

---

# **CS335 Class Notes**

Luis F. Uceta

2020-01-28

## Contents

<b>1</b>	<b>Lecture 1: C++ basics</b>	<b>6</b>
1.1	C++ classes . . . . .	6
1.2	Initialization list . . . . .	6
1.3	Constant member functions . . . . .	7
1.4	Interface and implementation . . . . .	7
1.4.1	IntCell interface . . . . .	8
1.4.2	IntCell implementation . . . . .	8
1.4.3	Using IntCell . . . . .	9
1.5	C++ Objects . . . . .	9
1.6	vectors and strings . . . . .	9
1.6.1	vector . . . . .	10
1.7	C++ details . . . . .	10
1.7.1	Pointers . . . . .	10
1.7.2	Dynamic object allocation and garbage collection . . . . .	11
1.7.3	Lvalues, and Rvalues . . . . .	11
1.7.4	References . . . . .	12
1.7.5	Uses of lvalue references . . . . .	12
<b>2</b>	<b>Lecture 2: C++ Basics (cont.)</b>	<b>14</b>
2.1	Parameter passing . . . . .	14
2.2	Return passing . . . . .	15
2.3	C++11 <code>std::swap</code> and <code>std::move</code> . . . . .	16
2.4	The big five . . . . .	17
2.4.1	Destructor . . . . .	17
2.4.2	Copy constructor and move constructor . . . . .	18
2.4.3	Copy assignment and move assignment (operator =) . . . . .	18
2.4.4	Defaults . . . . .	19
2.5	Templates . . . . .	20
2.5.1	Function templates . . . . .	20
2.6	Function objects . . . . .	21
2.6.1	Class templates . . . . .	25
2.7	Using matrices . . . . .	25
2.8	Class problems . . . . .	27
<b>3</b>	<b>Lecture 3: Algorithm analysis</b>	<b>29</b>
3.1	Mathematical background . . . . .	29

3.2	Notation . . . . .	29
3.2.1	Big-Oh notation . . . . .	29
3.2.2	Big-Omega notation . . . . .	30
3.2.3	Big-Theta notation . . . . .	31
3.3	Typical growth rates . . . . .	31
3.4	Rules . . . . .	32
3.5	Few points . . . . .	32
3.6	Model of computation . . . . .	33
3.7	What to analyze . . . . .	33
3.7.1	Types of performance . . . . .	33
3.8	Running-time calculations . . . . .	34
3.9	Sample problems . . . . .	35
3.9.1	Sum of cubes . . . . .	35
3.9.2	Factorial . . . . .	35
3.9.3	Maximum subsequence sum problem . . . . .	36
3.9.4	Algorithm 1 (brute force) . . . . .	37
3.9.5	Algorithm 2 (brute force) . . . . .	38
3.9.6	Algorithm 3 . . . . .	38
3.9.7	Algorithm 4 . . . . .	39
<b>4</b>	<b>Lecture 4: Algorithm analysis (cont.) and lists/stacks/queues</b>	<b>40</b>
4.1	Binary search . . . . .	40
4.1.1	Running time . . . . .	40
4.2	List/stacks/queue ADTs . . . . .	41
4.2.1	The List ADT . . . . .	41
4.2.2	Vector in the STL . . . . .	42
4.2.3	List in the STL . . . . .	42
4.3	Containers . . . . .	43
4.4	Iterators . . . . .	43
4.4.1	How to get an iterator . . . . .	43
4.4.2	Iterator methods . . . . .	43
4.4.3	Container operations that require iterators . . . . .	44
4.4.4	Example: Using erase on a list . . . . .	45
4.4.5	const_iterators . . . . .	45
4.4.6	Printing a container . . . . .	46
<b>5</b>	<b>Lecture 5: Lists, Stacks, Queues, and the STL</b>	<b>47</b>
5.1	Implementation of a vector . . . . .	47

5.2	Implementation of a list . . . . .	47
5.2.1	Insertion . . . . .	49
5.2.2	Erasure . . . . .	50
5.3	The Stack ADT . . . . .	50
5.3.1	Implementations . . . . .	51
5.3.2	Applications . . . . .	52
5.4	The Queue ADT . . . . .	52
5.4.1	Implementations . . . . .	53
5.4.2	Applications . . . . .	54
5.5	Trees . . . . .	54
<b>6</b>	<b>Lecture 6: Trees</b>	<b>54</b>
6.1	Preliminaries . . . . .	54
6.2	Implementation of trees . . . . .	56
6.3	Tree traversals . . . . .	57
6.3.1	Depth-first search . . . . .	58
6.4	Binary trees . . . . .	60
6.4.1	Implementation . . . . .	61
6.4.2	Expression trees . . . . .	61
6.4.3	Constructing an expression tree . . . . .	62
6.5	Binary search trees (BSTs) . . . . .	63
6.5.1	Operations on a BST . . . . .	63
6.5.2	Average-case analysis . . . . .	68
6.6	Balanced trees . . . . .	68
6.7	AVL trees . . . . .	69
6.7.1	Insertion cases . . . . .	69
6.7.2	When do you use rotations? Why? . . . . .	79
<b>7</b>	<b>Lecture 7: AVL trees (cont.), splay trees, and B-trees</b>	<b>79</b>
7.1	AVL tree implementation . . . . .	79
7.1.1	Insertion . . . . .	80
7.1.2	Balance . . . . .	80
7.1.3	Left-left single rotation (case 1) . . . . .	81
7.1.4	Left-right double rotation (case 2) . . . . .	81
7.1.5	Deletion . . . . .	81
7.2	Amortized cost . . . . .	81
7.3	Splay trees . . . . .	81
7.3.1	A simple idea (that doesn't work) . . . . .	82

7.3.2	Splaying . . . . .	84
7.3.3	Fundamental property of splay trees . . . . .	85
7.3.4	Summary . . . . .	85
7.4	B-trees . . . . .	86
7.4.1	Example . . . . .	88
7.4.2	Insertion . . . . .	88
7.4.3	Deletion . . . . .	89
<b>8</b>	<b>Lecture 08: Sets and maps</b>	<b>90</b>
8.1	STL containers . . . . .	90
8.2	The set container . . . . .	90
8.2.1	Insertion . . . . .	90
8.2.2	Insertion with hint . . . . .	91
8.2.3	erase . . . . .	91
8.2.4	find . . . . .	91
8.3	The map container . . . . .	92
8.3.1	Example . . . . .	93
8.4	Implementation of set/map in C++ . . . . .	93
8.4.1	Threaded binary trees . . . . .	93
8.4.2	Types of threaded binary trees . . . . .	94
8.4.3	An example . . . . .	95
8.5	Summary to avoid performance issues . . . . .	100
8.5.1	AVL trees . . . . .	100
8.5.2	Splay trees . . . . .	100
8.5.3	B-trees . . . . .	100

# 1 Lecture 1: C++ basics

## 1.1 C++ classes

```
class IntCell {
public:
    // This is a constructor
    IntCell() { _value = 0; }

    // This is another constructor. It takes an initial value.
    IntCell( int initialValue ) { _value = initialValue; }

    // This is an accessor (or getter) method.
    int read() { return _value; }

    // This is a setter method.
    void write( int newValue ) { _value = newValue; }

private:
    int _value;
};

int main() {
    // Using the first constructor w/o parameters
    IntCell x;
    int xValue = x.read(); // 0
    x.write(5);
    xValue = x.read();    // 5

    // Using the first constructor with a single parameter
    IntCell y(12);
    int yValue = y.read(); // 12
    y.write(y.read() * 2);
    yValue = y.read();    // 24
}
```

The keywords **public** and **private** determine the visibility of class members. A member that is **public** may be accessed by any method in any class while a member that is **private** may only be accessed by methods in the class it's declared. Data members are usually declared **private** to hide a class's internal details (i.e., **information hiding**).

## 1.2 Initialization list

Instead of initializing data members inside a constructor's body, an **initialization list** can be used to do so right in the constructor's signature:

```
class IntCell {
public:
    explicit IntCell( int initialValue = 0 ) : _value { initialValue }
    { }
    int read() const { return _value; }
    void write( int newValue ) { _value = newValue; }

private:
    int _value;
};
```

By default C++ does behind-the-scenes type conversions in all one-parameter constructors. Thus, to avoid this the keyword `explicit` is used.

With **implicit type conversion**,

```
IntCell obj;
obj = 37;
```

means

```
IntCell temp = 37;
obj = temp;
```

With **explicit type conversion** (by using `explicit`), a “type mismatch” compiler error is thrown instead. The use of `explicit` means a one-parameter constructor cannot be used to generate an implicit temporary object.

### 1.3 Constant member functions

A member function that examines but does not change the state of its object is an **accessor**. In the `IntCell`, `read` is an accessor and thus it's marked explicitly as an accessor by using the keyword `const` after the closing parenthesis that ends the parameter list. This signals that this method doesn't change the state of an `IntCell` and if it tries the compiler complains.

### 1.4 Interface and implementation

An **interface** lists the class and its members (both data and methods) and the **implementation** provides the implementations of the methods.

In C++, this is done by placing the interface in a `.h` file and the implementation in a `.cpp` file. Source code that requires knowledge of the interface must include the interface file (e.g., `#include Interf`

.h). To avoid including files multiple times, a few preprocessor commands are used in the interface file:

```
#ifndef INTCELL_H // if INTCELL_H is not defined...
#define INTCELL_H // ...define it.

// Here's the IntCell class's interface

#endif // end the definition
```

### 1.4.1 IntCell interface

```
// IntCell.h
#ifndef INTCELL_H
#define INTCELL_H

class IntCell {
public:
    explicit IntCell( int initialValue = 0 ) {}
    int read() const { return _value; }
    void write( int newValue ) { _value = newValue; }

private:
    int _value;
};

#endif
```

### 1.4.2 IntCell implementation

```
// IntCell.cpp
#include "IntCell.h" // we're making use IntCell's interf. so we include it.

IntCell::IntCell( int initialValue = ) : _value { initialValue } { }

// Return the stored value.
int IntCell::read() const { return _value; }

// Change the stored value.
void IntCell::write( int newValue ) { _value = newValue; }

};
```

Notice that we must use the **scope resolution operator** (`::`) to identify the class each member function belongs to. Otherwise, it's assumed that the function is in the **global scope**. Also note that the signature of an implemented member function *must* match exactly the signature in the class interface.



### 1.4.3 Using IntCell

```
// main.cpp
#include "IntCell.h"

int main() {
    IntCell m{};
    m.write(5);
    int val = m.read(); // 5

    return 0;
}
```

To compile the program, run:

```
g++ -o prog main.cpp IntCell.cpp
```

The header files (e.g., `IntCell.h`) aren't listed since they're included in the implementation files and thus are compiled too.

## 1.5 C++ Objects

In C++ objects are declared much like primitive types (e.g., `char`, `int`, `float`, etc.):

```
IntCell obj1;      // Zero parameter constructor
IntCell obj2(12);  // One parameter constructor
```

However, due to an effort to standardize the uniform initialization syntax using braces, it's recommended to declare objects as follows:

```
IntCell obj1;      // Zero parameter constructor
IntCell obj1{};    // same as before
IntCell obj2{12};  // One parameter constructor
```

## 1.6 vectors and strings

The C++ standard defines the classes `vector` and `string` that intend to replace the built-in C++ array. Unlike arrays, these classes behave like first-class objects. Thus, they can be copied with `=`, they remember how many items they can store, and their indexing operator check that the provided valid (to access some element) is valid.

### 1.6.1 vector

```
vector<int> numbers = {1, 2, 3, 4, 5}; // Or...
vector<int> numbers {1, 2, 3, 4, 5};
```

However, there's some ambiguity with the declaration. For instance, `vector<int> a(12)` declares a vector that stores 12 integers but `vector<int> a{12}` declares a vector of size 1 with the value 12 at position 0. This is because C++ gives precedence to the initializer list. If the intention is to declare a vector of size 12, use the old C++ syntax using parentheses: `vector<int> a(12)`.

Instead of using the regular for loop, vectors can be looped over using the following syntax:

```
int sum = 0;
vector<int> numbers = {1, 2, 3, 4, 5};
for( int num : numbers ) {
    sum += num;
}
```

The keyword `auto` can be used to let the compiler determine the type:

```
int sum = 0;
vector<int> numbers = {1, 2, 3, 4, 5};
for( auto num : numbers ) {
    sum += num;
}
```

However, keep in mind that this range syntax is only appropriate when 1) accessing elements sequentially and 2) when the index isn't needed.

## 1.7 C++ details

### 1.7.1 Pointers

A **pointer** is a variable that stores the address of a memory location (i.e., where an object resides). The syntax for declaring a pointer is as follows:

```
char* ptr_c;    // ptr_c is a pointer-to-char
int *ptr_i;     // ptr_i is a pointer-to-int
float *ptr_f;   // ptr_f is a pointer-to-float
void *u_ptr;    // ptr_u is a untyped pointer
```

Note that the position of the asterisk `*` (known as the **indirection operator** in this context) doesn't matter.

In order to obtain an object's address, the **address-of** operator & is used and a pointer variable is used to store it:

```
int num = 5;
int *ptr_num = &num; // get the address of num
int val = *ptr_num;   // get the value stored in the address ptr_num points
                      // to
*ptr_num = 6;         // now num is also 6. after all ptr_num to the same
                      // memory location
```

Thus, the indirection operator has two jobs: 1) it declares a pointer and 2) it dereferences a pointer.

### 1.7.2 Dynamic object allocation and garbage collection

Objects are created dynamically by allocating memory with **new** which returns a pointer to the newly created object:

```
IntCell *m;
m = new IntCell(); // OK
m = new IntCell{}; // OK
m = new IntCell;   // OK
```

If a pointer variable points at a class type, then a (visible) member of the object being pointed at can be accessed via the **->** operator:

```
int x = m->read();
m->write(6);
```

When an object that is allocated by **new** is no longer referenced, the **delete** operation must be applied to the object (through a pointer). Otherwise, the memory that it consumes is lost (until the program terminates). This is known as a **memory leak**.

```
delete m;
```

### 1.7.3 Lvalues, and Rvalues

An **lvalue** is an expression that identifies a non-temporary object:

```
vector<int> arr(3);
const int x = 2;
int y;
int z = x + y;
vector<string> *ptr = &arr;
```

The expressions `arr`, `x`, `y`, `z`, `ptr`, `*ptr`, and `z` are all lvalues.

An **rvalue** is an expression that identifies 1) a temporary object or 2) a value (e.g., literal constant) not associated with any object:

```
vector<int> arr(3);
const int x = 2;
int y;
int z = x + y;
vector<string> *ptr = &arr;
```

Here, `2`, `x+y`, and `&arr` are all rvalues.

### 1.7.4 References

A reference type allows us to define a new name for an existing value. In classic C++, a reference can generally only be a name for an lvalue, since having a reference to a temporary would lead to the ability to access an object that has theoretically been declared as no longer needed, and thus may have had its resources reclaimed for another object. However, in C++11, we can two types of references:

- An **lvalue reference** is declared by placing an `&` after some type. An lvalue reference then becomes a synonym for the object it references.

```
string str = "hello";
string &str_r = str;
str_r += " world"; // now str is "hello world"

string &lit_r = "hi"; // ILLEGAL: "hi" is not modifiable
string &bad_r = str + "!"; // ILLEGAL: str + "!" isn't an lvalue
```

- An **rvalue reference** is declared by placing an `&&` after some type. An rvalue reference has the same characteristics as an lvalue reference except that, unlike an lvalue reference, *an rvalue reference can also reference an rvalue* (i.e., a temporary).

```
string str = "hell";
string && bad1 = "hello"; // LEGAL
string && bad2 = str + ""; // LEGAL
string && sub = str.substr( 0, 4 ); // LEGAL
```

### 1.7.5 Uses of lvalue references

- **Aliasing complicated names.** We can rename objects that are too long and complicated to simpler names.

```

auto &whichList = theLists[ myhash( x, theLists.size( ) ) ];
if( find( begin( whichList ), end( whichList ), x ) != end( whichList ) )
{
    return false;
}
whichList.push_back( x );

```

Note that simply writing `auto whichList = theLists[ myhash( x, theLists.size( ) ) ];` wouldn't work because this would create a copy, and `whichList.push_back( x );` would be applied to the copy, not the original.

