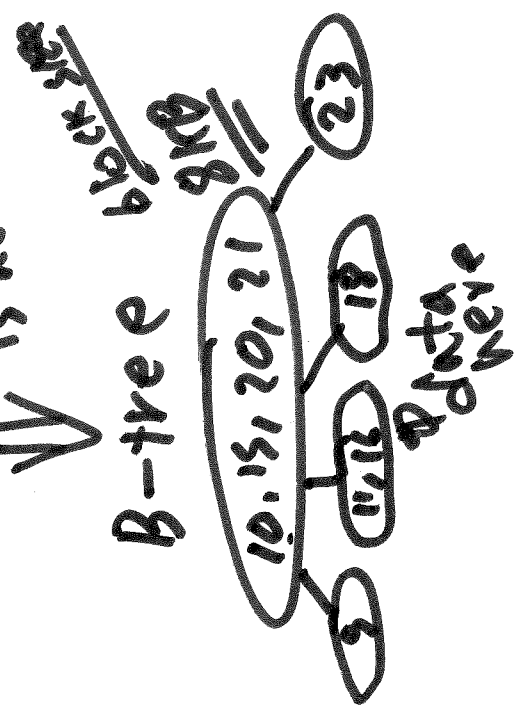
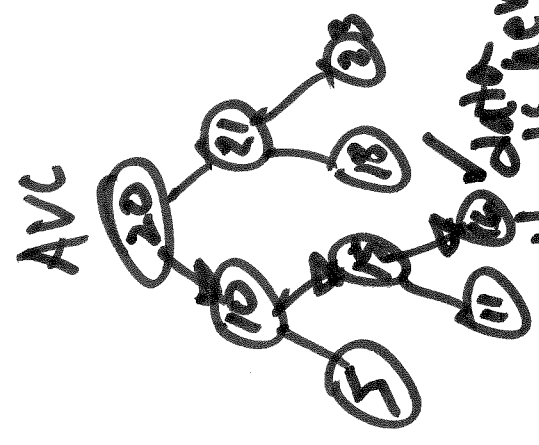
In B-tree shown, access to the bottom will take only 2 disk accesses

Assume each node is in a different block.

The maximum number of keys in a AVL node will be limited to the number of keys that you can fit in a block

B-trees are used to store indexes in data bases.

B-trees minimize the number of disk accesses



Merge Sort

- It is a sorting mechanism that uses divide and conquer

Divide and conquer

* ~~three~~ three steps

- Divide - divide data into two or more disjoint sets
- Recur - solve the problem in the smaller subsets recursively.
- Conquer - Take solutions of subproblems and merge them to create a solution of the original problem.

