

CS251

①

Summer 2016

Data Structures.

- Gustavo Rodriguez-Rivera
grr@purdue.edu

LWSN 1203

- www.cs.purdue.edu/homes/cs251

- Grades

25% - Mid-term

25% - Final

10% - Attendance.

40% - Labs + HW's

50.

- ~~Pre~~ Post questions in Piazza

Goal

②

- In this course you will learn how to represent different types of data in optimal way in computer data structures.
- It is also an introduction to computer algorithms.

Algorithm + Data Structure = Program
② * A. W. R. H.

- You will also improve your programming skills.

3

Content

- Analysis Tools
- Stacks, Queues
- Vectors, Lists, Sequences.

- Trees
- Priority Queues (Heap data structure) $\log n$

- Dictionaries
- Search Trees (AVL trees) $2/3$ trees, RB trees
- Sorting (Quick Sort, Merge Sort, Heap Sort, Bucket Sort, Radix Sort).
- Text processing / Strings
- Graphs / Graph Algorithms.

Math Review

(4)

$$a^{b+c} = a^b a^c$$

$$a^{bc} = (a^b)^c$$

$$\frac{a^b}{a^c} = a^{b-c}$$

$$\log_a b = c \iff a^c = b$$

$$\log_b xy = \log_b x + \log_b y$$

$$\log_b \frac{x}{y} = \log_b x - \log_b y$$

$$\log_b x^a = a \log_b x$$

— Floor $\lfloor x \rfloor =$ the largest integer $\leq x$
Ceiling $\lceil x \rceil =$ the smallest integer $\geq x$

Summations

⑤

$$\sum_{i=s}^t f(i) = f(s) + f(s+1) + f(s+2) \dots f(t)$$

$s \rightarrow$ start index

$t \rightarrow$ end index.

Geometric Progression

$$\sum_{i=0}^n a^i = 1 + a + a^2 + \dots + a^n = S(a, n)$$

$$aS(a, n) = a + a^2 + a^3 + \dots + a^{n+1}$$

$$aS(a, n) - S(a, n) = a + a^2 + a^3 + \dots + a^{n+1} - 1 - a - a^2 - a^3 - \dots - a^n$$

$$aS(a, n) - S(a, n) = a^{n+1} - 1$$

$$S(a, n)(a-1) = a^{n+1} - 1$$
$$S(a, n) = \frac{a^{n+1} - 1}{a - 1}$$

Arithmetic Progressions

⑥

$$A(i) = \sum_{i=1}^n i = 1+2+3+\dots+n = \frac{n^2+n}{2}$$

$$= \frac{n(n+1)}{2}$$

When measuring an algorithm we ~~use~~ measure two parameters

- Time

- Space

* Time - How long does the algorithm take.
(secs)

It depends on:

+ Algorithm (itself)

+ Hardware (machine).

+ Input size and type of input.

+ Programmer/Language/Compiler/Optimizer

* Space - Computer storage used
(bytes) by the data structure

⑦ In this course we put more emphasis in Time rather than Space.

* Space is important but not as important as time.

Types of running time.

- Worst case - Time that algorithm takes with most difficult input.
- Average case - Average time across all inputs.
- Best case - Time that algorithm takes with easiest input.
- Time is measured experimentally using inputs of different sizes and types

- You may use the time command in Unix
time . command → prints how long
(prog) command takes.

8

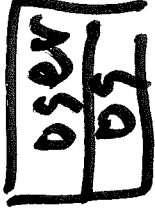
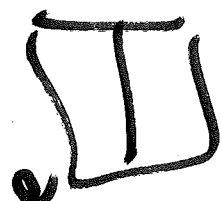
Eg. gcc myprog.c → How long compilation takes.

User space ← user: 2 secs

OS time (time in kernel) ← sys: 0.5 secs

← total: 3 secs

wall clock time



N CPU computer

user + sys ≤ N_{total}

single CPU

user + sys ≤ total

9

- Experimental studies of how long a program or algorithm takes have limitations.

- It is necessary to implement the algorithm.

- There is a limited set of inputs used.

It is impossible to measure execution for all inputs

- To compare algorithms they need to run in the same machine, same compiler, same input etc. same language etc.

- Pseudo code

- High-level description of an algorithm.

- Enough information to come up with a program that implements the algorithm.

(10)
We will develop a technique to compare algorithms that does not require an experimental analysis. This technique is called "Asymptotic Analysis".

Asymptotic Analysis

We will use:

- The worst case
- Very large input size $n \rightarrow \infty$

Steps:

1. We count the operations for an input size n .
2. We assume that n is very large ($n \rightarrow \infty$) and eliminate the components that become unimportant as n grows to infinite.

Asymptotic notation

(11)

Also called "Big-Oh" notation.

Given functions $f(n)$ and $g(n)$
we say that $f(n)$ is $O(g(n))$
if and only if $f(n) \leq c g(n)$ for $n > n_0$
for some c .



$n_0 \rightarrow f(n) \leq c g(n)$ for $n > n_0$
then $f(n)$ is $O(g(n))$

Example:

(12)

$$\text{Execution time } f(n) = .000001 n^2 + 100n + 1000$$

$$\text{Execution time } f(n) = O(?)$$

$$\text{Execution Time } (n=1) = .000001(1)^2 + 100(1) + 1000 \\ = .000001 + 100 + 1000 \sim 1100$$

↳ most significant

$$\text{Execution Time } (n=1000) = .000001(1000)^2 + 100(1000) + 1000$$

↑

$$1 + 100,000 + 1000$$

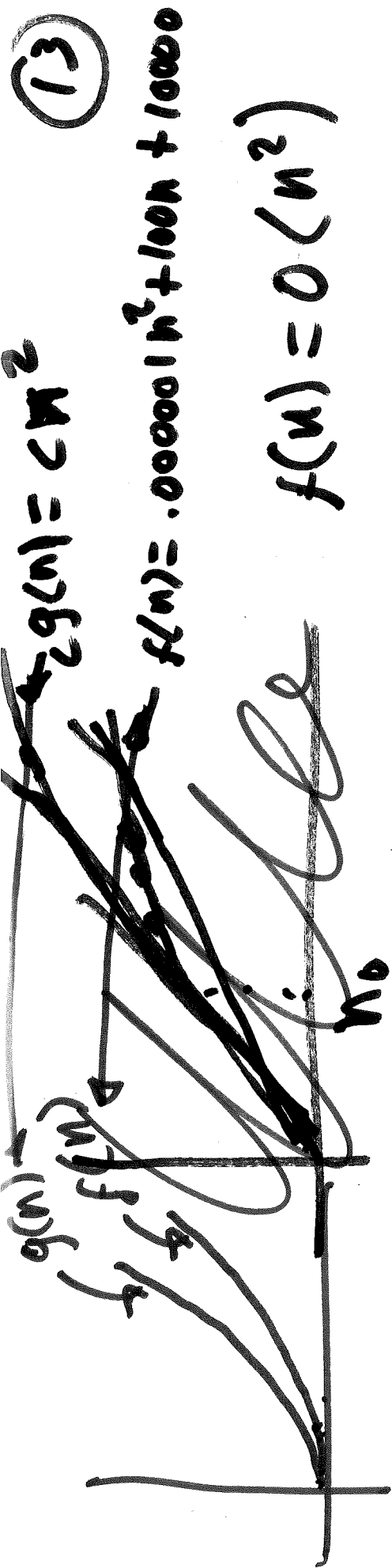
$$\text{Execution time } (n=10^{20}) = .000001(10^{20})^2 + 100(10^{20}) + 1000$$

$$10^{40} + 10^{22} + 1000$$

$$\text{Execution Time } = O(n^2)$$

$$10^{34} + 10^{22} + 1000$$

↳ most important.



$f(n)$ is ~~also~~ technically also $O(n^3)$ or $O(n^4)$ but we will use the $g(n)$ that most accurately describe how $f(n)$ grows.

Goal of Asymptotic analysis is to get rid of unneeded information.

- It is expected that the approximation should be as small in order as possible.

$f(n) = 7n - 3$ is also $O(n^5)$ but we choose instead $O(n)$.

- Drop lower order terms and constant factors

(14)

$$7n - 3 \text{ is } O(n)$$
$$8n^2 \log n + 5n^2 \text{ is } O(n^2 \log n)$$
$$7 \log n + n \text{ is } O(n)$$

Special cases of algorithms

$O(\log n)$ logarithmic

$O(n)$ linear

$O(n^2)$ quadratic

$O(n^k)$ polynomial

$O(a^n)$ exponential

$$5n^3 + (1.03)^n \rightarrow O(a^n) \text{ exponential}$$

more important $n \rightarrow \infty$

Example

Sort an array of n numbers

Algorithm bubbleSort(A, n)

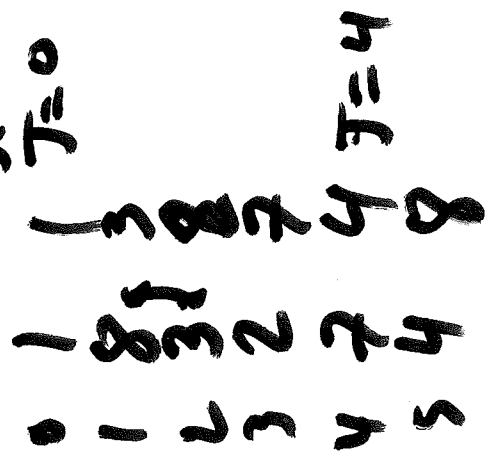
Input: array of n integers

Output: Sorted array.

```

for (i=0; i < n-1; i++) {
    for (j=0; j < n-i-1; j++) {

```



```

    if (A[j] > A[j+1]) {

```

```

        // Need to swap
        int tmp = A[j];
        A[j] = A[j+1];
        A[j+1] = tmp;
    }
}

```

$\rightarrow O(n^2)$

Time = $O(?)$

time = $\sum_{i=0}^{n-1} ((n-1) + (n-2) + (n-3) + \dots + 0) =$

$$\text{time} = c [n(n-1) - \sum_{i=1}^{n-1} i]$$

$$= c [n^2 - n - \frac{(n-1)(n)}{2}]$$

$$= c \left[\frac{2n^2 - 2n - n^2 + n}{2} \right]$$

$$= c \left[\frac{n^2 - n}{2} \right] = O(n^2)$$

$$\textcircled{16} \quad \sum_{i=1}^{n-1} i = \frac{(n-1)(n)}{2}$$

Stacks, Queues, Linked Lists (17)

Stack

It is a container of objects where they are inserted and removed according to the last-in-first-out principle (LIFO).