- **Range for loops.** By default, a range **for** loop cannot change the elements it iterates over, however taking a lvalue reference allows to modify those elements.

```

vector<int> numbers = {1, 2, 3};

for( auto &number : numbers ) {
    number++;
}

```

- **Avoiding a copy.** Given a function that returns an element of an array/vector, we could return a non-modifiable reference to that element instead of a copy. For instance, instead of

```

// @arr: A non-empty vector of strings.
// @return the maximum string in the @arr.
// Will abort() if @arr is empty.
string FindMax1(constvector<string>&arr) {
    if (arr.empty()) abort();
    int max_index=0;
    for (inti =1; i < arr.size(); ++i) {
        if(arr[max_index]<arr[i]) { max_index = i; }
    }
    return arr[max_index];
}

```

we could have

```

// @arr: A non-empty vector of strings.
// @return the maximum string in the @arr.
// Will abort() if @arr is empty.
const string &FindMax1(constvector<string>&arr) {
    if (arr.empty()) abort();
    int max_index=0;
    for (inti =1; i < arr.size(); ++i) {
        if(arr[max_index]<arr[i]) { max_index = i; }
    }
    return arr[max_index];
}

```

Syntax is needed in function declarations and returns to enable the passing and returning using references instead of copies. Notice the **const** keyword and & in the function's header.

## 2 Lecture 2: C++ Basics (cont.)

### 2.1 Parameter passing

Many languages, including C++, pass all parameters using **call-by-value**: the actual argument is copied into the formal parameter. However, this might be utterly inefficient if large complex objects are being passed since they're copied in their entirety. Historically C++ has had three ways to pass parameters:

- **call-by-value**: Useful to pass small objects that shouldn't be mutated by the function.

```
double average( double a, double b );
double x = 2.5, y = 3.5;
double z = average(x, y); // x and y remain unchanged

void swap( double a, double b ); // swap a and b
swap(x, y); // Not what's expected; x and y remain unchanged
```

- **call-by-reference (call-by-lvalue-reference)**: Useful for all type of objects that may be changed by the function.

```
double x = 2.5, y = 3.5;
swap( &a, &b );
swap(x, y); // now x = 3.5 and y = 2.5
```

- **call-by-constant-reference**: Useful for large objects that are expensive to copy and that must not be changed by the function. For this, the parameter is declared as a reference and the keyword **const** is used to signal that it cannot be modified.

```
string randomItem( const vector<string> &arr ); // return a random item in
arr
```

There's still another fourth way to pass parameters:

- **call-by-rvalue-reference**: Instead of copying a temporary object stored in an rvalue, a move is used; moving an object's state is easier than copying, as *it may involve just a simpler pointer change*. Functions know if a value is temporary or not based on their signature so the primary use for this type of parameter passing is overloading a function based on whether a parameter is an lvalue or rvalue.

```

string randomItem( const vector<string> &arr ); // return random item in
    lvalue arr (single &)
string randomItem(      vector<string> &&arr ); // return random item in
    rvalue arr (double &)

vector<string> v = {"hello", "world"};
string x = randomItem(v);                // invoke lvalue function
string x = randomItem({"hello", "world"}); // invoke rvalue function

```

This idiom is particularly useful for defining the behavior of = and in writing constructors.

## 2.2 Return passing

In C++, there are several mechanism for returning from a function:

- **return-by-value:** A copy of the object is returned which can be inefficient, however in C++11 return-by-value may still be efficient for large objects if the returned objects are rvalues.

```

double average( double a, double b );
LargeType randomItem( const vector<LargeType> &arr ); // potentially
    inefficient
vector<int> partialSum( const vector<int> &arr );      // efficient in C
    ++11

```

The following are two versions of the function `randomItem`. The second version avoids the creation of a temporary `LargeType` object, but only if the caller accesses it with a constant reference:

```

LargeType randomItem1( const vector<LargeType> &arr ) {
    return arr[ randomInt(0, arr.size() -1) ];
}

const LargeType & randomItem2( const vector<LargeType> &arr ) {
    return arr[ randomInt(0, arr.size() -1) ];
}

vector<LargeType> vec;
LargeType item1 = randomItem2(vec);           // copy
LargeType item2 = randomItem2(vec);           // copy
const LargeType &item3 = randomItem2(vec);    // no copy
auto &item4 = randomItem2(vec);               // no copy

```

- **return-by-constant-reference:** This avoid creating an immediate copy of an object. However, the caller must also use a constant reference to access the return value; otherwise, there will be still a copy. What does the constant reference mean? It means that we don't want to allow changes to be made by the caller by using the return value.

```
const LargeType & randomItem2( const vector<LargeType> &arr ) {  
    return arr[ randomInt(0, arr.size() -1) ];  
}  
  
vector<LargeType> vec;  
const LargeType &item3 = randomItem2(vec); // no copy
```

- **return-by-reference:** A reference is returned and the caller can modify the returned value. This is used in a few places to allow the caller of a function to have modifiable access to the internal data representation of a class.

## 2.3 C++11 `std::swap` and `std::move`

Given that copying large objects is expensive, C++11 allows the programmer to easily replace expensive copies with moves provided the object's class supports moves.

Take the following example of a `swap` function that swap its arguments by three copies:

```
void swap( vector<string> &x, vector<string> &y ) {  
    vector<string> tmp = x;  
    x = y;  
    y = tmp;  
}
```

In C++11, if the right-hand side of the assignment operator (or constructor) is an rvalue, then if the object supports moving, we can automatically avoid copies. In the example above, we know that `vector` supports moving so instead of copy operations we could do move operations. These could be done either by casting the right-hand side of an assignment to an rvalue reference or by using `std::move`.

```
// Using type-casting  
void swap( vector<string> &x, vector<string> &y ) {  
    vector<string> tmp = static_cast<vector<string>> &&>( x );  
    x = static_cast<vector<string> &&>( y );  
    y = static_cast<vector<string> &&>( tmp );  
}  
  
// Using std::move, equivalent to casting but more succinct  
void swap( vector<string> &x, vector<string> &y ) {  
    vector<string> tmp = std::move(x);  
    x = std::move(y);  
    y = std::move(tmp);  
}
```

**NOTE:** `std::move` doesn't move anything; rather, it makes a value (either lvalue or rvalue) subject to



be moved.

It's worth noting that `std::swap` is already part of STL and works for any part:

```
vector<string> x;  
vector<string> y;  
std::swap(x, y); // x contains y's contents and y contains x's contents
```

## 2.4 The big five

In C++, classes come with five special functions already written for each class. These are the **destructor**, **copy constructor**, **move constructor**, **copy assignment operator**, and **move assignment operator**. In many cases, you accept the default behavior provided by the compiler can be relied on. Sometimes you cannot.

Let's assume the following interface for the `IntCell` class:

```
#ifndef INTCELL_H  
#define INTCELL_H  
  
class IntCell {  
public:  
    explicit IntCell( int initialValue = 0 ) { value_ = new int{  
        initialValue}; }  
    int read() const {}  
    void write( int x ) {}  
private:  
    int *value_;  
};  
  
#endif
```

### 2.4.1 Destructor

This function is called whenever an object goes out of scope or is subjected to a `delete` operation. Typically, the only responsibility of the destructor is to free up any resources that were acquired during the use of the object. This includes calling `delete` for any corresponding `new`s, closing any files that were opened, and so on.

```
IntCell::~IntCell() {  
    delete value_;  
}
```

### 2.4.2 Copy constructor and move constructor

These two constructors are required to construct a new object, initialized to the same state as another object of the same type.

- If the existing object is an **lvalue**, then it's a **copy constructor**.
- If the existing object is an **rvalue**, then it's a **move constructor**.

```
// Copy constructor, the parameter is an lvalue of the same type.
IntCell::IntCell( const IntCell &rhs ) {
    stored_ = new int{*rhs.value_};
}

// Move constructor, the parameter is an rvalue of the same type.
IntCell::IntCell( IntCell &&rhs ) : value_{rhs.value_} {
    rhs.value_ = nullptr;
}
```

#### When is either constructor called?

- During a declaration with initialization.

```
IntCell B = C;    // Copy constructor if C is lvalue; Move constructor if C
                  is rvalue
IntCell B { C }; // same as above
```

- An object passed using call-by-value (rarely done).
- An object returned by value.

### 2.4.3 Copy assignment and move assignment (operator =)

The assignment operator is called when = is applied to two objects that have both been previously constructed. Given the expression `lhs = rhs`, then the state of `rhs` (right hand side) is copied into `lhs` (left hand side).

- If `rhs` is an **lvalue**, this is done using the copy assignment operator.
- If `rhs` is an **rvalue**, this is done using the move assignment operator.

```
// Copy assignment operator
IntCell & IntCell::operator=( const IntCell &rhs ) {
    if (this &= &rhs) {
        *value_ = *rhs.value_;
    }
    return *this;
}

// Move assignment operator
IntCell & IntCell::operator=( IntCell &&rhs ) {
    std::swap(value_, rhs.value_);
    return *this;
}
```

In C++11, the copy assignment operator could be implemented more idiomatically with a copy-and-swap idiom:

```
// Copy assignment operator
IntCell & IntCell::operator=( const IntCell &rhs ) {
    IntCell copy = rhs; // calls the copy constructor
    std::swap(*this, copy);
    return *this;
}
```

#### 2.4.4 Defaults

It's often the case that the defaults “big five” are perfectly acceptable, so nothing need to be done. If a class consists of data members that are exclusively primitive types and objects for which the defaults make sense, the class defaults will usually make sense. Thus a class whose data members are **int**, **double**, **string**, and even **vector<string>** can accept defaults.

**When do defaults fail?** The defaults fail in a class that contains a data member that is a pointer:

- The default destructor does nothing to data members that are pointers. It's imperative that we **delete** them ourselves.
- The copy constructor and copy assignment operator both copy the value of the pointer (i.e., a memory address) rather than the object being pointed at. This means we end up with two class instances that contain pointers that point to the same objects being pointed at. This is known as **shallow copy**. However, we would typically expect a **deep copy**, in which a clone of the entire object is created.

Thus, as a result, when a class contains pointers as data members, and deep semantics are important, we typically must implement the destructor, copy assignment, and copy constructors ourselves. Doing so removes the move defaults, so we also must implement move assignment and the move constructor. As a general rule, either you accept the default for all five operations, or you should declare all five, and explicitly define, default (use the keyword **default**), or disallow each (use the keyword **delete**). Generally we will defined all five.

**NOTE:** If you write any of the big-five, it would be good practice to explicitly consider all the others, as the defaults may be invalid or inappropriate. Changing a default method is an all or nothing situation.

**When don't the defaults work?** The most common situation in which the defaults do not work occurs when a data member is a pointer type and the pointer is allocated by some object member function (such as the constructor). A problem that might arise is the default copy assignment and copy constructor doing shallow copies. Another problem is **memory leak**; an object allocated by a constructor might remain allocated and it's never reclaimed.

## 2.5 Templates

An algorithm is **type independent** if the logic of the algorithm does not depend on the type of items it handles. When we write C++ code for a type-independent algorithm or data structure, it's preferable to write the code once rather than recode it for each different type. This is accomplished by using **templates**.

### 2.5.1 Function templates

A **function template** isn't an actual function, but instead a pattern for what could become a function. For example, the following is a function template:

```
/*
Return the maximum item in array a.
Assumes a.size() > 0.
Comparable object must provide operator< and operator=.
*/

template <typename Comparable>
const Comparable & findMax( const vector<Comparable> &a ) {
    int maxIndex = 0;
    for (int i = 1; i < a.size(); i++) {
        if (a[maxIndex] < a[i]) {
            maxIndex = i;
        }
    }
    return a[maxIndex];
}
```

Using the function template would look as follows:

```
vector<int> v1 = { 37, 12, 1, 89 };
vector<double> v2 = { 78.1, 89.8, 12.4, 1.1 };
vector<string> v3 = { 'hi', 'ha', 'ho', 'he' };
vector<IntCell> v4(75);

findMax(v1); // OK: Comparable = int
findMax(v2); // OK: Comparable = double
findMax(v3); // OK: Comparable = string
findMax(v4); // Illegal; IntCell doesn't implement operator< and thus not
               a Comparable
```

Thus, one of the limitations of function templates is that object passed to the function need to implement whatever operator the function uses internally. However, instead of relying on “hardcoded” conditions, a function that perform the required operation could be passed alongside the object. This type of functions are known as **function objects**.

## 2.6 Function objects

Previously we implemented `findMax` as a function template but it was only limited to objects that have an `operator<` function defined. Instead, we need to rewrite `findMax` to accept as parameters an array of object and a comparison function that explains how to decide which of two objects is the larger and which is the smaller. Instead of relying on the array objects knowing how to compare themselves, we completely decouple this information from the object in the arrays.

**How do we pass a functions as parameters though?** One way is to define a class with no data and one member function, and pass an instance of the class. In fact, the function is passed by placing it

inside an object, which is known as a **function object**.

The following is the simplest implementation of the function object idea:

```
#include <iostream>
#include <vector>
#include <string>

// Generic findMax, with a function object.
// Precondition: a.size() > 0.
// Comparator object is assumed to implement the isLessThan method.
template <typename Object, typename Comparator>
const Object & findMax( const std::vector<Object> &arr, Comparator cmp ) {
    int maxIndex = 0;

    for (int i = 1; i < arr.size(); i++) {
        if (cmp.isLessThan(arr[maxIndex], arr[i])) {
            maxIndex = i;
        }
    }

    return arr[maxIndex];
};

class StringComparisonByLength {
public:
    bool isLessThan( const std::string &lhs, const std::string &rhs )
    {
        return lhs.length() < rhs.length();
    }
};

class StringComparisonCaseInsensitive {
public:
    bool isLessThan( const std::string &lhs, const std::string &rhs )
    {
        return std::strcasecmp(lhs.c_str(), rhs.c_str()) < 0;
    }
};

int main() {
    std::vector<std::string> greetings = {"hi", "hello", "bonjour", "HOLA"};

    std::cout << findMax( greetings, StringComparisonByLength{} ) << "\n";
    std::cout << findMax( greetings, StringComparisonCaseInsensitive{} )
        << "\n";

    return 0;
}
```

C++ function objects are implemented using this basic idea, but with some fancy syntax. First, instead of using a function with a name, we use operator overloading. Instead of the function being `isLessThan`, it is `operator()`. Second, when invoking `operator()`, `cmp.operator()(x,y)` can be shortened to `cmp(x,y)`. Third, we can provide a version of `findMax` that works without a function object that uses a default ordering; the implementation uses the Standard Library function object template `less` (defined in header file `functional`) to generate a function object that imposes the normal default ordering.

The following is the more idiomatic implementation:

```
#include <iostream>
#include <vector>
#include <string>

// Generic findMax, with a function object.
// Precondition: a.size() > 0.
template <typename Object, typename Comparator>
const Object & findMax( const std::vector<Object> &arr, Comparator
    isLessThan ) {
    int maxIndex = 0;

    for (int i = 1; i < arr.size(); i++) {
        if (isLessThan(arr[maxIndex], arr[i])) {
            maxIndex = i;
        }
    }

    return arr[maxIndex];
};

// Generic findMax, using default ordering.
const Object & findMax( const vector<Object> &arr ) {
    return findMax(arr, less<Object>{} );
}

class StringComparisonByLength {
public:
    bool operator()( const std::string &lhs, const std::string &rhs )
    {
        return lhs.length() < rhs.length();
    }
};

class StringComparisonCaseInsensitive {
public:
    bool operator()( const std::string &lhs, const std::string &rhs )
    {
        return std::strcasecmp(lhs.c_str(), rhs.c_str()) < 0;
    }
};

int main() {
    std::vector<std::string> greetings = {"hi", "hello", "bonjour", "HOLA"};

    std::cout << findMax( greetings, StringComparisonByLength{} ) << "\n";
    std::cout << findMax( greetings, StringComparisonCaseInsensitive{} )
        << "\n";
    std::cout << findMax( greetings ) << "\n";

    return 0;
}
```



### 2.6.1 Class templates

In its most simplest form, a class template works much like a function template. The following `MultiCell` is an implementation that is like `IntCell`, but works for any type `Object`, provided that `Object` has a zero-parameter constructor, a copy constructor, and a copy assignment operator.

```
template <typename Object>
class MemoryCell {
public:
    explicit MemoryCell( const Object &initialValue = Object{} ) :
        _value{initialValue} {}
    const Object & read() const { return _value; }
    void write( const Object &x ) { _value = x; }

private:
    Object _value;
};
```

Notice that

- `Object` is passed by constant reference.
- the default parameter for the constructor is not 0, because 0 might not be a valid `Object`. Instead, the default parameter is the result of constructing an `Object` with its zero-parameter constructor.