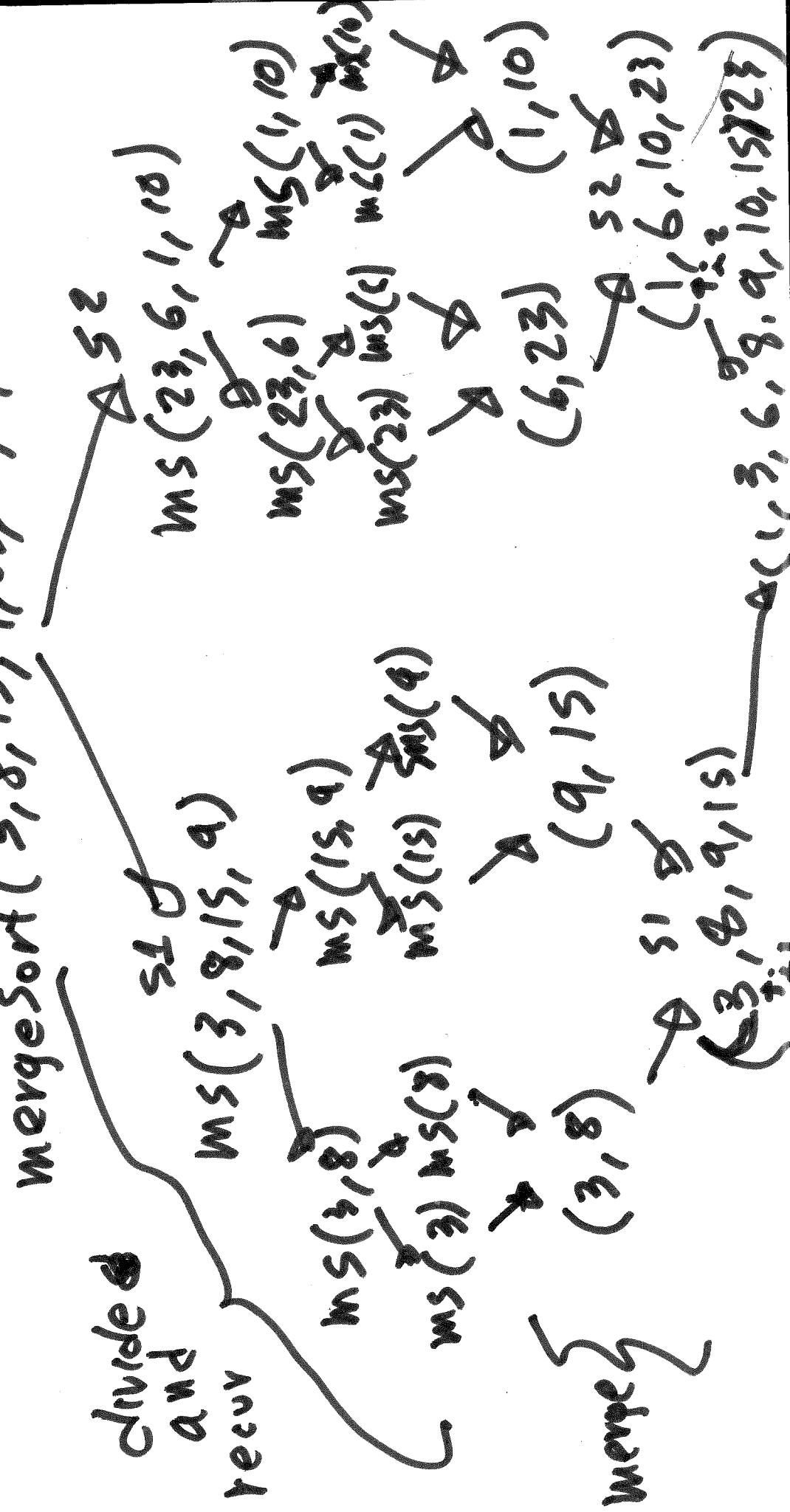
Merge sort uses divide and conquer

Algorithm mergesort(s)

Input: array s

Output: sorted array s

Example: mergesort(3, 8, 15, 9, 23, 6, 1, 10)



Quiz 17

751

What is wrong in the following code that computes the height of a node n

```
n->height = max(left, right) + 1 + max(n->left->height, n->right->height);
```

- a) The code is correct
- b) It should be just $\max(n->left->height, n->right->height) + 1$, no "+1" needed.
- c) It has to check if n is root
- d) left or right may be null.

Time complexity of Mergesort. (153)

We use "amortized time", that is we look at the work done for one element, and then we multiply the work by n .

- Every element (eg. 23) is copied $\log(n) + \log(n)$ (23 is copied 3+3 for $n=8$).

- You can see this from the drawing above since it looks like a tree and an inverted tree. The height of the tree is $\log(n)$.

- Work for one element is $O(\log n)$
Work for n elements is $nO(\log n)$
so the worst case takes $O(n \log n)$.

Implementation

451

```
void mergeSort(int *s, int n) {  
    // sorts array of type int  
    // with n elements  
    if (n <= 1) { return }  
    return; // End condition
```

```
}  
// Divide  
// create s1 and s2  
int n1 = n/2;  
int n2 = n - n1;  
int *s1 = new int [n1];  
int *s2 = new int [n2];  
// put elements into s1, s2  
for int i;  
for (i = 0; i < n1; i++) {  
    s1[i] = s[i];  
}
```

(12)

```

for (i = phi; i < n2; i++) {
    s2[i] = s[i+n1];
}

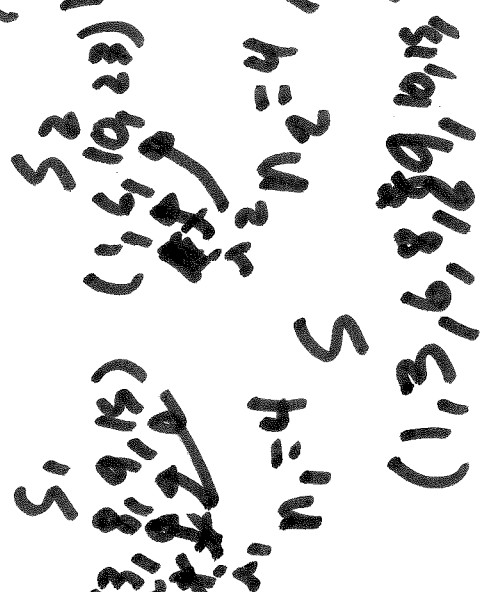
```

// call mergesort recursively

```

mergesort(s1, n1);
mergesort(s2, n2);

```



```

// merge them back into s
int i1 = phi; // index for s1
int i2 = phi; // index for s2
i = phi; // index for s

```

```

while (i1 < n1 && i2 < n2) {
    if (s1[i1] < s2[i2]) {
        s[i] = s1[i1]; // i, n1, s1 is smallest
        i1++;
    }
    else {
        s[i] = s2[i2]; i2++;
    }
}

```

```

↑ i++
} // while

```

```

// copy remaining elements
// copy remaining elements in s1, if any
for (; i1 < n1; i1++, i++) {
  s[i] = s1[i1];
}

```

```

} // copy remaining elements in s2, if any
for (; i2 < n2; i2++, i++) {
  s[i] = s2[i2];
}

```

```

} // s is sorted!!
// Delete temporal arrays s1, s2
delete [] s1;
delete [] s2;
} // End of merge sort.

```

Notes on mergesort

157

- Merge Sort is not an

"in-place" algorithm, that is

we cannot use the same array S to do the sorting. We need to allocate extra memory.

- However, we could reduce the number of times memory is allocated with "new" using a scratch array that is allocated at the beginning and deleted at the end.

Mergesort takes $O(n \log n)$

worst case

QuickSort

- Popular sorting algorithm that also uses divide-and-conquer.

Algorithm $quicksort(S)$

input: unsorted array S

output: sorted array S

1. Divide

- If S has only one element return

- Otherwise choose an element x in S called the "pivot".

- Divide S into three sets: L, E, G such that the elements in L are smaller than x , the elements in E are equal to x , and the elements in G are larger than x .

$x > y$ $y < x$

$x = y$ $y = x$

$x < y$ $y > x$

(159)
2. Recur
quicksort(L), E, quicksort(G)

3. Conquer
put elements in L, E, G back
into S.

In-place quicksort

- We use the same array S for storage of L, E, G at each step.
- We use subranges to represent each set.
- The last element is used as pivot X.
- To merge L and G we will scan L left to right ~~and G~~ and G right to left to look for element not in place and swap them.

- This is how we create L, E, G in place.

