# Data Structures Review

CS 225
Brad Solomon

December 9, 2024
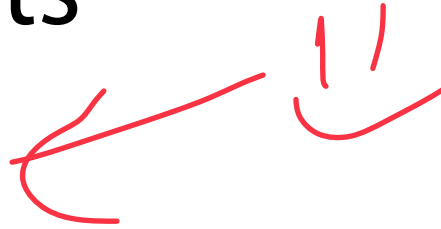
University of Illinois Urbana-Champaign
Department of Computer Science

:)
You got this!

# Announcements

Fill out ICES forms!

Interested in being a CA? Apply for CS 225 or CS 277!

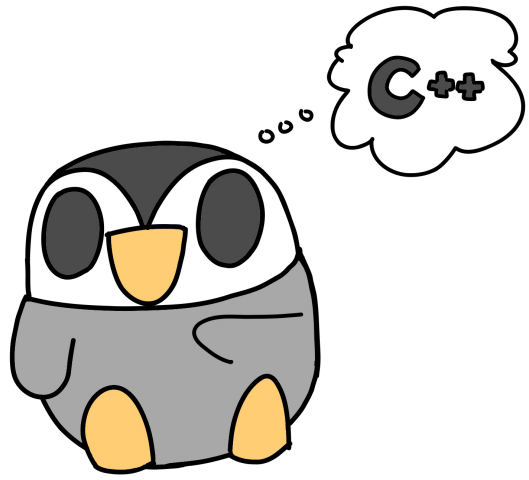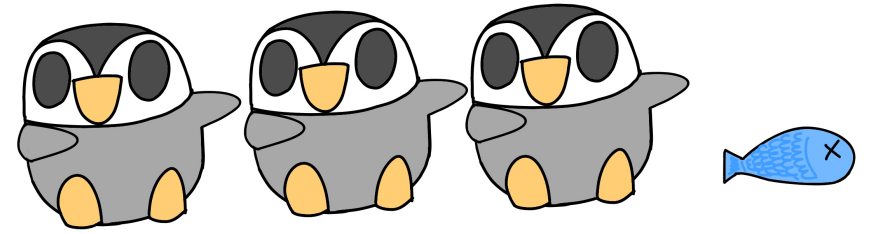https://opportunities.cs.illinois.edu/courses/positions/

# Material covered here is not only material in class!

Represents only an attempt to provide some helpful resources.

# Lists

# List Implementation  September 9 (Array List Lecture)

| | Singly Linked List | Array |
|---|---|---|
| Look up **arbitrary** location ↳ index | O(n) | O(1) ⌣ |
| Insert after **given** element ↳ ref pointer | O(1) ⌣ | O(n) |
| Remove after **given** element | O(1) ⌣ | O(n) |
| Insert at **arbitrary** location | Find O(n) change O(1)    O(n) | Find O(1) change O(n)    O(n) |
| Remove at **arbitrary** location | O(n) | O(n) |
| Search for an input **value** | O(n) | O(n) |

Special cases    insert/remove front O(1)    if array got full insert/remove is O(1)*

always amortized

# Lists

*The not-so-secret underlying implementation for many things*

|  | Singly Linked List | Array |
|---|---|---|
| Look up **arbitrary** location | O(n) | O(1) |
| Insert after **given** element | O(1) | O(n) |
| Remove after **given** element | O(1) | O(n) |
| Insert at **arbitrary** location | O(n) | O(n) |
| Remove at **arbitrary** location | O(n) | O(n) |
| Search for an input **value** | O(n) | O(n) |

## Special Cases:

Insert Front  O(1)     Insert Back  O(1)*

# Stack and Queue

*Taking advantage of special cases in lists / arrays*

Queue

Insert / remove Back

O(1) access

O(n) change value

O(1)*



Front

Don't move items!

O(1)*

Back

O(1)*

Front          Back

head   O(1)                              tail   O(1)



| C | | S | | 2 | | 2 | | 5 | | Ø |

# Stack ADT

- [Order]: Last in first out (LIFO)

- [Implementation]: Trivially as vector or LL ← student Q

Front is top

- [Runtime]: $O(1)$ *

\* if array is not full
if array is full, amortized still says $O(1)$

# Stack ADT

- [Order]: LIFO

- [Implementation]: Array (such as std::vector)

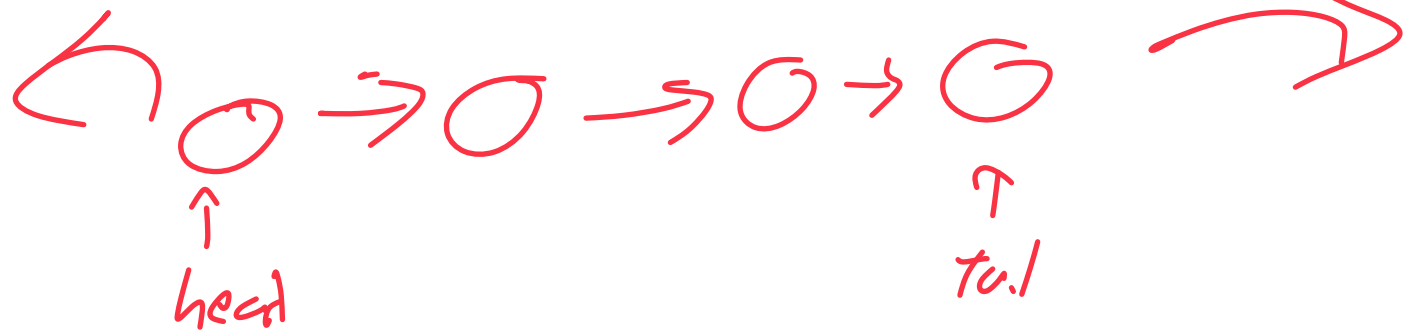- [Runtime]: O(1) Push and Pop

# Queue ADT

- [Order]: First in First out (FIFO)

- [Implementation]: Vector / dequeue → LL is possible easily

- [Runtime]: $O(1)$ *

# Queue ADT

- [Order]: FIFO

- [Implementation]: Circular Queue as Array

- [Runtime]: O(1)

# Iterators

The actual iterator is defined as a class **inside** the outer class:

1. It must be of base class `std::iterator`

2. It must implement at least the following operations:

`Iterator& operator ++()` — move to next item

`const T & operator *()` — return the data/value at current pos

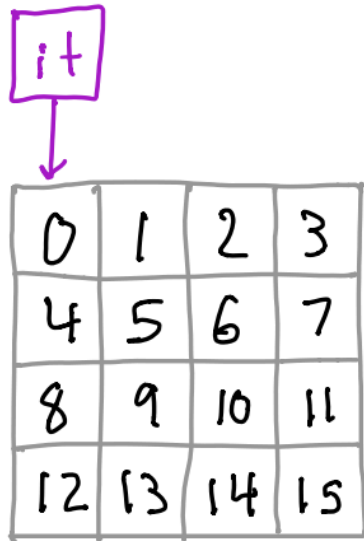`bool operator !=(const Iterator &)` — check if iterators are equal

# Iterators

Here is a (truncated) example of an iterator:

```cpp
template <class T>
class List {

    class ListIterator : public
std::iterator<std::bidirectional_iterator_tag, T> {
      public:

        ListIterator& operator++();

        ListIterator& operator--()

        bool operator!=(const ListIterator& rhs);

        const T& operator*();
    };

    ListIterator begin() const;

    ListIterator end() const;
};
```
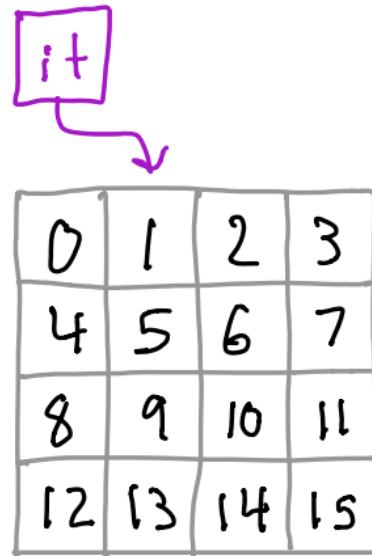
```cpp
std::vector<Animal> zoo;


/* Full text snippet */

  for ( std::vector<Animal>::iterator it = zoo.begin(); it != zoo.end(); ++it ) {
    std::cout << (*it).name << " " << (*it).food << std::endl;
  }


/* Auto Snippet */

  for ( auto it = zoo.begin(); it != zoo.end; ++it ) {
    std::cout << (*it).name << " " << (*it).food << std::endl;
  }

/* For Each Snippet */

  for ( const Animal & animal : zoo ) {
    std::cout << animal.name << " " << animal.food << std::endl;
  }
```

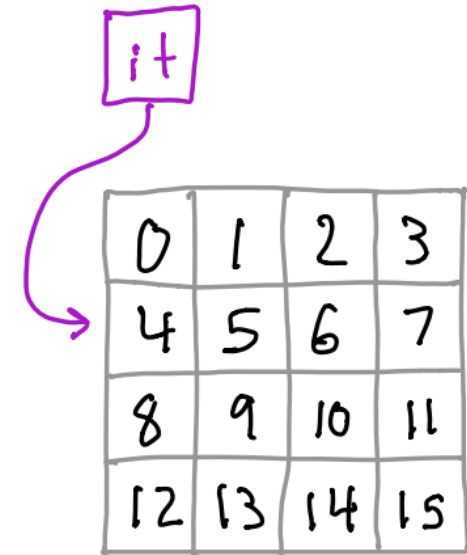# Iterators (225 Webpage Resources)



