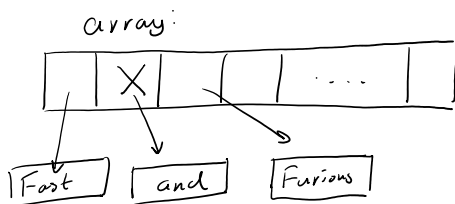


Data Structures

- A collection containing:
 - Data values
 - Relationship among data
 - Operations applied to data
- Describes exactly how the data are organized and how tasks are performed

Eg. ArrayList



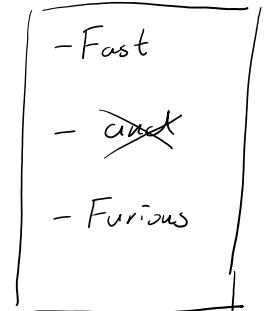
Counter: ~~X~~ ~~X~~ ~~X~~ ~~X~~ 2

Abstract Data Types

- Defines behavior relative to a user
 - what operations must it have?
- Describes only what needs to be done, not how it's done.

Eg. List

- add elements
- find elements
- remove elements
- size of list
- is ordered



C++ Reference Notes

Friday, September 24, 2021 4:15 PM

Data Types : Signed or Unsigned

- char (1 bytes)
 - short (2 bytes)
 - int (4 bytes)
 - long (8 bytes)
 - long long (16 bytes)
- } integers
- float (4 bytes)
 - double (8 bytes)
 - long double (16 bytes)
- } floating points
- bool (1 byte)

Strings

- string type.
 - mutable (modifiable)
 - can only concatenate strings
- s.substr(beginIndex, length)

Relational Operators

- ==, !=, <, <=, >, >=
- work for all types, even classes due to operator overloading.

Variables

Friday, September 24, 2021 4:10 PM

In C++

- variable initialization is not checked

```
int fast;  
int furious;  
int fastFurious = fast * furious; ← undefined behavior
```

- narrowing is not checked

```
int x = 40000;  
short y = x; ← overflow, truncation
```

- variables can be declared outside of a class:

```
int m = 42; ← global variable, poor practice except for constants  
class MyClass {  
    ...  
}
```

C++ Classes:

- public/private modifiers are grouped into regions, where all public/private fields/methods are grouped.
- methods do not require implementation at declaration.
- can use a Member Initializer List:

private: int x; int y;

Point::Point(int i, int j): x(i), y(j) {}

- header files and source files.

header files store class, field, method declaration

source files store field, method implementations

Example Class

Friday, September 24, 2021 4:25 PM

Watch later Sha

C++

```
class Student {  
    public:  
        static int numStudents;  
  
        Student(string n);  
  
        void setName(string n);  
        string getName() const;  
  
    private:  
        string name;  
};  
  
int Student::numStudents = 0;  
  
Student::Student(string n) { /* CODE */ }  
  
void Student::setName(string n) { /* CODE */ }  
string Student::getName() const { /* CODE */ }
```

declare

Example Member Initializer List

Friday, September 24, 2021 4:26 PM

Member Initializer List

```
class Point {  
    private:  
        int x;  
        int y;  
  
    public:  
        Point(int i, int j);  
};  
  
Point::Point(int i, int j) {  
    x = i;  
    y = j;  
}
```



```
class Point {  
    private:  
        int x;  
        int y;  
  
    public:  
        Point(int i, int j);  
};  
  
Point::Point(int i, int j) : x(i), y(j) {}
```

Example Header and Source Files

Friday, September 24, 2021 4:26 PM

Student.h header

```
class Student {  
    public:  
        static int numStudents;  
        Student(string n);  
  
    private:  
        string name;  
};
```

Student.cpp

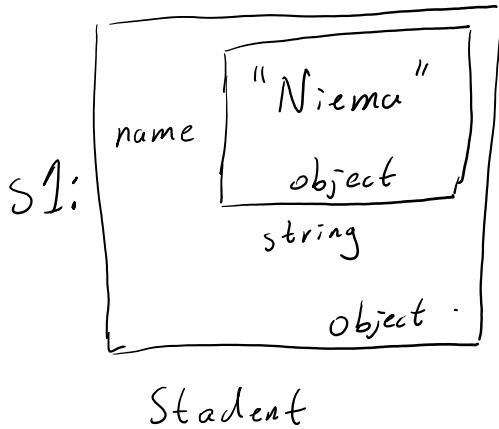
```
int Student::numStudents = 0;  
  
Student::Student(string n) : name(n) {  
    numStudents++;  
}
```

Source

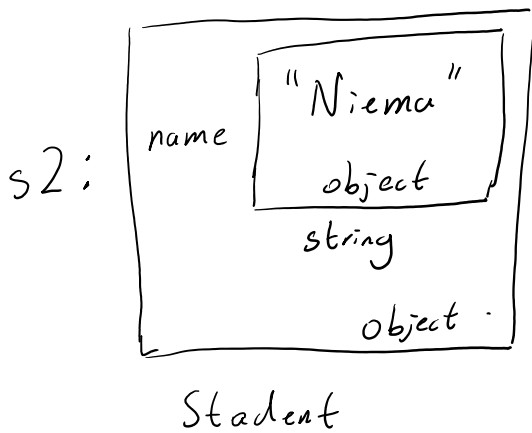
Memory Models (Classes), Reference

Friday, September 24, 2021 4:27 PM

```
Student s1 ("Niema");  
Student s2 = s1;
```



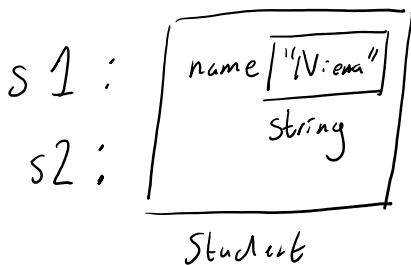
`s1` is not a reference to the `Student` object, it is the object.



`s2` is a copy of `s1`.
It is not a reference.

To create a reference use:

```
Student &s2 = s1;
```



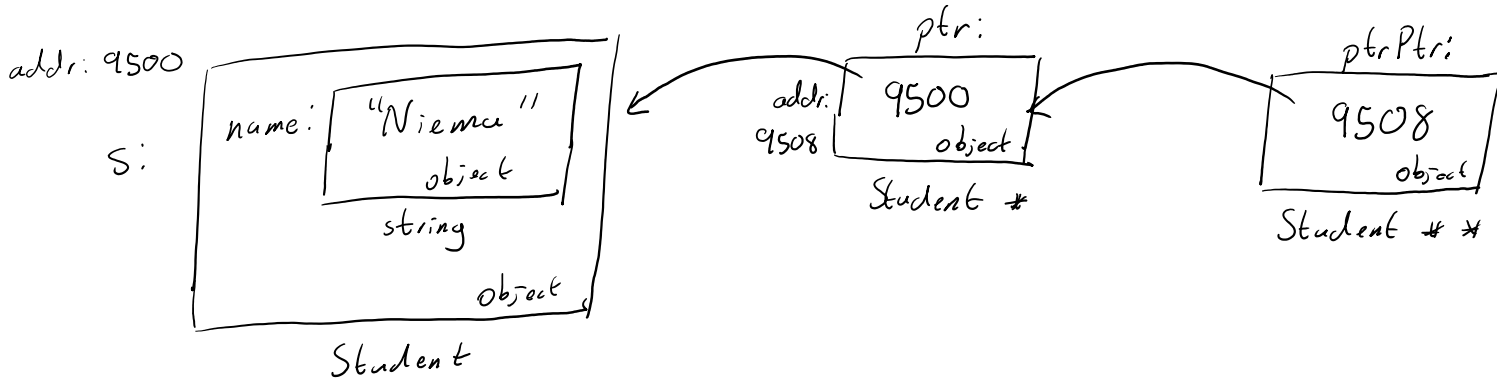
Pointers

Friday, September 24, 2021 4:35 PM

```
Student s = Student("Niema");
```

```
Student * ptr = &s;
```

```
Student ** ptrPtr = &ptr;
```