Two operations:
push(val) - inserts item at the top of stack

val = pop() - removes item at the top of stack and returns

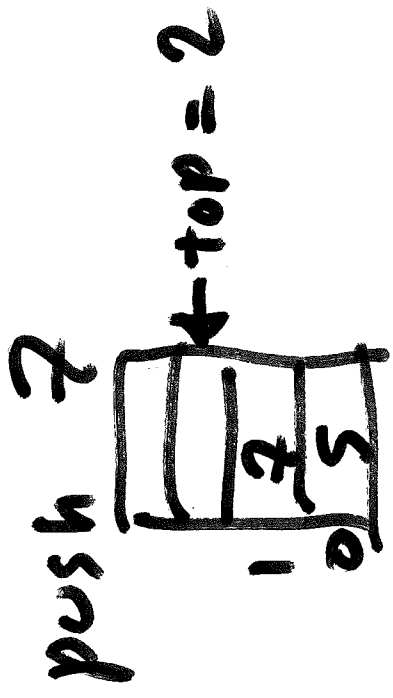
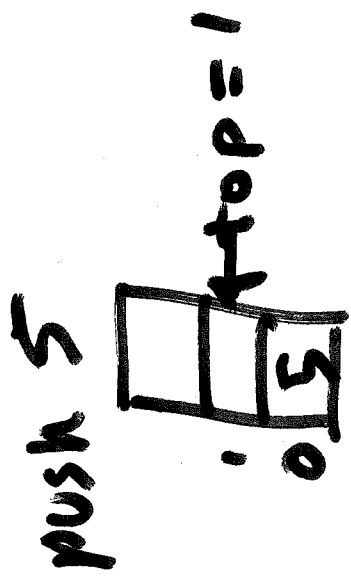
Additional methods: 17.

size() → # of elements in stack

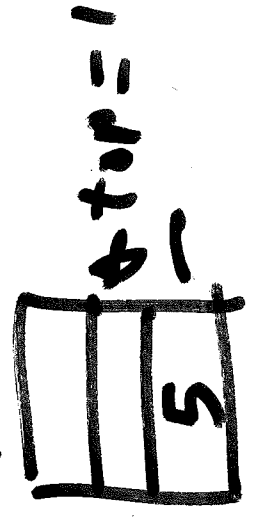
isEmpty()

top()

(8)

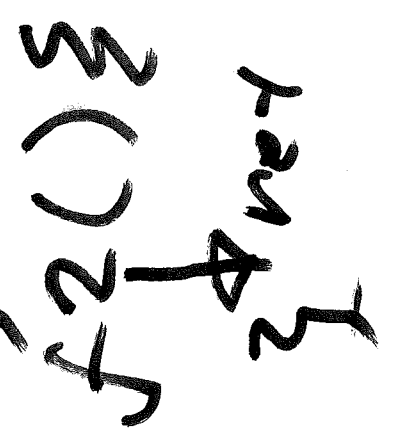
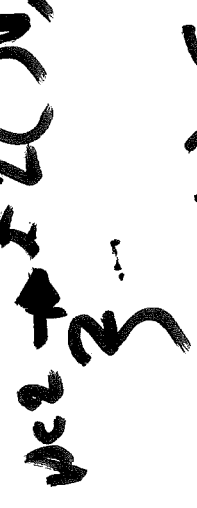
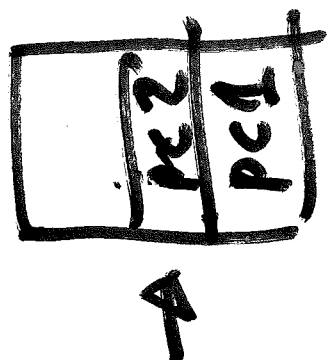
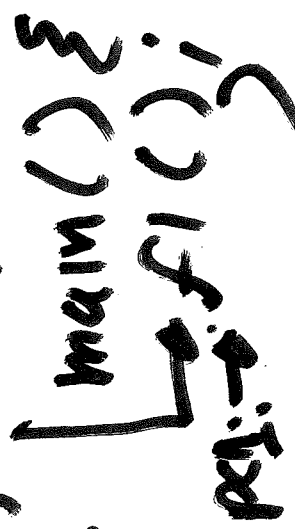


val = pop() → val == 7



where are stacks being queried? (19)

→ Execution Stack.



20

Quiz 1

The function

$$f(n) = 3 \log(n) + 50n + 3n^{1/2}$$

is

a) $O(n)$

b) $O(\log n)$

c) $O(n^{1/2})$

②

Hash Tables

A hash table is a data structure that stores pairs (key, data)

key(name) data(grade)

student, grade

George 90
Mary 92
Peter 92

struct

key(name) data (student record)

E.g.

George
Mary
Peter

address, SID, Maj
(2006 Grant, 3527, CS)
(34 Main, 2275, EE)
(43 7th, 753, Econ)

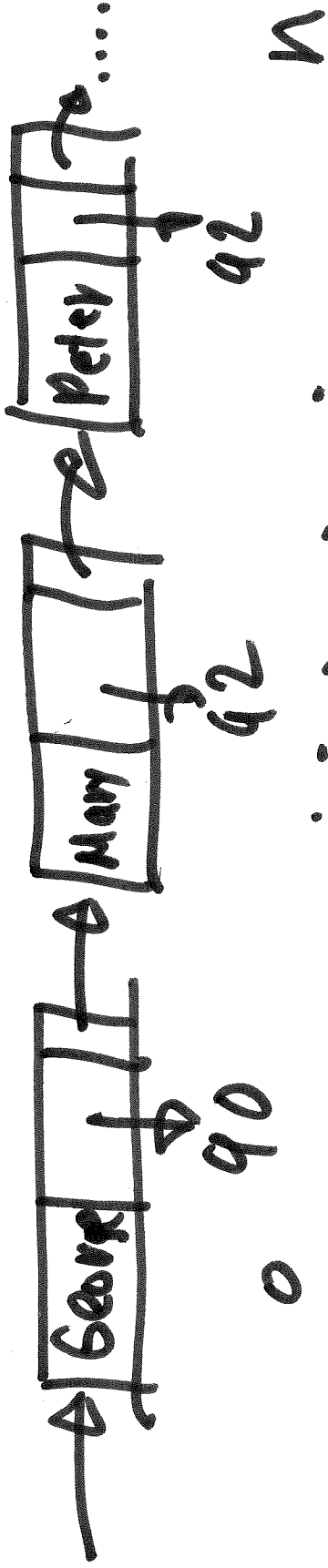
If we store table in an array, a lookup in n entries will be

0	George	90
1	Mary	92
2	Peter	92
⋮	⋮	⋮
n		

$O(n)$ worst case.

grade = lookup("Paul");

If we use linked list



also grade = lookup("Paul");

$O(n)$

Quiz 2

Assume we allocate an array as

~~array = new int [n];~~
array = new int [n];

To deallocate array we use:

- a) delete array
- b) free array
- c) delete [] array

as 23

23

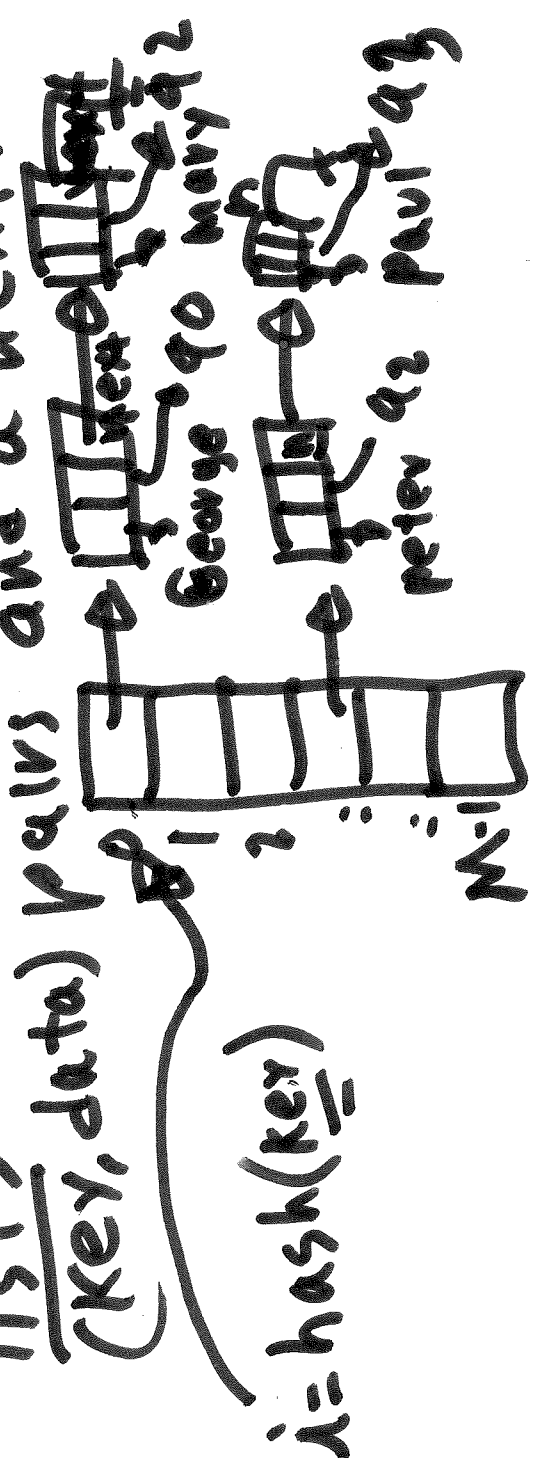
52

A hash table also implements a dictionary to store (key, data) pairs.

A hash table has an expected time of $O(1)$ for lookup operation

The worst case is $O(n)$ but it is unlikely to happen if we follow some rules we will talk about later.

A hash table is made of an array of lists where each entry in the list have



The hash function takes a key (25) as argument and generates an integer $0 \rightarrow M-1$ hash table

The hash function/needs to have the following properties:

1. $\hat{h} = \text{hash}(\text{key})$ has to generate numbers $0 \leq \hat{h} \leq M-1$

2. The hash function has to be evenly distributed, that is all numbers $0 \leq i \leq M-1$ ~~has~~ will have the same probability from occurring for the input keys

3. The number M of buckets (lists) in hash table have to be larger or comparable to N (# of keys) $M \geq N$

97
If the hash table has these 3 characteristics, the expected time will be $O(1)$.



Algorithm backup
data = lookup(key):

$\lambda = \text{hash}(key)$

list = bucket[λ];

while list != null && list->key != key:
list = list->next;

if list == null: //reached end
return null;

else: //found
return list->data

Algorithm insert

27

boolean insert(key, data):

i = hash(key)

list = bucket[i];

while list != null & list->key != key:

list = list->next;

if list == null: // key exists

overwrite data

~~entry = new Entry
entry->key = key
entry->data = data~~

list->data = data

return true

else: // key does not exist

entry = new Entry

entry->key = key

entry->data = data

entry->next = bucket[i];

```

bucket[i] = entry;
end
return false
end

```

(28)

— Hash functions:

```

// If key is an int
int hash(int key) {
    return key % M;
}

```

String key is a string, a hash function
 If could be implemented by adding the ASCII
 value of all characters in key

(29)

```
int hash(char *key) {  
    int sum = 0;  
    for (int i = 0; key[i] != '\0'; i++) {  
        sum += key[i];  
    }  
    return sum % M;  
}
```

One problem of this hash function will be that some ~~characters~~ keys with the same characters but permuted will have the same hash value. apple → pleap
In an improvement is to use the position of a character as part of the hash function.

(30)

```
int hash(char *key) {  
    int sum = 0;  
    for (int i = 0; key[i] != '\0'; i++) {  
        sum += i * key[i];  
    }  
    return sum % M;  
}
```

Also to improve the hash table performance and have the hash function be more evenly distributed is to choose M to be prime to prevent hash values from being periodic.

31

```
HashTable Word:: HashTableVoid() {  
    _buckets = new HashTableVoidEntry * [Table Size];
```

```
    for (int i = 0; i < Table Size; i++) {
```

```
        _buckets [i] = NULL;
```

```
    }
```

```
}
```



Quiz 3

32

Which of the following properties are necessary to have a lookup() expected time of $O(1)$ in a hash table;

a) hash function needs to be evenly distributed

b) $M \geq N$

c) $N \leq M$

d) a and b

e) a and c

~~$M = \text{size}$~~

$M = \# \text{ of buckets}$

$N = \# \text{ of keys}$

32

Another use of stacks

Matching enclosing { }, (), []

function isMatchOk(const char *s) {

// It returns true if s has nested

// { }, (), [], that match each other

// without interleaving if ([]) {

} for(

{

{ () () { [] } ok



{ [] ([]) { } ok



{ [] ([]) { } Not ok



{ { () } } Not ok



// there is an element in the
 // stack. pop <= stk[top-1];
 top--;

```

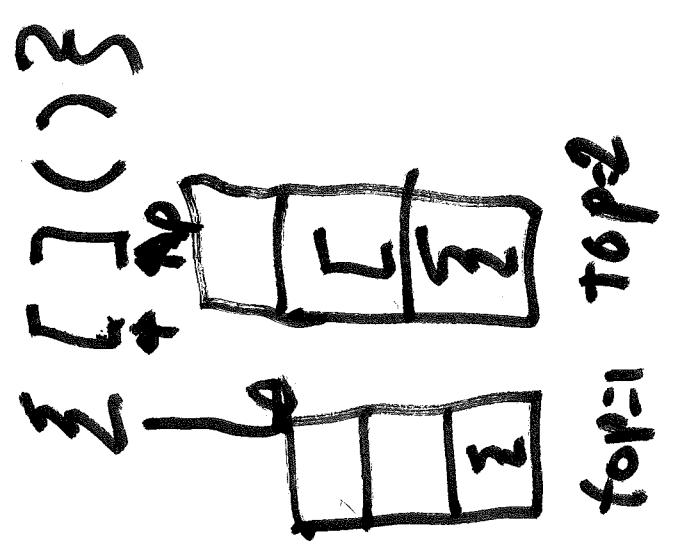
if (c == 'Σ') {
  if (*p == 'Σ') {
    delete stk;
    return false;
  }
}
else
  return false;
  
```

```

else if (c == 'Σ') {
  if (*p != 'Σ') {
    delete stk;
    return false;
  }
}
  
```

```

else if (c == '(') {
  if (*p != '(') {
    delete stk;
    return false;
  }
}
  
```



93

}

```
{  
  // we reached the end of string  
  // make sure nothing was left  
  if (top > 0) {  
    delete stk;  
    return false;  
  }  
  delete stk;  
  return true;  
}
```

}

Queues

37

They are data structures that follow the policy FIFO, that is the next item removed will be the oldest one in the queue.