(160)

(85, 24, 63, 45, 17, 31, 96, 50)

2 move 2 to right

Move r to the left until we find element out of place

all elements < 50

all elements > 50

2 (85 is not in G) 31 (31 is not in G)

swap

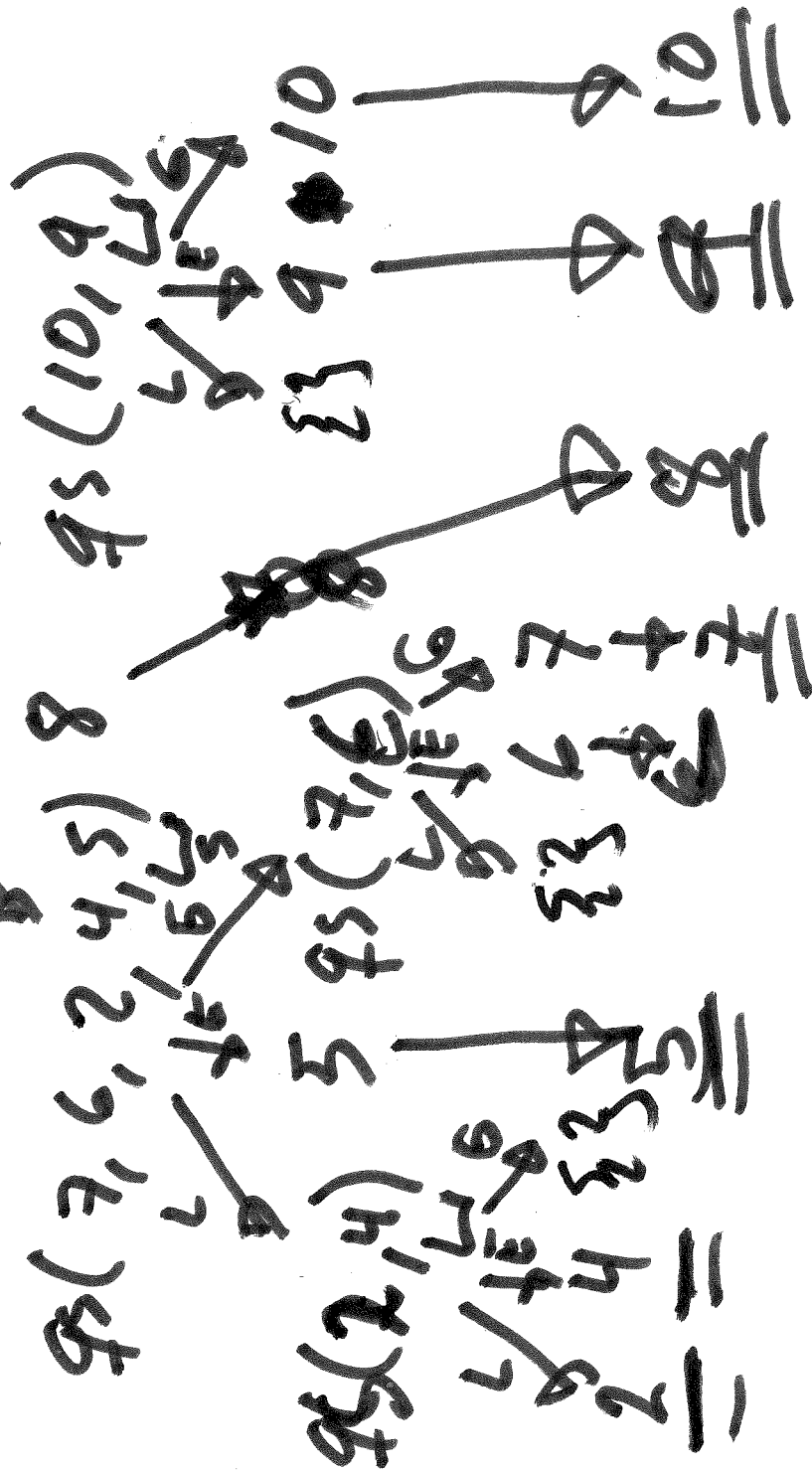
(31, 24, 63, 45, 17, 85, 96, 50)

(31, 24, 17, 45, 63, 85, 96, 50)

if not found instead of r

(31, 24, 17, 45, 50, 85, 96, 63)

Quicksort (7, 6, 2, 10, 4, 5, 9, 8)



In average if we choose a random pivot, L and R will have more less the same size. That means that every time we split S into L, E, R we will have half the elements