If we implement class templates as a single unit, then there is very little syntax baggage. Many class templates are, in fact, implemented this way because, currently, separate compilation of templates does not work well on many platforms. **Therefore, in many cases, the entire class, with its implementation, must be placed in a .h file.** Popular implementations of the STL follow this strategy.

### 2.7 Using matrices

This will be implemented by using a vector of vectors (e.g., a vector of `int` vectors).

```
#ifndef MATRIX_H
#define MATRIX_H

#include <vector>

template <typename Object>
class Matrix {
public:
    /*
     * Create _array as having rows entries each of type vector<Object>
     * >.
     * We have a rows zero-length vectors of Object so each row is
     * resized
     * to have cols columns. Thus this creates a two-dimensional array
     * .
     */
    Matrix( int rows, int cols ) : _array(rows) {
        for (auto &row : _array) {
            row.resize(cols);
        }
    }

    Matrix( std::vector<std::vector<Object>> v ) : _array{ v } { }

    Matrix( std::vector<std::vector<Object>> &&v ) : _array{ std::move
        (v) } { }

    /*
     * Array indexing operator. This returns the row (a vector<Object>
     * at
     * the index row.
     */
    const std::vector<Object> & operator[]( int row ) const {
        return _array[row]
    }

    int numrows() const {
        return _array.size();
    }

    int numcols() const {
        return numrows() ? _array[0].size() : 0;
    }

private:
    /*
     * A matrix is represented by _array, a vector of vector<Object>.
     */
    std::vector<std::vector<Object>> _array;
};

#endif
```

## 2.8 Class problems

1. What's  $1 + 2 + 3 + \dots + n$ ? Prove it.

$$1 + 2 + 3 + \dots + n = \frac{n(n+1)}{2}$$

**Proof by induction:**

- **Base case:** When  $n = 1$ , then  $1 = 1(1 + 1)/2 = 1$ .
- **Inductive step:** Since  $1 = (1 * 2)/2$ , the statement  $P(1)$  is true. Assume that  $P(k)$  is true for an arbitrary positive integer  $k$ . We show that  $P(k + 1)$  is true. In other words,

$$1 + 2 + 3 + \dots + (k + 1) = \frac{(k + 1)(k + 2)}{2}$$

Thus

$$\begin{aligned} 1 + 2 + 3 + \dots + (k + 1) &= (1 + 2 + 3 + \dots + k) + (k + 1) \\ &= \frac{k(k + 1)}{2} + (k + 1) \\ &= \frac{k(k + 1)}{2} + \frac{2(k + 1)}{2} \\ &= \frac{(k + 1)(k + 2)}{2} \end{aligned}$$

Therefore, by the principle of mathematical induction,  $P(n)$  is true for every positive integer  $n$ .

2. What's  $1 + 2 + 4 + \dots + 2^n$ ? Prove it.

$$\begin{aligned} n = 0 &\Rightarrow 1 \Rightarrow 1 \\ n = 1 &\Rightarrow 1 + 2 \Rightarrow 3 \\ n = 2 &\Rightarrow 1 + 2 + 4 \Rightarrow 7 \\ n = 3 &\Rightarrow 1 + 2 + 4 + 8 \Rightarrow 15 \\ &\dots \\ n = n &\Rightarrow 1 + 2 + \dots + 2^n \Rightarrow 2^{(n+1)} - 1 \end{aligned}$$

Thus,  $1 + 2 + 4 + \dots + 2^n = 2^{(n+1)} - 1$ .

**Proof by induction:**

- **Base case:** When  $n = 0$ , then  $1 = 2^{(0+1)} - 1 = 1$ . When  $n = 1$ , then  $2^0 + 2^1 = 3 = 2^{(1+1)} - 1 = 3$ .
- **Inductive step:** Since  $1 = 2^{(0+1)} - 1$ , the statement  $P(1)$  is true. Assume that  $P(k)$  is true for an arbitrary positive integer  $k$ . We show that  $P(k + 1)$  is true. In other words, we show that

$$1 + 2 + 4 + \dots + 2^{(k+1)} = 2^{(k+2)} - 1.$$

Thus,

$$\begin{aligned} 1 + 2 + 4 + \dots + 2^{(k+1)} &= (1 + 2 + 4 + \dots + 2^k) + 2^{(k+1)} \\ &= 2^{(k+1)} - 1 + 2^{(k+1)} \\ &= 2^{(k+1)} + 2^{(k+1)} - 1 \\ &= 2 \cdot 2^{(k+1)} - 1 \\ &= 2^{(k+2)} - 1 \end{aligned}$$

Therefore, by the principle of mathematical induction,  $P(n)$  is true for every non-negative integer  $n$ .

3. If  $A_0 = 1$  and  $A_n = 2A_{(n-1)} + 1$ , what's a closed formula for  $A_n$ ? Prove it.

$$A_1 = 2 \cdot A_0 + 1 = 3$$

$$A_2 = 2 \cdot A_1 + 1 = 7$$

$$A_3 = 2 \cdot A_2 + 1 = 15$$

$$A_4 = 2 \cdot A_3 + 1 = 31$$

...

$$A_n = 2 \cdot A_{(n-1)} + 1 = 2^{(n+1)} - 1$$

### 3 Lecture 3: Algorithm analysis

An **algorithm** is a clearly defined set of simple instructions which must be followed in order to solve a particular problem.

#### 3.1 Mathematical background

Algorithm analysis is grounded on mathematics and the definitions below set up a formal framework to study algorithms. These definitions try to establish a **relative order among functions**. Given two functions, there are usually points where one function is smaller than the other so it doesn't make sense to claim, for instance,  $f(N) < g(N)$ . Instead, we compare their **relative rates of growth**.

**Big-Oh**  $T(N) = O(f(N))$  if there are positive constants  $c$  and  $n_0$  such that  $T(N) \leq c \cdot f(N)$  when  $N \geq n_0$ .

Informally, the growth rate  $T(N)$  is less than or equal to that  $f(N)$ .

**Big-Omega**  $T(N) = \Omega(g(N))$  if there are positive constants  $c$  and  $n_0$  such that  $T(N) \geq c \cdot g(N)$  when  $N \geq n_0$ .

Informally, the growth rate  $T(N)$  is greater than or equal to that  $g(N)$ .

**Big-Theta**  $T(N) = \Theta(h(N))$  **if and only if**  $T(N) = O(h(N))$  and  $T(N) = \Omega(h(N))$ .

Informally, the growth rate  $T(N)$  equals the growth rate of  $h(N)$ . This means that  $T(N)$  is eventually *sandwiched* between two different constant multiples of  $h(N)$ .

**Little-Oh**  $T(N) = o(p(N))$  if, for all positive constants  $c$ , there exist an  $n_0$  such that  $T(N) < c \cdot p(N)$  when  $N > n_0$ .

Informally, the growth rate  $T(N)$  is less than the growth rate of  $f(N)$ . This is different from Big-Oh, given that Big-Oh allows the possibility that the growth rates are the same.

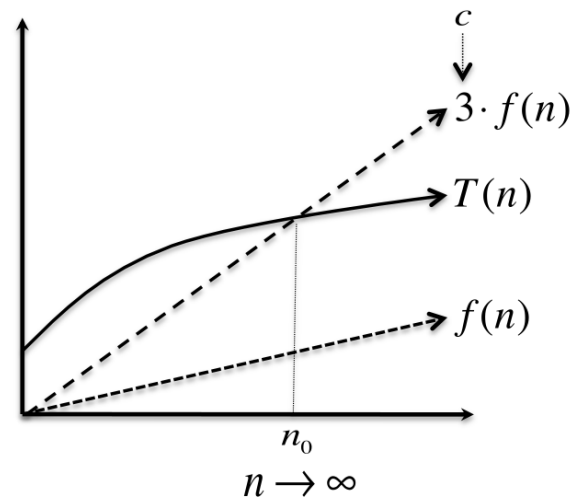
**NOTE:** In the above definitions, both  $c$  and  $n_0$  are constants and they cannot depend on  $N$ . If you see yourself saying "take  $n_0 = N$ " or "take  $c = \log_2 N$ " in an alleged Big-Oh proof, then you need to start with choices of  $c$  and  $n_0$  that are independent of  $N$ .

#### 3.2 Notation

##### 3.2.1 Big-Oh notation

When we say that  $T(N) = O(f(N))$ , we're guaranteeing that the function  $T(N)$  grows at a rate no faster than  $f(N)$ ; thus  $f(N)$  is an **upper bound** on  $T(N)$ . Alternatively, we could say that  $T(N)$  is

**bounded above** by a constant multiple of  $f(N)$ .



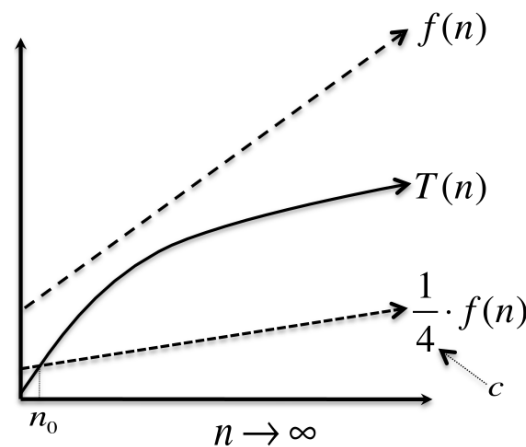
**Figure 2.1:** A picture illustrating when  $T(n) = O(f(n))$ . The constant  $c$  quantifies the “constant multiple” of  $f(n)$ , and the constant  $n_0$  quantifies “eventually.”

### Figure 1: Big-Oh

For example,  $T(N) = O(N^2)$  means that the function  $T(N)$  grows at a rate no faster than  $N^2$ ; its growth could equal that of  $N^2$  but it could never surpass it.

#### 3.2.2 Big-Omega notation

When we say that  $T(N) = \Omega(g(N))$ , we’re guaranteeing that the function  $T(N)$  grows at a rate no lower than  $g(N)$ ; thus  $g(N)$  is a **lower bound** on  $T(N)$ . Alternatively, we could say that  $T(N)$  is **bounded below** by constant mutiple of  $g(N)$ .



$T(n)$  again corresponds to the function with the solid line. The function  $f(n)$  is the upper dashed line. This function does not bound  $T(n)$  from below, but if we multiply it by the constant  $c = \frac{1}{4}$ , the result (the lower dashed line) does bound  $T(n)$  from below for all  $n$  past the crossover point at  $n_0$ . Thus  $T(n) = \Omega(f(n))$ .

**Figure 2:** Big-Omega

For example,  $T(N) = \Omega(N^2)$  means that the function  $T(N)$  grows at a rate no lower than  $N^2$ ; its growth could equal that of  $N^2$  but it could never drop below it.

Whenever talking about Big-Oh, there's always an implication about Big-Omega. For instance,  $T(N) = O(f(N))$  (i.e.,  $f(N)$  is an **upper bound** on  $T(N)$ ) implies that  $f(N) = \Omega(T(N))$  (i.e.,  $T(N)$  is a **lower bound** on  $f(N)$ ). As an example,  $N^3$  grows faster than  $N^2$ , so we can say that  $N^2 = O(N^3)$  or  $N^3 = \Omega(N^2)$ .

### 3.2.3 Big-Theta notation

When we say that  $T(N) = \Theta(h(N))$ , we're guaranteeing that the function  $T(N)$  grows at the same rate as  $h(N)$ . However, when two functions grow at the same rate, then the decision of whether or not to signify this with  $\Theta()$  can depend on the context.

## 3.3 Typical growth rates

Function	Name
$c$	Constant
$\log_2 N$	Logarithmic
$\log^2 N$	Log-squared
$N$	Linear
$N \log N$	
$N^2$	Quadratic
$N^3$	Cubic
$2^N$	Exponential

### 3.4 Rules

**Rule 1** If  $T_1(N) = O(f(N))$  and  $T_2(N) = O(g(N))$ , then

1.  $T_1(N) + T_2(N) = O(f(N) + g(N))$ . Less formally it is  $O(\max(f(N), g(N)))$ , and
2.  $T_1(N) \cdot T_2(N) = O(f(N) \cdot g(N))$

**Rule 2** If  $T(N)$  is a polynomial of degree  $k$ , then  $T(N) = \Theta(N^k)$ .

**Rule 3**  $\log^k N = O(N)$  for any constant  $k$ . In other words, logarithms run at a lower rate than linear functions which means logarithms grow very slowly.

### 3.5 Few points

- In simple terms, the goal of asymptotic notation is *to suppress constants factors and lower-order*. Thus, they aren't included in Big-Oh. For instance, don't say  $T(N) = O(2N^2)$  or  $T(N) = O(N^2 + N)$ . In both cases, the correct form is  $T(N) = O(N^2)$ .

When analyzing the running time of an algorithm, why would we want to throw away information like constant factors and lower-order terms? 1) Lower-order terms become increasingly irrelevant as you focus in large inputs, which are the inputs that require algorithmic ingenuity and 2) constant factors are generally highly dependent on the details of the environment (e.g., programming language, architecture, compiler, etc.) and thus ditching them allows to generalize by not committing ourselves to a specific programming language, architecture, etc.

- The relative growth rates of two functions  $f(N)$  and  $g(N)$  can always be determined by computing  $\lim_{N \rightarrow \infty} \frac{f(N)}{g(N)}$ , using [L'Hopital's rule](#) if necessary. The limit can have four possible values:



- The limit is 0, meaning that  $f(N) = O(g(N))$ .
  - The limit is  $c \neq 0$ , meaning that  $f(N) = \Theta(g(N))$ .
  - The limit is  $\infty$ , meaning that  $g(N) = o(f(N))$ .
  - The limit doesn't exist and thus no relation exist.
- Stylistically it's bad to say that  $f(N) \leq O(g(N))$  since the inequality is already implied by the definition of Big-Oh.
  - It's wrong to write  $f(N) \geq O(g(N))$ , because it doesn't make sense.

### 3.6 Model of computation

The model of computation is basically a normal computer with the following characteristics:

- Instructions are executed sequentially.
- The model has the standard repertoires of simple instructions, such as addition, multiplication, comparison, and assignment. Unlike real computers, this computer takes exactly one time unit to do anything.
- The computer has fixed-size (e.g., 32-bit) integers and no fancy operations (e.g., matrix inversion, sorting, etc.).
- The computer has infinite memory.

### 3.7 What to analyze

The most important resource to analyze is generally the *running time* and typically, the size of the input is the main consideration. We define two functions,  $T_{avg}(N)$  (for the average-case running time) and  $T_{worst}(N)$  (for the worst-case running time) used by an algorithm on input of size  $N$ . It's worth noting that  $T_{avg}(N) \leq T_{worst}(N)$  (i.e.,  $T_{avg}(N)$  has a lower rate growth than that of  $T_{worst}(N)$ ).

#### 3.7.1 Types of performance

**Best-case performance** Although it can be occasionally analyzed, it's often of little interest since it doesn't represent typical behavior.

**Average-case performance** It often reflects typical behavior, however it doesn't provide a bound for all input and it can be difficult to compute. Furthermore, it's also hard to define what's the average input

**Worst-case performance** It represents a guarantee for performance on any possible input. This is the quantity generally required because it provides a bound for all input.

### 3.8 Running-time calculations

When computing a Big-Oh running time, there are several general rules:

**For loops** The running time of a **for** loop is *at most* the running time of the statements inside the **for** loop times the number of iterations.