https://courses.grainger.illinois.edu/cs225/fa2024/resources/iterators/

# Trees

↳ Binary tree!
↳ complete / full / perfect

# Tree Traversals

**Pre-order:** 1 2 3 5 4 6 11 8 7 10 9

root

**In-order:** 3 2 4 5 6 1 8 7 11 9 10

← root

Left then right
Left most child

**Post-order:** 3 4 6 5 2 ⊗ 7 8 9 10 11

root

1

# Tree Traversals



**Pre-order:** 1, 2, 3, 5, 4, 6, 11, 8, 7, 10, 9

**In-order:** 3, 2, 4, 5, 6, 1, 8, 7, 11, 9, 10

**Post-order:** 3, 4, 6, 5, 2, 7, 8, 9, 10, 11, 1

# Depth First Search <span style="color:red">September 20 (BST Lecture)</span>

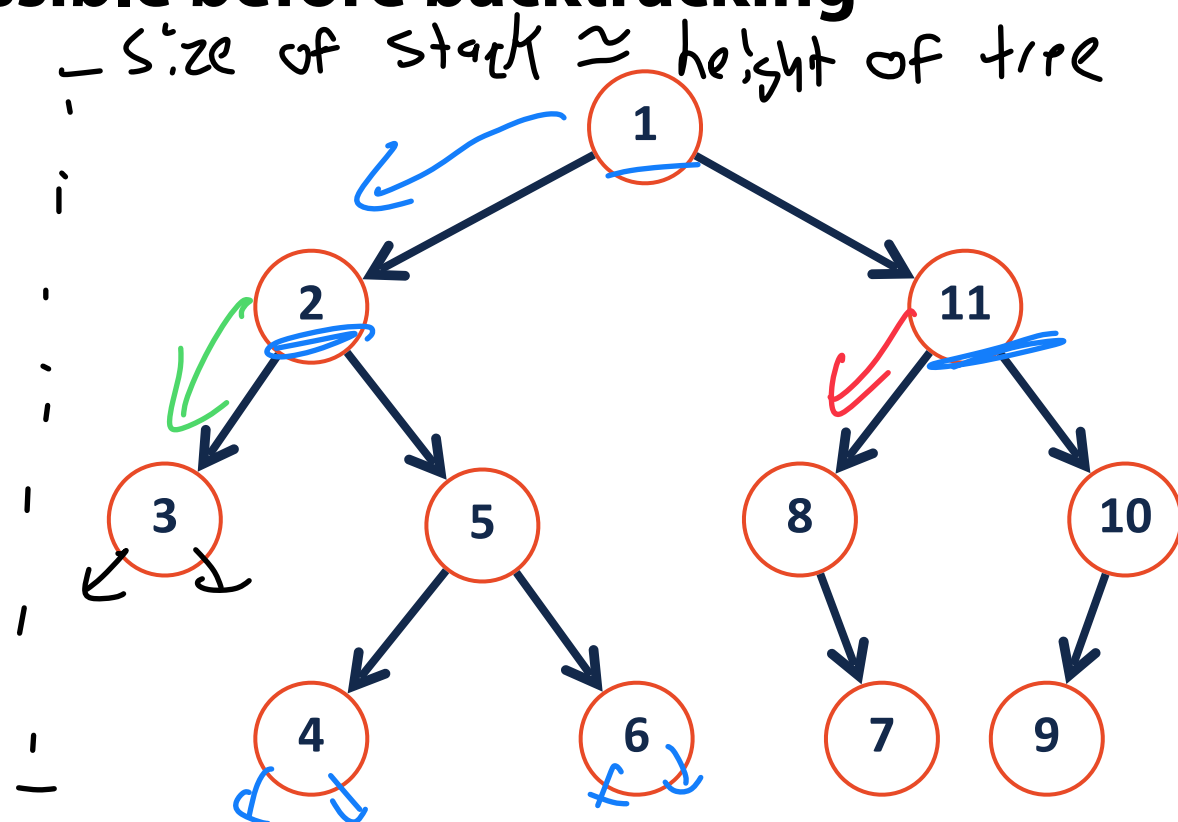**Explore as far along one path as possible before backtracking**

Make a stack, initialize to root — Size of stack ≃ height of tree

While stack not empty
   Pop the top element (as tmp)
   Print tmp
   Push tmp → right
   Push tmp → left



Stack: 1 , 11 , 2 , 5 , 3 , 6,4 , 10, 8 7 9 .

Print: 1 , 2 , 3 , 5 , 4 , 6 , 11 , 8 , 7 , 10 , 9

# Depth First Search

**Explore as far along one path as possible before backtracking**

Make a stack initialized with root
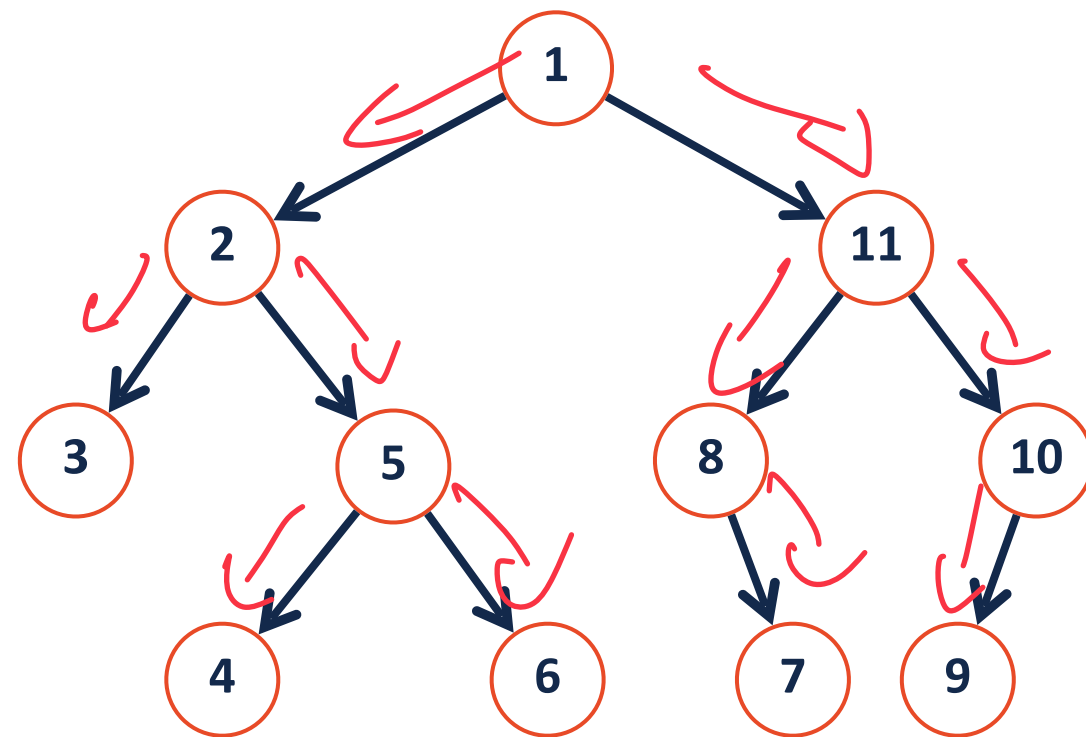
While stack isn't empty:

Pop top element (as `tmp`)

Print `tmp`

Push `tmp->right` to stack

Push `tmp->left` to stack

LIFO

Stack: 1, 11, 2, 5, 3, 6, 4, 10, 8, 7, 9

Print: 1, 2, 3, 5, 4, 6, 11, 8, 7, 10, 9 ← Pre order!

# Breadth First Search

**Fully explore depth i before exploring depth i+1**

Make a queue initialized with root

While queue isn't empty:

Dequeue front element (as tmp)

Print tmp

Enqueue tmp->left

Enqueue tmp->right

FIFO

Size of queue ≈ width of tree
height ≈ width
1
2    2
3    8    4
4    16

equal to width

Queue: 1, 2, 11, 3, 5, 8, 10, 4, 6, 7, 9

Print: 1, 2, 11, 3, 5, 8, 10, 4, 6, 7, 9

# Breadth First Search

**Fully explore depth i before exploring depth i+1**

Make a queue initialized with root
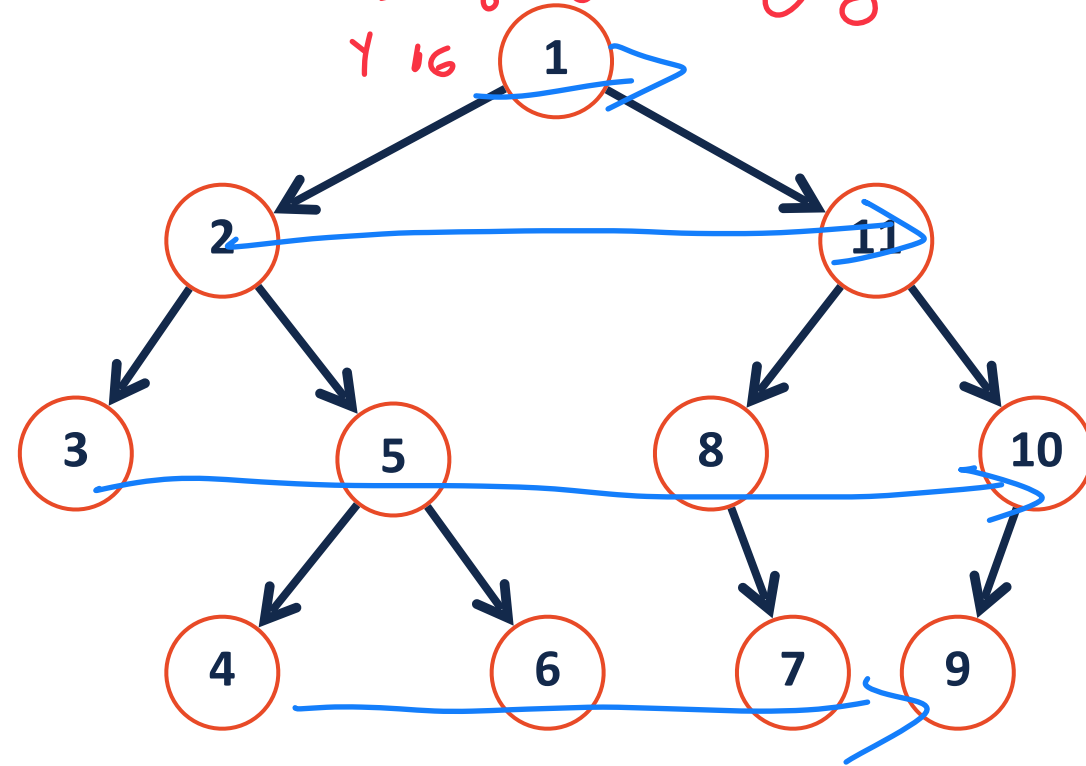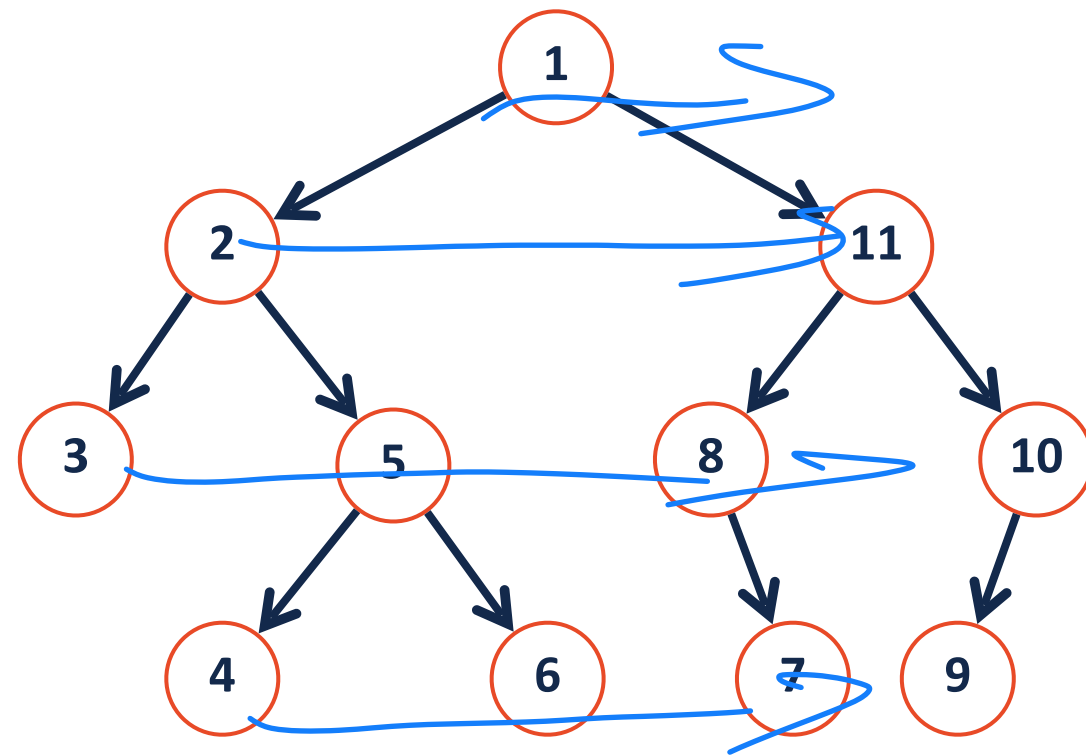
While queue isn't empty:

Dequeue front element (as `tmp`)

Print `tmp`

Enqueue `tmp->left`

Enqueue `tmp->right`

Queue: 1, 2, 11, 3, 5, 8, 10, 4, 6, 7, 9

Print: 1, 2, 3, 5, 4, 6, 11, 8, 7, 10, 9

# BST Find

root → nullptr ( )

root → (66)



find(66)

Start @ root

Recursive Problem!

Base Case:
↳ If tree empty, return root
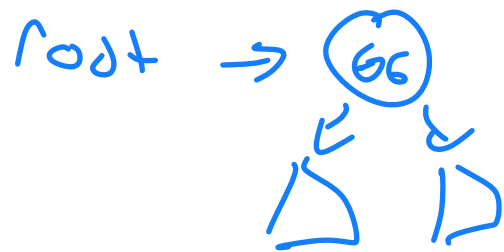↳ If root is query, return root

Recursive Step:
Compare root→key w/ query
$\underbrace{\qquad}_{tmp}$

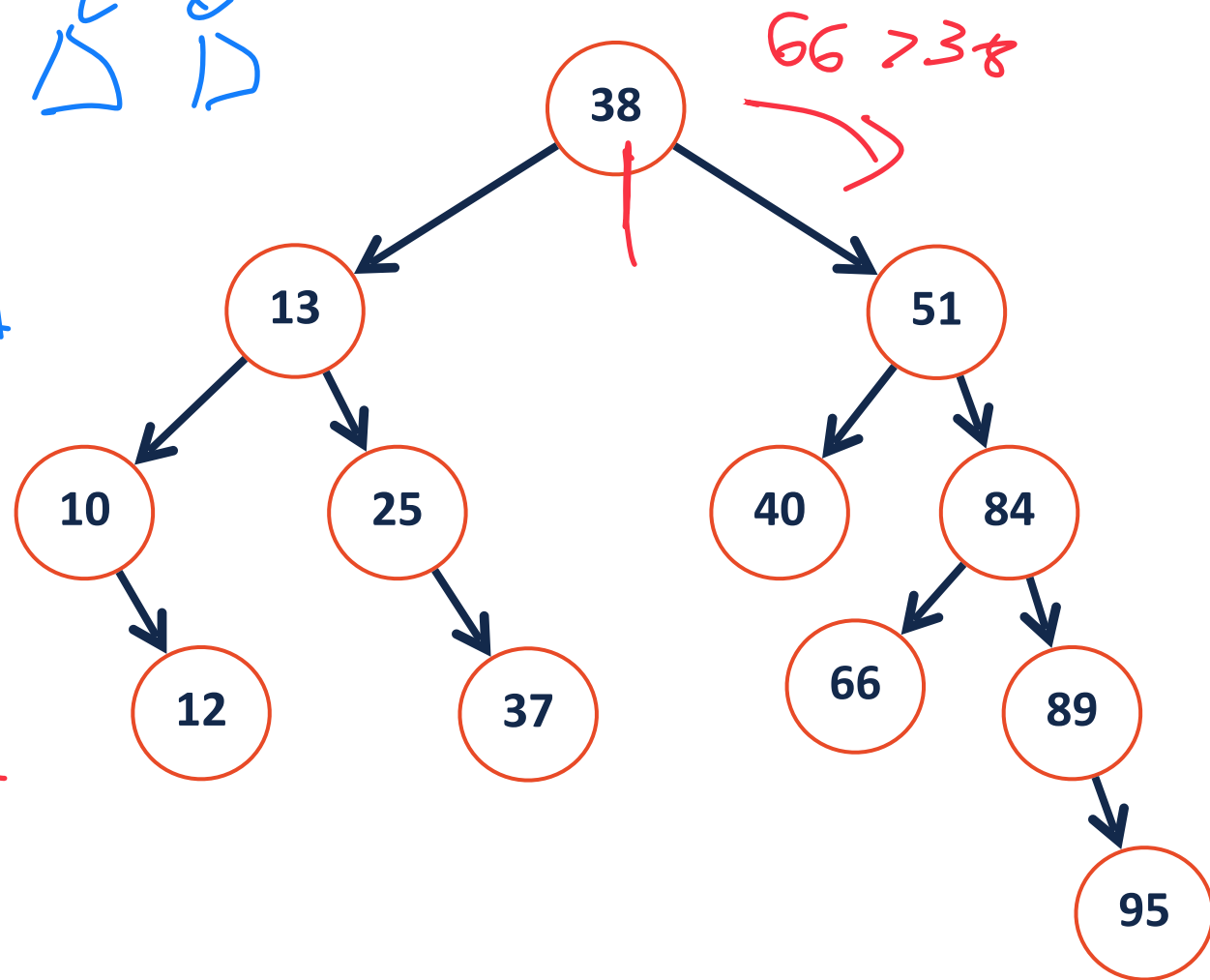if
tmp > query, recurse right
tmp < query, recurse left
==