Implementation of a queue using an array

insert 5

insert 6

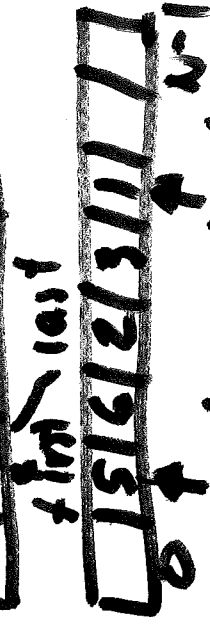
insert 2

remove \rightarrow 5

remove \rightarrow 6

remove \rightarrow 2.

5, 6, 2



We use an array and first and last.

two indexes: first and last.

- We add at the "last" position 38 and remove from "first" position and increment first and last after every operation.
- When last or first reach the end they will start over at ϕ .
- The maximum number of elements in the queue will be the size of the array (n).

Quiz 4

39

What is the advantage of using templates vs. using void* for data
(Answer best option)

- a) Type checking
- b) More clear code
- c) Runs faster
- d) All of the above.

Implementation:

QueueInt.h

```
class QueueInt {
```

```
    int first;
```

```
    int last;
```

```
    int n; // # of elems in queue
```

```
    int maxN;
```

```
    int * arr;
```

```
public:
```

```
    QueueInt(int maxN);
```

```
    bool enqueue(int val); // returns false  
                            // if full
```

```
    bool dequeue(int &val); // return false  
                            // if empty
```

```
    ~QueueInt();
```

```
}
```

QueueInt.cpp

```

QueueInt::QueueInt(int maxn) {
  first = 0;
  last = 0;
  n = 0;
  this->maxn = maxn;
  array = new int [maxn];
}

```

```

QueueInt::~QueueInt() {
  delete[] array;
}

```

body

41

```
QueueInt::dequeue(int &val) {
```

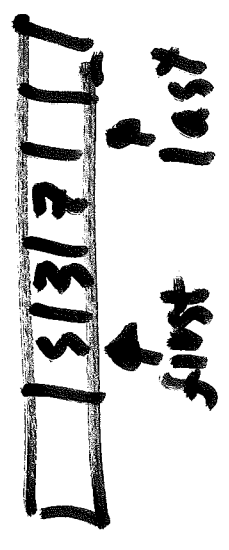
```
    // check if queue is empty  
    // dequeue will
```

```
    if (n == 0) {  
        // queue is empty  
        return false;
```

```
    }  
    // get oldest value  
    val = array[first];  
    first = (first + 1) % maxn;
```

```
    n--;  
    return true
```

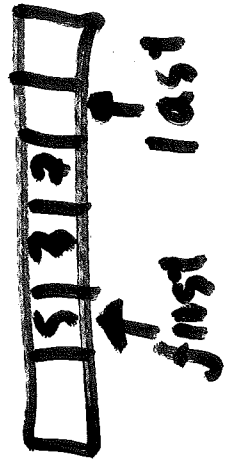
```
}
```



42 (34)

```
bool
QueueInt::enqueue(int val) {
```

```
// last is the index of
// the empty space in
// queue where we will
// insert
```



```
if (n == maxn) {
    // no space
    return false;
```

```
}
// there is space
```

```
array [last] = val;
```

```
last = (last+1) % maxn;
return true;
```

```
// increment and
// wrap around if
// necessary
```

```
}
```

Resizable Array Int.h

```

class ResizableArrayInt {
    int maxN; // size of array before
              // having to resize
    int n; // # of elements used in
           // array
    int *array;

```

```

public:
    ResizableArray(int initialSize = 100);
    ~ResizableArray();
    bool setAt(int pos, int val); // set val
    bool getAt(int pos, int &val); // at pos.
    void insertAtEnd(int val); // pos 0
                                // of bound

```

Σ

Resizable Array Int.cpp

46

```
#include "ResizableArrayInt.h"
```

```
ResizableArrayInt::ResizableArrayInt(  
    int initialSize) {}
```

```
    maxn = initialSize;
```

```
    n = 0;
```

```
    array = new int [initialSize];
```

```
}
```

```
· ResizableArray::~ResizableArrayInt() {  
    delete [] array;
```

```
}
```

void
ResizableArrayInt :: InsertAtEnd(int val) (45)

// check if there is space

if (n == maxn) {

// Expand array

~~int newarray;~~

maxn = 2 * maxn;

int newarray = new int[maxn];

// copy old array to new array

for (int i = 0; i < n; i++) {

newarray[i] = array[i];

}

delete [] array;

array = new array;

}

```
// there is space in array
array[n] = val;
n++;
```

```
bool
ResizableArray:: setAt (int pos, int val) {
  // check if pos is between bounds.
  if (pos < 0 || pos >= n) {
    return false;
  }
}
```

```
array[pos] = val;
return true;
}
```

```
bool
ResizableArray:: getAt (int pos, int &val) {
  if (pos < 0 || pos >= n) {
    return false;
  }
}
```

```
val = array[pos];
return true;
}
```

Analysis

Inserting an item will be most of the time $O(1)$ except when we have to resize. Resizing is $O(n)$.

For n inserts

Cost of insert: $n O(1)$

Cost of resize: # of resizes $\neq O(n)$.

of resize for n elements: $O(\log n)$

$$n = \max n = 16 m 2^r$$

$$n < 16 m 2^r$$

$$\frac{n}{16m} < 2^r \Rightarrow \log_2 \frac{n}{16m} < r$$

$$r = O(\log n)$$

of resize

Total cost for n inserts = $n O(1) + O(\log n)$
 $= O(n \log n)$

48

If instead of duplicating $\max n$

$$\max n = 2 * \max n$$

during resize, we increase by a constant.

$$\max n = \max n + \text{increment};$$

then the number of resizes for n elements will be:

$$\# \text{ of resize} = \frac{n}{\text{increment}} = O\left(\frac{n}{\text{increment}}\right)$$

$$= O(n)$$

So the total cost for n inserts

if we increment $\max n$ by a constant

$$\text{is: } nO(1) + O(n) = O(n) = O(n^2)$$

inserts # resize cost of ~~resize~~

So it is better to double the size of the array for inserts

Amortized cost = cost of n inserts
Average cost for n inserts
Worst case = $O(n)$ for n inserts
when insert resizes = $O(n \log n)$

Quiz 5

619

What is the ^{total} cost of n inserts in a resizable array assuming you double the max size at every resize.

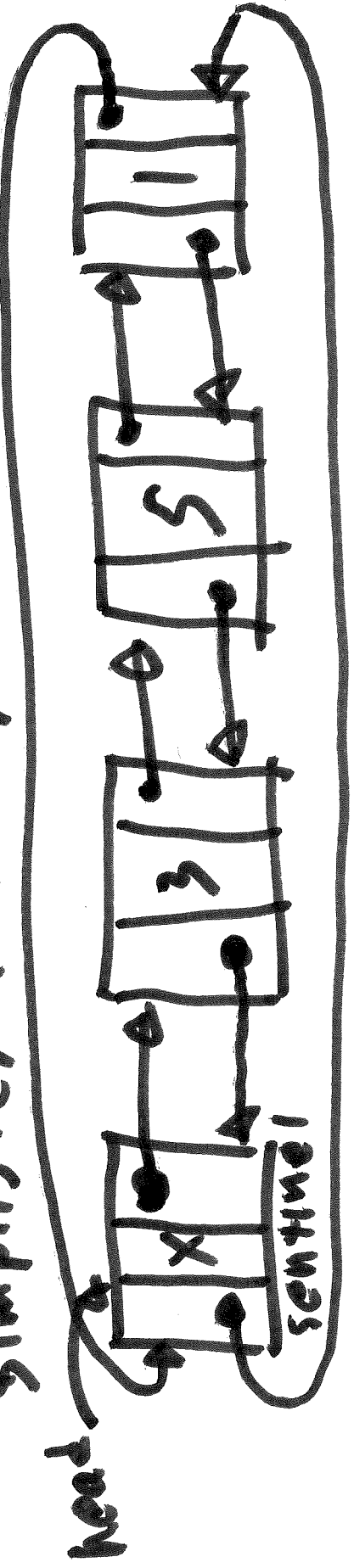
- a) $O(n)$
- b) $O(n^2)$
- c) $O(n \log n)$
- d) $O(\log n)$

51

Double Linked lists

→ ~~Prepared~~ data structure that stores elements in a list where each node has a previous and next pointer.

→ There is a node at the beginning of the list called "sentinel" or "head" that does not store any elements but simplifies the implementation.



Double linked lists are popular because

52

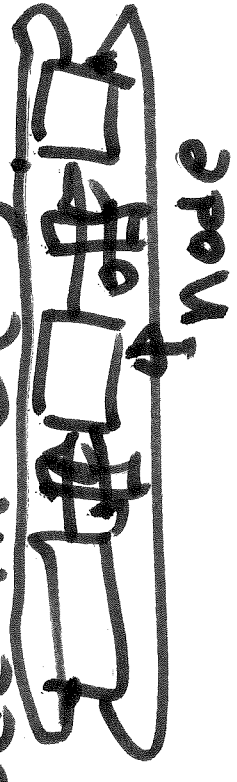
Advantages → No resizing is needed like in arrays

insert is $O(1)$ worst case.

insert Front (val) $O(1)$
insert End (val) $O(1)$

→ You can insert at the end in $O(1)$ compared to single linked lists that insertion at end is $O(n)$.

→ Given a pointer to a Node, we can insert before or after the node in $O(1)$



Disadvantage

getElemAt(i) set an element at a position in a double linked list is $O(n)$. In an array this operation is $O(1)$.

DLList.h

```
class DLLNode {
```

```
    int val;
```

```
    DLLNode *next;
```

```
    DLLNode *prev;
```

```
}
```

```
class DLLList {
```

```
    DLLNode *head;
```

```
public:
```

```
    DLLList();
```

```
    ~DLLList();
```

```
    void insertFront(int val);
```

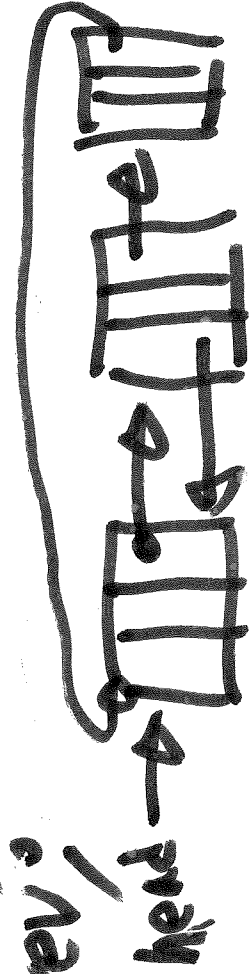
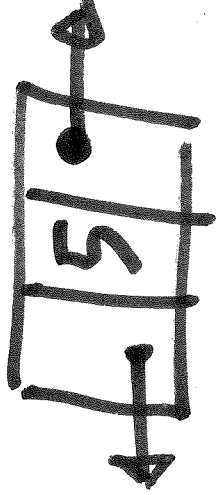
```
    bool removeFront(int val);
```

```
    void insertEnd(int val);
```

```
    bool removeEnd(int val);
```

53

DLLNode



// puts val at
 // in val and
 // ret false
 // if err.