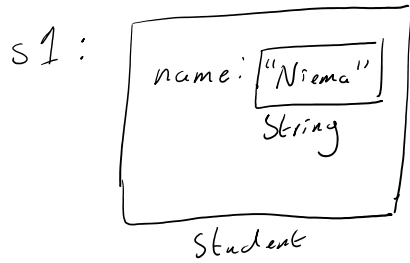
To use the value at a pointer:

- 1) Dereference $(*ptr).name$
 - 2) Arrow $ptr \rightarrow name$
- } equivalent

Memory Management (new, delete)

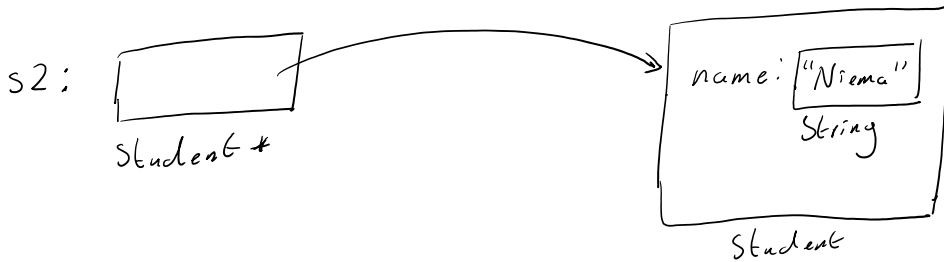
Friday, September 24, 2021 4:40 PM

```
Student s1 = Student ("Niema");  
Student *s2 = new Student ("Ryan");
```



s1 resides on the runtime stack, meaning it is destroyed once outside the scope

• The new keyword is dynamic memory allocation



s2 resides on the heap and will not be automatically deallocated

To deallocate memory, use the delete keyword.

Const(ant) Variables

Friday, September 24, 2021 4:45 PM

- `const` variables can't be modified
trying to do so will result in a compiler error.

```
const int a = 42; } identical  
int const b = 42;
```

Const References

Friday, September 24, 2021 4:54 PM

```
int a = 42;
```

```
const int & ref1 = a;
```

```
int const & ref2 = a;
```

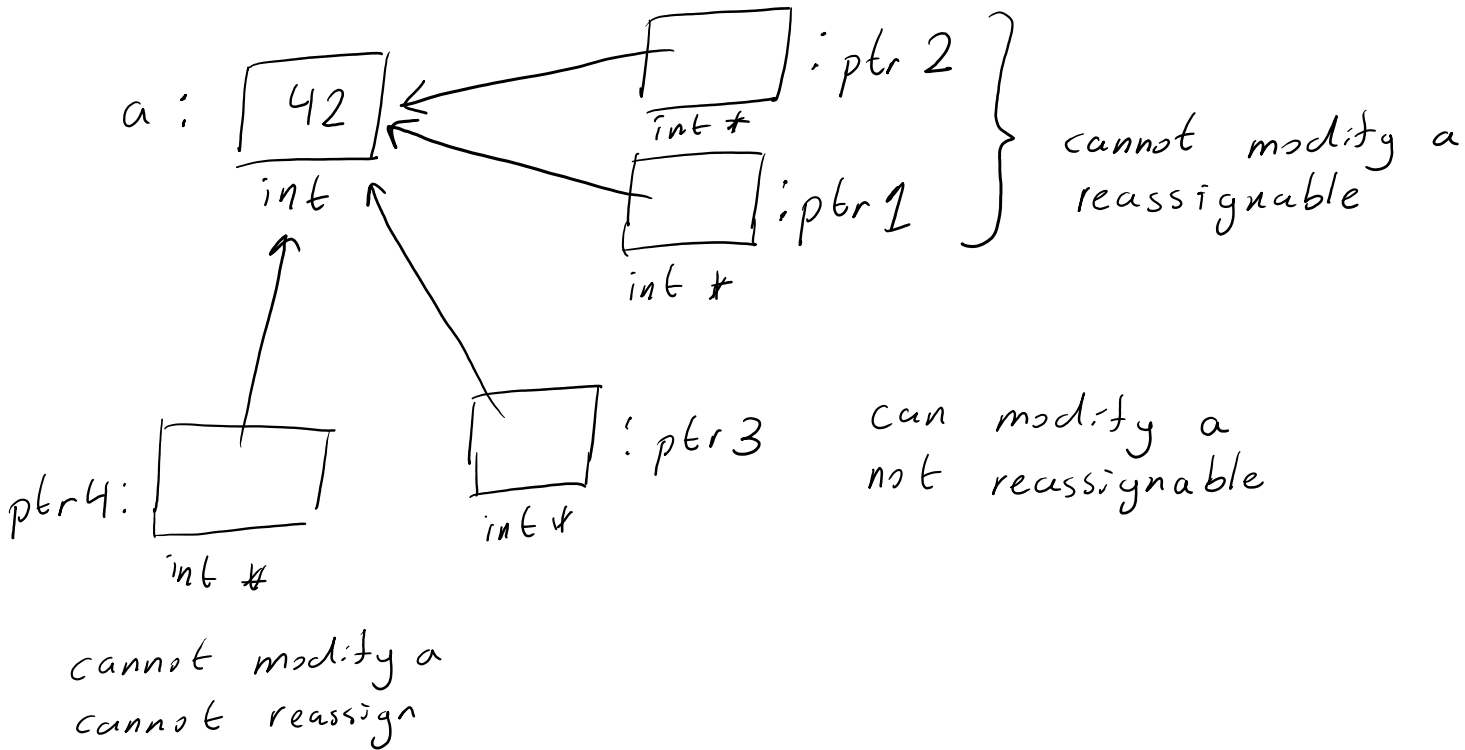
} equivalent,

cannot modify the variable through
the reference
can be reassigned.

Const Pointers

Friday, September 24, 2021 4:49 PM

```
int a = 42;  
const int * ptr1 = &a;  
int const * ptr2 = &a; } equivalent  
int * const ptr3 = &a;  
const int * const ptr4 = &a;
```



Const Functions

Friday, September 24, 2021 4:56 PM

const functions cannot modify the object

- cannot modify instance variables
- can only call other const functions.

eg: `public : string getName () const ;`

Functions, Global Functions

Friday, September 24, 2021 7:16 PM

- Functions may be defined outside any class, these are global functions
 - main function is a global function, returns an exit code
ie:

```
int main() {  
    ...  
}  
class MyClass {  
    ...  
};
```


Pass by Value

```
void swap (int a, int b) {  
    int tmp = a;  
    a = b;  
    b = tmp;  
}
```

Note because a, b are passed by value, this method creates no change

Pass by Reference

```
void swap (int &a, int &b) {  
    int tmp = a;  
    a = b;  
    b = tmp;  
}
```

Note because a, b are passed by reference, this method swaps the values of a, b

Vectors

Friday, September 24, 2021 7:22 PM

C++ Vector

- is a dynamic array

vector <int> a;

- add elements with:

a.push_back(42);

a.push_back(21);

- remove elements with:

a.pop_back();

- access specific index with:

a[0];

Note accessing out of bounds
does not throw an error

- assigning one vector to another copies all elements.

vector <int> b = a;

b is a separate vector copy of a.