**Nested loops** Analyze these inside out. The total running time of a statement inside a group of nested loops is the running time of the statement multiplied by the product of the sizes of all the loops.

For example, the following program is  $O(N^2)$  because the inner loop does  $N$  iterations and the outer loops does  $N$  too. This amounts to  $N \times N = N^2$  iterations. The first assignment counts as one operation (thus  $O(N)$ ) while the innermost statement counts for 2 operations (1 multiplication and 1 assignment). Thus, to be more precise the program's running time is  $1 + 2N^2$ , however Big-Oh suppresses constant factors and lower-order terms.

```
k ← 0
for [1, n] → i:
  for [1, n] → j:
    k ← i * j
```

**Consecutive statements** They just add (which means that the maximum is the one that counts).

For example, the following program fragment, which has  $O(N)$  work followed by  $O(N^2)$  work, is ultimately  $O(N^2)$  which dominates  $O(N)$ :

```
for [0, n) → i:
  a[i] = 0

for [0, n) → i:
  for [0, n) → j:
    a[i] += a[j] + i + j
```

**If/else** For the fragment **if** condition { S1 } **else** { S2 }, the running time of an **if/else** statement is never more than the running time of the test plus the larger of the running times of S1 and S2.

## 3.9 Sample problems

### 3.9.1 Sum of cubes

Input: a positive integer  $n$ .

Output: the sum of all cubes from 1 to  $n^3$ .

```
SumOfCubes(n):  
    sum ← 0                # 0(1)  
    for [1, n] → i:        # 0(n)  
        sum += i * i * i    # 0(4), 1 assignment, 1 addition and 2 products  
    return sum             # 0(1)
```

Thus, for  $\text{SumOfCubes}(N) = 1 + 4n + 1 = 4n + 2 = O(n)$ . We discard both the constant factors and the lower-order terms.

### 3.9.2 Factorial

The factorial is defined as

$$n! = n \times (n - 1) \times (n - 2) \times \cdots \times 2 \times 1$$

Implemented recursively, the algorithm is as follows:

Input: a non-negative integer  $n$

Output: the factorial of  $n$

```
Factorial(n):  
    if n ≤ 1:  
        return 1  
    else:  
        return n * Factorial(n-1)
```

However, this is a thinly veiled **for** loop. In this case, the analysis involves the recurrence relation  $T(n) = 1 + T(n - 1)$  for  $n > 1$ ,  $T(1) = 2$  that needs to be solved:

$$\begin{aligned}
T(N) &= 1 + T(n-1) \\
&= 1 + (1 + T(n-2)) \\
&= 1 + (1 + (1 + T(n-3))) \\
&= \dots \\
&= 1 + (1 + (1 + T(n-k))) \text{ } k \text{ 1's}
\end{aligned}$$


---

$$\begin{aligned}
T(N) &= k + T(n-k) \\
&= (n-1) + T(n-(n-1)) \\
&= (n-1) + T(1) = (n-1) + 2 = n + 1
\end{aligned}$$

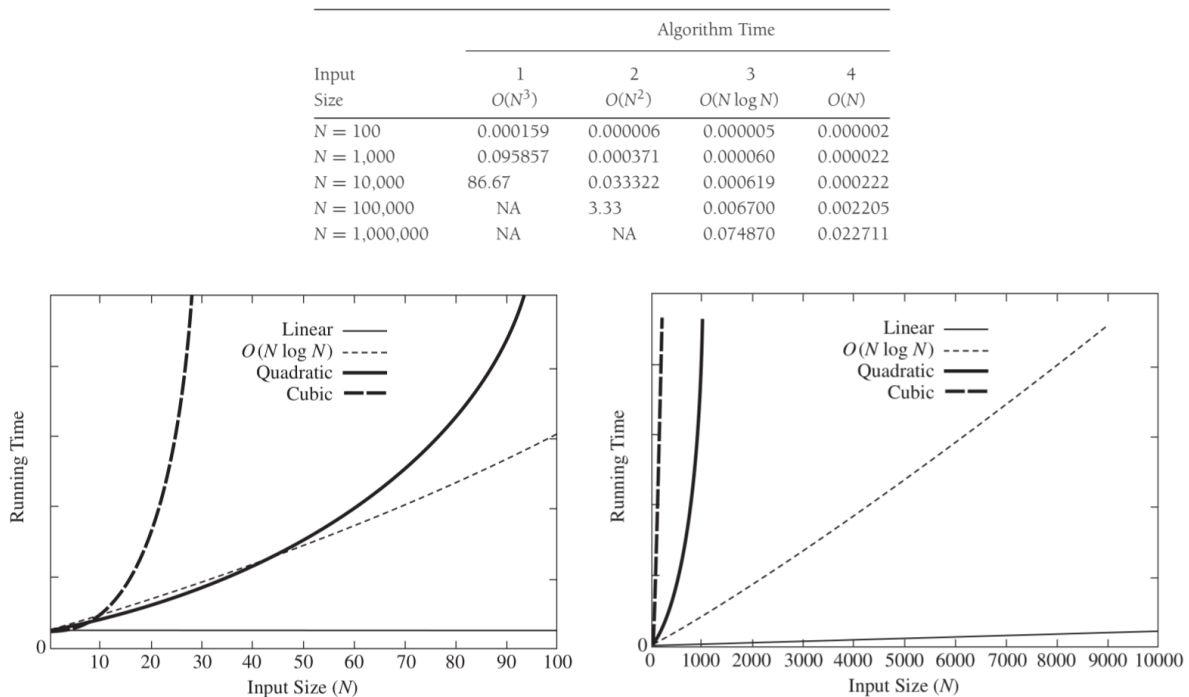
Thus,  $T(n) = O(n)$ .

### 3.9.3 Maximum subsequence sum problem

Given (possibly negative) integers  $A_1, A_2, \dots, A_N$ , find the maximum value of  $\sum_{k=1}^i A_k$ .

In other words, given a one-dimensional array of numbers, find a contiguous subarray with the largest sum. For example, with the input  $-2, 11, -4, 13, -5, -2$  the answer is  $11 - 4 + 13 = 20$ .

This problem is interesting mainly because there are many algorithms to solve it, and the performance of these algorithms varies drastically.



**Figure 3:** Plot (N vs. time) of various algorithms

For algorithm 4 (linear), as the problem size increases by a factor of 10, so does the running time. For algorithm 3 (quadratic), a tenfold increase in input size yields roughly hundredfold ( $10^2$ ) increase in running time.

### 3.9.4 Algorithm 1 (brute force)

implementations/max-subsequence-sum:

Input: array A of integers.  
Output: the maximum positive subsequence sum.

```

MaxSubsequenceSum(A):
    maxSum ← 0
    for [0, A.size) → i:
        for [i, A.size) → j:
            currentSum ← 0
            for [i, j] → k:
                currentSum += A[k]
            if currentSum > maxSum:
                maxSum = currentSum
    return maxSum

```

For this algorithm, there's  $N$  starting places, average  $\frac{N}{2}$  lengths to check, and average  $\frac{N}{4}$  numbers to add. Thus we have  $O(N^3)$ .

### 3.9.5 Algorithm 2 (brute force)

[implementations/max-subsequence-sum:](#)

```
Input: array A of integers.
Output: the maximum positive subsequence sum.

MaxSubsequenceSum( A ):
    maxSum ← 0
    for [0, A.size) → i:
        currentSum ← 0
        for [i, A.size) → j:
            currentSum += A[j]
            if currentSum > maxSum:
                maxSum = currentSum
    return maxSum
```

For this algorithm, the innermost **for** loop has been removed. There's  $N$  starting places and an average  $\frac{N}{4}$  numbers to add. Thus we have  $O(N^2)$ .

### 3.9.6 Algorithm 3

[implementations/max-subsequence-sum:](#)

Input: array A of integers.

Output: find the maximum sum **in** subarray spanning A[LEFT..RIGHT]. It does not attempt to maintain actual best sequence.

```
MaxSumRec( A, LEFT, RIGHT ):
    # Base case
    if LEFT = RIGHT:
        if A[LEFT] > 0:
            return A[LEFT]
        else
            return 0

    center ← (LEFT + RIGHT) div 2
    maxLeftSum ← MaxSumRec(A, LEFT, center)
    maxRightSum ← MaxSumRec(A, center + 1, RIGHT)

    maxLeftBorderSum ← 0
    leftBorderSum ← 0

    for [center, LEFT] → i:
        leftBorderSum += A[i]
        if leftBorderSum > maxLeftBorderSum:
            maxLeftBorderSum ← leftBorderSum

    maxRightBorderSum ← 0
    rightBorderSum ← 0

    for [center + 1, RIGHT] → j:
        rightBorderSum += A[j]
        if rightBorderSum > maxRightBorderSum:
            maxRightBorderSum ← rightBorderSum

    return max(maxLeftSum, maxRightSum, maxLeftBorderSum +
              maxRightBorderSum)

# Driver for divide-and-conquer maximum contiguous subsequence sum
algorithm.
MaxSubsequenceSum( A ):
    return MaxSumRec(A, 0, A.size - 1)
```

### 3.9.7 Algorithm 4

implementations/max-subsequence-sum:

Input: array A of integers.  
Output: the maximum positive subsequence sum.

```
MaxSubsequenceSum( A ):  
    maxSum ← 0  
    currentSum ← 0  
  
    for [0, A.size] → j:  
        currentSum += A[j]  
  
        if currentSum > maxSum:  
            maxSum ← currentSum  
        else if currentSum < 0:  
            currentSum = 0  
  
    return maxSum
```

## 4 Lecture 4: Algorithm analysis (cont.) and lists/stacks/queues

### 4.1 Binary search

Given an integer  $X$  and integers  $A_0, A_1, \dots, A_{N-1}$ , which are **presorted** already in memory, find index  $i$  such that  $A_i = X$ , or return  $-1$  if  $X$  is not in the group of integers.

implementations/binary-search

Input: array A of elements and item x which we are searching **for**.  
Output: item where item is found or **-1** **if** not found.

```
BinarySearch( A, x ):  
    low ← 0  
    high ← A.size - 1  
    while low ≤ high:  
        mid ← (low + high) / 2  
        if A[mid] < x:  
            low = mid + 1  
        else if A[mid] > x:  
            high = mid - 1  
        else:  
            return mid  
    return -1
```

#### 4.1.1 Running time

We have the following recurrence relation  $T(N) = 1 + T(N/2)$  and  $T(1) = 1$ . We must solve it to find out the algorithm's running time:



$$\begin{aligned}T(N) &= 1 + T(N/2) \\&= 1 + 1 + T(N/2^2) \\&= 1 + 1 + 1 + T(N/2^3) \\&\dots \\&= k + T(N/2^k)\end{aligned}$$

If  $N$  is a power of 2 (i.e.,  $N = 2^k$  with  $k = \log(N)$ ) we'll have:

$$\begin{aligned}T(N) &= k + T(N/2^k) \\&= k + T(1) = \log(N) + 1\end{aligned}$$

This means that  $T(N) = O(\log N)$ .

## 4.2 List/stacks/queue ADTs

An **abstract data type** (ADT) is a set of objects (lists, sets, graphs, etc.) together with a set of related operations (add, remove, size, etc.). They provide a template for the objects, not how the objects and their respective set of operations are implemented.

### 4.2.1 The List ADT

We usually deal with a general list of the form  $A_0, A_1, A_2, \dots, A_{N-1}$ . We say that the size of this list is  $N$ . We will call the special list of size 0 an **empty list**. The **position** of the element  $A_i$  in a list is  $i$ .

The List ADT could describe the following operations:

**find(x)** Return the position of the occurrence of item  $x$ .

**insert(x, i)** Insert some element  $x$  at position  $i$ .

**remove(i)** Remove some element at position  $i$ .

### 4.2.2 Vector in the STL

**Pros:**

- Constant time indexing
- Fast to add data at the end (not the front).

**Cons:**

- Slow to add data in the middle.
- Inefficient for searches.

**Operations**

- `void push_back(const T &x)` - add `x` at the end of the list.
- `void pop_back()` - remove element at the end of the list.
- `const T &back() const` - return the element at the end of the list.
- `const T &front() const` - return the element at the front of the list.
- `void push_front(const T &x)` - add `x` to the front of the list.
- `void pop_front()` - remove the element at the front of the list.

### 4.2.3 List in the STL

**Pros:**

- Implemented as a doubly linked list.
- Fast insertion/removal of items in any position.

**Cons:**

- No indexing.
- Inefficient for searches.

**Operations**

- `T & operator[] (int idx)` - return element at index `idx` with no bounds checking.
- `T &at(int idx)` - return element at index `idx` with bounds checking.
- `int capacity() const` - return internal capacity of vector.
- `void reserve(int new_capacity)` - set new capacity and possibly void expansion of vector.

### 4.3 Containers

A **container** is a holder object that stores a collection of other objects. This is usually implemented as a class templates which provides it with great flexibility for the types supported as elements.

A container

- manages the storage space for its elements and provides member functions to access them either directly or through **iterators**.
- replicates structures commonly used in programming such as dynamic arrays (e.g., vectors), queues, stacks, heaps (e.g., priority queues), linked lists, trees (e.g., sets), associative arrays (e.g., maps), etc.

### 4.4 Iterators

Some operations on lists require the notion of a position. In the STL, a position is represented by some nested type known as an **iterator**. In particular, for a `list<string>`, the position is represented by the type `list<string>::iterator`; for a `vector<int>`, `Vector<int>::iterator`.

#### 4.4.1 How to get an iterator

The STL lists (and all other STL containers) define a pair of methods:

- `iterator begin()`: returns an appropriate iterator representing the first item in the container.
- `iterator end()`: returns an appropriate iterator representing the endmarker in the container. This endmarker is “out-of-bounds”.

#### 4.4.2 Iterator methods

- `itr++` and `++itr`: advances the iterator `itr` to the next location.
- `*itr`: returns a reference, which may or may not be modifiable, to the object stored at iterator `itr`'s location.
- `itr1 == itr2`: returns true if iterators `itr1` and `itr2` refer to the same location and false otherwise.
- `itr1 != itr2`: returns true if iterators `itr1` and `itr2` refer to a different location and false otherwise.

For example the code:

```
for (int i = 0; i <= v.size(); i++) {  
    std::cout << v[i] << "\n";  
}
```

could be rewritten as follows using iterators:

```
for (vector<int>::iterator itr = v.begin(); itr != v.end(); itr++) {  
    std::cout << *itr << "\n";  
}
```

Alternatively:

```
vector<int>::iterator itr = v.begin();  
while (itr != v.end()) {  
    std::cout << *itr++ << "\n";  
}
```

#### 4.4.3 Container operations that require iterators

- `iterator insert( iterator pos, const T & x )`: adds `x` into the list, prior to the position by the iterator `pos`. This is a constant operator for `list`, but not for `vector`. The return value is an iterator representing the position of the inserted item.
- `iterator erase( iterator pos )`: removes the object at position and return the position of the element that followed `pos` prior to the call. It invalidates `pos`, making it stale.
- `iterator erase( iterator start, iterator pos )`: removes all items beginning at position up to, but not including `end`. An entire list `c` can be erased by `c.erase( c.begin(), c.end() )`.

#### 4.4.4 Example: Using erase on a list

```

/*
 * @brief Deletes every other element from lst, starting from the first
 *        item.
 * @param lst a list, or any object that supports iterators and erase.
 */
template <typename Container>
void removeEveryOtherItem( Container &lst ) {
    typename Container::iterator itr = lst.begin(); // could be shortened
    to                                              // auto itr = lst.
                                                    // begin();

    while (itr != lst.end()) {
        itr = lst.erase(itr);
        if (itr != lst.end()) itr++;
    }
}

```

#### 4.4.5 const\_iterators

The result of `*itr` is both the value of the item that the iterator is viewing but also the item itself. This is very powerful but also introduces some complications.