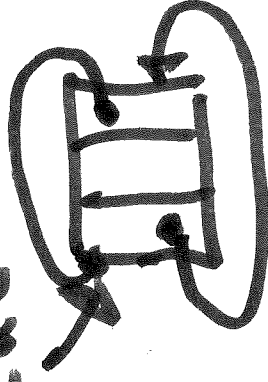
bool isEmpty();

void insertAfter(DLLNode *node, int val);
void insertBefore(DLLNode *node, int val);
void remove(DLLNode *node);

}

__DLList.cpp
#include "DLList.h"

head



DLList::DLList() {
head = new DLLNode();
head->next = head;
head->prev = head;
}

}

55

```
void
DLLisf:: insertFront (int val) {
```

```
    DLLNode *n = new DLLNode
        DLLNode();
```

```
    ① n->next = head->next;
```

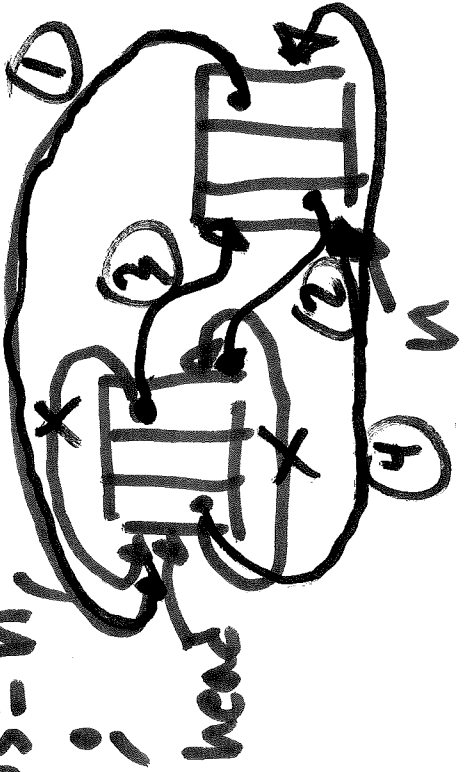
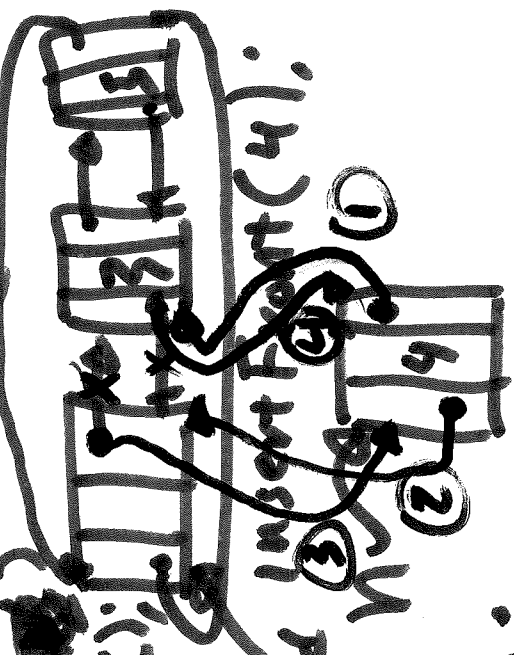
```
    ② n->previous = head; head = n;
```

```
    ③ head->next = n;
```

```
    ④ n->next->previous = n;
```

```
    n->value = val;
```

```
}
```



O(1)

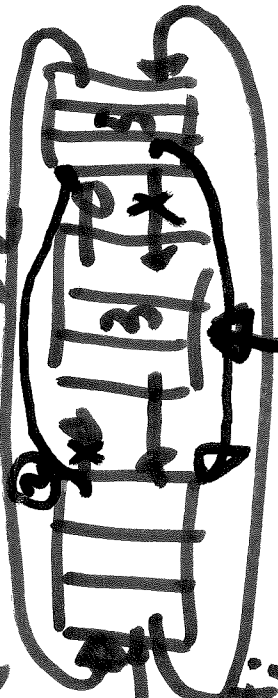
56

bool

Dllist::remove Front (int &val) {

if (isEmpty()) {

return false; head



} DllNode *n = head->next;

② head->next = n->next; n

③ n->next->previous = head;

val = n->val;

delete n;

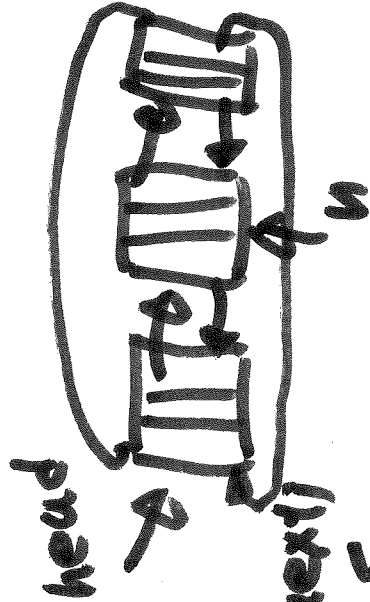
}

57

```
bool  
Dllist::IsEmpty() {
```

```
    return head->next == head;
```

```
}
```



```
Dllist::~Dllist() {
```

```
    DllNode *n = head->next;
```

```
    while (n != head) {
```

```
        // save next pointer
```

```
        DllNode next = n->next;
```

```
        delete n;
```

```
        n = next;
```

```
    } delete head;
```

```
}
```

void

DLList::insertEnd(int val) {

head = prev

DLLNode *n = new DLLNode(val, insertEnd);

n->val = val;

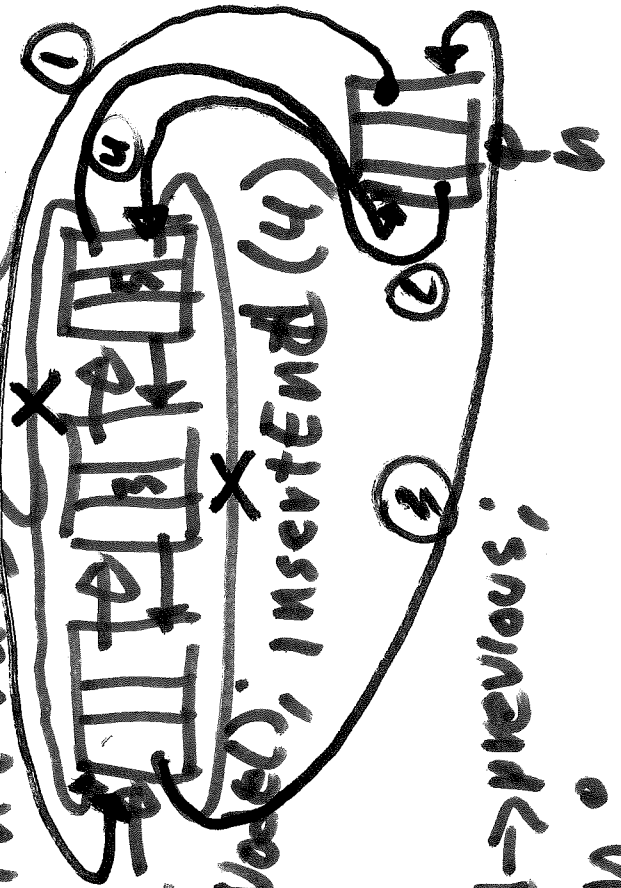
① n->next = head;

② n->previous = head->previous;

③ head->previous = n;

④ n->previous->next = n;

}



85

Quiz 6

~~Three following~~

What is the missing line in the following method

```

void insertFront(int val) {
    DLLNode * n = new DLLNode();

```

```

n->next = head->next;
n->previous = head

```

MISSING LINE

```

n->next->previous = n;
n->val = val;

```

- 3
- a) head->previous = n;
 - b) head->next = n;
 - c) n->previous = head
 - d) head->next->previous = n;

void

DLLList::insertAfter(DLLNode *node, int val) {

DLLNode *n = new DLLNode(val);

n->val = val;

① n->next = node->next;

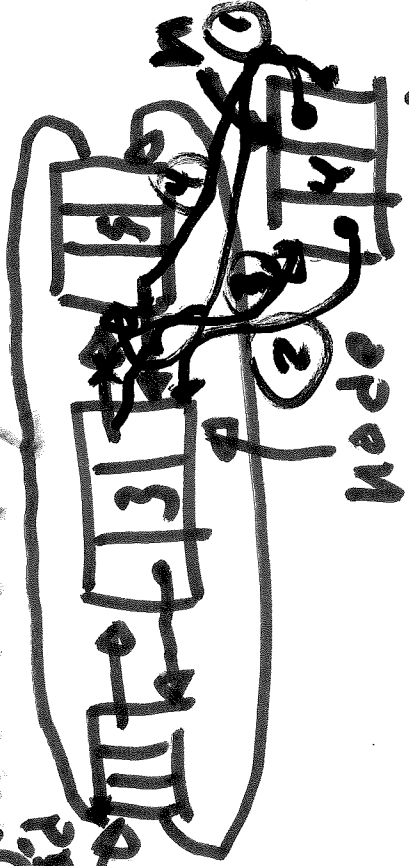
② n->previous = node;

③ node->next = n;

④ n->next->previous = n;

O(1)

insertAfter(node, 4);



}

void

DLList::remove(DLLNode *node) {

// removes node from list

// node is a valid node in the list

① node->previous->next =

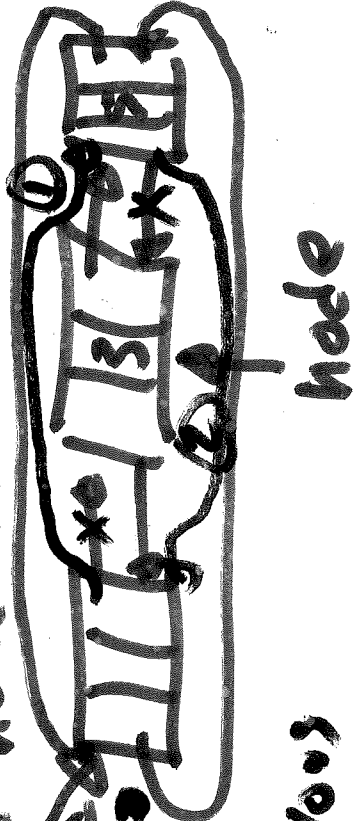
= node->next;

② node->next->previous =

= node->previous;

delete node;

61



Σ

Comparison in array and DLList implement.

- isEmpty()

- insertAtEnd

fixed worst about.