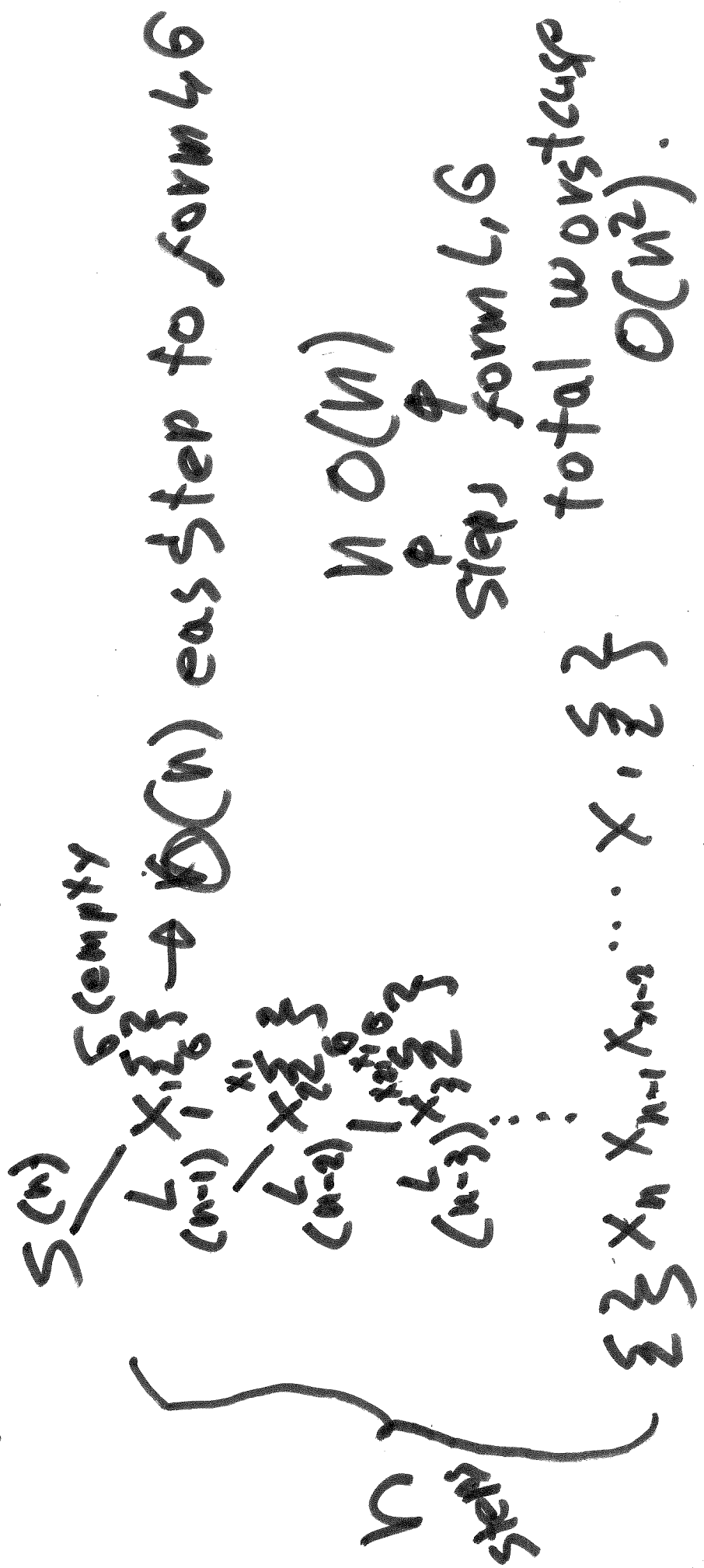
work for 1 element in average is $O(\log n)$.

$n \rightarrow n/2 \rightarrow n/4 \rightarrow n/8 \dots$

For n element in average total work is

$n O(\log n)$. average case $\rightarrow O(n \log n)$.

However in the worst case if the pivot is either the largest or the smallest then L or R will be empty.



In practice quicksort
works well. The $O(n^2)$

time usually will not happen.

(163)

Quiz 18

164

What algorithm has a worst-case time of $O(n^2)$

- a) heap sort
- b) merge sort
- c) quick sort

Implementation

165

In place quicksort. No need to allocate memory. All sorting is done in the array itself.

```
// Sort an array of int's using quicksort
void quicksort(int *s, int n) {
    quicksortSubrange(s, 0, n-1);
}
```

```
// sorts a subrange of integers using
```

```
// quicksort
```

```
void quickSortSubrange(int *s, int left, int right) {
```

```
    // base case. Return if s Subrange has no  
    // elements
```

```
    if (left >= right) {
```

```
        return;
```

```
    }
```

Improve next // Use the last element as pivot

```

int i = rand(left, right);
int tmp = s[right];
s[right] = s[i];
s[i] = tmp;
int x = s[right];
// Divide s into L, E, G
int l = left; // left index
int r = right - 1; // right index

```

```

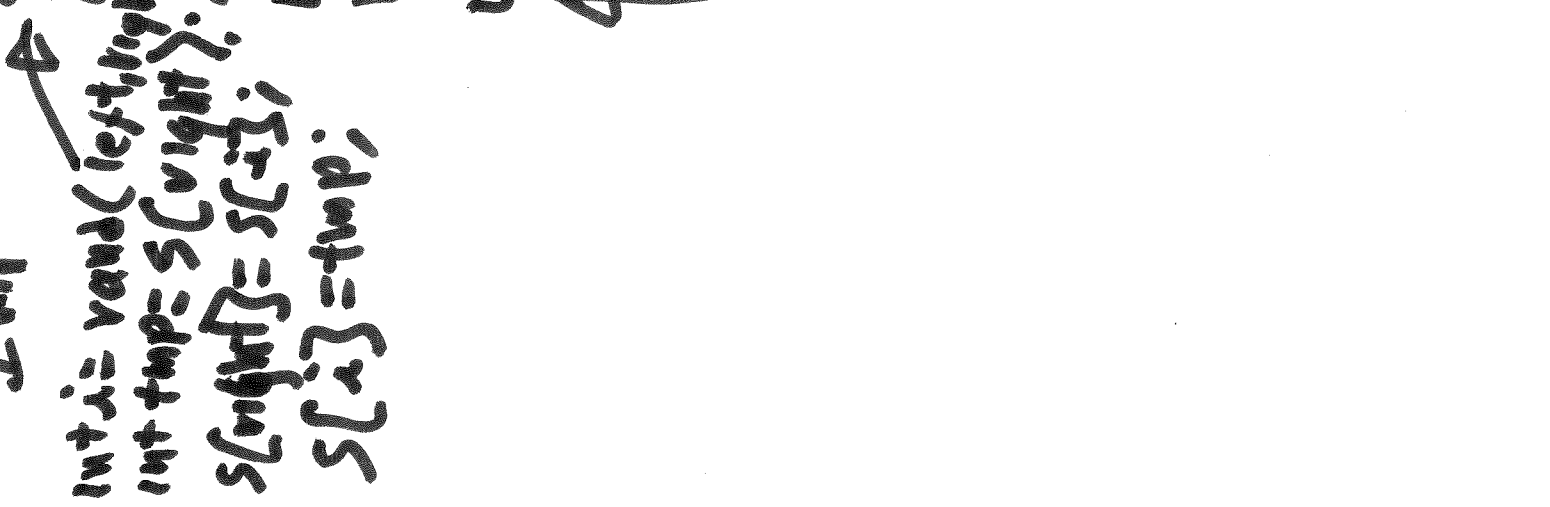
while (l < r) { // loop until l and r cross
  // Move l to the right until
  // we find a number that is larger
  // than pivot x since it is not in L.
  while (l < r && s[l] < x) {
    l++;
  }