Let's analyze the following routine that works for both `vector` and `list` and runs in linear time:

```

template <typename Container, typename Object>
void change( Container &c, const Object &newValue ) {
    typename Container::iterator itr = c.begin();
    while (itr != c.end()) {
        *itr = newValue; // set current iterator's value to newValue
        itr++;           // advance the pointer to next item
    }
}

```

The potential problem arises if `Container c` was passed to a routine using call-by-constant reference, meaning we would expect that not changes would be allowed to `c`, and the compiler would ensure this by not allowing calls to any `c`'s mutators. For example, consider the following code that prints a `list` of integers but also tries to do some changes:

```

void print( const list<int> &lst, ostream &out = cout ) {
    typename Container::iterator itr = lst.begin();
    while (itr != lst.end()) {
        out << *itr << "\n";
        *itr = 0; // <= This is the suspect.
        itr++;
    }
}

```

The solution provided by the STL is that every collection contains not only an `iterator` nested type but also `const_iterator` nested type. The main difference between them is that `operator*` for `const_iterator` returns a constant reference, and thus `itr*` for a `const_iterator` cannot appear on the left-hand side of an assignment statement.

Further the compiler will force you to use a `const_iterator` to traverse a *constant* collection and does so by providing two versions of `begin` and two versions of `end`:

- `iterator begin()`
- `const_iterator begin() const`
- `iterator end()`
- `const_iterator end() const`

**Note:** The two versions of `begin` can be in the same class only because of the const-ness of a method (whether an accessor or mutator) is considered to be part of the signature. The trick is overloading `operator[]`.

Using `auto` to declare your iterators means the compiler will deduce for you whether an `iterator` or `const_iterator` is substituted. This exempt the programmer from having to keep track of the correct iterator type and is precisely one of the intended uses of `auto`.

#### 4.4.6 Printing a container

```
/*
 * @brief Prints out the container on the output stream.
 * @param c a given container.
 * @param out an output stream.
 */
template <typename Container>
void print( const Container &c, ostream &out = cout ) {
    if (c.empty()) {
        out << "(empty)";
    }
    else {
        typename Container::iterator itr = begin(c);
        out << "[" << *itr++; // print first item
        while (itr != end(c)) {
            out << ", " << *itr++;
        }
        out << "]" << "\n";
    }
}
```

## 5 Lecture 5: Lists, Stacks, Queues, and the STL

### 5.1 Implementation of a vector

The main details of the implementation of the `Vector` class is as follows:

- It'll maintain the primitive array (via a pointer variable to the block of allocated memory), the array capacity, the current number of stored items.
- It'll implement the “big-five” to provide **deep-copy** semantics.
- It'll provide the `resize` and `reserve` methods that will change the size of the `Vector` and the capacity of the `Vector` respectively.
- It'll overload the `operator[]`.
- It'll provide several basic methods, such as `size`, `empty`, `clear`, `back`, `pop_back`, and `push_back`.
- It'll provide support for the nested types `iterator` and `const_iterator`, and associated `begin` and `end` methods.

**Implementation:** [implementations/Vector](#)

**Why create a class `Vector`? Is it better than a simple array?** There are few reasons why creating `Vector` class template makes sense. The main one is that `Vector` will be a first-class type, which means that 1) it can be copied, and 2) the memory it uses can be reclaimed back (via its destructor). This is not the case for primitive array in C++.

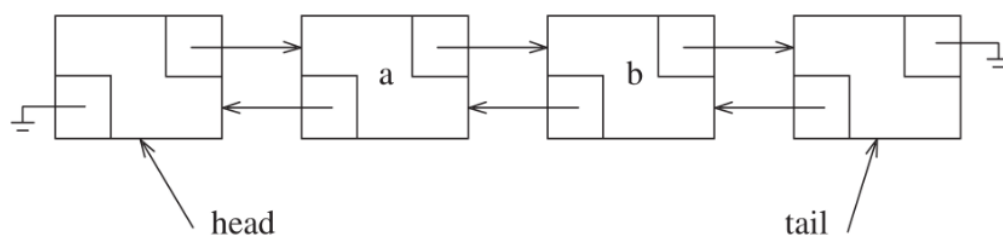
### 5.2 Implementation of a list

The main details of the implementation of the `List` class template is as follows:

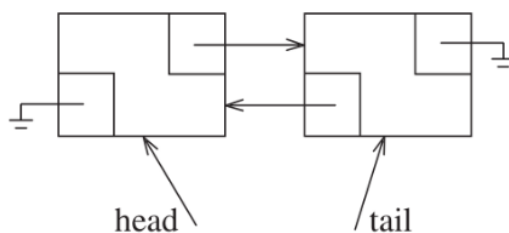
- The underlying data structure is a **doubly linked list** and thus we must maintain pointers to both ends of the list. This is in order to maintain constant time cost per operation, so as long as the operation occurs at a known position. The known position can be either end or at a position specified by an iterator.
- Four classes will be provided to better compose the `List` class template:
  - The `List` class itself, which contains links to both ends, the size of the list, and several methods.
  - The `Node` class, which is likely to be a private nested class. A node contains the data and pointers to the previous and next nodes, along with appropriate constructors.

- The `const_iterator` class, which abstracts the notion of a position, and is a public nested class. It stores a pointer to the current node, and provides implementation of the basic iterator operations.
- The `iterator` class, which abstracts the notion of a position, and is a public nested class. It has similar functionality to `const_iterator`, except that `operator*` returns a **reference** to the item being viewed, rather than a **constant reference**.
- Extra nodes are at both the front and end of the list. The node at the front represents the beginning marker while the one at the end represents the endmarker. These are the **sentinel nodes**; the one at the front is the **header node** and the one at the end is the **tail node**. Using the sentinels simplify the coding by removing a host of special cases (e.g., removing the first node).

**Implementation:** [implementations/List](#)



A doubly linked list with header and tail nodes



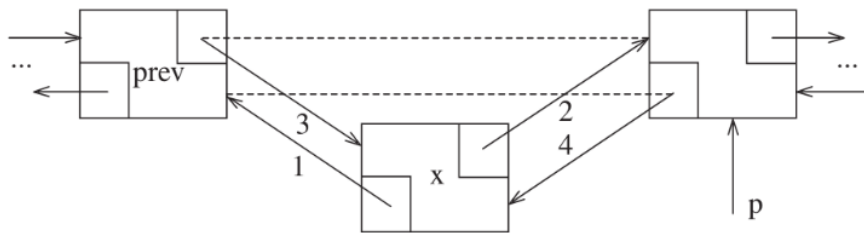
An empty doubly linked list with header and tail nodes

**Figure 4:** Doubly linked list



### 5.2.1 Insertion

```
Node *newNode = new Node{ x, p->prev, p }; // Steps 1 and 2
p->prev->next = newNode;                  // Step 3
p->prev = newNode;                         // Step 4
```

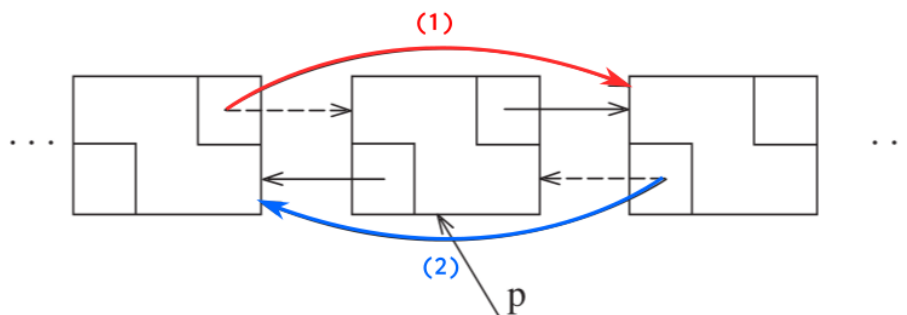


Insertion in a doubly linked list by getting a new node and then changing pointers in the order indicated

**Figure 5:** Adding a node

### 5.2.2 Erasure

```
p->prev->next = p->next; (1)  
p->next->prev = p->prev; (2)  
delete p;
```



Removing node specified by *p* from a doubly linked list

**Figure 6:** Removing a node

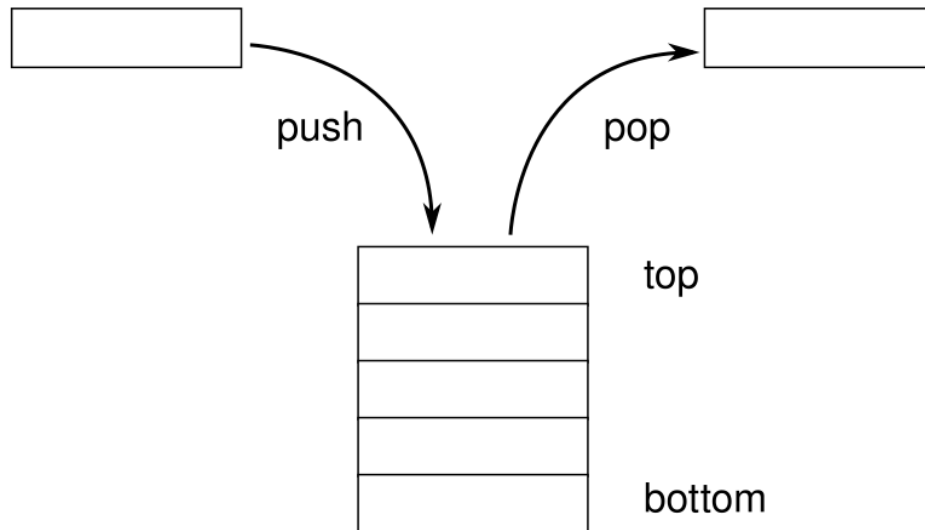
### 5.3 The Stack ADT

A **stack** (sometimes known as a LIFO (*last in, first out*) list) is a special list where insertions and deletions can be performed only from the end of the list (a.k.a., the **top**).

The fundamental operations on a stack are:

**push** Push an item to the top of the stack. Equivalent to *insert*.

**pop** Pop the most recently inserted item. The topmost item can be examined by using the *top* routine.



**Figure 7:** Stack model

### 5.3.1 Implementations

Two popular implementations of stacks are:

- **Linked list implementation**

This implementation of a stack uses a singly linked list. A **push** operation is performed by inserting at the front of the list, a **pop** operation is performed by deleting the element at the front of the list, and a **top** operation merely examines the element at the front of the list by returning its value.

**Implementation:** [implementations/LLStack](#)

- **Array implementation**

This implementation of a stack avoids pointers and is usually the most popular solution. It uses the **back**, **push\_back**, and **pop\_back** implementation from **vector**. Associated with each stack are **array** to store the items and **top\_of\_stack** which is  $-1$  for an empty stack.

**Implementation:** [implementations/ArrayStack](#)

### 5.3.2 Applications

After the array, the stack is probably the most fundamental data structure in computer science and for this reason, the application of stacks can be found in many places. Some of them are:

- **Balancing symbols**

Given that compiler check programs for syntax error, a useful tool to check for the lack of symbols that might cause errors is a program that checks whether everything is balanced (i.e., every right brace, bracket, and parenthesis must correspond to its left counterpart). For instance, the sequence `[ ( ) ]` is legal, but `[ ( ] )` is wrong. This algorithm uses a stack and could be described as follows:

Make an empty stack. Read characters until end of file. If the character is an opening symbol, push it onto the stack. If it is a closing symbol and the stack is empty, report an error. Otherwise, pop the stack. If the symbol popped is not the corresponding opening symbol, then report an error. At end of file, if the stack is not empty, report an error.

- **Postfix expressions**

This notation, also known as **reverse Polish notation**, is a mathematical notation in which operators follow their operands, in contrast to Polish notation, in which operators precede their operands. For example, the infix expression  $((15 \div (7 - (1 + 1))) \times 3) - (2 + (1 + 1))$  can be written in reverse Polish notation as `15 7 1 1 + - ÷ 3 × 2 1 1 + + -` which yields 5 regardless of the direction of the evaluation. An algorithm that uses a stack and that processes the expression from left to right is described below:

Input: `EXPR` is an expression in reverse Polish notation.

Output: `EXPR`

```
EvaluatePostfixExpression( EXPR ):
    stack ← []
    for EXPR → token:
        if token = operator:
            operand_2 ← stack.pop()
            operand_1 ← stack.pop()
            result ← evaluate operator_1 and operand_2 with operator
            stack.push(result)
        else if token = operand:
            stack.push(token)
    return stack.pop()
```

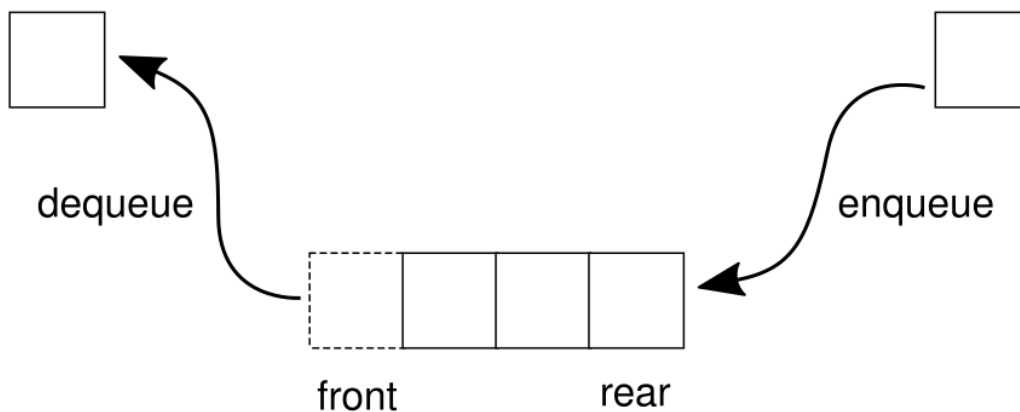
### 5.4 The Queue ADT

A **queue** (sometimes known as a FIFO (*first in, first out*) list) is a specialist where insertions are performed at one end and deletions are performed at the other end.

The fundamental operations on a queue are:

**enqueue** Insert an element at the end of the list (called the rear).

**dequeue** Delete (and return) the element at the start of the list (called the front).



**Figure 8:** Queue model

### 5.4.1 Implementations

Two popular implementations of stacks are:

- **Linked list implementation.**

This implementation is similar to its counterpart stack implementation.

**Implementation:** [implementations/LLQueue](#)

- **Array implementation**

For this implementation, we keep an array, `array`, and the positions `front` and `back`, which represents the ends of the queue. Likewise, we keep track of the number of items currently in the queue, `current_size`. To `enqueue` an element `x`, we increment `current_size` and `back`, then

set `_array[back] = x`. To `dequeue` an element, we set the return value to `array[front]`, decrement `current_size`, and increment `front`.

However there's one potential problem with this implementation. After a certain number of operations, a queue might appear full since `back` is now at the last array index, and the next `enqueue` would be in a nonexistent position. A simple solution would be to wrap around to the beginning whenever `front` or `back` gets to the end of the array. This is known as a **circular array** implementation for which modular arithmetic is used.

**Implementation:** [implementations/ArrayQueue](#)

### 5.4.2 Applications

There are many algorithms that used queues to give efficient running times. Some examples are:

- Jobs submitted to a printer are arranged in order of arrival. Thus, they're essentially placed in a queue.
- Virtually every real-life line is supposed to be a queue. For instance, lines at ticket counters are queues, because the service is first-come first-served.
- There are many network setups of personal computers in which the disk is attached to one machine, known as the file server. Users on other machines are given access to files on a first-come first-served basis, so the data structure is a queue.

## 5.5 Trees

Recursively, a **tree** is either empty or consists a root node  $r$  and zero or more non-empty subtrees whose roots are connected by an edge.

- Used to implement file systems of several popular OSs
- Used to evaluate arithmetic expressions
- Support  $O(\log N)$  search
- Set and map classes.