- insertFront

- getValAtPos(i-th)

array O(1)

O(n)

O(n)

O(1)

DLList

O(1)

O(1)

O(1)

O(n)

Trees

62

Data structures used to represent a

hierarchy

Living Things

Animals — Plants

Vertebrates

Invertebrates

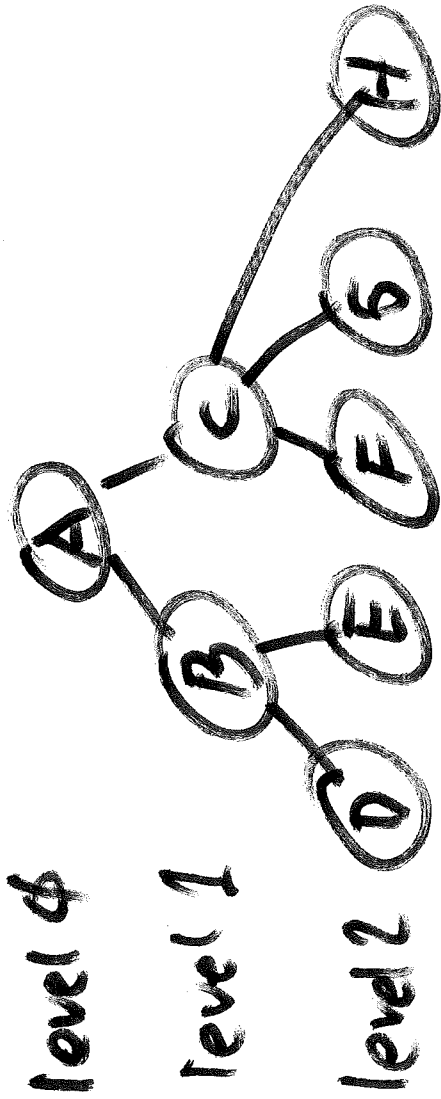
Fish — Amphibians — Reptiles — Birds — Mammals

File System

usr — bin etc

include — lib

Vocabulary



63

A is in the root node
B is the parent of D, E
C is the sibling of B
D, E are children of B
D, E, F, G, H are external nodes

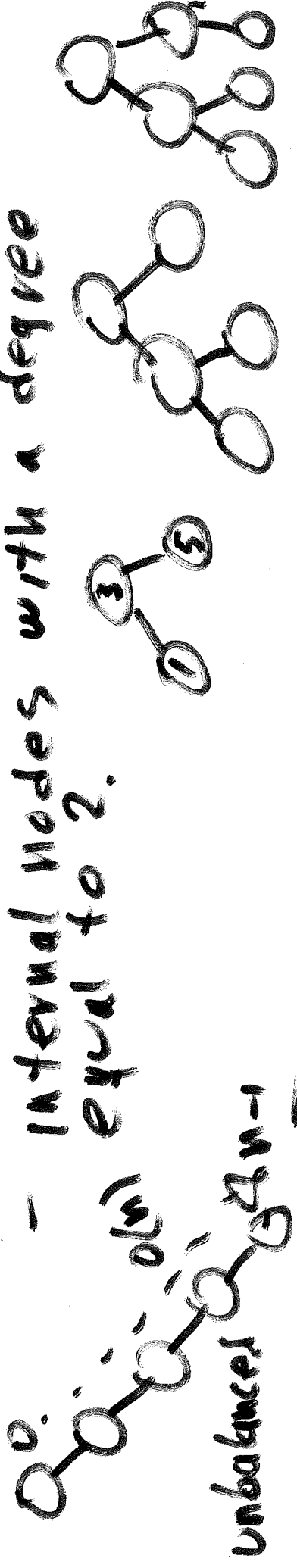
A, B, C are internal nodes
The depth (level) of E is 2.

The height of the tree is 3
The degree of B is 2 (# of children)

Ordered Tree - The children of each node are ordered.

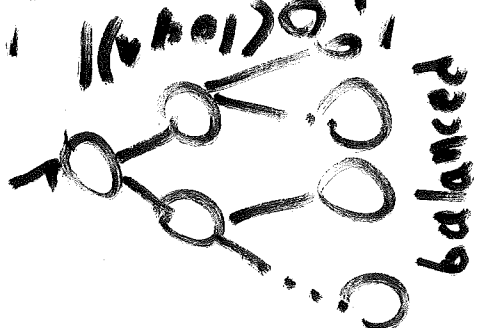
(there will be different types of ordering).

Binary Tree - ordered tree with all internal nodes with a degree equal to 2.



Balanced Tree -

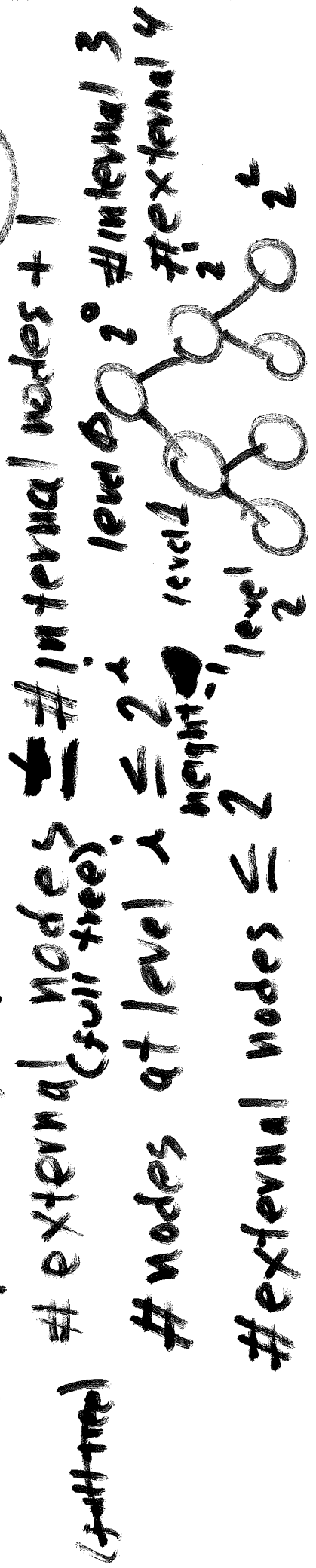
It is a tree where the distance from the root to any external node is $O(\log n)$ where n is the total number of nodes in the tree.



A balanced tree will allow a traversal from the root to any node to be $O(\log n)$

Properties of Trees

59



$\text{height} = 3$
 $\# \text{ external nodes} = 2^3 = 8$

Traversal of trees

66

Appendix

- Preorder Traversal
- Postorder Traversal
- Inorder Traversal

Preorder Traversal

Visit parent before children.

preorder(v) {

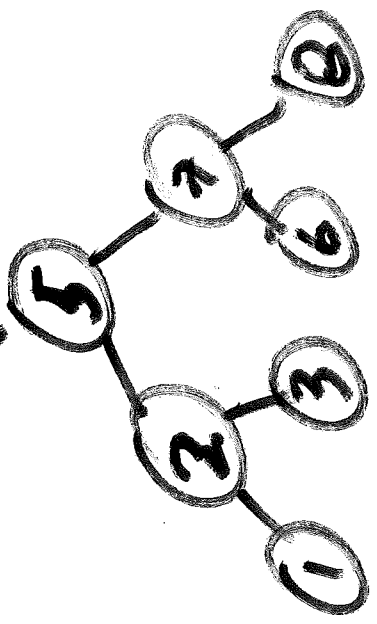
 visit parent

 for each child of v

 recursively call preorder(c)

}

root



visit \Rightarrow print value

preorder(root)

5, 2, 1, 3, 7, 6, 8

67

postorder traversal

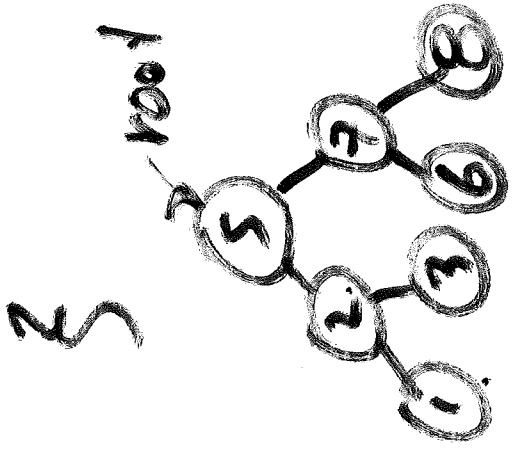
Visit the children before parent.

postorder(v) {
 for each child c of v
 recursively call postorder(c)

visit v

visit → print value
postorder(root)

1, 3, 2, 6, 8, 7, 5



inorder traversal

69

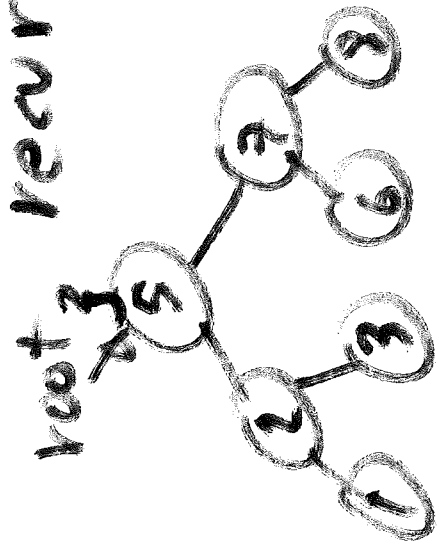
→ only for a binary tree
→ visit left child, parent, right child

$inorder(v) \{$

~~visit~~ recursively call $inorder(v \rightarrow left)$
visit v
visit $v \rightarrow right$

visit v

recursively call $inorder(v \rightarrow right)$



$inorder(root)$

1, 2, 3, 5, 6, 7, 8

Implementation of a binary tree

70

BinTree.h

```
struct BTreeNode {
```

```
    int val;
```

```
    BTreeNode *left;
```

```
    BTreeNode *right;
```

```
}
```

```
class BinTree {
```

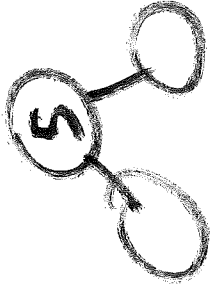
```
    BTreeNode *root;
```

```
public:
```

```
    BinTree();
```

```
    void insert(int val);
```

```
}
```



17

void
BinTree :: InOrderPrint (BTNode *node) {

if (node == NULL) return;
InOrderPrint (node->left);
printf ("%d", node->val);
InOrderPrint (node->right);

}

void
BinTree :: preOrderPrint (BTNode *node) {

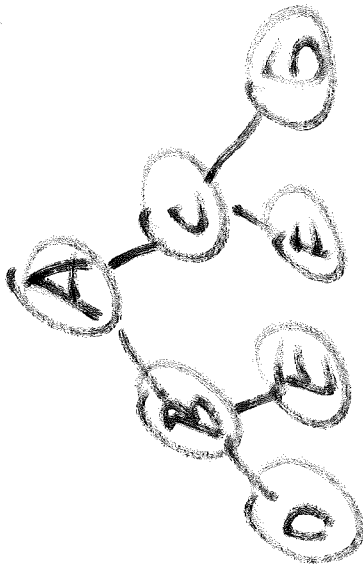
if (node == NULL) return;
~~preOrderPrint (node~~
printf ("%d", node->val);
preOrderPrint (node->left);
preOrderPrint (node->right);

}

Quiz 7

72

What is the post order traversal
of the following tree



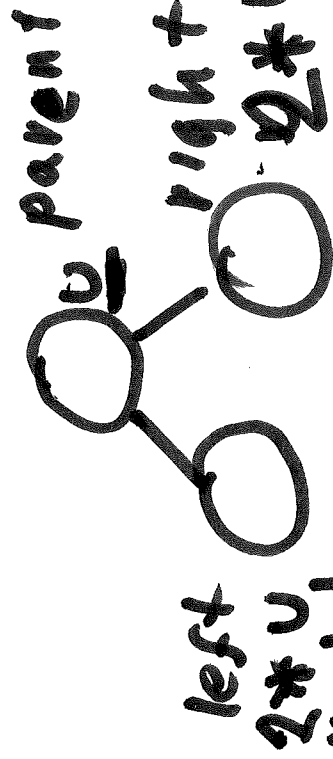
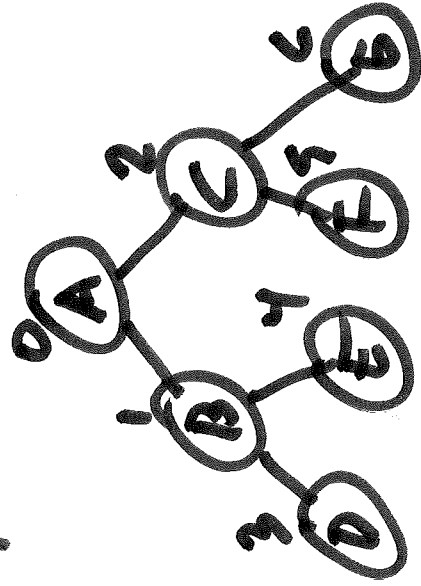
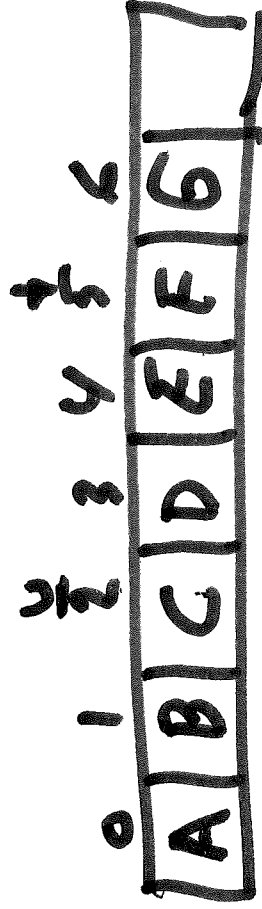
- a) A, B, C, D, E, F, G
- b) D, B, E, A, F, C, G
- c) D, E, B, F, G, C, A
- d) G, C, F, E, B, A

Representation of Binary Trees 73

We have seen already a pointer representation of a binary tree.

```
struct BTreeNode {  
    BTreeNode *left;  
    BTreeNode *right;  
    int val;  
};
```

We can also represent a binary tree using an array



We assume tree is full. +1 if not leave nodes empty if not used.