```

```

  // Move r to the left until
  // we find a number that is smaller
  // than pivot x since it is not in G
  while (l < r && s[r] > x) {
    r--;
  }
}

```



168

```

// If still l < r and l and r
// are indexes of elements that
// should not be in l and r then
// swap them

```

```

if (l < r) {
  // swap s[l] and s[r]
  int tmp = s[l];
  s[l] = s[r];
  s[r] = tmp;
}

```

```

} // while

```

```

// swap now pivot with l

```

```

int tmp = s[l];
s[l] = s[right];
s[right] = tmp;

```

// L and R have been formed.

// Recur on L and R

quicksort Subrange(S , left, $l-1$);

quicksort Subrange(S , $l+1$, right);

} // End quicksort Subrange

In average, the pivot will be chosen in such a way that L and R will have approximately the same size. If this is the case and looking at the execution tree, each number will be copied at most $\log(n)$ times. Therefore the work done for n elements will be $n \cdot \log(n)$ or $O(n \log(n))$

average case

69

However, in the worst case if we choose the pivot x every time so L or R are empty then quicksort will degenerate and take $O(n^2)$. For example in the algorithm above we try to sort a sequence that is already sorted.

(5, 8, 10, 12, 16, 20, 22, 23)

5 8 10 12 16 20 22 23
Pivot 10
L: 5 8 R: 12 16 20 22 23

5 8 10 12 16 20 22 23
Pivot 10
L: 5 8 R: 12 16 20 22 23

This effect can be prevented by choosing a random pivot in the range and swapping it with the right pivot
pivot = rand(left, right);
swap [left, pivot] = swap [pivot, right];

5 8 10 12 16 20 22 23
Pivot 10
L: 5 8 R: 12 16 20 22 23

Quiz 19

If at every step of quicksort the pivot is chosen to be the minimum in the subarray, the algorithm will take

- a) $O(n \log n)$
- b) $O(n^2)$
- c) $O(n^3)$
- d) $O(1)$

Radix Sort

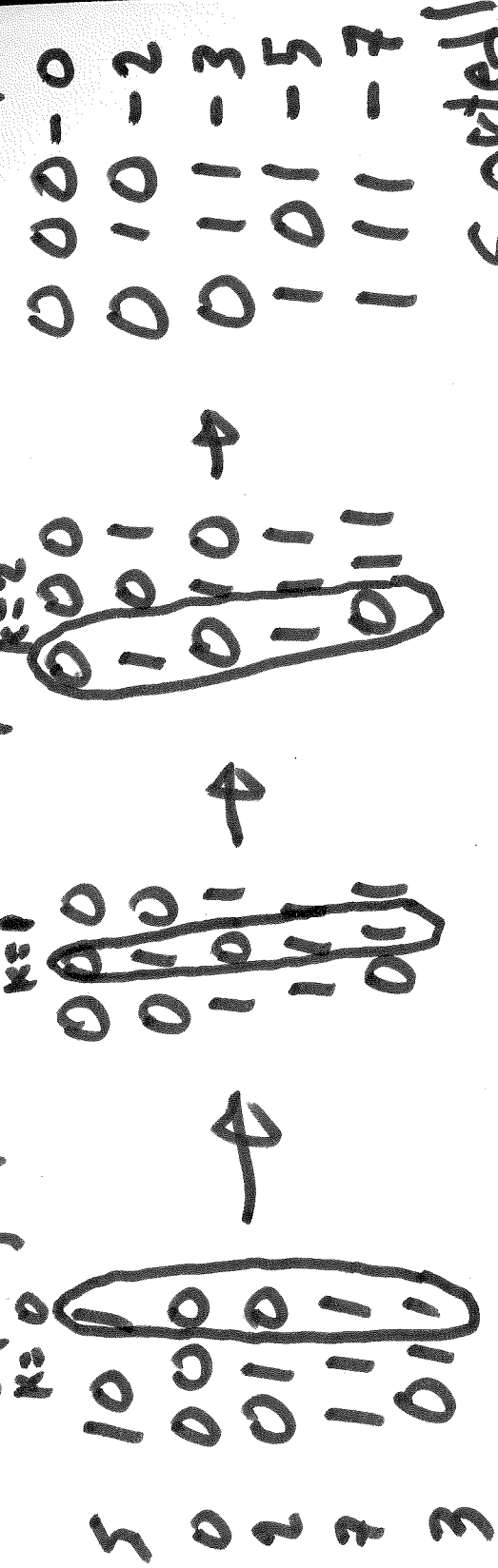
(171)

- Considers the structure of the keys.
- It assumes that the keys are represented in a number system.
- Algorithm Radix Sort