Input, Output

Friday, September 24, 2021 7:26 PM

Standard output : `cout << "message" << endl;`

Standard input : `cin >> n;`

Standard error : `cerr << "error" << endl;`

To read an entire line from stdin:

```
getline(cin, message)
```

// get from cin, store to message

Note `cout`, `cin`, `cerr`, `endl` are part of the `std` namespace.
to use a namespace:

```
using namespace std;
```

Templates

Friday, September 24, 2021 7:30 PM

- Similar to generics in Java.
- We can specify a template in C++ with:

```
template <typename T> class Node {  
    public:  
        T const data;  
        Node(const T &d): data(d) {}  
};
```

- We instantiate a template object with:

```
Node <string> a(s);  
Node <int> b(n);
```

Iterators

Friday, September 24, 2021 7:37 PM

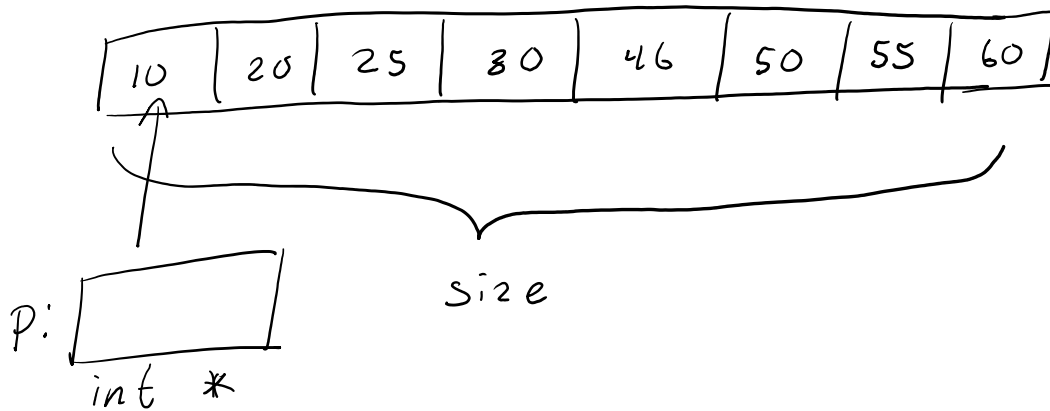
Motivation: given a for each loop like:

```
for (string name : names) {  
    cout << name << endl;  
}
```

names could be a vector, set, unordered_set, but with iterators, we don't need to know what specific type names is.

Iterating over an Array with Pointer

Friday, September 24, 2021 7:40 PM



We can iterate over the array with:

```
for (int i = 0; i < size; ++i) {  
    cout << *p << endl;  
    ++p;  
}
```

Annotations: An arrow points from the text "dereference" to the asterisk in `*p`. Another arrow points from the text "pointer arithmetic" to the `++p` line.

Steps:

- 1) Dereference pointer
 - 2) Used data to print
 - 3) Incremented pointer
- Arrows indicate the flow from step 1 to 2, and from 2 to 3. A large curved arrow on the left side encompasses all three steps.

Iterating through a Vector using Iterators

Friday, September 24, 2021 7:46 PM

- Given some vector:

```
vector<string> names;
```

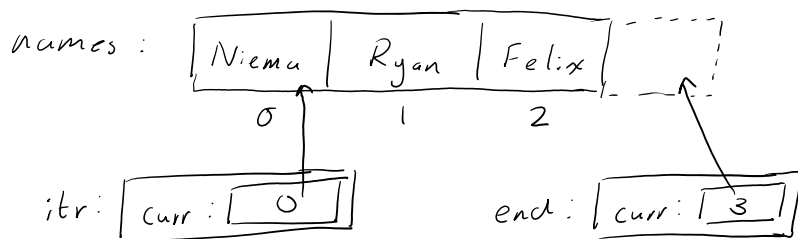
- We can get an iterator to the beginning of the vector with:

```
vector<string>::iterator itr = names.begin(); ← first item
```

and an iterator to the end with:

```
vector<string>::iterator end = names.end(); ← just after the last item
```

- Visualization:



- We can then iterate through the vector with:

```
while (itr != end) {  
    cout << *itr << endl;  
    ++itr;  
}
```

Iterating through a LinkedList using Iterators

Friday, September 24, 2021 7:51 PM

- Similar to iterators of vectors:

```
LinkedList <string> names;
```

```
LinkedList <string>::iterator itr = names.begin();
```

```
LinkedList <string>::iterator end = names.end();
```

```
while (itr != end) {
```

```
    cout << *itr << endl;
```

```
    ++itr;
```

```
}
```


Custom Iterators

Friday, September 24, 2021 7:59 PM

General idea of any Iterator class:

Overload operators in the Iterator class:

- $==$ true if iterators are pointing to the same item
- $!=$ true if iterators are not pointing to the same item
- $*$ (dereference) return a reference to the current data value
- $++$ (pre/post increment) move the iterator to the next item

New functions in the Data Structure class:

- `begin()` return an iterator to the first element
- `end()` returns iterator to just after the last element

Time Complexity

Sunday, September 26, 2021 6:07 PM

- How fast is my program?
 - ↳ hours?
 - ↳ minutes?
 - ↳ nanoseconds?} human time

Problem: human time depends on hardware

- How fast is my algorithm
 - given an input size (n)

Notations of Time Complexity

Sunday, September 26, 2021 6:11 PM

- Big-O : upper bound $f(n) \in O(g(n))$ iff $A \cdot g(n) \geq f(n) \forall n$
- Big- Ω : lower bound $f(n) \in \Omega(g(n))$ iff $B \cdot g(n) \leq f(n) \forall n$
- Big- Θ : upper & lower bound $f(n) \in \Theta(g(n))$ iff $f(n) \in O(g(n))$ & $f(n) \in \Omega(g(n)) \forall n$
 $f(n) \in \Theta(g(n))$ iff $B \cdot g(n) \leq f(n) \leq A \cdot g(n) \forall n$

Finding Big-O, common Big-O

Sunday, September 26, 2021 6:18 PM

- 1) Determine $f(n)$, number of operations given input size n
- 2) Drop all terms of n except highest term
- 3) Drop all constant coefficients.

Common Time Complexities.

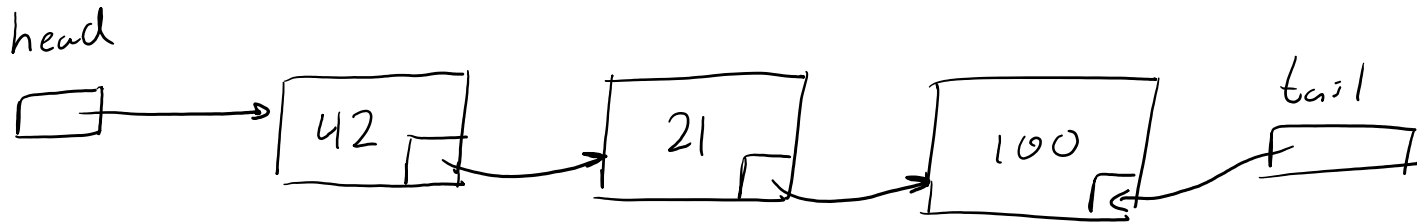
• $O(1)$	constant	time	} $O(n^c) \forall c$ polynomial time
• $O(\log n)$	logarithmic	time	
• $O(n)$	linear	time	
• $O(n \log n)$			
• $O(n^2)$	quadratic	time	
• $O(n^3)$	cubic	time	
• $O(k^n)$	exponential	time	} bad time complexities
• $O(n!)$	factorial	time	

Space Complexity

Sunday, September 26, 2021 7:54 PM

- We can use Big-O notation to generalize.

Eg: Linked List



space used: $O(n)$

Graphs, Tree

Saturday, October 2, 2021 12:53 PM

Def: Graphs are a collection of Nodes and Edges

- Edges can be directed or undirected:

undirected



directed



Def: Trees are graphs with:

- no undirected cycles (no cycles of undirected edges)
- connected (must exist a path between any two nodes)

Special Trees

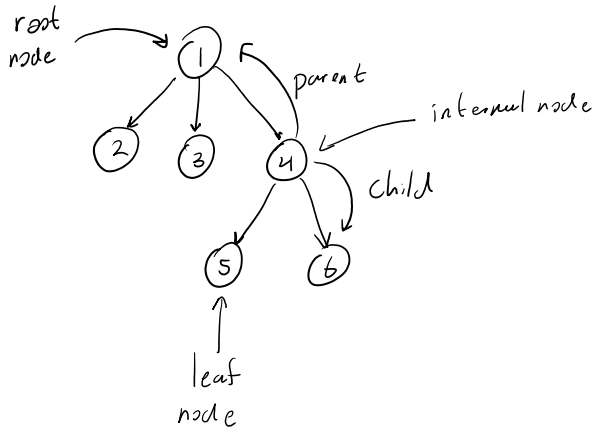
Saturday, October 2, 2021 1:04 PM

- Empty / Null Tree:
 - 0 nodes
 - 0 edges
- Single Node Tree:
 - 1 Node
 - 0 Edges