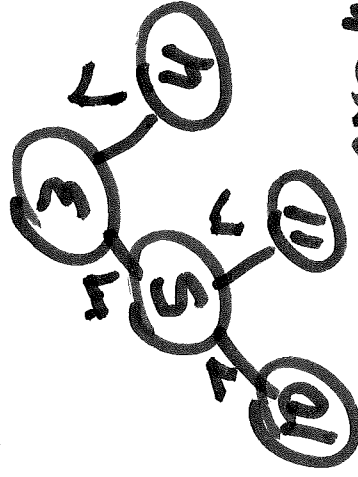
Priority Queue

- A priority queue is a data structure where we can insert items in any order but only the element with the minimum key can be removed.

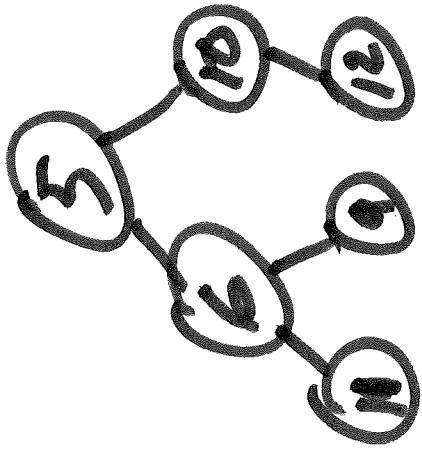
Heaps

- A heap is a data structure that implements a priority queue using a binary tree.

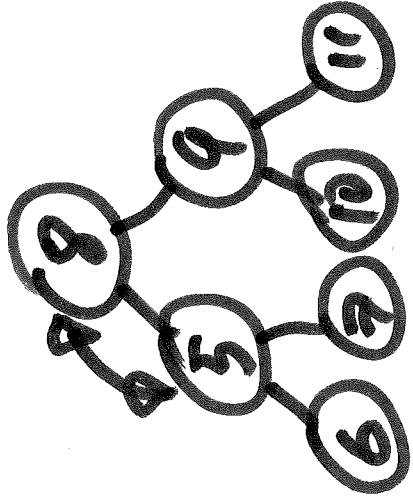
Example:



- It satisfies two properties
 Order Property: $\text{key}(\text{parent}) \leq \text{key}(\text{child})$
 Structural Property: All levels are full nodes except the last one where all nodes are left filled.



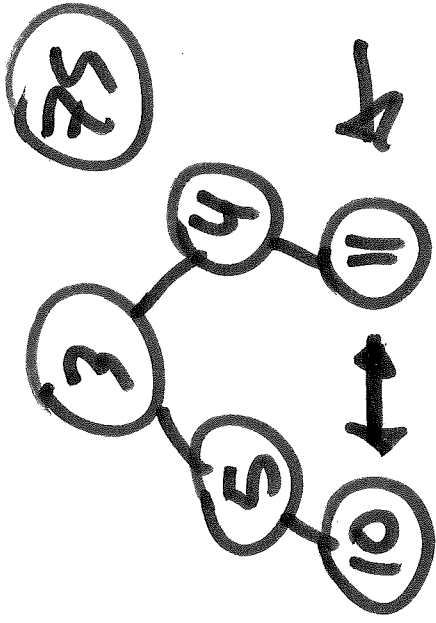
Heap



Not a heap

5 < 8
parent larger than child

Does not preserve order property



Not a heap

Breaks

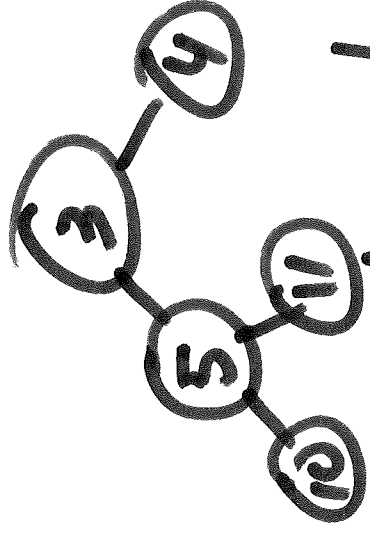
structural

property

Not left

filled

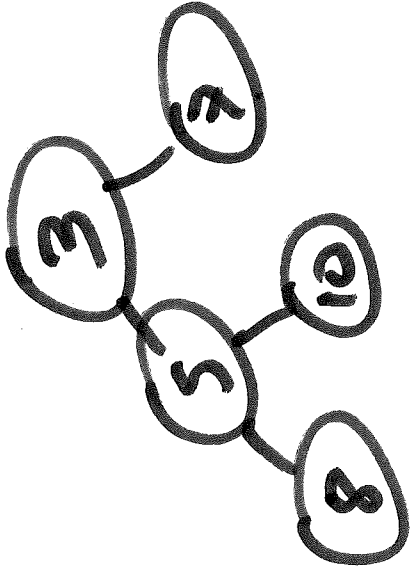
11 has to be child of 4



Heap!

Operations of a Heap

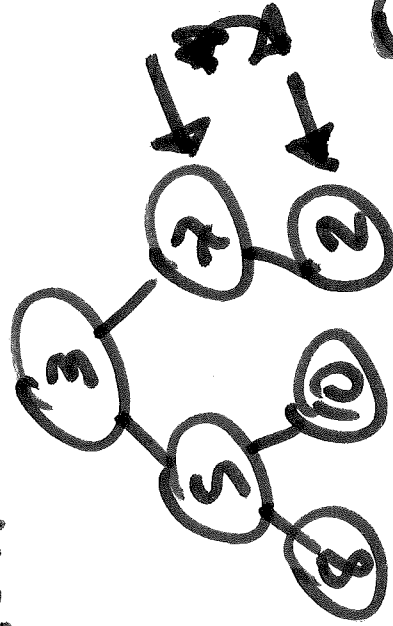
Initial Heap



Insertion

Insert (2)

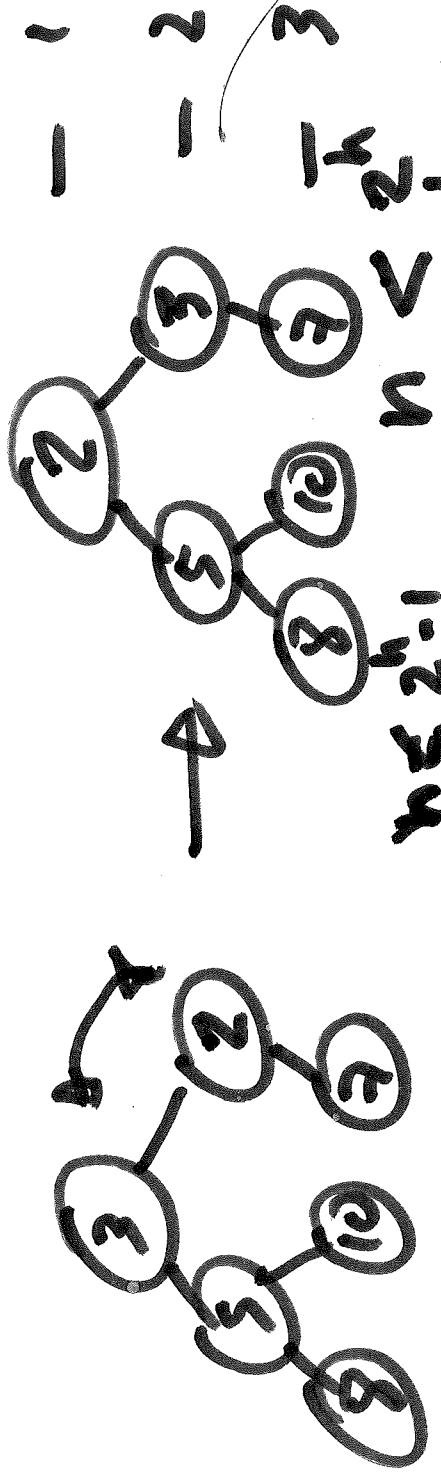
Add key in the first available position



This is not a heap $2 < 7$
child < parent

77

We do an "up heap" step
 swapping the child and parent
 and going up until the heap
 is ordered or we reach the root. h



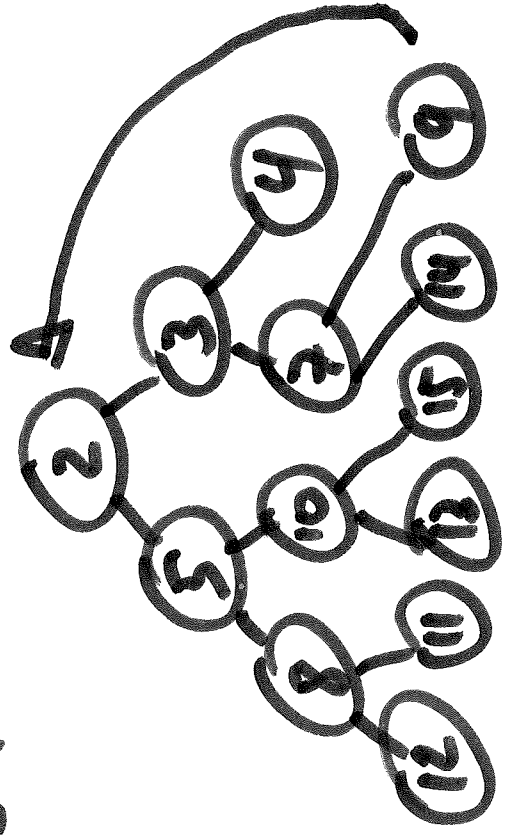
total # of swaps operations = $h-1$

height = $O(\log n)$

→ # of swaps is $O(\log n)$
 so insert in a heap takes
 $O(\log n)$.

Remove Min

- Remove element at the root
- Removal of root leaves a "hole"
- Replace root with last key of the heap.
- Do a "down heap" step swapping elements parent/children until heap is sorted or we reach bottom.