1. Assuming binary representation of the key, examine ~~from~~ bits from right to left

5-	101
0-	000
2-	011
7-	1011
3-	011

2. At every step k , sort the array by (172) looking only at bit k and keeping the order of the keys unchanged.



Complexity:

- b - # of bits used to represent a key
- n - Number of entries sorted

Radix Sort takes = b steps
 $* O(n)$ → work at every step
 $= O(bn)$

Implementation

173

```
// sorts an array of unsigned ints
// using radix sort
void radixSort(unsigned *array, int n) {
    // allocate scratch memory
    unsigned int *array2 = new int[n];
    assert(array2 != NULL);
    int nbits = sizeof(unsigned) * 8;
```

```
// for all bit positions
for (int i = 0; i < nbits; i++) {
    int j = 0; // index of array2
```

001 ← 0
010 ← 1
100 ← 2

```
// sort 0's first
for (int k = 0; k < n; k++) {
    // check if column k of array[k] is 0
    if ((array[k] & (1 << i)) == 0) {
        // array[k] has 0 at bit i
        array[j++] = array[k];
    }
}
```

i ↓

a[k] = 101
2 001 ← (1)
001 ≠ 0
a[k] 010
2 001 := 02 for bit

481

(i < j) & (i > j)

0 1 1 0
 1 0 1 0

 0 1 0 0

// Sort A's next

// check if bit i of array[k] is 1
 for (int k = 0; k < n; k++) {

if ((array[k] & (1 << i)) != 0) {

array[k] has 1 at bit i
~~array[j]~~
 array2[j] = array[k];
 j++;

// copy array 2 to array

memcpy (array, array2, n * sizeof(int))

// for i = 0 ...

delete [] array 2;

}

Bucket Sort

175

- It is a sorting algorithm useful to sort large sequence of numbers that have a small range of values. This means that the numbers will be repeated.

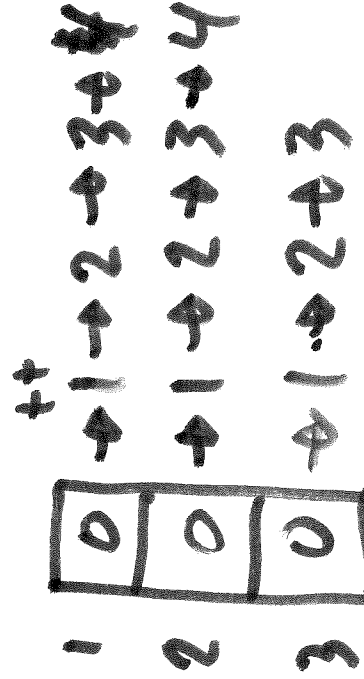
E.g. $n=1,000,000$ numbers

range of values is 1..10

2, 1, 3, 1, 2, 2, 3, 1, 1, 3, 2

$n=11$

$m=1..3=3$ (size of range)



4

\Rightarrow ~~00000~~

1, 1, 1, 1, 2, 2, 2, 2, 3, 3, 3

array of counters
(buckets)

Complexity of bucket sort

$O(n+m)$. But since $n \gg m$
then

$\approx O(n)$

Exercise:

write ~~an~~ implementation for

bucketSort (int *a, int n) {

// implement it. You may compute
min and max in
 $O(n)$ time.

}

String matching

Find the position of a substring in a text.

Brute Force

haystack Hello world
needle 0100

0 1 2 3 4 5 6 matches!
pos == 3

```

char * strstr(char * haystack, char * needle) {
  // Returns pointer to first occurrence
  // of 'needle' in 'haystack' or NULL if
  // it is not there.
  char * s = haystack;

```

```

while (*s != '\0') {

```

```

  char * h = needle;

```

```

  while (*h != '\0' && *h == *s) { // compare
    h++; s++; // needle
    h++; s++; // and s

```

```

}

```

```

if (*h == '\0') { // Found match.
  return s; // return position.
}

```

```

↓ { // while (s != 'a')
    // Match not found
    return NULL;
} // End strstr

```

Complexity: N - Length of string (haystack)
 M - Length of substring (needle)

In the worst case we do N steps
~~for~~ searching a substring that is M in length

$O(MN)$

Example of worst case.

haystack eeeeeeeeeeeeeeeee
 needle eeeeeeh

Forces to do all iterations

Quiz 20

179

Can radix sort be used for strings.

~~Q~~

a) Yes

b) No

c) Yes, only if strings ~~are~~ are made of '0' and '1' ascii characters.

Rabin-Karp String Matching

180

It uses a hash function that is applied to the m characters of the substring and the string matched hash₁

AAAAA AAAA A → N string
AAH → M substring.

compute

If $\text{hash}[i] \neq \text{hash}[i+1]$ → go to next char. Recompute hash₂ in $O(1)$

If $\text{hash}[i] = \text{hash}[i+1]$ then perform a bruteforce ~~match~~
comparison

An example of hash function in a string is adding the ASCII values

$$\text{hash}(s) = \sum_{i=0}^{n-1} s[i]$$

If we need to recompute hash₂ every single time, it will take $O(NM)$, that is the same as before.