Rooted vs Unrooted

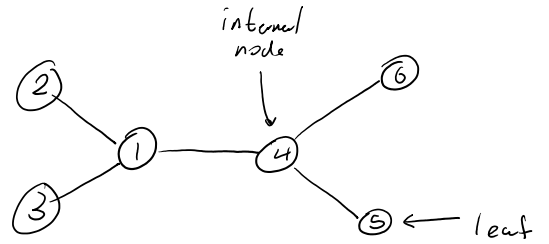
Saturday, October 2, 2021 1:08 PM

Rooted: Hierarchical Structure



root node: has no parent
internal node: has parent and children
leaf node: has no children

Unrooted: No structure



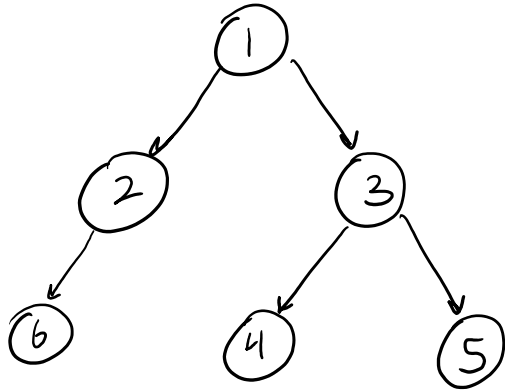
internal node: more than one neighbor
leaf node: only one neighbor

Rooted Binary Trees

Saturday, October 2, 2021 1:12 PM

Rooted: parent-child relationships

Binary: each node has at most 2 children



Tree Traversals

Saturday, October 2, 2021 1:17 PM

Preorder Traversal: Visit, Left, Right
In-Order Traversal: Left, Visit, Right
Post-Order Traversal: Left, Right, Visit

} DFS

Level-Order Traversal: 1st level (left to right), 2nd level (left to right) ... } BFS

Binary Search Tree

Saturday, October 2, 2021 1:46 PM

• Binary Search Trees are:

- Rooted Binary Tree

- Every node is larger than all nodes in its left subtree

Every node is smaller than all nodes in its right subtree

• BST Find Algorithm:

1) Start at the root

2) If $query == current$, return

3) If $query > current$, traverse right goto #2

4) If $query < current$, traverse left goto #2

If there are no nodes to traverse left/right, fail

• BST Insert Algorithm:

1) Perform Find Algorithm: duplicate if successful

2) If find doesn't succeed, insert new element at site of failure

• BST Successor Algorithm: Successor: next largest node

1) If the node has a right tree, traverse right once, then all the way left

2) Otherwise, traverse up until the node was the parent's left child:
the parent is the successor

• BST Remove Algorithm:

1) Perform Find Algorithm: exists if successful

If Node has no children

- delete the Node

If Node has one child

- replace Node with Node's child

If Node has two children

- swap node with successor, apply remove algorithm to Node

BST Height

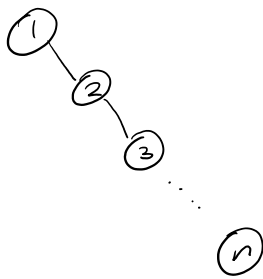
Saturday, October 2, 2021 2:11 PM

Def: Node height: Longest distance from a node to a leaf

Def: Tree height: height of the root of the tree

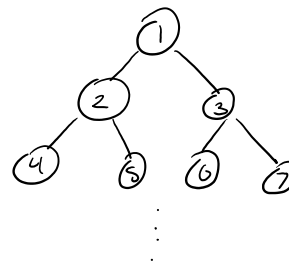
Def: Balance: ratio between tree height and number of elements

Perfectly Unbalanced:



$$\text{height} = n - 1$$

Perfectly Balanced:



$$\text{height} = \log_2(n+1) - 1$$

• Best vs. Worst vs. Average for Find Algorithm

Best: query is the root $\longrightarrow O(1)$

Worst: perfectly unbalanced, query not found $\longrightarrow O(n)$

Average: theoretical expected value $\longrightarrow O(\log(n))$

• Average Case

- Assumptions:

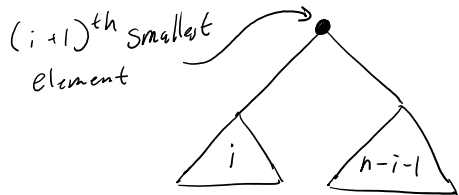
- 1) All n elements are equally likely to be searched
- 2) All $n!$ possible insertion orders are equally likely

Def: Node depth: number of nodes in the path from Node to root

avg case time complexity = expected # operations to find query
= expected depth

Total expected depth: sum of all node depths, $D(n)$

Def: Every tree can be represented with:



$D(n|i)$ = expected total depth w/ n nodes and i nodes in left subtree

$D(i)$ = expected total depth of left subtree

$$D(n|i) = D(i) + D(n-i-1) + i + (n-i-1) + 1$$

$$= D(i) + D(n-i-1) + n$$

Def: The expected number of operations for some BST is: $\left| \frac{1}{n} D(n) \right|$

$$D(n) = \sum_{i=0}^{n-1} D(n|i) \underbrace{P(I=i)}_{\substack{\text{probability } i \text{ nodes in left subtree} \\ = (i+1)^{\text{th}} \text{ smallest node was inserted first}}} \left. \vphantom{\sum} \right\} \frac{1}{n}$$

$$D(n) = \frac{1}{n} \sum_{i=0}^{n-1} D(i) + D(n-i-1) + n$$

Thus:

$$D(n) = \frac{1}{n} \sum_{i=0}^{n-1} D(i) + \sum_{i=0}^{n-1} D(n-i-1) + n$$

$$D(n) = \frac{1}{n} \sum_{i=0}^{n-1} D(i) + \sum_{i=0}^{n-1} D(n-i-1) + n$$

$$= \frac{2}{n} \sum_{i=0}^{n-1} D(i) + n$$

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

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

$$n D(n) - (n-1) D(n-1) = 2 D(n-1) + n^2 - (n-1)^2$$

$$n D(n) = (n+1) D(n-1) + 2n - 1$$

$$\frac{D(n)}{n+1} = \frac{D(n-1)}{n} + \frac{2n-1}{n(n+1)} \xrightarrow{n=1} \frac{D(1)}{2} = \frac{D(0)}{1} + \frac{2-1}{(1)(2)} \rightarrow D(1) = 1$$

$D(0) = 0$

$$\frac{D(n)}{n+1} = \sum_{i=1}^n \frac{2i-1}{i(i+1)} = \sum_{i=1}^n \frac{2}{i+1} - \sum_{i=1}^n \frac{1}{i(i+1)} \rightarrow (\text{expanded})$$

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

We can approximate this to $2 \frac{n+1}{n} \ln(n) - 3 \approx 1.386 \log_2(n)$

which is $O(\log(n))$

Treeps

Saturday, October 9, 2021 3:32 PM

• Treep = Tree + Heap stores objects of (key, priority)

1) BST properties wrt keys for any given node:

- larger than all keys in left subtree
smaller than all keys in right subtree

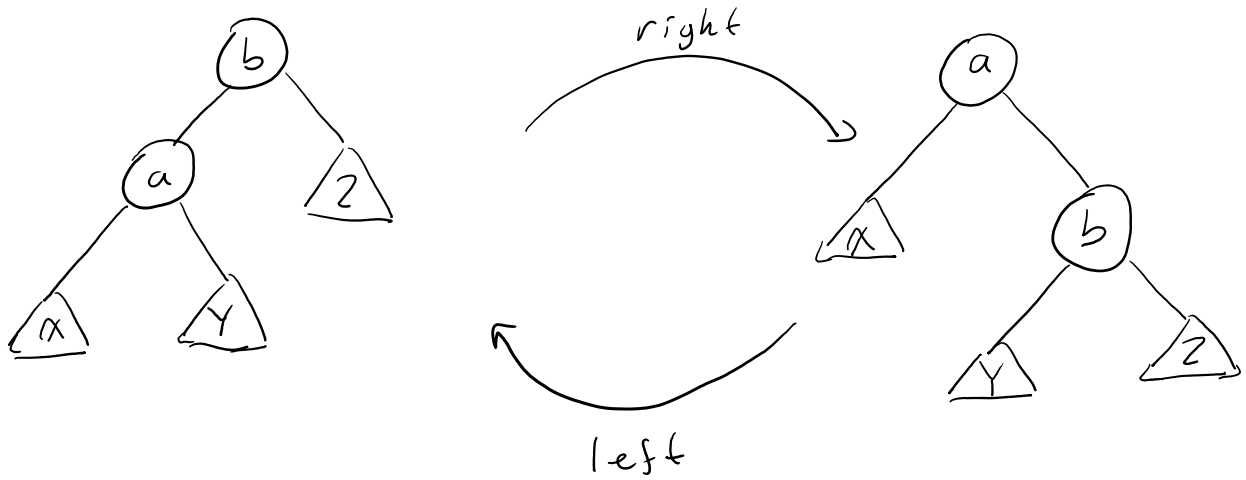
2) Heap property wrt priorities for any given node:

- larger than all priorities below

AVL Rotation

Saturday, October 9, 2021 3:36 PM

• Given some tree: AVL right/left rotations are:



Note: BST properties are maintained

Treep Insertion

Saturday, October 9, 2021 3:47 PM

Treep Insertion Algorithm:

1) Insert via BST insert algorithm wrt keys

2) Use AVL Rotations to "bubble up" to fix priorities

- If the inserted node has a higher priority than parent

- If inserted node is left child \rightarrow AVL right

- If inserted node is right child \rightarrow AVL left

- repeat until inserted node has lower priority than parent

Random Search Trees

Saturday, October 9, 2021 3:57 PM

RSTs are implemented with Treaps:

- use elements (values) as keys
- Randomly generate priorities maintaining Heap property

Advantage: addresses the assumption of BSTs that each element has an equal chance of being searched, creates a more balanced tree.

Runtime: Worst: $O(n)$ Average: $O(\log(n))$

AVL Trees

Saturday, October 9, 2021 4:12 PM

- AVL Tree : BST with balancing
 - Balance factor = (Height of right subtree) - (Height of left subtree)
 - AVL tree: BST where every node has BF of $(-1, 0, 1)$

AVL Time Complexities

Saturday, October 9, 2021 4:18 PM

Find algorithm: Best = $O(1)$ Average = $O(\log(n))$ Worst = $O(\log(n))$

Note: AVL Trees have better worst time complexity than BST

AVL Insertion Algorithm

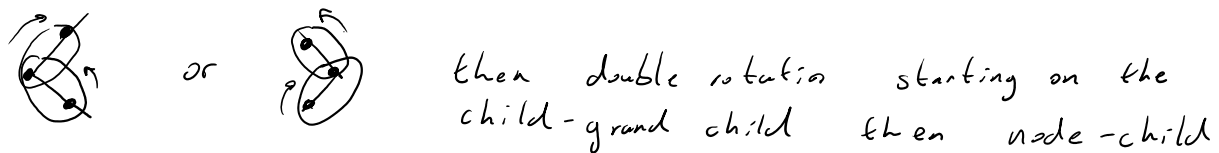
Saturday, October 9, 2021 4:31 PM

• AVL Insertion Algorithm:

- Perform regular BST insertion
- Update Balance Factors
- Fix broken Balance Factors using AVL Rotations
 - If out of balance node, child, grand child are:



but if



Red-Black Tree, RB vs AVL

Saturday, October 9, 2021 5:00 PM

- RB Tree = BST with balancing
 - All nodes must be red or black
 - The root must be black
 - If a node is red, all of its children must be black
 - For every node, every possible path to a null reference must have the same # of black nodes
 - null references are black
- RB vs AVL:
 - RB Trees are not always AVL Trees, all AVL Trees could be RB Trees

RB Time Complexities

Saturday, October 9, 2021 5:09 PM

- Like AVL Trees, RB Trees have worst case time complexity for find algorithm of $O(\log(n))$

RB Insertion Algorithm

Saturday, October 9, 2021 5:11 PM

- Case 1: Empty Tree

- Insert the new node as root
- Color it black

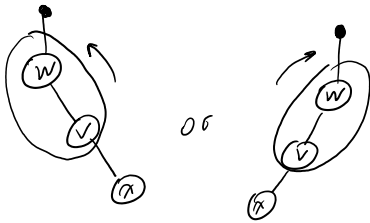
- General steps for inserting in non-empty tree

- Insert using regular BST
 - during traversal, if there is ever a black node w/ two red children, swap all colors red \leftrightarrow black for node and two children
 - insert node and color new node red
- potentially fix tree for Red-Black tree properties

- Case 2: Child of Black Node

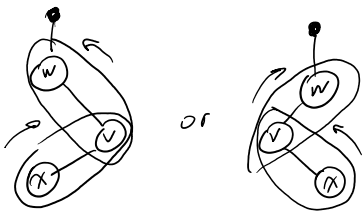
- no conflicts, done!

- Case 3: Child of Red Node, Straight Line



- perform insertion
- rotate parent-grandparent
- swap color of parent-grandparent red \leftrightarrow black

- Case 4: Child of Red Node, Crossed



- perform insertion
- rotate node-parent then (parent-grandparent) \leftarrow case 3
- swap color of node-parent

Set, Map ADTs

Friday, October 15, 2021 11:21 PM

- Set: Store multiple elements (keys)
 - find(x) : true if x exists in the set, false otherwise
 - insert(x) : add x to the set
 - remove(x) : remove x from the set
- Map: Store multiple (key, value) pairs
 - get(k) : return value associated with key k
 - put(k, v) : add the key value pair
 - remove(k) : remove the key and associated value

Implementing Set/Map ADTs

Friday, October 15, 2021 11:25 PM

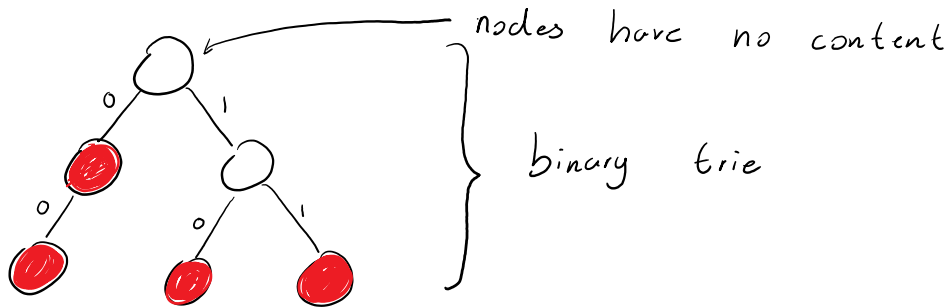
	Find/Remove	Insert	Notes
Unsorted Linked List :	$O(n)$	$O(1)$	
Sorted Linked List :	$O(n)$	$O(n)$	Can iterate in sorted order
Unsorted Array List :	$O(n)$	$O(1)$	
Sorted Array List :	$O(\log n)$	$O(n)$	can iterate in sorted order
Self Balancing BST :	$O(\log n)$	$O(\log n)$	can iterate in sorted order
Hash Table :	$O(1)$	$O(1)$	need to calculate $O(k)$ hash

Tries, Binary Tries

Friday, October 15, 2021 11:31 PM

Def Trie is a tree structure in which elements are represented by paths

Ex Binary Trie:



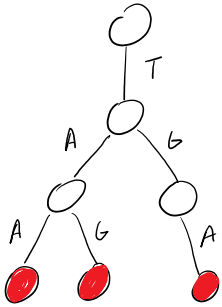
red nodes are word nodes: 0, 00, 10, 11

Multiway Tries, Algorithms (Find, Insert, Remove)

Friday, October 15, 2021 11:35 PM

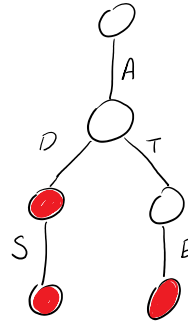
Def Multiway Tries are Tries which nodes can have more than 2 children

DNA Trie



words: TAA, TAG, TGA

English Trie



words: AD, ADS, ATE

Algorithms:

Initialize: 1) Create root node non-word

Insert : 1) Start at the root

2) If there exists an edge for each letter, traverse the edge
Else create the edge and next node

3) Repeat complete, mark last letter's node as word node

Find : 1) Follow edges corresponding with each letter until word node or run out of edges

Remove : 1) Perform find algorithm

2) Mark node as non-word if it exists

Multiway Tries Time, Space Complexity

Friday, October 15, 2021 11:42 PM

Given a MWT and: n = number of words, k = length of longest word

Time Complexity:

Insert $O(k)$

Find $O(k)$

Remove $O(k)$

Given a MWT and: $|\Sigma|$ = length of alphabet, n = number of words, k = length of longest word

For each node, we can have $|\Sigma|$ children

There can be n^k possible words

Space occupied is $O(|\Sigma|^{k+1})$

Significant Algorithms on Multiway Tries

Friday, October 15, 2021 11:58 PM

- Iterate in ascending order \rightarrow preorder traversal
- Iterate in descending order \rightarrow postorder traversal
- Iterate in order of length \rightarrow level order traversal
- Autocomplete: Given some prefix, return all words in subtree \rightarrow find, traverse from found node

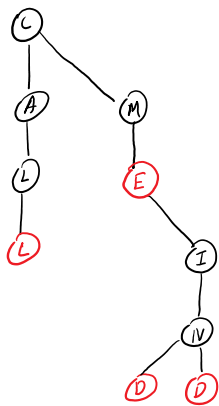
Ternary Search Trees

Saturday, October 16, 2021 12:06 AM

- BST: $O(k \log n)$, memory efficient
- MWT: $O(k)$, memory inefficient
- TST: Somewhere in between

Def: Ternary Search Trees are trees with each node having some value, and also is or is not a word node, when traversing a node is only used if the traversal brings it down.

Ex:



words: call
me
mind
mid

TST Find

Saturday, October 16, 2021 12:13 AM

Algorithm:

- 1) current node = root
current letter = first letter of query
- 2) If $c >$ current node:
traverse right
- If $c <$ current node:
traverse left
- If c is last letter and current node is word node
success!
- Else:
traverse down the middle child and go to next letter in query

TST Insert, Remove

Saturday, October 16, 2021 12:17 AM

Insert Algorithm:

- Traverse using comparisons with the query, if it does not exist create the node
- Make the last node ended at a word node

Remove Algorithm:

- Perform find algorithm
- Make the last node not a word node

TST Time Complexity

Saturday, October 16, 2021 12:30 AM

Given a TST with n words, k longest word, $|\Sigma|$ alphabet size

$O(n)$ worst case

$O(\log(n))$ avg. case

Hash Functions, Tables, Maps

Saturday, October 23, 2021 9:43 PM

Hash Function:

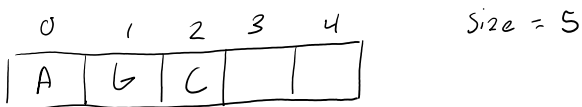
Input: Object

Output: Integer representation of $x \leftarrow$ hash value

Property of equality: if $x = y$, then $h(x) = h(y)$

Property of inequality: if $x \neq y$, then $h(x) \neq h(y)$ \leftarrow not necessary

Hash Table: Stores elements by hash value: implements Set ADT



$$A: h(A) = 65 \% 5 = 0$$

$$C: h(C) = 67 \% 5 = 2$$

$$G: h(G) = 71 \% 5 = 1$$

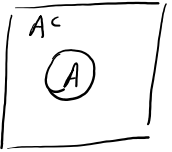
Hash Map: Hash table where each element is a key-value pair: implements Map ADT

Collisions, Probability

Saturday, October 23, 2021 10:02 PM

• Basic Probability Review

$$P(A) + P(A^c) = 1 \rightarrow P(A) = 1 - P(A^c)$$



$$P(A \text{ and } B) = P(A) \cdot P(B) \text{ for independent events } A, B$$

$$P(A \text{ or } B) = P(A) + P(B) \text{ for mutually exclusive } A, B$$

• Probability of At Least 1 collision given N elements, M slots:

$$P_{N,M}(\geq \text{collision}) = 1 - P_{N,M}(0 \text{ collisions})$$

$$P_{N,M}(0 \text{ collisions}) = 1 \cdot \frac{M-1}{M} \cdot \frac{M-2}{M} \cdot \dots \cdot \frac{M-N+1}{M}$$

• Load Factor:

$$LF = \frac{N}{M} = \alpha$$

$$\text{expected collisions} = \sum_{i=1}^{N-1} \frac{1}{m} = \frac{N(N-1)}{2M} = 1 \rightarrow N = \sqrt{2M}, M = O(N^2)$$

$$= \frac{1}{2} \left(1 + \frac{1}{1-\alpha} \right)$$

Collision Resolution

Saturday, October 23, 2021 10:47 PM

- Open Addressing (Linear Probing) → chance of future collisions get more likely
- Attempt to insert as normal
If collision: iterate through array until an empty slot is found
- Attempt to find as normal
If occupied: iterate through array until query is found or empty slot is found
- Double Hashing:
 - Create a new hash function, so that when we try to insert, we use the next hash function
 - Finding an element: hash once, if not query, hash again, if not query then it does not exist
- Closed Addressing (Separate Chaining) → chance of future collisions are always the same
 - Make each array index a linked list. During each insertion search through to check
 - Find same as insertion.

Bloom Filters

Saturday, October 23, 2021 11:24 PM

Def Probabilistic datastructure, memory-efficient

- ↳ no false negatives
- ↳ possible false positives

Bloom Filters are hash tables of booleans, multiple hash functions for double hashing

For each item:

insert: set each index from every hash function to true

find: if all index from every hash function is true, then it likely exists

Design:

- Assume 1) Each hash function uniformly distributes across the array
2) Each insertion is independent

Given: $k = \#$ hash functions $m = \#$ bits $n = \#$ insertions

Probability of false-positive in the Bloom Filter

Probability that a specific bit is set to true: $\approx 1 - e^{-\frac{kn}{m}}$

$$P(\text{false-positive}) = \epsilon \approx (1 - e^{-\frac{kn}{m}})^k$$

$$\text{Thus: } k = \frac{m}{n} \ln(2) = \lceil -\log_2(\epsilon) \rceil, \quad m = \frac{-n \ln(\epsilon)}{\ln(2)^2}$$

- 1) Estimate n
- 2) Pick ϵ that is appropriate
- 3) $k = -\log_2(\epsilon)$
- 4) $m = \frac{-n \ln(\epsilon)}{\ln(2)^2}$

Count-Min Sketch

Saturday, October 23, 2021 11:46 PM

Def Probabilistic data structure, memory efficient

Used to count number of occurrences, provides upper-bound on counts

	0	1	2	3	4
h_1	0	0	0	0	0
h_2	0	0	0	0	0
h_3	0	0	0	0	0

$k = \begin{cases} h_1 \\ h_2 \\ h_3 \end{cases}$

insert: For each element, compute the hash functions, increment corresponding row's count at the corresponding index.

find: compute all hash functions, the maximum possible is the minimum individual count

Design:

Given $n = \#$ elements in stream, c_x the true count of element x , \hat{c}_x the estimated count $c_x \leq \hat{c}_x$
 $\hat{c}_x \leq c_x + \epsilon n$ with probability $1 - \delta$

1) Guess value of n

2) Pick upper bound $\hat{c}_x - c_x \leq \epsilon n$ for ϵ

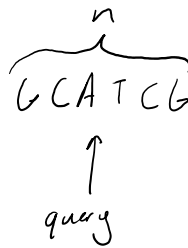
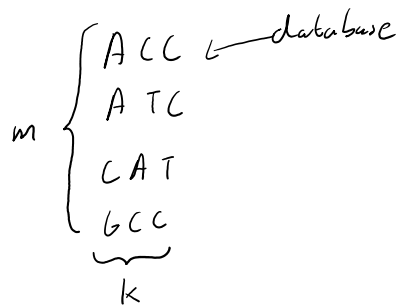
3) Pick probability of being in range δ

4) $m = \lceil \frac{e}{\epsilon} \rceil$ where $m = \text{num columns}$ and $k = \lceil \ln(\frac{1}{\delta}) \rceil$

String Sequencing

Saturday, October 30, 2021 6:37 PM

Problem: Matching database of short strings to a long string query



Time complexity: $O(n \cdot k \cdot m)$

Aho-Corasick Automaton

Saturday, October 30, 2021 4:34 PM

Solution: Aho - Corasick Automaton

- 1) Construct the multinary tree relating to the database
- 2) For every node, connect a failure link to the longest suffix that exists in MWT excluding itself. CATG has suffixes ATG, TG, G, \emptyset
root node links only to itself
- 3) For every node, connect a dictionary link to the first word node encountered after traversing the failure links until reaching a root node

For each query: traverse the MWT given the query. If we reach a word node, then it exists in the query

If there is no further traversal, traverse the failure link

If there is a dictionary link on any node, then we encountered that word.