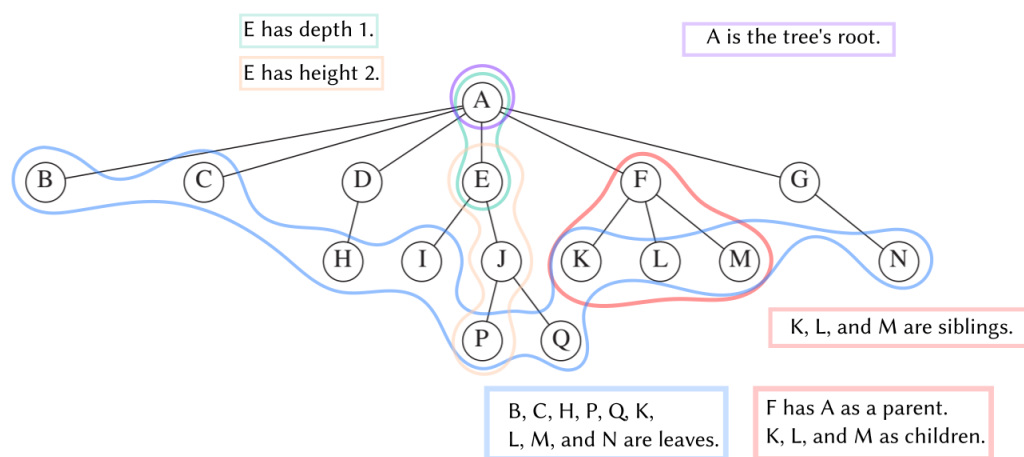
## 6 Lecture 6: Trees

### 6.1 Preliminaries

A **tree** is a widely used abstract data type (ADT) that simulates a hierarchical tree structure, with a root value and subtrees of children with a parent node, represented as a set of linked nodes.

Trees (more specifically binary trees) are really important so a certain terminology has developed for them:

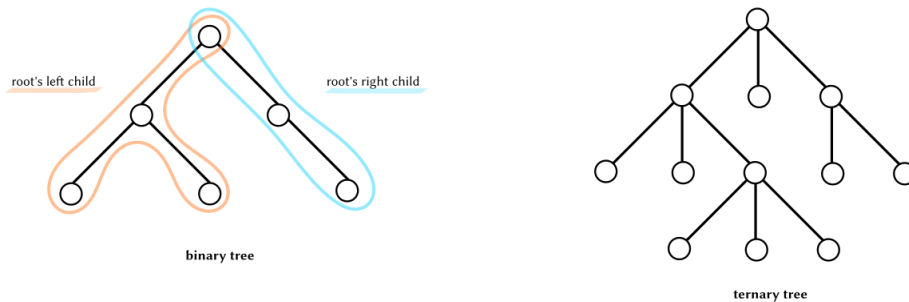
- A **node** is a structure which may contain a value or condition, or represent a separate data structure.
- An **edge** is connection between one node and another.
- A **path** is a sequence of nodes and edges connecting a node with a descendant.
- An **ancestor** is a node reachable by repeated proceeding from child to parent. A **descendant** is a node reachable by repeated proceeding from parent to child. For example, if there's a path from node  $u$  to node  $v$ , then  $u$  is an ancestor of  $v$  and  $v$  is a descendant of  $u$ . If  $u \neq v$ , then  $u$  is a **proper ancestor** of  $v$  and  $v$  is a **proper descendant** of  $u$ .
- The **root**  $r$  of a tree is its top node known as the prime **ancestor**.
- The root of each subtree is said to be a **child** of the root  $r$ , and  $r$  is the **parent** of each subtree root.
- Two nodes  $u$  and  $v$  are **siblings** if they have the same parent.
- The **subtree** of a node,  $u$ , is the tree rooted at  $u$  and contains all of  $u$ 's descendants.
- The **depth** of a node,  $u$ , is the length of the path from  $u$  to the root of the tree. This translates into the number of edges from node to the tree's root node. Thus, a root node will have a depth 0.
- The **height** of a node,  $u$ , is the length of the longest path from  $u$  to one of its descendants. It's defined as  $h = \max(\text{height}(u.\text{left}), \text{height}(u.\text{right})) + 1$ . Thus, the height of a tree is equal to the height of the root. Since leaves have no children, all leaves have height 0.
- A node,  $u$ , is a **leaf** if it has no children.



**Figure 9:** A tree

An  $m$ -ary tree is a rooted tree in which each node has no more than  $m$  children. When  $m = 2$ , the tree

is known as a **binary tree** (i.e., each node has at most two children, the left child and the right child). When  $m = 3$ , the tree is a **ternary tree**.



**Figure 10:** Binary and ternary tree

## 6.2 Implementation of trees

An obvious way to implement a tree would be to have in each node, besides its data, a link to each child of the node. However, the number of children per node can vary and is not known in advance which might translate into too much wasted space. The solution is simple: Keep the children of each in a linked list of nodes. In this way, each node can have as many children as the tree it composes allows.

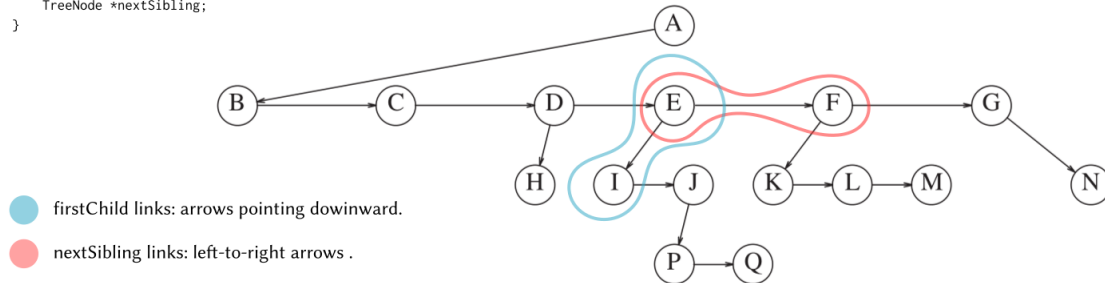
```
template<typename Object>
struct TreeNode {
    Object element;
    TreeNode* firstChild;
    TreeNode* nextSibling;
}
```



```

struct TreeNode {
    Object element;
    TreeNode *firstChild;
    TreeNode *nextSibling;
}

```



First child/next sibling representation of the tree

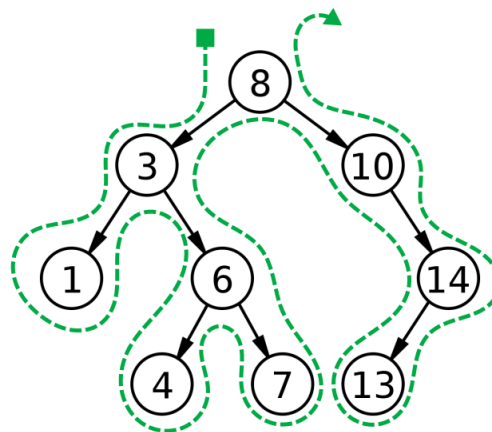
**Figure 11:** First child/next-sibling representation

### 6.3 Tree traversals

A **tree traversal** refers to the process of visiting each node in a tree data structure, exactly once. Such traversals are classified by the order in which the nodes are visited.

Trees can be traversed either:

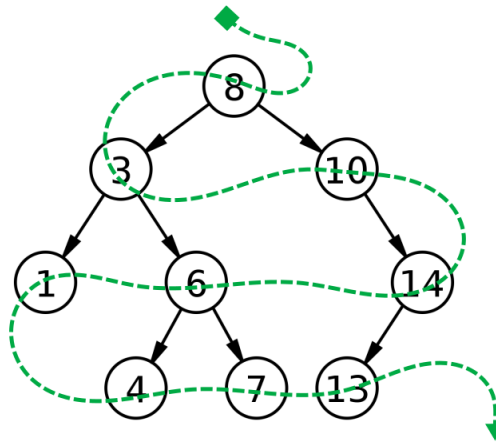
- by deepening the search as much as possible on each child before going to the next sibling. This is known as **depth-first search** (DFS).



**Figure 12:** Depth-first search

- by visiting every node on a level before going to a lower level. This is known as **breadth-first**

**search** (BFS).



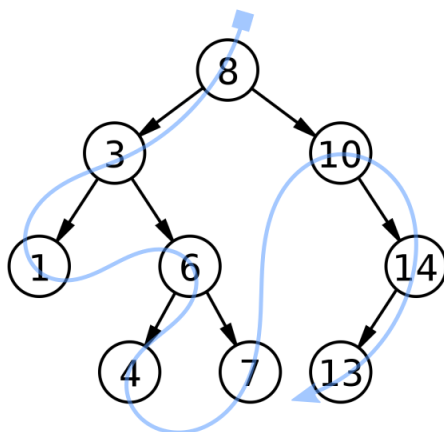
**Figure 13:** Breadth-first search

### 6.3.1 Depth-first search

Depending on when the work at a node is performed, a traversal of this type can be classified into:

- **pre-order traversal.** In a pre-order traversal, work at a node is performed before (*pre*) its children are processed.

1. Check if the current node is empty or null.
2. Display the data part of the root (or current node).
3. Traverse the left subtree by recursively calling the pre-order function.
4. Traverse the right subtree by recursively calling the pre-order function.



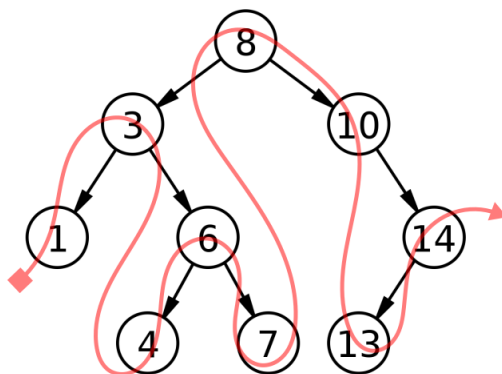
```
preorder(node)
  if (node == null)
    return
  visit(node)
  preorder(node.left)
  preorder(node.right)
```

8, 3, 1, 6, 4, 7, 10, 14, 13

**Figure 14:** Pre-order traversal

- **in-order traversal.** In an in-order traversal, work at a node is performed after its left child is visited, followed by the right child.

1. Check if the current node is empty or null.
2. Traverse the left subtree by recursively calling the in-order function.
3. Display the data part of the root (or current node).
4. Traverse the right subtree by recursively calling the in-order function.



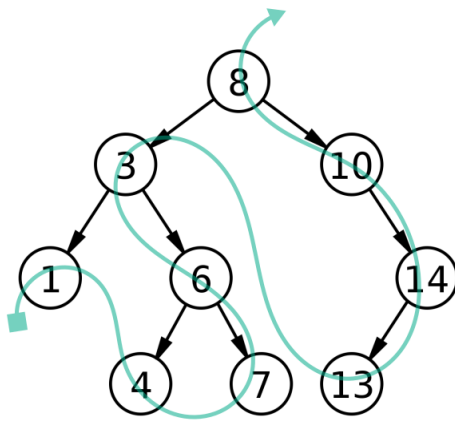
```
inorder(node)
  if (node == null)
    return
  inorder(node.left)
  visit(node)
  inorder(node.right)
```

1, 3, 4, 6, 7, 8, 10, 13, 14

**Figure 15:** In-order traversal

- **post-order traversal.** In a post-order traversal, work at a node is performed after (*post*) its children are evaluated.

1. Check if the current node is empty or null.
2. Traverse the left subtree by recursively calling the post-order function.
3. Traverse the right subtree by recursively calling the post-order function.
4. Display the data part of the root (or current node).



```
postorder(node)
    if (node == null)
        return
    postorder(node.left)
    postorder(node.right)
    visit(node)
```

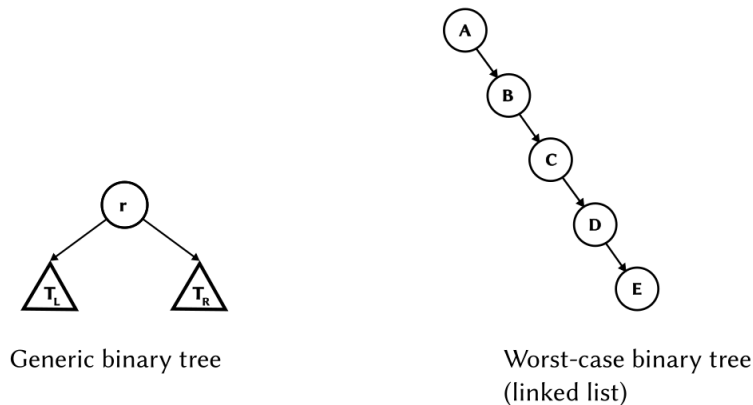
1, 4, 7, 6, 3, 13, 14, 10, 8

**Figure 16:** Post-order traversal

## 6.4 Binary trees

As stated earlier, a **binary tree** is a tree in which no node can have more than two children.

A property of a binary tree that is sometimes important is that the depth of an average binary tree is considerably smaller than  $N$ . An analysis shows that the average depth is  $O(\sqrt{N})$ , and that for a special type of binary tree, namely the **binary search tree**, the average value of the depth is  $O(\log N)$ . Unfortunately, the depth can be as large as  $N - 1$ . In its worst case, we end up with a linked list (which is also a type of tree).



**Figure 17:** Generic binary tree and linked list

### 6.4.1 Implementation

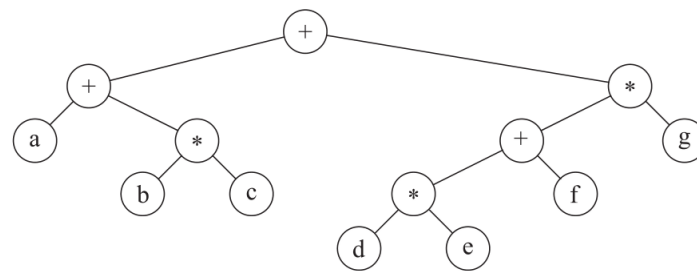
Given that a binary tree has at most two children per node, we can keep direct links to them. Thus, a binary node is just a structure consisting of the data plus two pointers to current node's left and right nodes. For example:

```
struct BinaryNode {  
    Object element;    // data in the node  
    BinaryNode* left; // left child  
    BinaryNode* right; // right children  
}
```

Binary trees have many important uses not associated with searching. One of the principal uses of binary trees is in the area of compiler design (e.g., expression trees).

### 6.4.2 Expression trees

A **binary expression tree** is a specific kind of a binary tree used to represent expressions. The leaves of an expression are **operands**, such as constants or variable names, and the other nodes contain **operators**. We can evaluate an expression,  $T$ , by applying the operator at the root to the values obtained by recursively evaluating the left and right subtrees.



Expression tree for  $(a + b * c) + ((d * e + f) * g)$

**Figure 18:** Expression tree

- To produce an (overly parenthesized) *infix expression*, recursively produce parenthesized left expression, then printing out the operator at the root, and finally recursively producing a parenthesized right expression. This is an **in-order traversal**.
- To produce a *postfix expression*, recursively print out the left subtree, the right subtree, and then the operator. This is an **post-order traversal**.
- To produce a *prefix expression*, print out the operator first and then recursively print out the left and right subtrees. This is an **pre-order traversal**.

### 6.4.3 Constructing an expression tree

This algorithm describes the steps to convert a postfix expression into an expression tree:

```

ConvertPostfixToTree( EXPR ):
    stack ← []
    for EXPR → token:
        if token = operand:
            T ← Tree(root ← operand)
            stack.push(T)
        if token = operator:
            T1 ← stack.pop()
            T2 ← stack.pop()
            T ← Tree(root ← operator, left ← T1, right ← T2)
            stack.push(T)
    return stack
  
```

**Implementation:** [implementations/ExpressionTree](#)

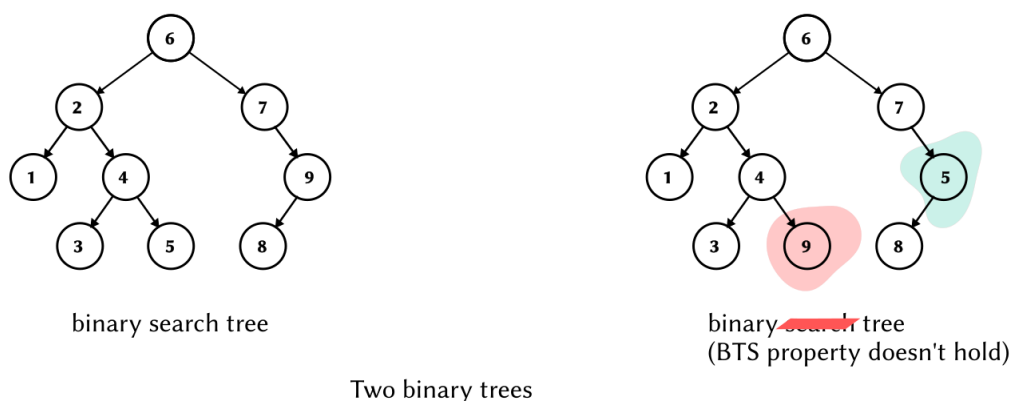
## 6.5 Binary search trees (BSTs)

A **binary search tree** is a special type of binary trees with the following property:

For every node,  $X$ , in the tree, the values of all the items in its left subtree are smaller than the item in  $X$ , and the value of all the items in its right are larger than the item in  $X$ .