However, if we compute hash2 using the previous value, like subtracting the ~~new~~ oldest character and adding

new one

$$\text{hash}(s, i, m) = \sum_{j=i}^{i+m-1} s[j] \rightarrow O(M)$$

takes

this can be rewritten as previous

$$\text{hash}(s, i, m) = \text{hash}(s, i-1, m) - s[i-1] + s[i+m-1]$$

remove oldest - add new.

takes $O(1)$ assuming that we have computed the hash $(s, i-1, m)$ already.

s-string
i-start position
m-length of substring

"Hello world"
i m l

Rabin Karp algorithm (haystack, needle)

- Compute hash 1 from needle
- compute hash 2 from first M chars in haystack.

while we have not reach end ($i=0..N-m$)

if hash 1 == hash 2

- brute force match of needle and haystack at i .
- If they match return i

end

Recompute hash 2. Subtract oldest character in hash 2 and add new one $i++$;

end

Return Null (No match found)

Do

In the typical case (average case),
~~that~~ hash 1 and hash 2 will be equal
where

only when the strings match. So

in average:

$O(N)$

(compare hash $O(1)$
recompute $O(1)$)

$* N$

Worst Case:

If we ~~are~~ have a false match where
the hash 1 and hash 2 are equal but
the string and substring do not match
at every position.

worst case:

$O(N^2)$

(compare hash $O(1)$
recompute hash $O(1)$
but force match $O(N)$)
 $* N$

Improving hash function

(184)

We can improve the hash function by computing it using a polynomial like:

$$S[i] \rightarrow \text{character at position } i$$

b - a constant $X(i)$

\rightarrow hash value starting at position i

$$\text{Eg. } b = 5 \quad X(i) = 5[S[i]]b^{M-1} + 5[S[i+1]]b^{M-2} + \dots + 5[S[i+M-1]]b^0$$

Note: If $b=1$ then this hash function is the same as adding the ASCII values as before.

Hello world

\downarrow hash =

$$X(i+1) = \underbrace{5[S[i+1]]b^{M-1} + 5[S[i+2]]b^{M-2} + \dots + 5[S[i+M]]b^0}_{5[S[i+M-1]]b}$$

$\text{ascii}(h) \neq 5$
 $\text{+ ascii}(e)$

$$X(i+1) = X(i)b - \underbrace{5[S[i]]b^M}_{\text{removes oldest char}} + \underbrace{5[S[i+M]]b^0}_{\text{add new character}}$$

Recompute $X(i+1)$ from $X(i)$ takes $O(1)$.

Implementation

181

```
char * haystack, * needle;  
int main() {  
    char * haystack, * needle;  
}
```

// Rabin Karp matching algorithm.
// We will use the simple hash function
// where we add the ASCII values of
// the substring (b == 1).

```
int hash1 = 0;  
int hash2 = 0;
```

```
int M = strlen(needle);  
int N = strlen(haystack);
```

```
char * hay = haystack;  
char * need = needle;
```

```
while (*hay && *needle) {  
    hash1 = hash1 + *hay;  
    hash2 = hash2 + *hay;  
    need++; hay++;  
}
```

```
if (M > N) {  
    return 0; // no match  
}
```

```
// compute  
// hash1 and  
// hash2
```

```

char * hay-end = haystack + N - M - 1;
hay = haystack;

```

```

while (hay <= hay-end) {

```

```

    // compare hash1 and hash2

```

```

    if (hash1 == hash2) {
        // likely match. Test it.

```

```

return haystack; needle;

```

```

        char *p = hay;

```

```

        char *q = needle;

```

```

        while (*p && *q && *p == *q) {
            p++; q++;

```

```

        }
        if (*q == '\0') {

```

```

            // real match

```

```

            return hay;

```

```

        }
        // no match

```

```

    }

```

(187)

// recompute hash2

~~hash2 = hash2 - (*hay) +~~ ~~hash2~~ + ~~hash2~~

[m] key
hay [m]

old char

hash2 = hash2 - hay[0] + hay[m];

hay++; // go to next char

} // while

// No match
return NULL;

}

Data Compression

189

Fixed-length encoding

- Each character is represented with the same number of bits Example: ASCII. (8bits)

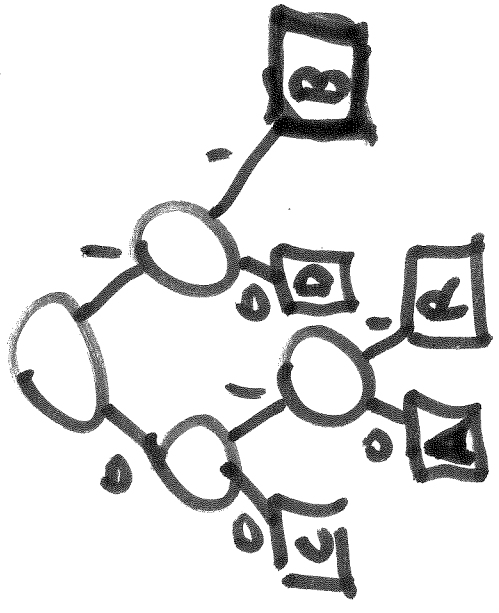
Variable length encoding

- Each character uses binary codes of different length.

Prefix Rule

To prevent ambiguities in variable length encoding, each encoding needs to satisfy the prefix rule, that is, no code is a prefix of another code.

We use an encoding ⁽¹⁹⁰⁾ to define an encoding that satisfies the prefix rule:



A = 010
 B = 11
 C = 00
 D = 10
 R = 011

010110110100010100101101011010
 A B R A C A D A B R A

11 chars = 29 bits

Compare to the ASCII representation
 8 bits/char

11 chars x 8 bits = 88 bits

Can we do better than 29 bits?

- Yes, we can do better.

A is used 5 times and has a representation of 010 (3 bits) and D is used 1 time and has a representation of 10 (2 bits).

- We could have a shorter representation of A and that will give us a smaller representation of the string.

Huffman Encoding Tree

- It builds an optimal encoding tree for a given string. Characters that happen more often than others get shorter representation.