Shortcut for dictionary links: subwords of a word which are also in the automaton get connected to the word. The last letter of longest subword connects to last letter of word in automaton.

Suffix Arrays

Saturday, October 30, 2021 6:37 PM

Solution: swap database and queries

database:

GCATCGC
n

queries:

ACC }
ATC } k
GC }

suffix array:

0: GCATCGC
1: CATCGC
2: ATCGC
3: TCGC
4: CGC
5: GC
6: C

} sort

2: ATCGC
6: C
1: CATCGC
4: CGC
5: GC
0: GCATCGC
3: TCGC

for each query: we search the suffix array by

- 1) Binary search for the lowest index occurrence of the query, prefix matches are counted
- 2) Binary search for the highest index occurrence of the query, prefix matches are counted
- 3) every index between low and high are instances of query

time complexity: sort: $O(n \log n)$ → closer to $O(n^2 \log n)$

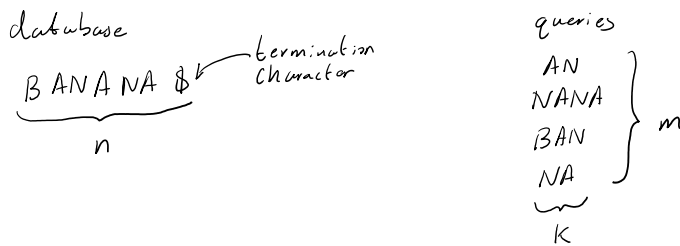
search: $O(k \cdot \log n)$ per query sequence

space complexity: $O(n^2)$

Burrows - Wheeler Transform

Saturday, October 30, 2021 7:01 PM

Solution: Burrows - Wheeler Transform



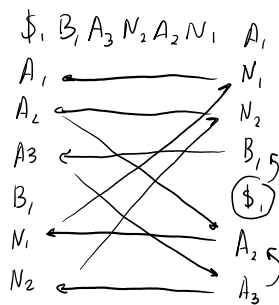
To create the BWT:

- 1) create all rotations of the database
- 2) sort rotations alphabetically, assume termination character is the smallest letter
- 3) BWT is the last column of all the rotations

Note: BWT is invertible:

To create string from BWT:

- 1) the first character one the characters in the BWT in alphabetical order
- 2) the number of times each character in each column are the same across all columns
- 3) Starting from the termination character the character above each character must have followed it, jumping from column to column:



Searching a BWT for pattern matching:

Given a BWT: we can match queries:

i	first	last	Last → First
0	\$ ₁	A ₁	1
1	A ₁	N ₁	5
2	A ₂	N ₂	6
3	A ₂	B ₁	4
4	B ₁	\$	0
5	N ₁	A ₂	2
6	N ₂	A ₃	3

- 1) start with top and bottom rows
- 2) search for first and last instances of each character looking at last column and using the last → first mapping
- 3) The top bottom pointers contain all instances of the query
- 4) convert positions to

b

10_2

10_3

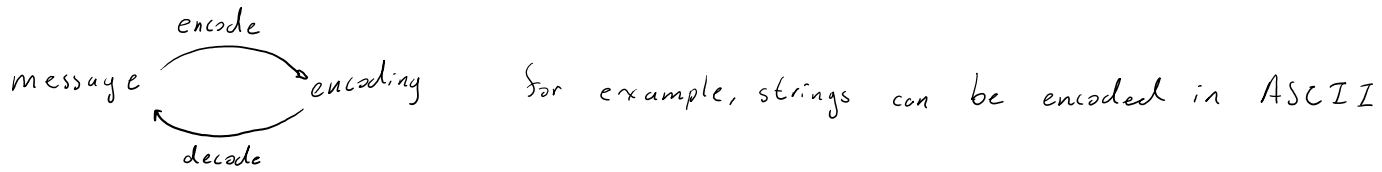
3

4) convert positions by

Encoding/Decoding, Entropy

Saturday, November 6, 2021 5:44 PM

Goal Given some message, we can convert it between raw and encoded formats:



Def Information \neq Data :

Information : content of some message

Tells us details about some system

Data: Raw unit of information

Representation (or encoding) of information

Def Entropy is the measure of disorder or non-uniformness of a system

Shannon Entropy: Expected value of the information contained in some data

Ex: A coin flip has two outcomes H or T, which is 1 bit of information

A biased coin flip that has only H, then 0 bits of information

Def when there are n outcomes then it contains $\lceil \log_2(n) \rceil$ bits of information

Def Uniformity: Lack of variation among symbols in a message

Ex: uniform
AAAAAAA

variable
ACGTACGT

generally, the more uniform the data, the less entropy and thus less information

Coding Trees, Prefix Codes

Saturday, November 6, 2021 5:48 PM

Idea We can encode a message using a tree. These are called encoding trees:

Ex Given the encoding map:

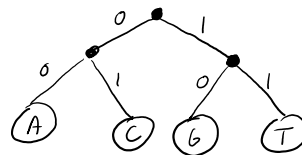
A \rightarrow 00

C \rightarrow 01

G \rightarrow 10

T \rightarrow 11

we can create the tree:



Def An encoding tree is a rooted tree where each leaf represents a character and the path to the node is its encoding.

Def A Prefix Code is an encoding which no symbol is represented by a code that is a prefix of another symbol's code.

A Prefix Code can always be represented by a coding tree

Data Compression

Saturday, November 6, 2021 6:01 PM

Idea We can try to compress some data

lossless: retains all data intact

lossy: smaller compression but may lose some data

Ex Given the two encoding methods:

ASCII

A \rightarrow 65

C \rightarrow 67

G \rightarrow 71

T \rightarrow 84

Two-Bit

A \rightarrow 0

C \rightarrow 1

G \rightarrow 2

T \rightarrow 3

A compression method could convert ASCII to the Two-Bit method and save 6 bits per character. It is a lossless compression.

Def The lowest bound on Data Compression:

Shannon Entropy = $\sum_x p_x \log_2 \left(\frac{1}{p_x} \right)$ where x are the symbols in a message

the Shannon Entropy is the lowest bound on Data Compression

Huffman Trees

Saturday, November 6, 2021 6:06 PM

Def A Huffman Tree is the optimal lossless compression

Alg: Construction:

- 1) Compute frequencies of symbols
- 2) Start with forest of individual trees with single symbol
- 3) While there are more than 1 tree in the forest
 - remove the 2 trees with lowest frequencies
 - create new node with combined frequencies of both trees
 - make new node parent of both trees, lower frequency on left
 - insert new tree into forest.

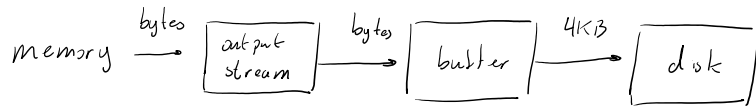
Alg: Encoding: Given a message, encode each symbol using the Huffman tree as a coding tree.

Alg: Decoding: Given an encoding, decode each symbol by traversing the tree according to the encoding, when a leaf is reached that is the next decoded symbol and restart at the root.