66 > 38

# BST Find

**find(66)**

**A recursive function based around value of root:**

**Base Case:** If root is `null`, return root
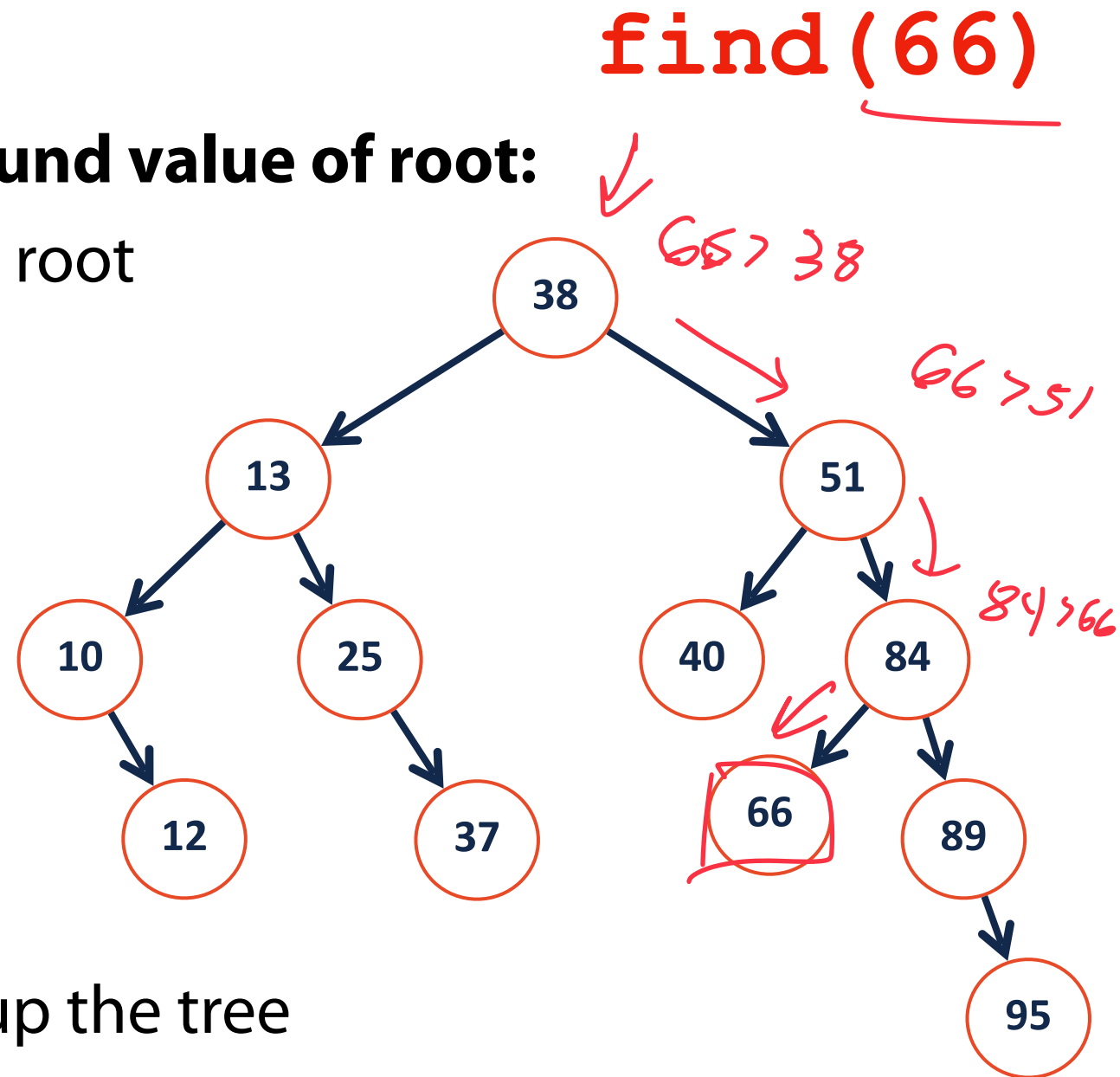
Let `tmp = root->key()`

`tmp == query`, return root

**Recursion:**

`tmp <   query`, recurse right

`tmp >   query`, recurse left

**Combining:**
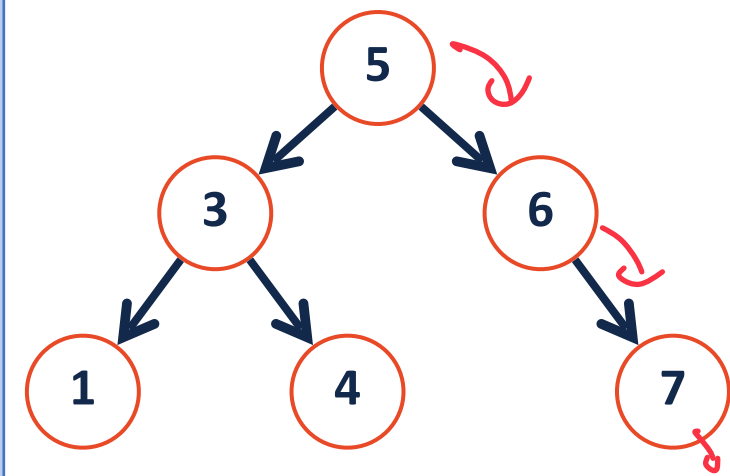
Return the recursive value back up the tree

```
template<typename K, typename V>

_____TreeNode *&_____ _find(TreeNode *& root, const K & key) {


// Base Case
if(root == nullptr || root->key == key){
    return root;
}


// Recursive Step ("Combining step" is 'return')
if (root->key > key){
    return _find(root->left, key);
}

return _find(root->right, key);



}
```

*(annotations)*
query
No const here
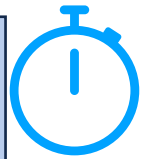No const here
Find(8) returns pointer by ref (7->right)
Not nullptr!
"else"
Smaller go left
larger go right

5
3    6
1    4    7
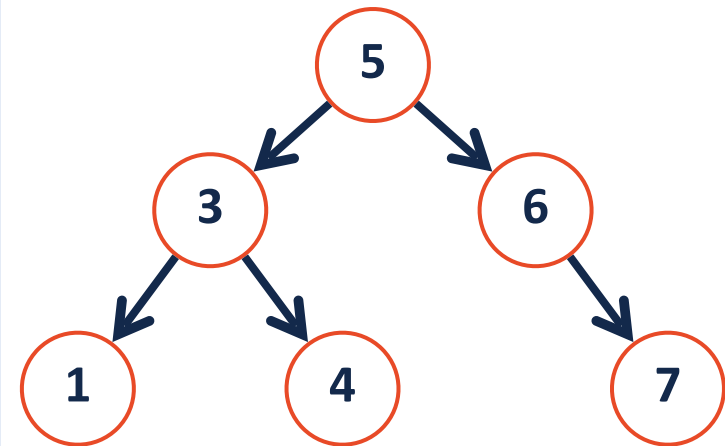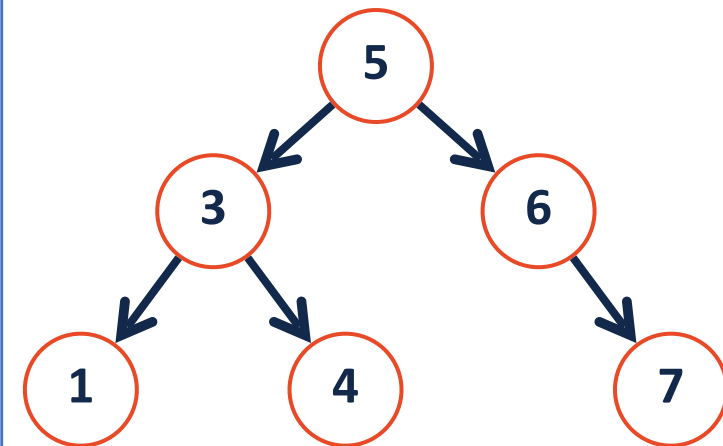
```
1  template<typename K, typename V>
2
3  void _insert(const K & key, const V & val) {
4
5      return _insert(root, key, val);
6  }
7
```

```
1  template<typename K, typename V>
2
3  void _insert(TreeNode *& root, const K & key, const V & val) {
4
5  TreeNode *& tmp = _find(root, key);
6
7
8  tmp = new treeNode(key, val);
9
10
11
12
13 }
14
15
16
```
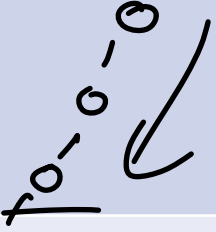
find is key!

```
template<typename K, typename V>

void _remove(TreeNode *& root, const K & key) {
```

This weeks lab!

0 - child

1 - child

2 - child → find IOP/IOS
  ↳ swap IOP w/ target
  ↳ remove target

```
}
```

# BST Analysis – Running Time

| Operation | BST Worst Case |
|-----------|----------------|
| find | $O(h)$ |
| insert | $O(h)$  b/c |
| remove | $O(h)$ + $O(h)$ + $O(h)$ = $O(h)$ <br> find    find (IOP)    remove() |
| traverse | $O(n)$ |

# BST Analysis – Running Time

| Operation | BST Worst Case |
|---|---|
| find | O(h) = O(n) |
| insert | O(h) = O(n) |
| remove | O(h) = O(n) |
| traverse | O(n) |

# Left Rotation

1) Create a tmp pointer to root

2) Update root to point to mid

3) tmp->right = root->left

4) root->left = tmp

# AVL Rotations

Four kinds of rotations: (L, R, LR, RL)

1. All rotations are local (subtrees are not impacted)

2. The running time of rotations are constant

3. The rotations maintain BST property

**Goal:** AVL tree will be balanced
↳ This will make height bounded by log(n)

# AVL Tree Analysis

For an AVL tree of height h:

Find runs in: _____ $O(h)$ _____.

Insert runs in: _____ $O(h)$ _____.

Remove runs in: _____ $O(h)$ _____.

**Guarantee:**

1) Tree is balanced

**Claim:** The height of the AVL tree with n nodes is: _____ $O(\log n)$ _____.

# Nearest Neighbor: k-d tree

A **k-d tree** is similar but splits on points:

$$(7,2), \ (5,4), \ (9,6), \ (4,7), \ (2,3), \ (8,1), \ (9,8)$$

Median of all items in X dim

$(7,2)$

X

$(9,6)$

$(5,4)$

med of all items $x < 7$
in the $y$ dim

# Nearest Neighbor: k-d tree

Search by comparing query and node in single **alternating** dimension

# Nearest Neighbor: k-d tree

**Backtracking:** start recursing backwards -- store "best" possibility as you trace back

(2,3) or (5,4) better nearest point?

query = (6,3)
cur best = (5,4)

(7,2)

(5,4)          (9,6)

(2,3)    (4,7)  (8,1)    (9,8)

↳ Defined first search radius

# Nearest Neighbor: k-d tree

May have to recursively check other branches of tree — **why?**

# Nearest Neighbor: k-d tree

# BTree Properties

A **BTrees** of order **m** is an m-ary tree and by definition:

- All keys within a node are ordered

- All nodes contain no more than **m-1** keys.

- All internal nodes have exactly **one more child than keys**

> 0 children

Root nodes can be a leaf or have $[2, m]$ children.

All non-root, internal nodes have $[\lceil \frac{m}{2} \rceil, m]$ children.

If $\frac{\text{int}(\frac{m}{2})}{+1}$ is keys is children

All leaves in the tree are at the same level.

# BTree Find

Note edge case if no larger value   **Find(7)**

Find(43)

Base Case:

If root is empty, return

If leaf, do array find() and return

Recursive Step:

Array find() for match or first greater value

Recurse on appropriate child *

**Tip:** Index of first greater value is index of child we want to visit!

7

Size=1

23

-3

42 | 55

-11 | 8

25 | 31

43

60

1

B is 0?

O(m) x Nodes
↳ O(log n)

this is constant

# BTree Insertion

When we hit **M** items, split into three nodes!

1) Find median

2) "Raise median up"

⤷ Cut array in half
as 2 new BTree Nodes

BTree Node

3

BTree Node

BTree Node

| 1 | 2 | ~~3~~ | 4 | 5 |
|---|---|---|---|---|

3

1 2

4 5

M-1 items

MAX

Insert(1)
Insert(2)
Insert(3)
Insert(4)
Insert(5)
**Insert(6)**
**Insert(7)**
**Insert(8)**

# BTree Recursive Insert

Insert always starts at a leaf but can propagate up repeatedly.

# Final thoughts on Trees

Trees have a large space of **possible coding questions**

We hit **tree iterators** multiple times…

You saw **tree constructors of unusual shapes**…

$[0,1,2,3,4,5]$

Even indices    left
Odd indices    right

$[2,4]$

$[1,3,5]$

# Heap

*Taking advantage of special cases in lists / arrays*

## Array List (Pointer implementation)

insert Back
O(1)*

T* Start

O(1) lookup

Swap

T* Size

T* Capacity

| | 4 | 5 | 6 | 15 | 9 | 7 | 20 | 16 | 25 | 14 | 12 | 11 | | | | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| 0 | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9 | 10 | 11 | 12 | 13 | 14 | 15 | 16 |

size_t Start

size_t Size

size_t Capacity

## Array List (Index implementation)

# (min)Heap (Priority Queue)

By storing as a complete tree, can avoid using pointers at all!