- The algorithm builds a trie bottom-up using the characters and subtrees of smaller frequency first.

A B R A C A D A B R A

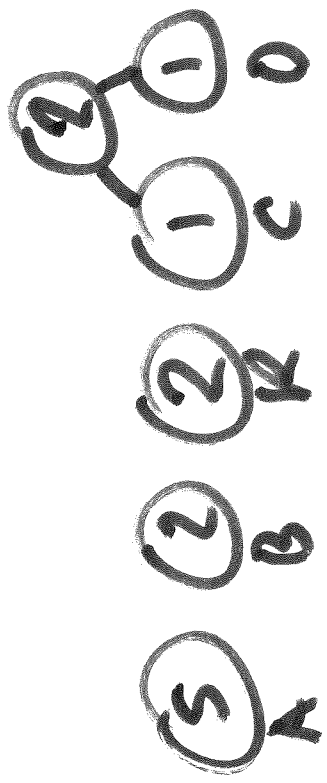
A	B	R	C	D	Frequency	$O(n)$
5	2	2	1	1		

Step 1



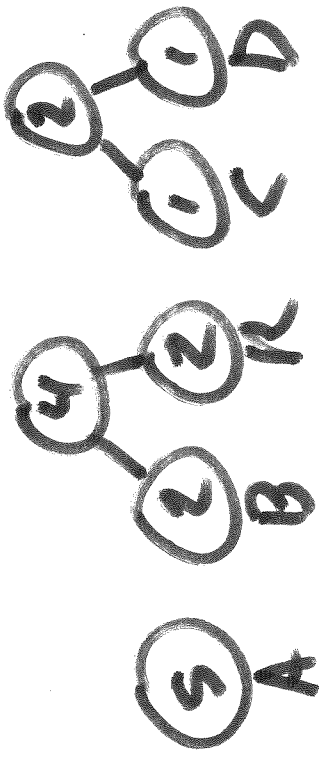
Step 2

Choose the 2 subtrees with minimum frequency.

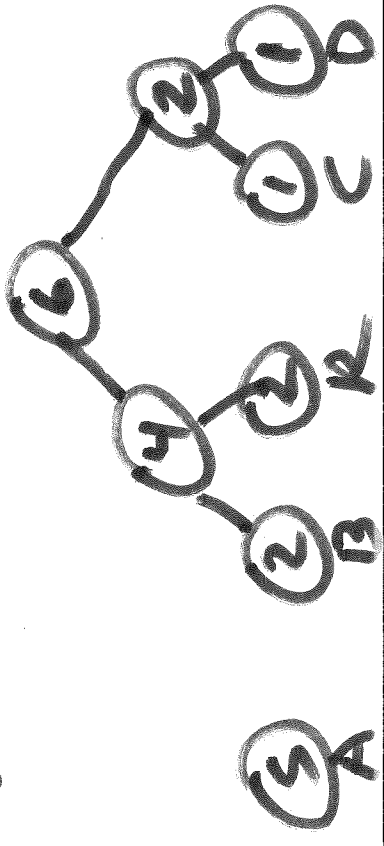


Step 3 →

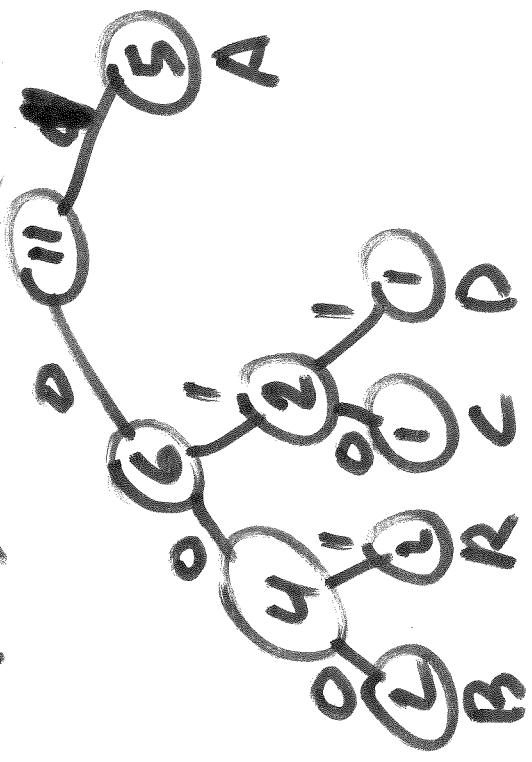
Again choose 2 subtrees with min frequency.
If there are more than 2 min subtrees choose any.



Step 4 Again choose 2 min subtrees



Step 5 choose min subtrees again



Optimal subtree

- A-1
- B-000
- R-001
- C-010
- D-011

ABRACADABRA

$\frac{1}{A} \frac{000}{B} \frac{001}{R} \frac{1}{A} \frac{010}{C} \frac{1}{A} \frac{011}{D} \frac{1}{A} \frac{000}{B} \frac{001}{R} \frac{1}{A}$

23 bits

Complexity of Huffman Encoding Optimal!

$O(N)$ - Build Frequency Histogram

$O(K \log K)$ - Build trees assuming K characters.

- Huffman encoding will not do well if

- We use all characters.
- All characters appear in the string with the same frequency.

- However with normal text some characters appear more often than others so Huffman encoding will give a good compression.

Quiz 22

196

In what situation Huffman encoding will not give a smaller representation?

- a) All characters are used
- b) All characters have the same number of occurrences.
- c) a ab d b
- d) a or b

Quiz 21

188

If Rabin-Karp string matching the worst case time is

a) $O(N)$

b) $O(M)$

c) $O(MN)$

d) $O(N^2)$