Byte Buffers, Bit Buffers (Read/Write)

Saturday, November 6, 2021 6:45 PM

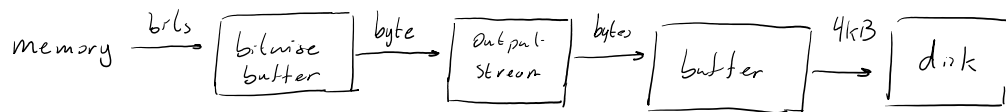
Idea We can use a buffer to write bytes to I/O efficiently:



We can use a buffer to read bytes from I/O efficiently:



Idea Instead, if we want to write/read bits instead of bytes, we can use a bitwise buffer:



Def Reading from a bitwise buffer: Given a byte:



we can extract a specific bit using bit masks and shifts

thus: to extract the bit at position c :

$$(buff \gg (7-c)) \& 1;$$

Def Writing to a bitwise buffer: Given some bits:

we can shift the bit and OR with the buffer.

Graphs, Directed/Undirected, Weighted/Unweighted

Saturday, November 13, 2021 8:54 PM

Def Graphs are a collection of nodes and edges

Def Nodes are a single entity

Def Edges are a relationship between nodes

Def Undirected graphs: an edge (u, v) is $u \rightarrow v$ and $v \rightarrow u$

Directed graph: an edge (u, v) is only $u \rightarrow v$

Def Unweighted graph: edges do not have a weight

Weighted graph: edges have a cost/weight value

Cycles

Saturday, November 13, 2021 8:59 PM

Def A cycle is a path starting from and ending at the same node

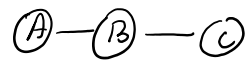
Unstructured, Sequential, Hierarchical, Structured

Saturday, November 13, 2021 9:01 PM

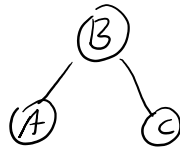
Def Unstructured: A disconnected collection of nodes:



Def Sequential: An ordered collection of connected nodes:



Def Hierarchical: A ranked collection of connected nodes:



Def Structured: A collection of connected and unconnected nodes:

Multigraphs

Saturday, November 13, 2021 9:04 PM

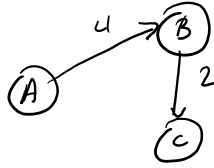
Def A multigraph is a graph in which a pair of nodes can have multiple distinct edges:



Graph Representations in Memory

Saturday, November 13, 2021 9:05 PM

Given some graph:



thus $V = \text{set of vertices}$
 $E = \text{set of edges}$

Def Adjacency matrix:

	A	B	C
A	X	4	X
B	X	X	2
C	X	X	X

size: $O(|V|^2)$

Def Adjacency list:

A: $\{(B, 4)\}$

B: $\{(C, 2)\}$

C: $\{\}$

size: $O(|V| + |E|)$

Graph Traversals, BFS, DFS, Complexities

Saturday, November 13, 2021 9:10 PM

BFS: Add $(0, \text{start})$ to the queue

While the queue is not empty:

pop (d, curr) from the queue

if curr has not been visited:

mark curr as visited with distance d

for all outgoing edges (curr, w) :

if w has not been visited, add $(d+1, w)$ to queue

Time complexity:

$$O(|V| + |E|)$$

DFS: Add start to the stack

While the stack is not empty:

pop curr from the stack

if curr has not been visited:

mark curr as visited

for all edges (curr, w) :

if w has not been visited, add w to stack

$$O(|V| + |E|)$$

Dijkstra's Algorithm

Saturday, November 13, 2021 9:38 PM

- 0) Set all node distances to ∞ and start to 0
- 1) Add $(0, \text{start})$ to priority queue
- 2) While the priority queue is not empty:
 - pop (d, curr) from priority queue
 - if curr is not visited:
 - mark curr as visited
 - for all edges (curr, w, e) :
 - if $d + e < w$'s current distance:
 - w 's distance = $d + e$
 - w 's previous node is curr
 - add $(d + e, w)$ to priority queue

Time complexity:

$$O(|V| + |E|)$$

Spanning Trees, Minimum Spanning Trees

Sunday, November 21, 2021 5:47 PM

Def a Spanning Tree is a subtree of any graph G :

- 1) contains all nodes of G
- 2) contains a subset of the edges of G
- 3) Has no cycles
- 4) Is connected

Def a Minimum Spanning Tree minimizes the edge cost of creating the spanning tree

Def Prim's Algorithm: $O(|E| \log |E|)$

Start at any node

for $|V|-1$ times:

Find the smallest weight edge such that one node on edge is in MST and the other is not
add edge to the MST

Def Kruskal's Algorithm: $O(|E| \log |E|)$

Repeat $|E|$ times:

Find the smallest edge such that adding it would not cause a cycle

Add edge to the MST

Disjoint Set ADT

Sunday, November 21, 2021 6:18 PM

Def Disjoint set ADT:

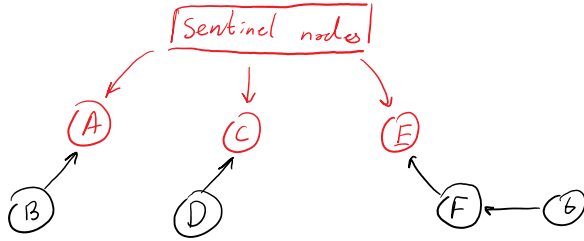
Union : Given two elements u & v , merge the sets in which they belong

Find : Given an element u , returns the set in which it belongs.

Up-Trees

Sunday, November 21, 2021 6:22 PM

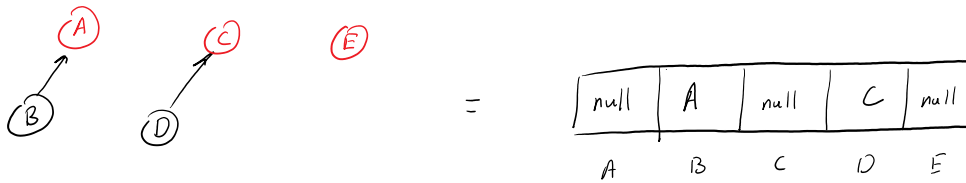
Def Up-Tree is a graph where each element is a single node, Each set is represented as a sentinel node. Sentinel nodes are nodes without a parent.



Find: Return the element's sentinel node

Union: Find elements' sentinel nodes and make one the child of the other.

Idea We can implement an Up-Tree using an array:



Def Union operation By Size: which ever sentinel node has more children should be the parent

Def Union operation By Height: which ever sentinel node has greater height should be the parent

Def Path compression: For every node under a sentinel, call find and point each node directly to sentinel.

Def Time complexity: Given an Up-Tree with n nodes

Union: $O(1)$ By Size, $O(n)$ By Height.

Find: Uncompressed: $O(n)$, Compressed: $O(1)$

Problem Complexity, P = NP

Sunday, November 28, 2021 9:57 PM

Def We can define a problem by describing its Inputs and Outputs

Ex Find maximum of list:

Input: A list of n numbers x_n

Output: A number $v \in x$ such that v is the maximum

Def We can group problems by complexity:

P: Problems can be solved in $O(n^c)$

NP: Problems can be verified in $O(n^c)$

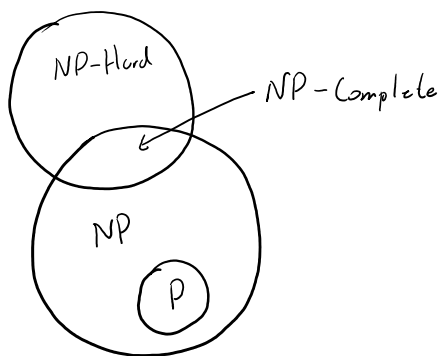
- P is a subset of NP

NP-Hard: Problems at least as hard as the hardest problem in NP

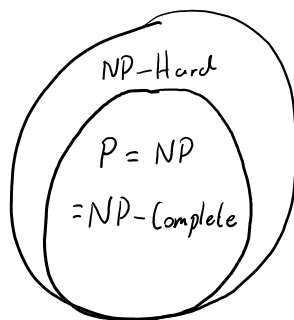
NP-Complete: The intersection of NP and NP-Hard

- The hardest problems in NP

Def P = NP problem:



$P \neq NP$



$P = NP$