This implies that all the elements in the tree can be ordered in some consistent manner.

Binary search trees keep their keys in sorted order, so that lookup and other operations can use the principle of binary search: when looking for a key in a tree (or a place to insert a new key), they traverse the tree from root to leaf, making comparisons to keys stored in the nodes of the tree and deciding, on the basis of the comparison, to continue searching in the left or right subtrees. On average, this means that *each comparison allows the operations to skip about half of the tree*, so that each lookup, insertion or deletion takes time proportional to the logarithm of the number of items stored in the tree.



Two binary trees

**Figure 19:** Binary search tree vs Non-binary search tree

**Complete implementation:** [implementations/BinarySearchTree](#)

### 6.5.1 Operations on a BTS

**contains** For this operation, we're required to return **true** if there is a node in the tree  $T$  that has item  $X$ . Otherwise, return **false**. If  $T$  is empty, then we can return **false** right away. Otherwise, if the item stored at  $T$  is  $X$ , we can return **true**. Otherwise, we can make a recursive call on a subtree of  $X$ , either left or right, depending on the relationship between  $X$  and the item stored in  $T$ . It's implied that the objects being compared have implemented the used relational operators appropriately.

Input:  $X$  is an item to search **for** and  $T$  is the node that roots the subtree.

Output: Return **true** or **false** depending on whether the tree contains the node storing the item.

```
contains( X, T ):
    if T = null:
        return false
    else if X < T.element():
        return contains X, T.left()
    else if X > T.element():
        return contains X, T.right()
    else:
        return true # Match
```

**findMin and findMax** These operations can either return the node (or its value) containing the smallest and largest element in the tree, respectively. To perform a **findMin**, start at the root and go left as long as there's a left child. The stopping point is the smallest element. The **findMax** routine is similar, except that branching is to the right child.

Recursive algorithm of **findMin**:

Input: the root of the tree.  
Output: the smallest element **in** the tree.

```
findMin( T ):
    if T = null:
        return T
    if T.left() = null:
        return T;
    return findMin T.left()
```

Iterative algorithm of **findMax**:

Input: the root of the tree.  
Output: the largest element **in** the tree.

```
findMax( T ):
    if T ≠ null:
        until T.left() = null:
            T ← T.left()
    return T
```

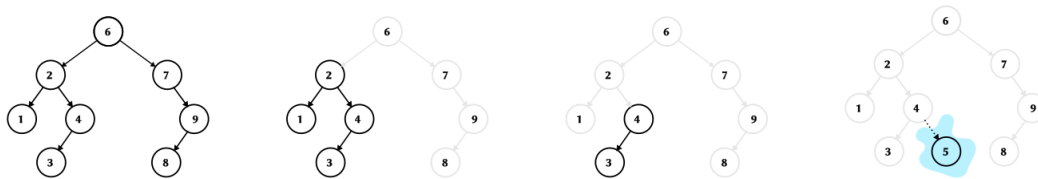
**insert** To insert  $X$  into tree  $T$ , proceed down the tree as you would with the **contains** operation. If  $X$  is found, do nothing. Otherwise, insert  $X$  at the last spot on the path traversed.



Duplicates can be handled by keeping an extra field in the node object indicating the frequency of occurrence for a particular element  $X$ . This adds some space extra space to the entire tree but is better than putting duplicates in the tree (which tends to make the tree very deep).

Input:  $X$  is an item to inserted in the tree  $T$ .  
Output: void.

```
insert( X, T ):
    if T ← null:
        T ← BinaryNode(X, null, null)
    else if X < T.element():
        insert X, T.left()
    else if X > T.element():
        insert X, T.right()
    else:
        # duplicate; do nothing
```

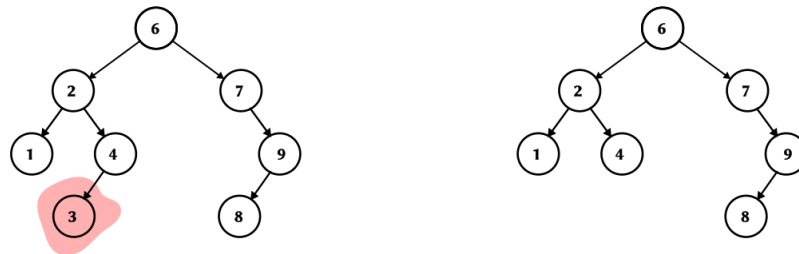


Progression of inserting 5 into the tree

**Figure 20:** Progression of an insertion

**remove** As is common with many data structures, the hardest operation is deletion. After finding the item to be deleted several possibilities must be considered:

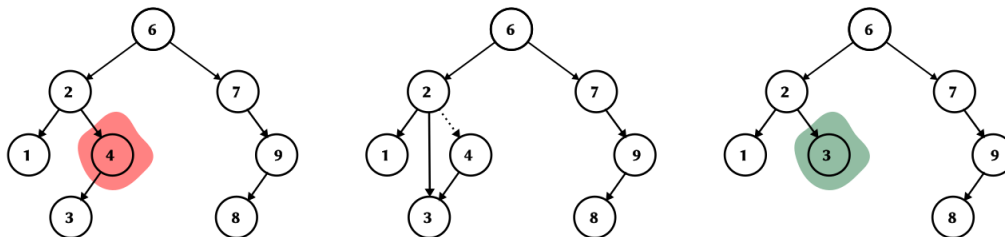
- If the node is a leaf, it can be deleted immediately.



Removing a childless node (3) from the tree

**Figure 21:** Removing childless node from tree

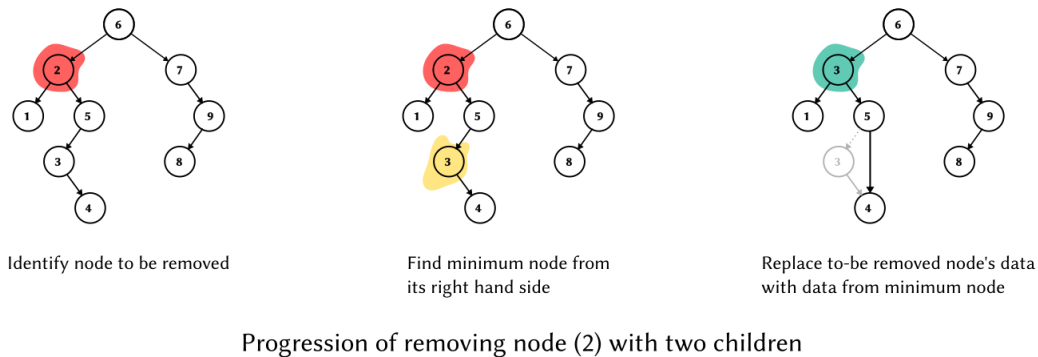
- If the node has one child, the node can be deleted after its parent adjusts a link to bypass the node.



Progression of removing node (4) with one child

**Figure 22:** Removing single child node from tree

- If the node has two children, replace the data of this node with the smallest data of the right subtree and recursively delete that node (which is now empty). Because the smallest node in the right subtree cannot have a left child, the second `remove` is an easy one.



**Figure 23:** Removing node with two children from tree

Input:  $X$  is an item to be removed and  $T$  is the tree.  
Output: void

```
remove( X, T ):
    if T = null:
        return; # Item not found; do nothing
    if X < T.element():
        remove X, T.left()
    else if X > T.element():
        remove X, T.right()
    else if T.left() ≠ null AND T.right() ≠ null:
        min ← findMin T.right()
        T.element ← min.element()
        remove T.element(), T.right()
    else:
        old-node ← T
        T ← T.left ≠ null ? T.left() : T.right()
```

The algorithm above is inefficient because it makes two passes down the tree to find and delete the smallest node in the right subtree when this is appropriate. It's easy to remove this inefficiency by writing a special `removeMin` method.

If the number of deletions is expected to be small, then a popular strategy to use is **lazy evaluation**: When an element is to be deleted, it's left in the tree and merely *marked* as being deleted.

**Pros of lazy evaluation (in this context):**

- Easy to handle if a deleted item is reinserted, do not need to allocate a new node.
- Do not need to handle finding a replacement node.

**Cons of lazy evaluation (in this context):**

- The depth of the tree will increase. However this increase is usually a small amount relative to the size of the tree.

**6.5.2 Average-case analysis**

The running time of all the operations (except `makeEmpty` and copying) is  $O(d)$ , where  $d$  is the depth of the node containing the accessed item.

The sum of the depths of all nodes in a tree is known as the **internal path length**.

Let  $D(N)$  be the internal path length for some tree  $T$  of  $N$  nodes. An  $N$ -node tree consists of an  $i$ -node left subtree and an  $(N - i - 1)$ -node right subtree, plus a root at depth zero for  $0 \leq i \leq N$ .  $D(i)$  is the internal path length of the left subtree with respect to its root. The same holds for the right subtree. We get the following recurrence  $D(N) = D(i) + D(N - i - 1) + N - 1$ .

In a BST all subtree sizes are equally likely, so

$$D(N) = \frac{2}{N} \left[ \sum_{j=0}^{N-1} D(j) \right] + N - 1$$

By solving the above recurrence we show that  $D(N) = O(N \log N)$ . Thus the expected depth of any node is  $O(\log N)$ .

Although it's tempting to say that this result implies the average running time of all the operations is  $\log N$ , but this is not entirely true. The reason for this is that because of deletions, it's not clear that all BSTs are equally likely. In particular, the deletion algorithm favors making the left subtrees deeper than the right, because we're always replacing a deleted node with a node from the right subtree.

**6.6 Balanced trees**

If the input comes into a tree presorted, then a series of inserts will take *quadratic time* and give a very expensive implementation of a linked list, since the tree will consist only of nodes with no left children. One solution to the problem is to insist on an extra structural condition called **balance** where no node is allowed to get too deep.

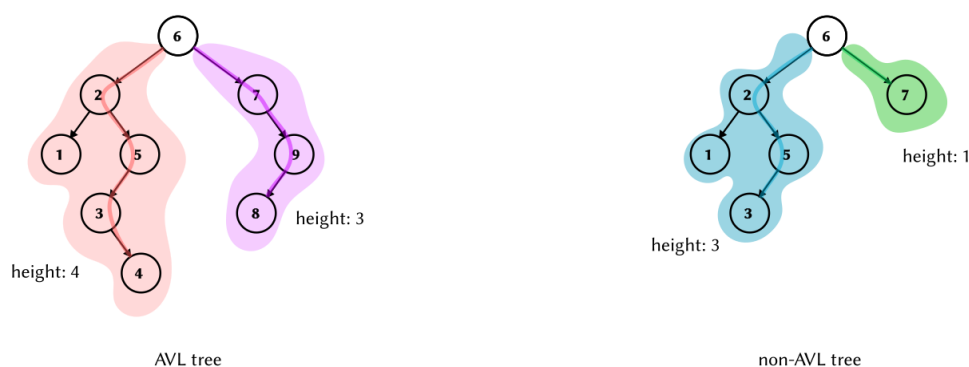
There are quite a few general algorithms to implement **balanced trees**. Most are quite a bit more complicated than a standard binary search tree, and all take longer on average for updates. They do, however, provide protection against the embarrassingly simple cases. Below, we will sketch one of the oldest forms of balanced search trees, the AVL tree.

## 6.7 AVL trees

An AVL (Adelson-Velskii-Landis) tree is a binary search tree that is *height-balanced*: At each node  $u$ , the height of the subtree rooted at  $u$ .left and the subtree rooted at  $u$ .right differ by at most 1. The balance condition must be easy to maintain, and it ensures that the depth of the tree is  $O(\log N)$ . The simplest idea is to require that the *left and right subtrees have the same height*.

It follows immediately that, if  $F(h)$  is the minimum number of leaves in a tree of height  $h$ , then  $F(h)$  obeys the Fibonacci recurrence  $F(h) = F(h-1) + F(h-2)$  with base cases  $F(0) = 1$  and  $F(1) = 1$ . This means  $F(h)$  is approximately  $\frac{\phi^h}{\sqrt{5}}$ , where  $\phi = (1 + \sqrt{5})/2 \approx 1.61803399$  is the *golden ratio* (More precisely,  $|\phi^h/\sqrt{5} - F(h)| \leq 1/2$ .) This implies  $h \leq \log_\phi n \approx 1.44042008 \log n$ , although in practice is only slightly more than  $\log N$ .

An **AVL tree** is identical to a binary search tree, except that for every node in the tree, the height of the left and right subtrees can differ by at most 1. The height of an empty tree is defined to be  $-1$ .



**Figure 24:** AVL tree and non-AVL tree

### 6.7.1 Insertion cases

After an insertion, only nodes that are on the path from the insertion point to the root might have their balance altered because only those nodes have their subtrees altered. As we follow the path up to the root and update the *balancing information*, we may find a node whose new balance violates the AVL condition.

Suppose that  $\alpha$  is the first node on the path that needs to be rebalanced. Since any node has at most two children, and a height imbalance requires that the two subtrees's heights of  $\alpha$  differ by two, it's easy to see that a violation might occur in four cases:

1. An insertion into the left subtree of the left child of  $\alpha$ .
2. An insertion into the right subtree of the left child of  $\alpha$ .
3. An insertion into the left subtree of the right child of  $\alpha$ .
4. An insertion into the right subtree of the right child of  $\alpha$ .

Cases 1 and 4 are mirror image symmetries with respect to  $\alpha$ , as are cases 2 and 4. Thus, as a matter of theory, there are only two basic cases. Although programmatically, there are still four cases.

A **single rotation** of the tree fixes case 1 (and 4) in which the insertion occurs on the *outside* (i.e., left-left or right-right). A **double rotation** (which is slightly more complex) fixes case 2 (and 3) in which the insertion occurs on the *inside* (i.e., left-right or right-left). The implementation of a double rotation simply involves two calls to the routine implementing a single rotation, although conceptually it may be easier to consider them two separate and different operations.