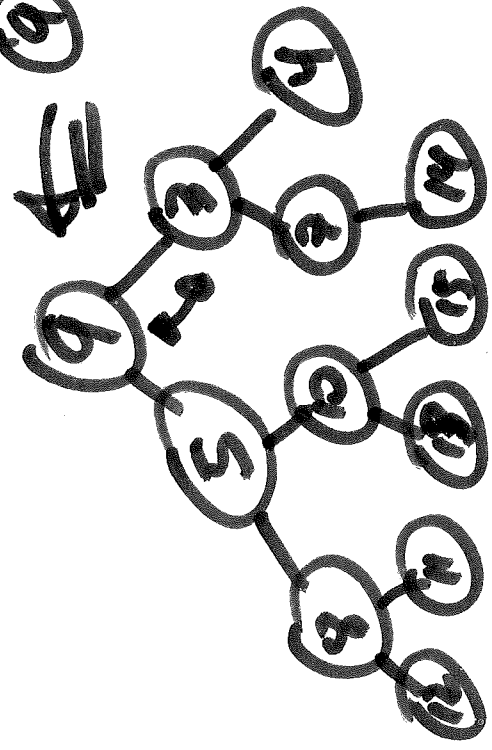
79

val = remove Min()

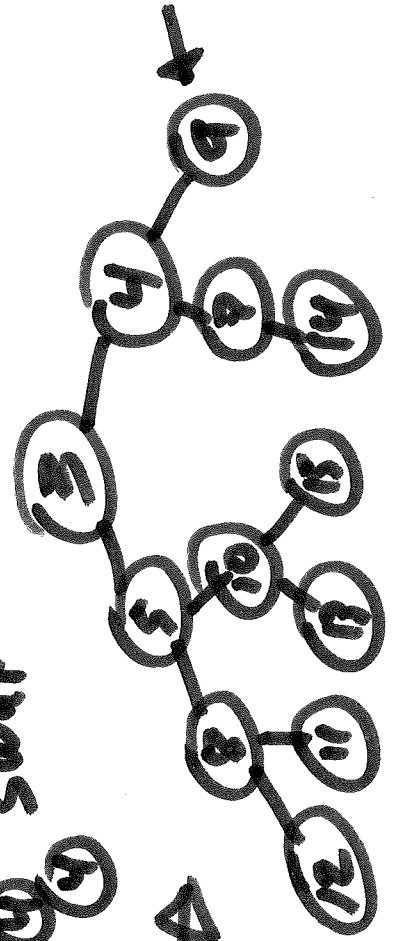
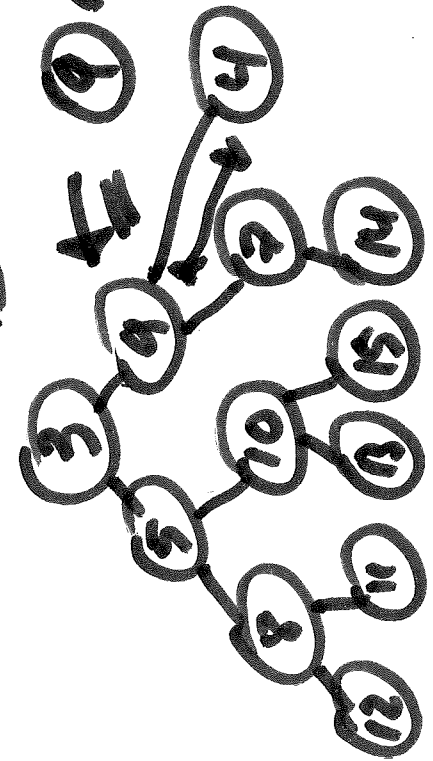
val = 2
place 9 instead of 2

9 is larger than both

5 and 3
Swap the smallest of 5 and 3 with 9



9 is larger than both
5 and 4
Swap with 4

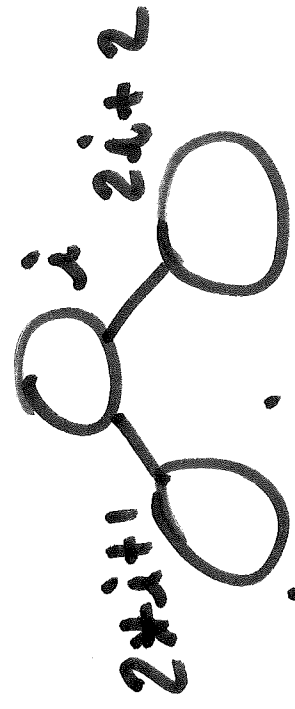
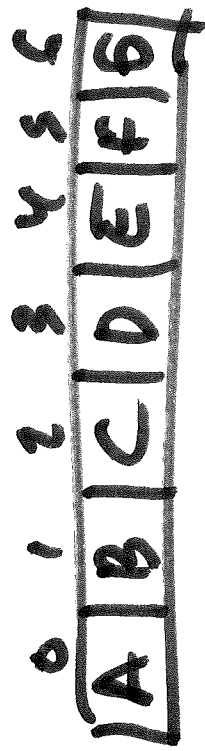
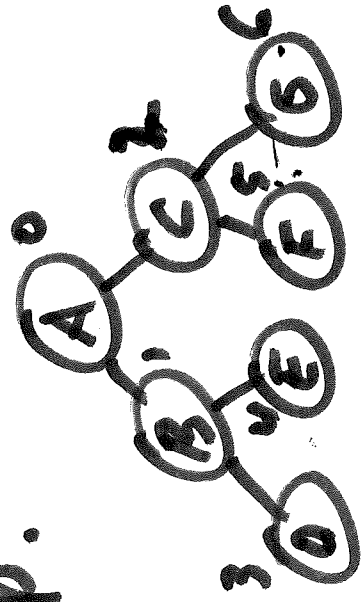


Heap!

Total # of swaps is $h-1$ (80) $h = O(\log n)$

Since h is $O(\log n)$ then the cost of remove Min is $O(\log n)$

We can use the array representation of a binary tree to represent a heap.



$parent(i) = (i-1)/2;$

If a parent is in location i then left child is in $2i+1$ and right child is in $2i+2$.

Implementation

(18)

Heap of type int.

HeapInt.h

```
class HeapInt {
```

```
    int n; // current number of elements in heap.  
    int nmax; // maximum number of elements  
    int *array; // array that stores heap.
```

```
public:
```

```
    HeapInt(int maxsize); // Inserts key.
```

```
    void insert(int key); // returns min key
```

```
    int removeMin(); // returns min key
```

```
};  
~HeapInt();
```

```
#define left-child(i) (2*(i)+1)
```

```
#define right-child(i) (2*(i)+2)
```

```
#define parent(i) ((i)-1)/2
```

HeapInt.cc

```

HeapInt::HeapInt (int maxSize) {
    n = 0;
    nmax = maxSize;
    array = new int [maxSize];
}

```

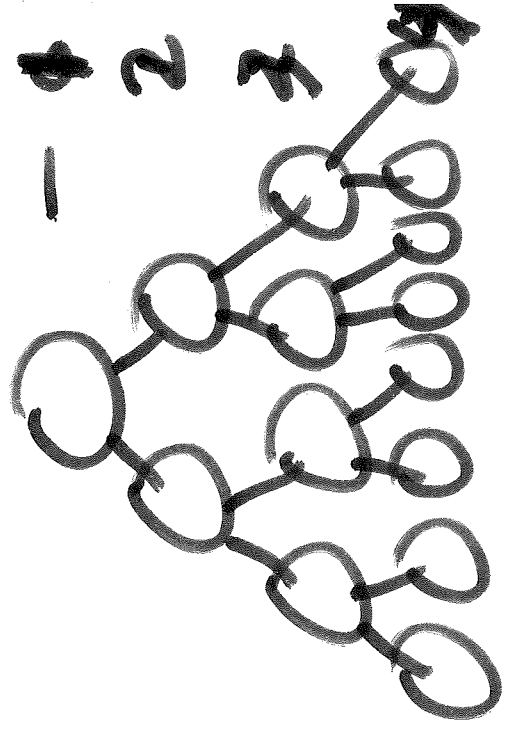
```

HeapInt::~HeapInt() {
    delete array;
}

```


Quiz 8

What is the maximum number of keys in a heap with height h



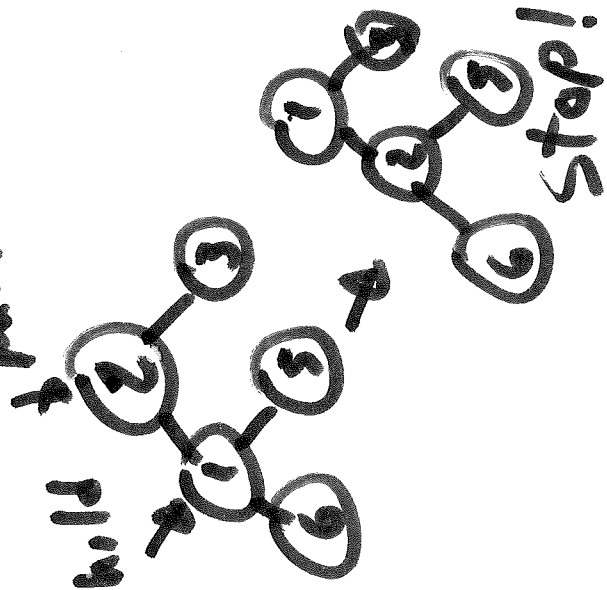
- a) $n_{\max} = h$
- b) $n_{\max} = h^2$
- c) $n_{\max} = 2^h$
- d) $n_{\max} = 2^h - 1$

// we need to swap parent and child

```
int tmp = array[child];  
array[child] = array[parent];  
array[parent] = tmp;  
child = parent; // go up  
parent = &parent[child];
```

} // while

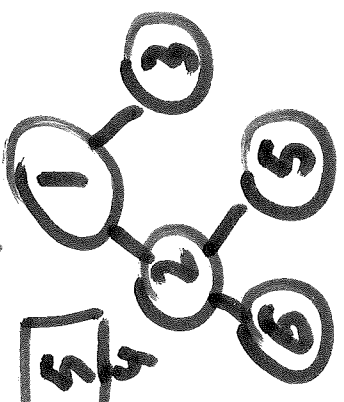
} // insert



98

```
int
HeapInt::removeMin() {
```

// removes and returns the minimum key.



```
assert(n > 0);
```

```
int minkey = array[0]; n =
```

```
n--;
```

```
if (n == 0) {
// heap empty. No need to fix.
return minkey;
}
```

parent

```
// Move last element in heap to top
```



```
array[0] = array[n]; // swap
```

```
// Fix heap (down heap) parent
```

```
// starting at the top
```

```
int parent = 0;
```

```
int left = left_child(parent);
```

```
int right = right_child(parent);
```

88

```
while (left < n) { // repeat until bottom is reached.
```

```
    // Determine smallest child
    int minChild = left;
```

```
    if (right < n && // Make sure there is a right child
```

```
        array[right] < array[left]) {
        minChild = right; // right child's minIndex.
```

```
    // check if we need to swap
    if (array[parent] > array[minChild])
```

```
        // No need to swap.
        break;
```

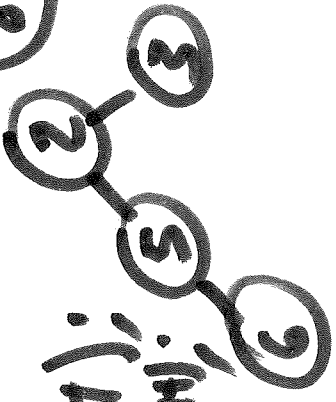
```
    // we need to swap parent and minChild
    int tmp = array[minChild];
    array[minChild] = array[parent];
    array[parent] = tmp;
```

Bounded by $\log(n)$
 n right nodes
 n left nodes
 $2n$ total nodes
 $O(\log n)$

```

// go down
parent = minchild;
left = left-child(parent);
right = right-child(parent);
} // while

```



```

// heap has been fixed
return minkey;
}

```

Heap insert — $O(\log n)$
 Heap removeMin — $O(\log n)$.

Heap Sort

89

- It is a sorting algorithm that uses a Heap.

- It runs in $O(n \log n)$.

```
#include "HeapInt.h"
void heapSort(int *array, int n) {
    // sorts array with n elements
    // in ascending order using heap sort.
```

```
    HeapInt heap(n);
```

```
    for (int i=0; i < n; i++) { //insert elements
        heap.insert(array[i]); //in heap.
```

```
    }
```

```
    // Now put elements back into array from heap
```

```
    // using removeMin
```

```
    for (int i=0; i < n; i++) {
        heap.removeMin();
        array[i] = heap.removeMin();
    }
```

```
}
```

Time Complexity of Heapsort (90)

$$N - \text{heap Inserts} = n O(\log n) \\ O(\log n) = O(n \log n)$$

$$N - \text{remove Min ops.} = n O(\log n) \\ O(\log n) = O(n \log n)$$

$$\text{Total time is } O(n \log n) + O(n \log n) \\ = O(n \log n).$$

worst case!