**If index starts at 1:**

`leftChild(i): 2i`

`rightChild(i): 2i+1`

`parent(i): floor(i/2)`

$2^{h+1} - 1$

→ = extra space

min value
is parent
of children



| | 4 | 5 | 6 | 15 | 9 | 7 | 20 |
|---|---|---|---|---|---|---|---|
| 0 | 1 | 2 | 3 | 4 | 5 | 6 | 7 |

# insert

O(1)*

1) Insert at end of array

2) Check if minHeap still valid

3) Swap with parent if needed

O(1)

**Steps 2 and 3 are recursive!**



$P(i) = i/2$

$i = 7$

$i$

# removeMin

1) Swap root w/ last item
   ↳ Delete last item
   ↳ Size --;

2) heapify Down ()
   ↳ Repeated swaps w/ min child
   until leaf of smaller than both children

# removeMin

1) Swap root with last item

    (and remove)

    (and modify size)

2) HeapifyDown( ) root

# Disjoint Sets

# Disjoint Set Implementation

*Taking advantage of array lookup operations*

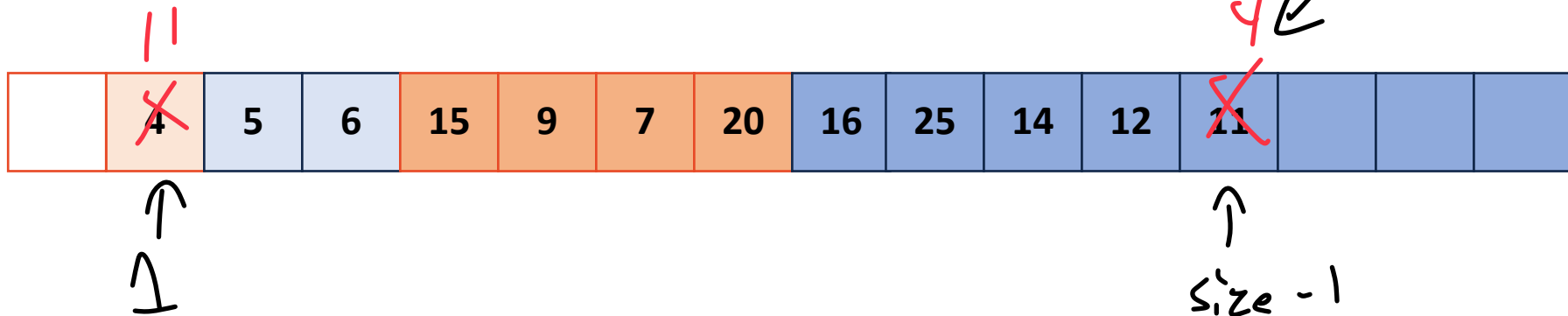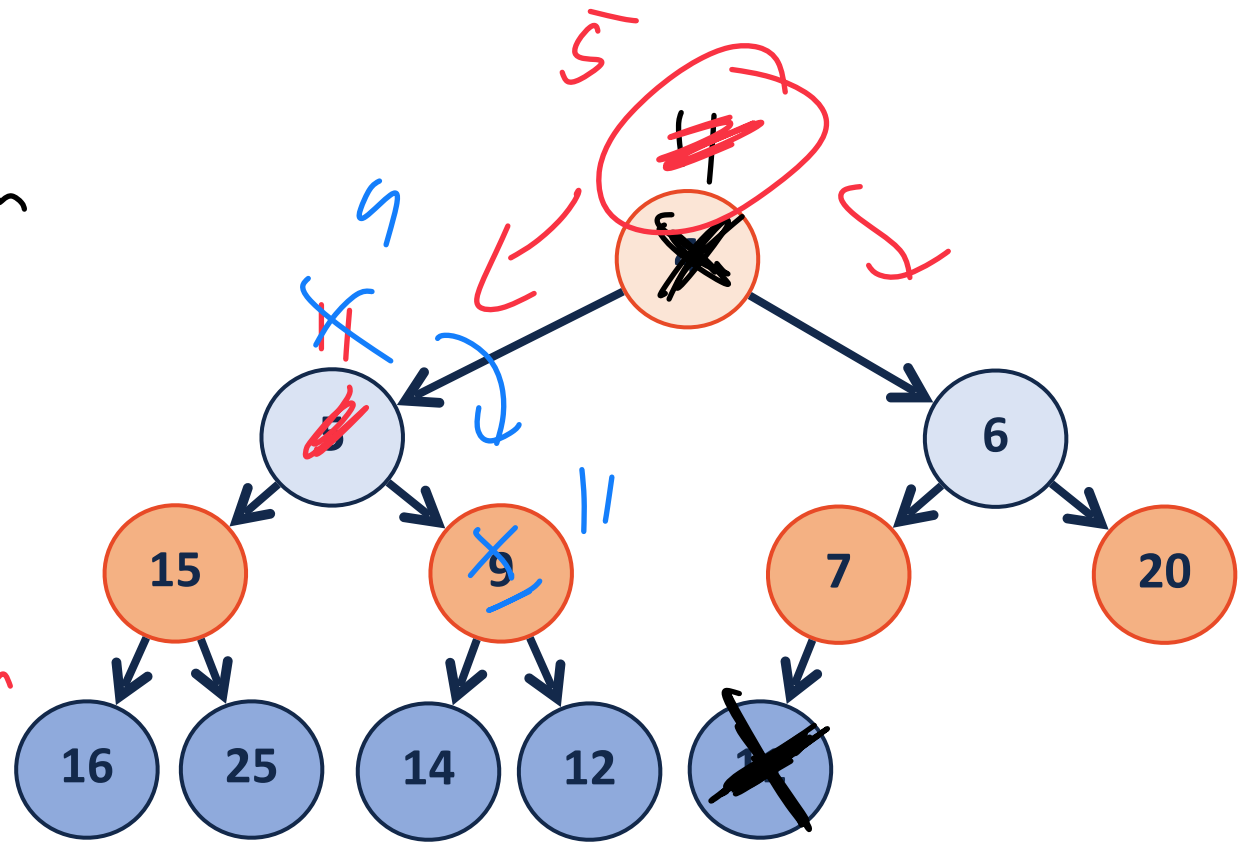Store an UpTree as an array, canonical items store **height** / **size**

| | 0 | 1 | 2 | 3 | 4 | 5 | 6 | 7 |
|---|---|---|---|---|---|---|---|---|
| | -2 | 0 | -2 | -2 | 0 | 3 | 3 | 2 |
| | -3 | | -2 | -3 | | | | |

Sets: 0 1 4    2 7    3 5 6

**Find(k):** Repeatedly look up values until **negative value**

**Union(k₁, k₂):** Update *smaller* canonical item to point to larger

Update value of remaining canonical item

$O(1)$

# Disjoint Sets – Smart Union

Two O(1) methods of combining two sets

Claim: Both limit height to: O(log n).

**Before Union**

**After Union**

**Union by height**

| 4 | ... | 7 |
|---|-----|---|
| -4 | | -3 |

| 4 | ... | 7 |
|---|-----|---|
| -4 | | 4 |

***Idea*: *Keep the height of the tree as small as possible.***

**Union by size**

| 4 | ... | 7 |
|---|-----|---|
| -4 | | -8 |

| 4 | ... | 7 |
|---|-----|---|
| 7 | | -12 |

***Idea*: *Minimize the number of nodes that increase in height***

# Disjoint Sets Path Compression

*Minimizing number of O(1) operations*

$O(1)$ kinda sorta

```
1  int DisjointSets::find(int i) {
2    if ( s[i] < 0 ) { return i; }
3    else {
4      int root = find( s[i] );
5      s[i] = root;
6      return root;
7    }
8  }
```

Taking advantage of arrays

# Graphs

# Graph Implementation: Edge List $|V| = n, |E| = m$

*The equivalent of an 'unordered' data structure*



**Vertex Storage:**

An optional list of vertices

**Edge Storage:**

A list storing edges as (V1, V2, Weight)

**Most graphs are stored as just an edge list!**

# Graph Implementation: Adjacency Matrix

$|V| = n, |E| = m$



## Vertex Storage:

A hash table of vertices

Implicitly or explicitly store index

|   |   |
|---|---|
| u | 0 |
| v | 1 |
| w | 2 |
| z | 3 |

|   | 0 | 1 | 2 | 3 |
|---|---|---|---|---|
| 0 | - | a | c | 0 |
| 1 |   | - | b | 0 |
| 2 |   |   | - | d |
| 3 |   |   |   | - |

## Edge Storage:

A $|V|$ x $|V|$ matrix of edges

Weight is stored at position (u, v)

# Adjacency List



**Vertex Storage:**

A bidirectional linked list with size variable

Each node is a pointer to edge in edge list

**Edge Storage:**

A list of (v1, v2, weight) edges

Also store pointers back to nodes

# Adjacency List

$|V| = n, |E| = m$



## Adj List Node:

| *Prev | *Edge | *Next |
|-------|-------|-------|

## Edge List:

| V1 | V2 | Weight |
|----|----|--------|
| *V1 | *V2 | |

$$|V| = n, |E| = m$$

| Expressed as O(f) | Edge List | Adjacency Matrix | Adjacency List |
|---|---|---|---|
| Space | n+m | n² | n+m |
| insertVertex(v) | 1* | n* | 1* |
| removeVertex(v) | n+m | n | deg(v) |
| insertEdge(u, v) | 1 | 1 | 1* |
| removeEdge(u, v) | m | 1 | min( deg(u), deg(v) ) |
| incidentEdges(v) | m | n | deg(v) |
| areAdjacent(u, v) | m | 1 | min( deg(u), deg(v) ) |

# Traversal: BFS

Initialize queue / depth / predecessor

or array!
or map!