1b

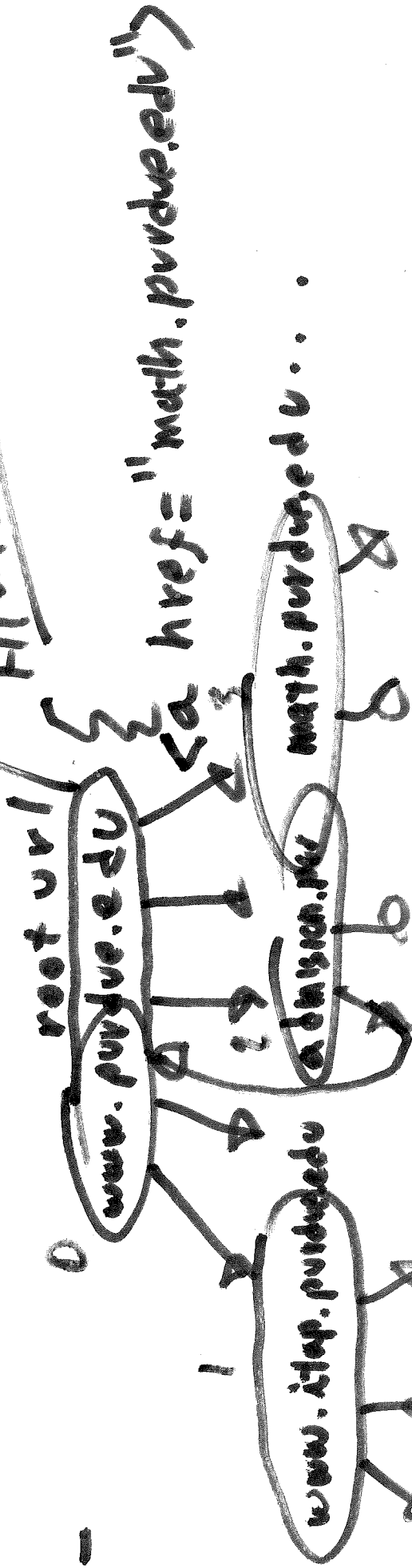
- The problem with the above implementation is that we need extra space in the heap to store array. This implementation is not an "in place" sorting implementation.
- We could modify the implementation to use the same array to store the heap but ~~we~~ will leave it to you to do at home.
- If is possible to make heapsort to be an in-place algorithm.

Lab 2

Web crawler.

(42)

HTML docs



www.purdue

web crawler will find all/max. docs pointed directly or indirectly from root url.

94

Quiz 9 instead of

In heap sort, if the heap

we use an array where insert

inserts at the end and removals

finds the minimum and shifts elements,

what is the complexity of sorting

a) $O(n)$

b) $O(n^2)$

c) $O(n \log n)$

d) $O(n^3)$

Dictionaries

It is an abstract model of a data base. Given a key, it finds the data that corresponds to the key.

Methods:

- Data findElement(key)
- insertElement(key, data)
- removeElement(key).

We can implement a dictionary using:

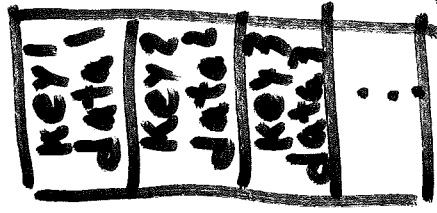
- Hash table
- Array
- Linked list
- AVL Tree
- Binary Search Dictionary ✓

Later on also we will cover

Binary Search Dictionary

98

- It uses an array to store entries with key, data.



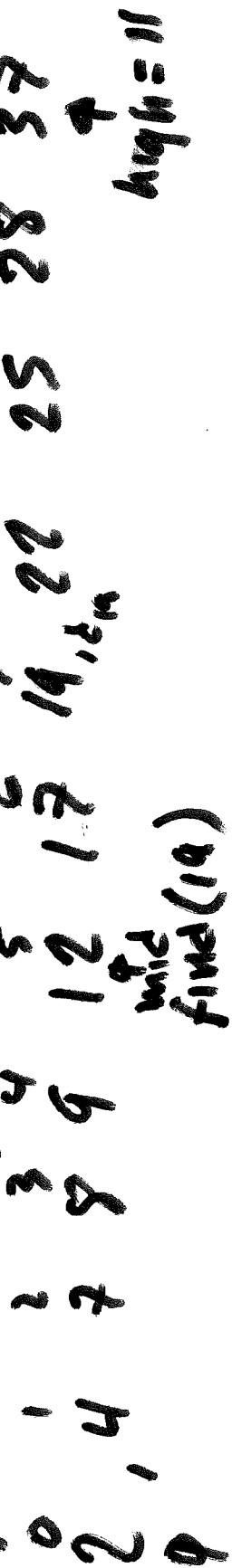
- The array is sorted by the key.

- The lookup uses binary search.

Binary Search

- Narrow down the range of the search by half at every iteration.

- High - low game.



low = 0

mid = $\frac{low + high}{2} = \frac{0 + 11}{2} = 5$

Time Complexity

At every step we divide the range by half

iteration range

- 0 n
- 1 n/2
- 2 n/4
- ...
- n

we stop when range = 1 $\rightarrow \frac{n}{2^x} = 1$

low = 6

high = 11

mid = $\frac{6 + 11}{2} = \frac{17}{2} = 8$

mid = 8

high = 7

low = 6

mid = $\frac{6 + 7}{2} = 6$

low = 7 high = 7

mid = $\frac{7 + 7}{2} = 7$

Found!

88

$$\frac{n}{2^i} = 1$$

$$n = 2^i$$

$$\log_2 n = \underline{\underline{i}}$$

$$\# \text{ iterations} = O(\log n)$$

Binary Search Find takes $O(\log n)$

Quiz 10

~~99~~ 99

What is the time complexity of `find(key)` using binary search

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

Implementation of Binary Search Table

100

BinSearchTable.h

```
struct BinSearchEnt {  
    char *key;  
    void *data;  
};
```

```
class BinSearchTable {  
    BinSearchEnt *array;
```

```
    int n;
```

```
    int nmax; // sorted;
```

```
public BinSearchTable();
```

```
    void insert(char *key, void *data);
```

```
    void *lookup(char *key);
```

```
    ~BinSearchTable();
```

```
}; // It can use any sorting algorithm.
```

```
BinSearchEnt() {  
    key = NULL;  
    data = NULL;  
};
```

```
BinSearchEnt(char *key, void *data) {  
    this->key = strdup(key);  
    this->data = data;  
};
```

```
BinSearchEnt() {  
    ~ BinSearchEnt();  
    if (key != NULL) {  
        free(key);  
    }  
};
```

```
BinSearchEnt() {  
    ~ BinSearchEnt();  
};
```

Not needed.

BinSearchTable.cc

(101)

```
BinSearchTable::BinSearchTable() {
```

```
    n = 0;
```

```
    nmax = 100;
```

```
    array = new BinSearchEnt [nmax];  
    isSorted = false;
```

```
}
```

```
void
```

```
BinSearchTable::insert(char *key, void *data) {
```

```
    // check if there is space.
```

```
    if (n == nmax) { // enlarge table
```

```
        BinSearchEnt * newarray;
```

```
        nmax = 2 * nmax;
```

```
        newarray = new BinSearchEnt [nmax];
```

```
        for (int i = 0; i < n; i++) {
```

```
            newarray [i] = array [i];
```

```
        }
```

```
delete [] array;
array = new array;
```

{ //if

```
// There is space in array
```

```
array[n].key = strdup(key);
```

```
array[n].data = data;
```

```
n++;
```

```
isSorted = false; // we no longer know if array is sorted.
```

}

103

void * p10
Bin Search Table :: lookup(char *key) {

Make sure table is sorted.
if (!isSorted) {
sortTable();
isSorted = true;
}

// Use binary search

int low = 0;

int high = n-1;

while (high >= low) {

int mid = (high + low) / 2;

int res = strcmp(key, array[mid].key);
"kiwi" "banana" > 0

if (res > 0) {

low = mid + 1;

else if (res < 0) {

high = mid - 1;

return NULL; // key not found
} else { // key found
return array[mid].data;
}

- 0 apple
- 1 banana
- 2 kiwi
- 3 orange

lookup("kiwi")

low = 0 high = 3
mid = (0+3)/2 = 1

Time Complexity

104

Bin Search Table insert

worst $O(n)$
case

best $O(1)$ #expansion of list
case

average case (amortized) = $O(n)$
cost of 1 insertion \times n insertions $\approx D(\log n)$
amort

lookup

worst $O(\log n)$
case

(assuming table is sorted)

worst case table is unsorted $O(n \log n)$
case

105

However binary search tables are used in cases where tables are "Read Only", that is we create the table once with n insertions and then we do m lookups where $m \gg n$.

Cost

n { insert } n insertions = $n O(\log n)$
 m { lookup } m lookups = $m O(\log n)$
 m { lookup } m lookups = $m O(\log n)$

Total: $O(n \log n) + O(m \log n)$
 $O(m \log n)$
 $= O(n \log n) + O(m \log n)$

if $m \gg n$ then ~~Total~~ Total = $O(m \log n)$

amortized case

Quiz 11 =

108

In a binary search table that

has resizing. If we do a lookup

after every insertion what will be the

cost ~~if we do not do any lookup for~~ ~~insertions and~~ lookups.

- insert
- lookup
- insert
- lookup
-
-

n

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

Todo:
→ destructor in table
→ Time complexity.

Midterm Review

107

- Math Review
- Types of Running time
 - worst case
 - best case
 - Average case
- Asymptotic Analysis and big-O notation $O(g(n))$ if $f(n) \leq c \cdot g(n)$
- Stacks (arrays) ^{push} pop LIFO \leftarrow enqueue
- Queues (arrays) \rightarrow (circular queue) \leftarrow dequeue FIFO
- Double linked lists.
- Arrays vs linked lists
- Hash Tables

	Avg	Worst
insert	$O(1)$	$O(n)$
lookup	$O(1)$	$O(n)$

~~extra~~ Requirements or Properties that make insert, lookup Avg time $O(1)$

1. — Hash function has to be Evenly Distributed
2. — Number of buckets $\geq N$

- Implementation

- Resizable array
 - Double size in every expansion
- Cost of insertion:

Best Case	Worst Case	Avg Case (amortized)
$O(1)$	$O(n)$	$O(\log n)$
- Implementation

- Trees

- Traversal
 - preorder
 - postorder
 - Inorder

- Tree representation

- Pointers left, right
- Array representation



```

struct TEntry {
  int val;
  TEntry * left;
  TEntry * right;
}
  
```

}

109

- Heap data structure

+ Implementation of priority queue

+ Two operations

- insert

- remove Min

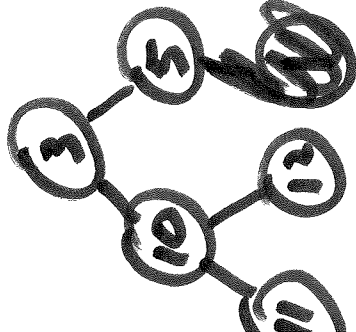
+ How to perform insert, remove Min in heap

+ Implementation using an array

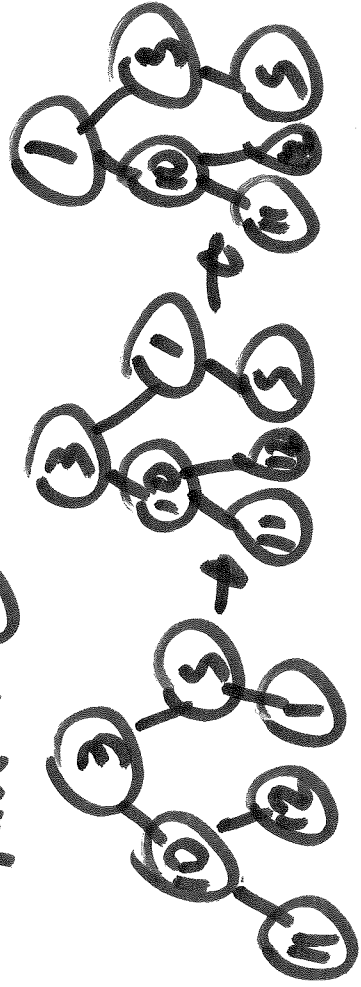
+ Analysis worst case

insert $O(\log n)$

remove Min $O(\log n)$



insert 1



- Heapsort

(110)

worst case - $O(n \log n)$
implementation

- Binary Search Dictionary

- <u>Implementation</u>	best case	worst case	Avg. n lks.
- lookup	$O(\log n)$	$O(n \log n)$	$O(n \log n)$
- insertion	$O(1)$	$O(n)$	$O(\log n)$

- Study templates.

To Study

- Class Notes (1)
- Book (4)
- Labs (3)
- Homework (2)