While queue not empty:

   Remove front vertex of queue

   Check if edge connects to new vertex

     Set dist / pred if new vertex

   Add unvisited edges to queue

Graph implementation Stores table
Vertex Node has
member variable (depth, pred)

| v | d | P | Adjacent Edges |
|---|---|---|---|
| A | 0 | - | B C D |
| B | 1 | A | A C E |
| C | 1 | A | A B D E F |
| D | 1 | A | A C F H |
| E | 2 | B | B C G |
| F | 2 | C | C D G |
| G | 3 | E | E F H |
| H | 2 | D | D G |

Cross edges have Meaning
↳ We already saw that vertex through
   a Shorter path
↳ Dist between vertices linked by cross's ≤ 1

A B C D E F H G

# Traversal: BFS

$O(n+m)$     $|V| = n$    $|E| = m$

Initialize queue / depth / predecessor

While queue not empty:

$O(mn)$

Remove front vertex of queue

Check if edge connects to new vertex

Set dist / pred if new vertex   $+O(m)$

Add unvisited edges to queue

| v | d | P | Adjacent Edges |
|---|---|---|---|
| A | 0 | - | B C D |
| B | 1 | A | A C E |
| C | 1 | A | A B D E F |
| D | 1 | A | A C F H |
| E | 2 | B | B C G |
| F | 2 | C | C D G |
| G | 3 | E | E F H |
| H | 2 | D | D G |

Given vertex, get all edges

$O(m)$ edgs
$O(n)$ matrix
$O(\deg v)$ adj;

$n \times$

sparse    dense

$\sum_V \deg(v) = 2|E| = m = O(m)$    $n-1 \le m \le n^2$

A B C D E F H G

# Traversal: DFS



0) Initialize dist/pred

1) Init Stack
   ↳ init w/ root

2) While stack not empty
   ↳ peek A & get 1 unvisited child
   ↳ Add child to stack
   ↳ If no children unvisited
      pop from stack

H
F
I
E
~~I~~
E
G
B
C
D
Bottom A

# Efficiency: DFS vs BFS

(Traversal)    $|V| = n, |E| = m$

**BFS:** $O(n+m)$



v vertices

A **B** **C** **D** E F H G

All neighbors

↳ width of graph

each deg(v)

$\sum_v \deg(v) = 2|E|$

**DFS:** $O(n+m)$



v vertices

**A B C D F G E** H

each deg(v)

↳ longest path

# Summary: DFS and BFS

$|V| = n, |E| = m$

Both are **O(n+m)** traversals! They label every edge and every node

**BFS**

Solves unweighted MST

Solves shortest path

Solves cycle detection

Memory bounded by width

**DFS**

Solves unweighted MST

Solves cycle detection

Memory bounded by longest path

↳ considered better in memory

# Kruskal's Algorithm

(A, D) ✓
(E, H) ✓
(F, G) ✓
(A, B) ✓
(B, D) ✗
(G, E) ✓
(G, H) ✗
(E, C) ✓
(C, H) ✗
(E, F) ✗
(F, C) ✗
(D, E) ✓
(B, C)
(C, D)
(A, F)
(D, F)

1) Build a **priority queue** on edges
   ↳ min heap
   ↳ Sorted list

2) Build a **disjoint set** on vertices
   ↳ All vertices start as own set

3) Repeat take min edge
   ↳ If connect two sets
      ↳ union sets
      ↳ record edge

4) Stop when:
   — n-1 nodes recorded
   — I have one disjoint set

# Kruskal's Algorithm

**1) Build a priority queue on edges**

```
 1  KruskalMST(G):
 2    DisjointSets forest
 3    foreach (Vertex v : G.vertices()):
 4      forest.makeSet(v)
 5
 6    PriorityQueue Q    // min edge weight
 7    Q.buildFromGraph(G.edges())
 8
 9    Graph T = (V, {})
10
11    while |T.edges()| < n-1:
12      Vertex (u, v) = Q.removeMin()
13      if forest.find(u) != forest.find(v):
14        T.addEdge(u, v)
15        forest.union( forest.find(u),
16                      forest.find(v) )
17
18    return T
19
```

**2) Build a disjoint set on vertices**

**3) Repeatedly find min edge**

If edge connects two sets

Union and record edge

**4) Stop after n-1 edges recorded**

# Kruskal's Algorithm

$|V| = n$     $|E| = m$

| Priority Queue: | Heap | Sorted Array |
|---|---|---|
| **Building** :7 | $O(m)$ | $O(m \log m)$ |
| **Each removeMin** :12 | $O(\log m)$ | $O(1)$ |

$m \times$

$M + m\log m$  vs  $m\log m + m$

why heap good?
↳ What if edge weight changes?

why sorted array good?
↳ sorted array **not** destroyed when used ≡ if we could use array later, this is better!

```
1   KruskalMST(G):
2     DisjointSets forest
3     foreach (Vertex v : G.vertices()):
4       forest.makeSet(v)
5
6     PriorityQueue Q     // min edge weight
7     Q.buildFromGraph(G.edges())
8
9     Graph T = (V, {})
10
11    while |T.edges()| < n-1:
12      Vertex (u, v) = Q.removeMin()
13      if forest.find(u) != forest.find(v):
14        T.addEdge(u, v)
15        forest.union( forest.find(u),
16                      forest.find(v) )
17
18    return T
19
```

$O(n)$

$m \times$

$O(1)$

# Prim's Algorithm



```
1   PrimMST(G, s):
2      Input: G, Graph;
3             s, vertex in G, starting vertex
4      Output: T, a minimum spanning tree (MST) of G
5
6      foreach (Vertex v : G.vertices()):        Init
7         d[v] = +inf
8         p[v] = NULL
9      d[s] = 0
10
11     PriorityQueue Q    // min distance, defined by d[v]
12     Q.buildHeap(G.vertices())
13     Graph T               // "labeled set"
14
15     repeat n times:
16        Vertex m = Q.removeMin()
17        T.add(m)
18        foreach (Vertex v : neighbors of m not in T):
19           if cost(v, m) < d[v]:
20              d[v] = cost(v, m)
21              p[v] = m
22
23     return T
```

update all neighbors
if new smaller edge

| A | B | C | D | E | F |
|---|---|---|---|---|---|
| 0 | 2 | 5 | 17 | 8 | 16 |

# Prim's Algorithm

Sparse Graph: $n \sim m$

    ↳ heap is better

Dense Graph: $m \sim n^2$

    ↳ unsorted array

        better

```
 6   PrimMST(G, s):
 7     foreach (Vertex v : G.vertices()):
 8       d[v] = +inf
 9       p[v] = NULL
10     d[s] = 0
11
12     PriorityQueue Q // min distance, defined by d[v]
13     Q.buildHeap(G.vertices())
14     Graph T          // "labeled set"
15
16     repeat n times:
17       Vertex m = Q.removeMin()   ←
18       T.add(m)
19       foreach (Vertex v : neighbors of m not in T):
20         if cost(v, m) < d[v]:            This is updating
21           d[v] = cost(v, m)   ]
22           p[v] = m
23
```

$$n-1 \leq m \leq n^2$$

$$m = n^2$$

| | Adj. Matrix | Adj. List |
|---|---|---|
| **Heap** | $O(n^2 + m \lg(n))$  → $n^2 \log n$ | Sparse  $O(n \log n)$  $O(n \lg(n) + m \lg(n))$  Dense  $n^2 \log n$ |
| **Unsorted Array** | $O(n^2)$ | $m = n$  $O(n^2)$ |

# MST Algorithm Runtime:

Kruskal's Algorithm:
**O(n + m log (n) )**

Prim's Algorithm:
**O(n log(n) + m log (n) )**

Sparse Graph: m ~ n

Dense Graph: m ~ $n^2$

# Dijkstra's Algorithm (SSSP)

Assume heap
Fib

## What is the running time of Dijkstra's Algorithm?

↳ This is Prim!

$$O\left(n + n\log n + m\right)$$

Find min — Dom term — update

@15 + @18: $\sum_v deg(v) = 2M$

Total # edge updates is M

```
DijkstraSSSP(G, s):
6      foreach (Vertex v : G):
7          d[v] = +inf                    O(n)
8          p[v] = NULL
9      d[s] = 0
10
11     PriorityQueue Q // min distance, defined by d[v]
12     Q.buildHeap(G.vertices())
13     Graph T              // "labeled set"
14
15     repeat n times:        nx
16        Vertex u = Q.removeMin()    O(log n)
17        T.add(u)
18        foreach (Vertex v : neighbors of u not in T):
19           if cost(u, v) + d[u] < d[v]:
20              d[v] = cost(u, v) + d[u]          O(1*)
21              p[v] = m
22
23     return T
```

# Dijkstra's Algorithm (SSSP)

Dijkstras Algorithm works only on non-negative weights

**Optimal implementation:**

Fibonacci Heap

If dense, unsorted list ties

**Optimal runtime:**

Sparse: O(m + n log n)

Dense: O(n²)

```
DijkstraSSSP(G, s):                      (Basically Prim)
6     foreach (Vertex v : G):
7         d[v] = +inf
8         p[v] = NULL
9     d[s] = 0
10
11    PriorityQueue Q // min distance, defined by d[v]
12    Q.buildHeap(G.vertices())
13    Graph T          // "labeled set"
14
15    repeat n times:
16        Vertex u = Q.removeMin()
17        T.add(u)
18        foreach (Vertex v : neighbors of u not in T):
19            if cost(u, v) + d[u] < d[v]:
20                d[v] = cost(u, v) + d[u]        ← This
21                p[v] = m                           changes
22
23    return T
```

# Floyd-Warshall Algorithm

Running time? $O(n^3)$ operation!

$\hookrightarrow$ Easy to code ( Mult! threadable!)

$\hookrightarrow$ Handles neg weight (but not cycle)

```
FloydWarshall(G):
 6    Let d be a adj. matrix initialized to +inf
 7    foreach (Vertex v : G):
 8      d[v][v] = 0
 9    foreach (Edge (u, v) : G):
10      d[u][v] = cost(u, v)
11
12    foreach (Vertex u : G):        nx
13      foreach (Vertex v : G):      nx
14        foreach (Vertex w : G):    ny
15          if d[u, v] > d[u, w] + d[w, v]:
16            d[u, v] = d[u, w] + d[w, v]    ] O(1)
```

matrix

# Final thoughts on Graphs

Graphs have a large space of **possible coding questions**

You should be able to solve common graph questions

  - Make sure you can use graphs to find all neighbors

  - Make sure you can use graphs to solve path questions

Consider how these fundamental skills can be challenged

  - What if I had labels on nodes and I need to find specific ones?

  - What if I need to label nodes or edges with specific properties?

  - Can I handle weights? Directions?

# Probability in CS

# Fundamentals of Probability

Imagine you roll a pair of six-sided dice. What is the expected value?

A **random variable** is a function from events to numeric values.

D1 is value of first dice $\rightarrow$

DBoth is value of D1 + D2

The **expectation** of a (discrete) random variable is:

$$E[X] = \sum_{x \in \Omega} Pr\{X = x\} \cdot x$$

$E[D1] = \frac{1}{6} \cdot 1 + \frac{1}{6} \cdot 2 + \dots \approx 3.5$

$E[DBoth] = \frac{1}{36} \cdot 2 + \frac{1}{36} \cdot (1+2) + \dots \approx 7$

# Next Class: Randomized Data Structures

Sometimes a data structure can be **too ordered / too structured**

Randomized data structures rely on **expected** performance

Randomized data structures 'cheat' tradeoffs!

↳ Add inaccuracy for speed gains

# Probabilistic Data Structures

# Randomized Algorithms

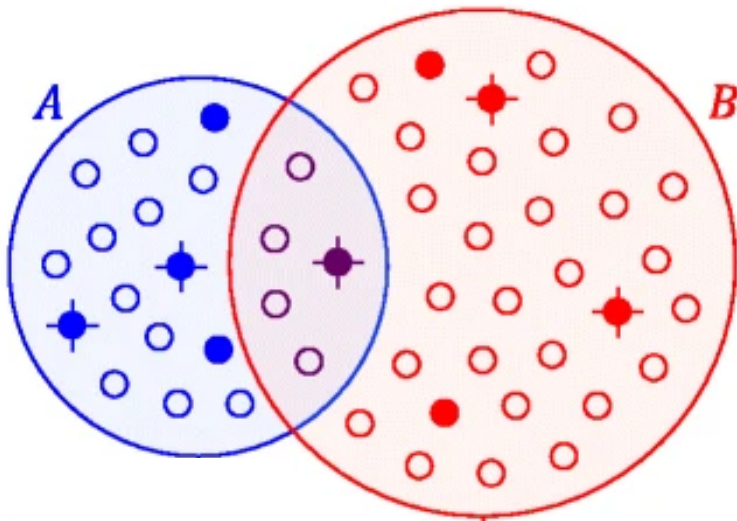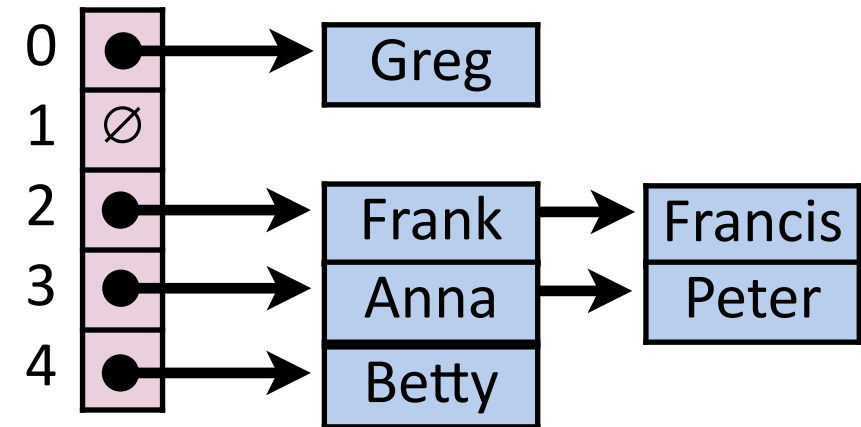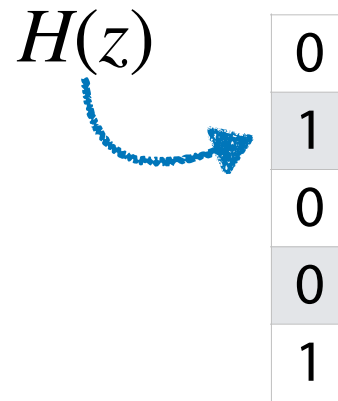A **randomized algorithm** is one which uses a source of randomness somewhere in its implementation.

Figure from Ondov et al 2016

| $H(z)$ | |
|---|---|
| | 0 |
| | 1 |
| | 0 |
| | 0 |
| | 1 |

| | | | |
|---|---|---|---|
| 0 | ● | → | Greg |
| 1 | ∅ | | |
| 2 | ● | → | Frank → Francis |
| 3 | ● | → | Anna → Peter |
| 4 | ● | → | Betty |

| $H(x)$ | 0 | 2 | 1 | 0 | 0 | 4 | 0 | 2 | 0 | 6 |
|---|---|---|---|---|---|---|---|---|---|---|
| $H(y)$ | 1 | 0 | 2 | 3 | 1 | 0 | 3 | 4 | 0 | 1 |
| $H(z)$ | 2 | 1 | 0 | 2 | 0 | 1 | 0 | 0 | 7 | 2 |

# A Hash Table based Dictionary

**User Code (is a map):**

```
1 Dictionary<KeyType, ValueType> d;
2 d[k] = v;
```

A **Hash Table** consists of three things:

1. A hash function    Assigns numeric (positive int) address to any key

Key -> Hash Value (Address)

2. A data storage structure    Array — very good at lookup given **index**

Hash Value (Address) is an index!

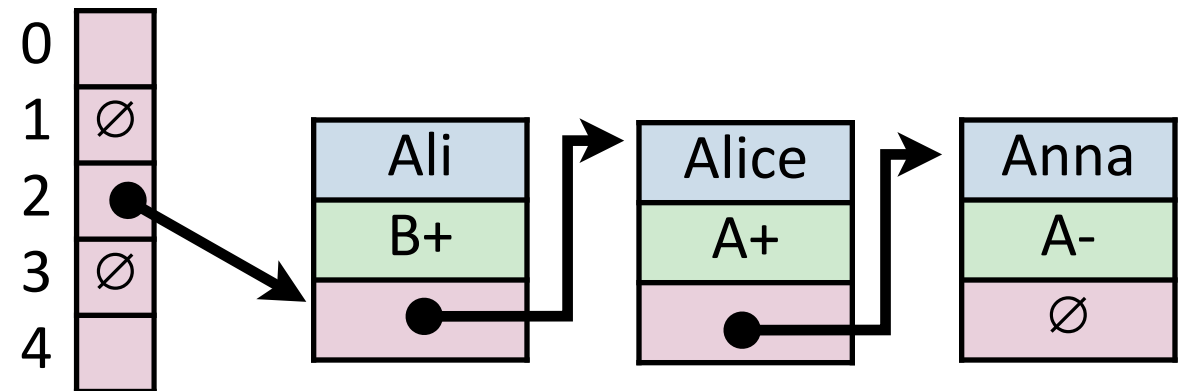**3. A method of addressing *hash collisions***

Two different keys, same hash value

# Open vs Closed Hashing

Addressing hash collisions depends on your storage structure.

- **Open Hashing:** store *k,v* pairs externally

  Such as a linked list

  Resolve collisions by adding to list

  

- **Closed Hashing:** store *k,v* pairs in the hash table

  Everything stored in one list

  How to store collisions? Unclear!

  

# Simple Uniform Hashing Assumption

Given table of size $m$, a simple uniform hash, $h$, implies

$$\forall k_1, k_2 \in U \text{ where } k_1 \neq k_2 , \; Pr(h[k_1] = h[k_2]) = \frac{1}{m}$$

**Uniform:** All keys equally likely to hash to any position

$$Pr(h[k_1]) = \frac{1}{m}$$

**Independent:** All key's hash values are independent of other keys

# Separate Chaining Under SUHA

**Under SUHA, a hash table of size _m_ and _n_ elements:**

Find runs in: $O(1+\alpha)$.

$$\alpha = \frac{n}{m}$$

$a$ — constant we control

Insert runs in: $O(1)$.

Remove runs in: $O(1+\alpha)$.

expected length at every position

0
1
2
3
4
5
6
7
8
9
10

Load Factor $(\alpha)$

$n/m$

# Running Times (Expectation under SUHA)

**Open Hashing:** $0 \leq \alpha \leq \infty$ (Length of chain)

insert: $\underline{\phantom{xxx}1\phantom{xxx}}$.

find/ remove: $\underline{\phantom{xxx}1 + \alpha\phantom{xxx}}$.

**Closed Hashing:** $0 \leq \alpha < 1$ (fraction full)

insert: $\underline{\dfrac{1}{1 - \alpha}}$.

find/ remove: $\underline{\dfrac{1}{1 - \alpha}}$.

**Observe:**

**- As α increases:**

OH: $\alpha \to \infty$, runtime $\to \infty$

CH: $\alpha \to 1$, runtime $\to \infty$

**- If α is constant:**

OH is constant
CH is constant $\Big\}$ $O(1)^*$

# Running Times    *(Don't memorize these equations, no need.)*

*The expected number of probes for find(key) under SUHA*

## Linear Probing:

- Successful:  ½(1 + 1/(1-α))
- Unsuccessful: ½(1 + 1/(1-α))²

## Double Hashing:

- Successful:  1/α * ln(1/(1-α))
- Unsuccessful: 1/(1-α)

## When do we resize?

Linear   ~ 0.7 - (0.8)

Double   ~ 0.7 - (0.9)

# Running Times

| | Hash Table | | AVL | Linked List |
|---|---|---|---|---|
| **Find** | Expectation*: $O(1)$*** <br><br> Worst Case: $O(n)$ | | $O(\log n)$ | $O(n)$ |
| **Insert** | Expectation*: $O(1)$*** <br><br> Worst Case: $O(n)$ | | $O(\log n)$ | $O(1)$ |
| **Storage Space** | $O(n)$ | | $O(n)$ | $O(n)$ |

# Bloom Filter

A probabilistic data structure storing a set of values

Built from a bit vector of length $m$ and $k$ hash functions

Insert / Find runs in: $\underline{\qquad O(k) \quad / \quad O(1) \qquad}$

Delete is not possible (yet)!

$$H = \{h_1, h_2, \ldots, h_k\}$$

| 0 |
|---|
| 0 |
| 1 |
| 0 |
| 0 |
| 1 |
| 0 |
| 1 |
| 0 |
| 0 |

# Probabilistic Accuracy in a Bloom Filter

# Bloom Filter: Error Rate

$m/n = 10$

FPR $\left(1 - e^{\frac{-nk}{m}}\right)^k$

Not enough random trials!

$k* = \ln 2 \cdot 10 = 6.93$

BF becomes too saturated w/ 1s

k small :(

$k$ # hashes

Figure by Ben Langmead

# Cardinality Estimation

Let min $= 95$. Can we estimate $N$, the cardinality of the set?



Conceptually: If we scatter $N$ points randomly across the interval, we end up with $N + 1$ partitions, each about $1000/(N + 1)$ long

Assuming our first 'partition' is about average:
$$95 \approx 1000/(N + 1)$$
$$N + 1 \approx 10.5$$
$$N \approx 9.5$$

# Set Similarity Review

To measure **similarity** of $A$ & $B$, we need both a measure of how similar the sets are but also the total size of both sets.

$$J = \frac{|A \cap B|}{|A \cup B|}$$

$J$ is the ***Jaccard coefficient***

# MinHash Sketch

**Claim:** Under SUHA, set similarity can be estimated by sketch similarity!

# MinHash Sketch

We can convert any hashable dataset into a **MinHash sketch**



hash

$k$

$\{1, 3, 5\}$

$\{2, 3, 7\}$

⅕

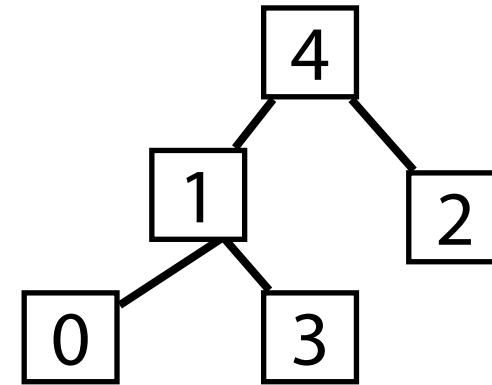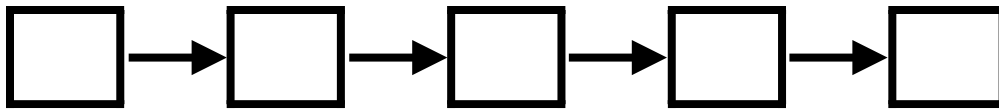We lose our original dataset, but we can still estimate two things:
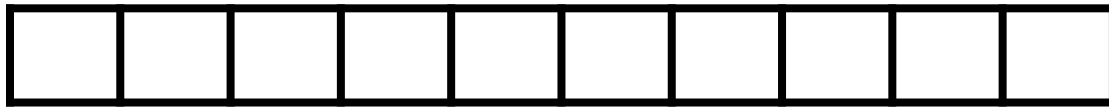
1. Cardinality (# of items)

2. Set Similarit

# Questions?

# CS 225 — Course Goals

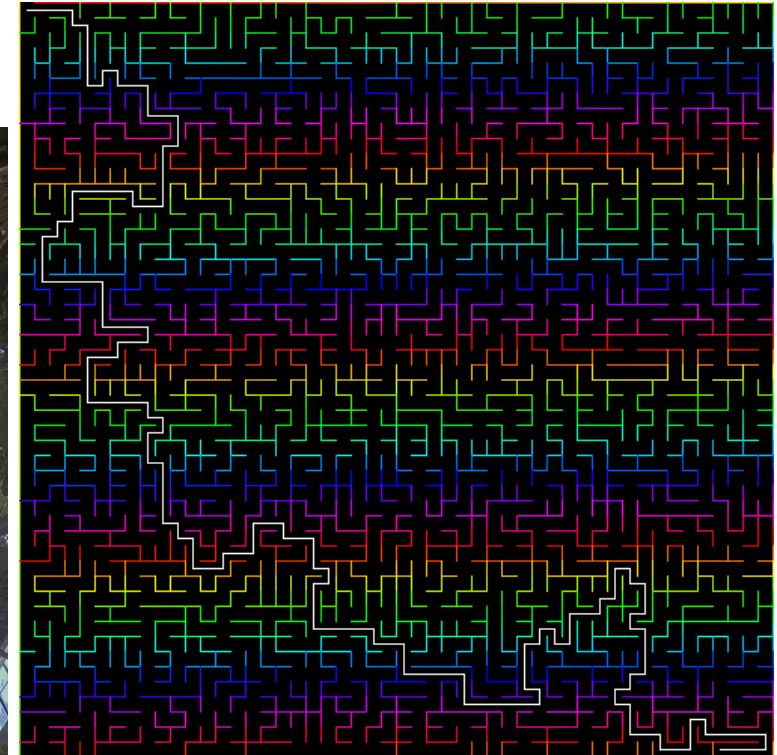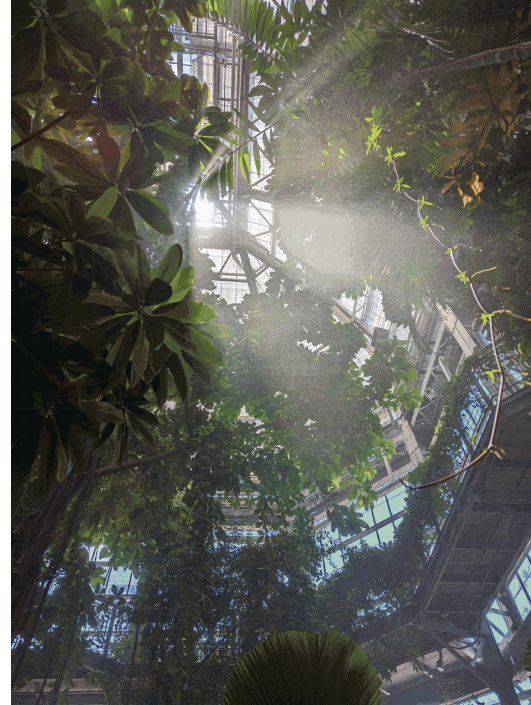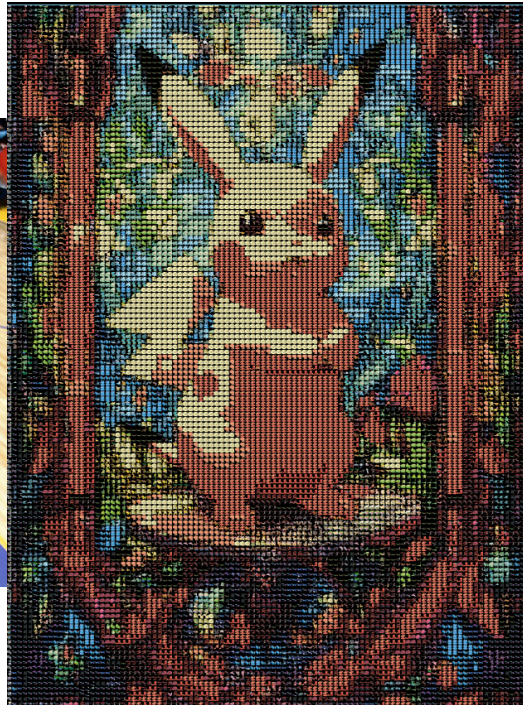## Understand foundational data structures and algorithms

# CS 225 — Course Goals

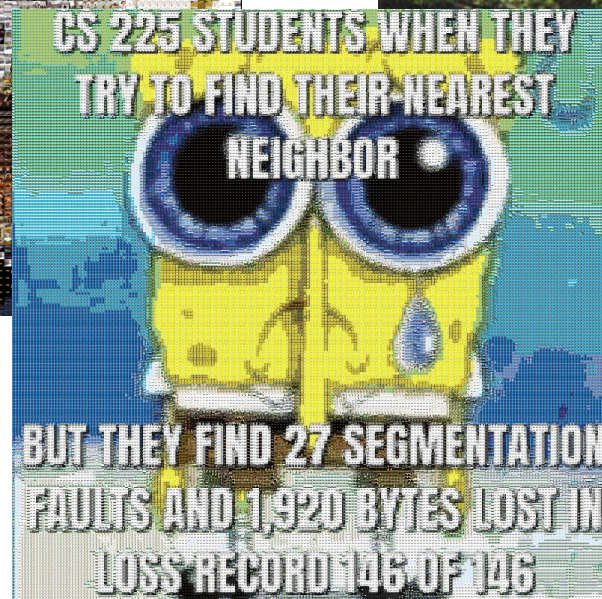## Justify appropriate algorithms for complex problems
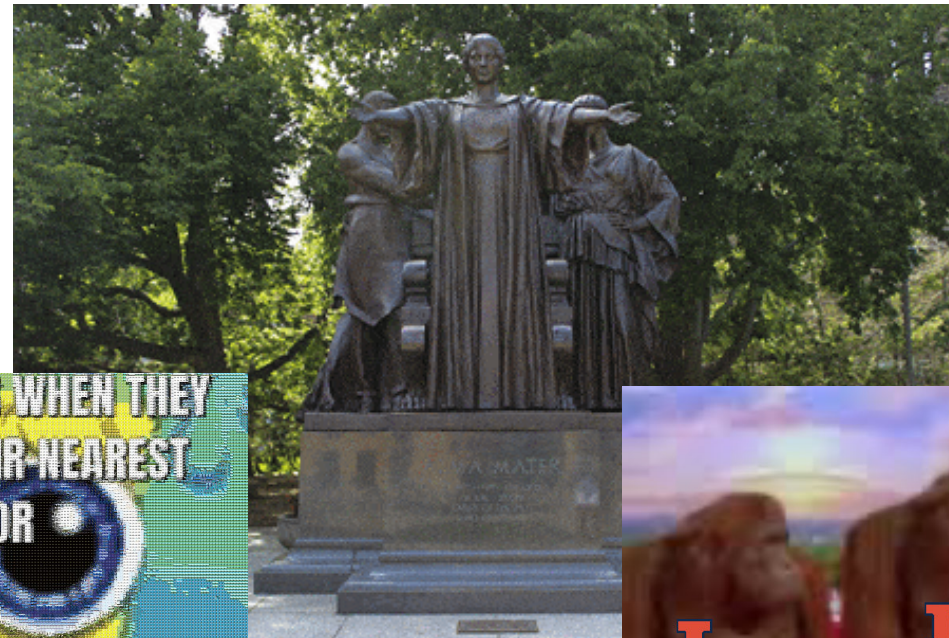
*Decompose problem into supporting data structures*

*Analyze efficiency of implementation choices*

# CS 225 — Course Goals

Implement intermediate difficulty problems in C++



CS 225 STUDENTS WHEN THEY TRY TO FIND THEIR NEAREST NEIGHBOR

BUT THEY FIND 27 SEGMENTATION FAULTS AND 1,920 BYTES LOST IN LOSS RECORD 146 OF 146

new mp where

# CS 225 — Course Goals

Understand foundational data structures and algorithms

Justify appropriate algorithms for complex problems

Implement intermediate difficulty problems in C++

Improve your foundation of CS theory

# Good luck on your finals!