### Single rotation

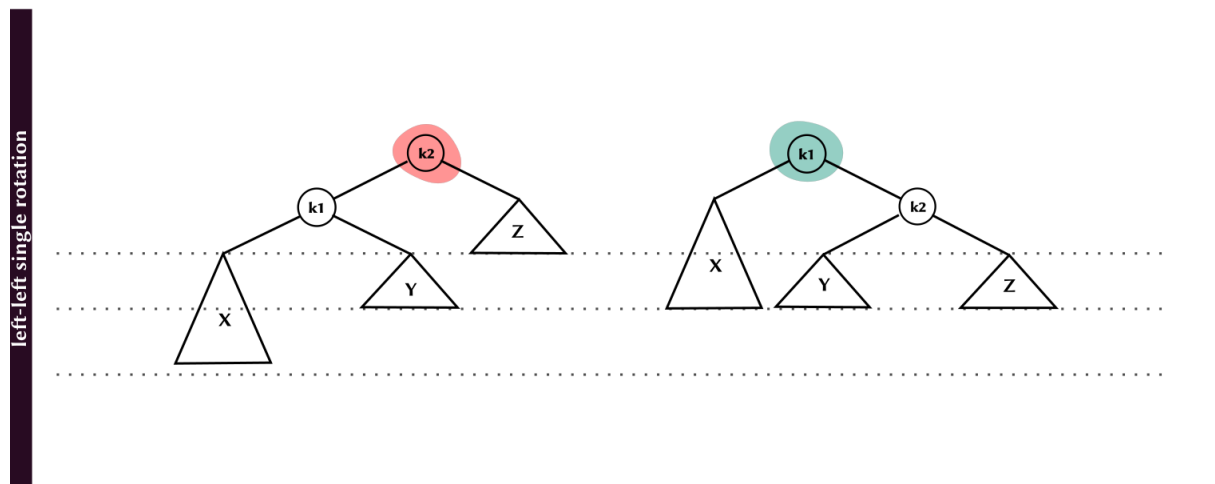
---

**NOTE:**  $k_1 < k_2$

#### Single rotation to fix case 1:

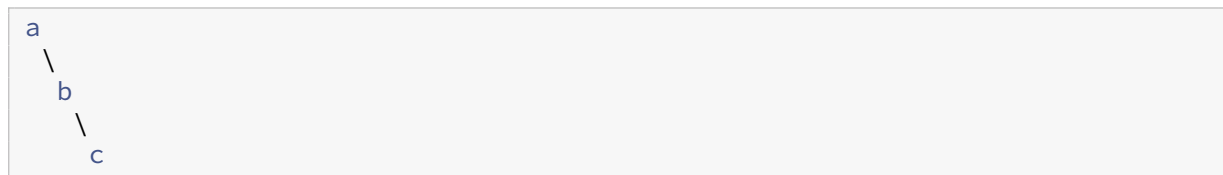
In the “Single rotation to fix case 1” figure, subtree  $X$  has grown to an extra level, causing it to be exactly two levels deeper than  $Z$ . The subtree  $Y$  cannot be at the same level as the  $X$  because then  $k_2$  would have been out of balance *before* the insertion, and  $Y$  cannot be at the same level as  $Z$  because then  $k_1$  would be the first node on the path toward the root that was in violation of the AVL balancing property.

To ideally rebalance the tree, we would like to move  $X$  up a level and  $Z$  down a level. By grabbing child node  $k_1$  and letting the other nodes hang, the result is that  $k_1$  will be the new root. The binary search tree property tells us that in the original tree  $k_2 > k_1$ , so  $k_2$  becomes the right child of  $k_1$  in the new tree.  $X$  and  $Z$  remain as the left child of  $k_1$  and right child of  $k_2$ , respectively. Subtree  $Y$ , which holds items that are between  $k_1$  and  $k_2$  in the original tree, can be placed as  $k_2$ 's left child in the new tree and satisfy all the ordering requirements.



**Figure 25:** Single rotation to fix case 1

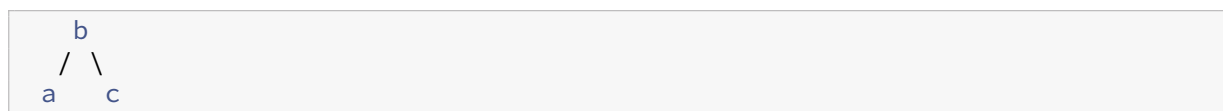
**Simplified example** Imagine we have this situation:



There's a height balance (i.e., the difference between his left and right subtree is greater than the balance factor, usually 1) in the node **a**. To fix this, we must perform a *left rotation* rooted at **a**. The steps are the following:

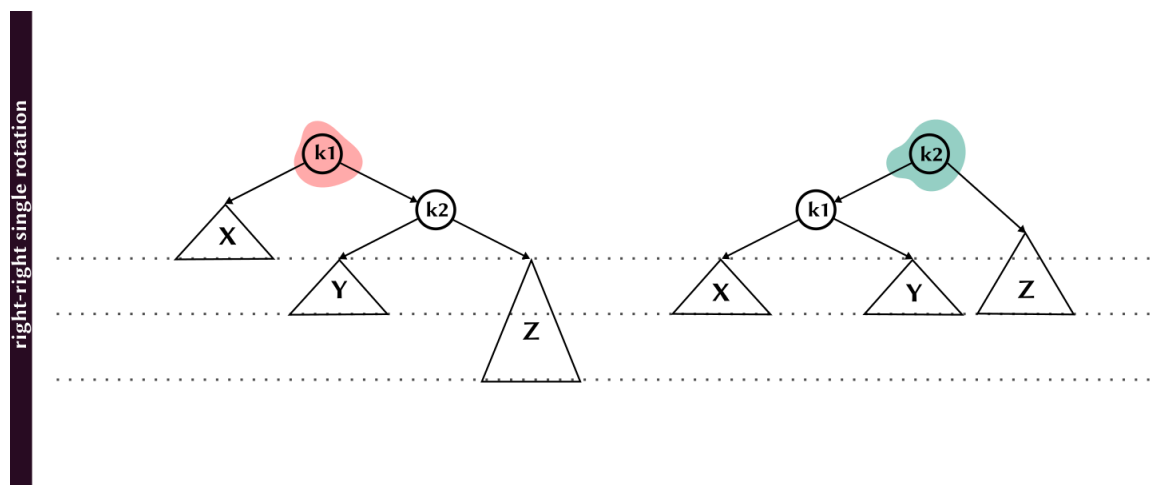
1. **b** becomes the new root.
2. **a** takes ownership of **b**'s left child as its right child. In this case, **null**.
3. **b** takes ownership of **a** as its left child.

The tree end up looking as follows:



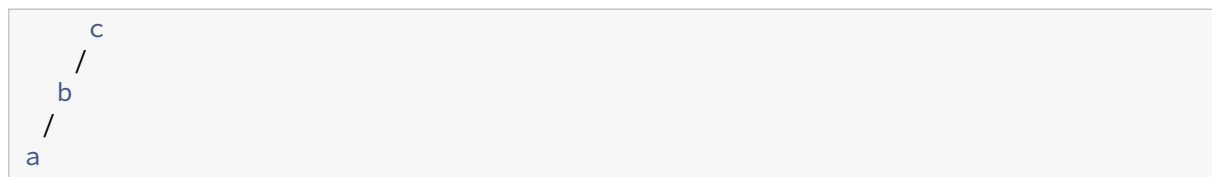
#### Single rotation to fix case 4:

Case 4 represents a symmetric case and the “Single rotation to fix case 4” figure shows how a single rotation is applied.



**Figure 26:** Single rotation to fix case 4

**Simplified example** Imagine we've this situation:



There's a height imbalance at node **c**, thus perform a *right rotation* rooted at **c**. The steps are the following:

1. **b** becomes the new root.
2. **c** takes ownership of **b**'s right child as its left child. In this case, **null**.
3. **b** takes ownership of **a** as its right child.

The tree ends up looking as follows:



**Example** Let's suppose we start with an empty AVL tree and insert the items 3, 2, 1, and then 4 through 7 sequentially. As we insert certain items, we must do some rotations along the ways.

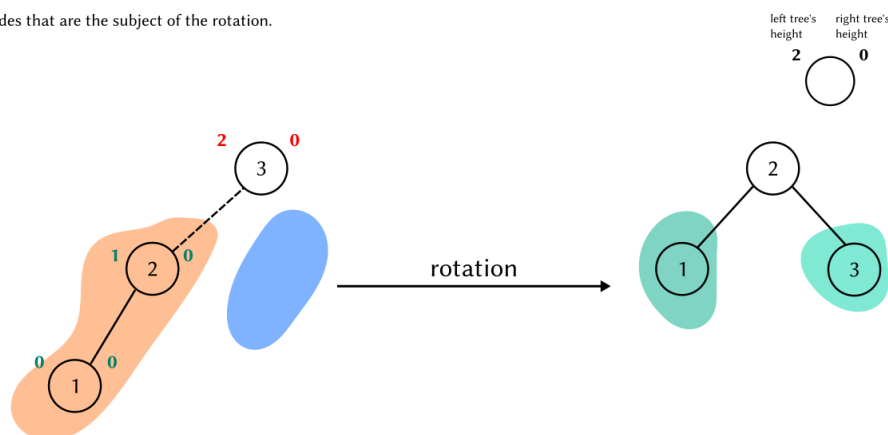
The first problem occurs when it's time to insert the item 1 because the AVL tree property is violated at the root. Thus, we perform a single rotation between the root and its left child to fix the problem. The "First rotation after violation of AVL property" figure shows the before and after the rotation.



----- A dashed line joins the nodes that are the subject of the rotation.

The insertion of 1 into the tree, causes a height imbalance at node 3, which then needs to be rebalanced. Here onwards this node is referred as  $\alpha$ .

The insertion occurs at the left subtree of the left child of  $\alpha$ . In other words, an insertion into the left subtree of node 2. This is case 1 and the rotation applied to it is known as **left-left single rotation**.



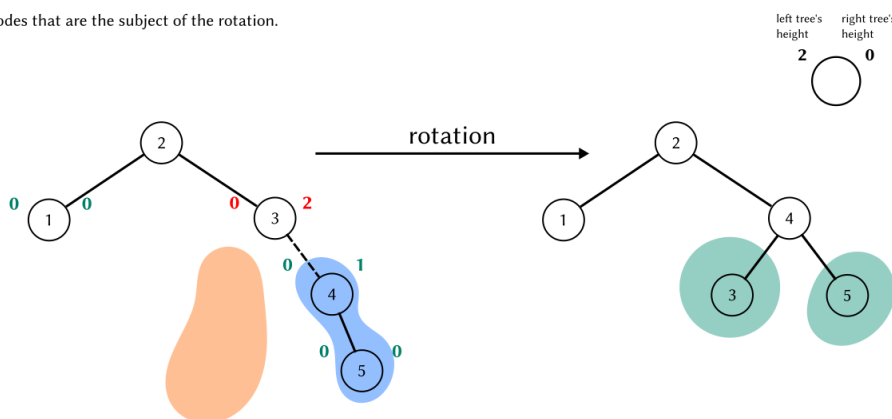
**Figure 27:** First rotation after violation of AVL property

Next we insert 4 but it causes no problems. However the insertion of 5 creates a violation at node 3 that is fixed by a single rotation. The rest of the tree has to be informed of this change so 2's right child must be reset to link to 4 instead of 3. The "Second rotation after violation of AVL property" figure shows the before and after the rotation.

----- A dashed line joins the nodes that are the subject of the rotation.

The insertion of 5 into the tree, causes a height imbalance at node 3, which then needs to be rebalanced. Here onwards this node is referred as  $\alpha$ .

The insertion occurs at the right subtree of the right child of  $\alpha$ . In other words, an insertion into the right subtree of node 4. This is case 4 and the rotation applied to it is known as **right-right single rotation**.



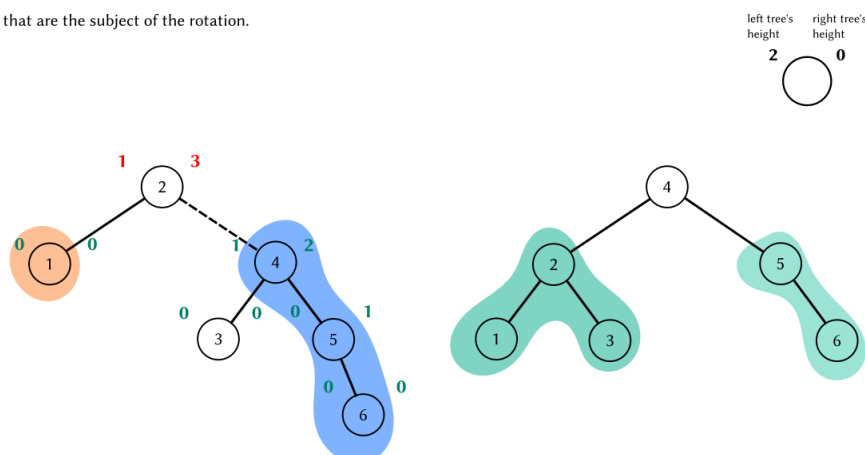
**Figure 28:** Second rotation after violation of AVL property

Next we insert 6 that causes a balance problem at the root, since its left subtree is of height 0 and its right subtree would be of height 2. Therefore, we perform a single rotation at the root between 2 and 4. The rotation is performed by making 2 a child of 4 and 4's original left subtree the new right subtree of 2. Every item in this subtree must lie between 2 and 4, so this transformation makes sense. The "Third rotation after violation of AVL property" figure shows the before and after the rotation.

----- A dashed line joins the nodes that are the subject of the rotation.

The insertion of 6 into the tree, causes a height imbalance at node 4, which then needs to be rebalanced. Here onwards this node is referred as  $\alpha$ .

The insertion occurs at the right subtree of the right child of  $\alpha$ . In other words, an insertion into the right subtree of node 5. This is case 4 and the rotation applied to it is known as **right-right single rotation**.



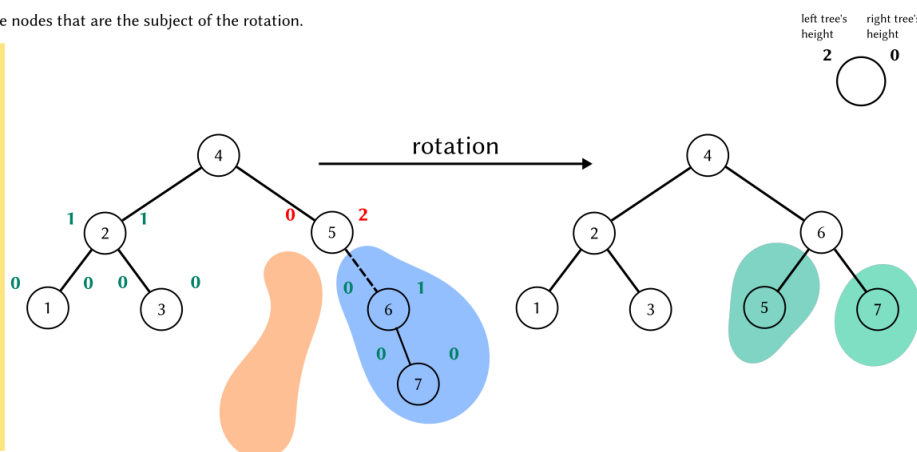
**Figure 29:** Third rotation after violation of AVL property

The next item we insert is 7, which causes another rotation. The “Fourth rotation after violation of AVL property” figure shows the before and after the rotation.

----- A dashed line joins the nodes that are the subject of the rotation.

The insertion of 7 into the tree, causes a height imbalance at node 5, which then needs to be rebalanced. Here onwards this node is referred as  $\alpha$ .

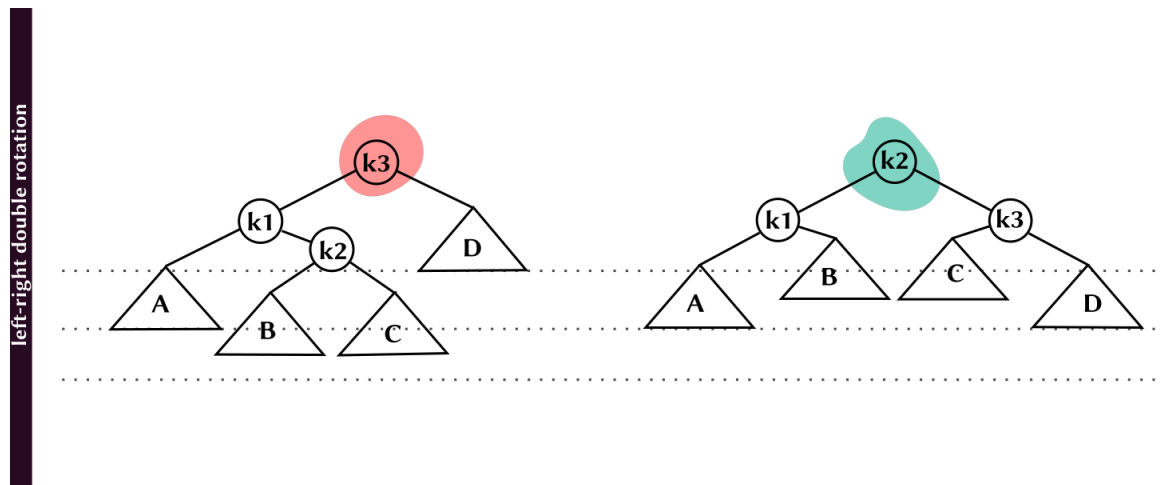
The insertion occurs at the right subtree of the right child of  $\alpha$ . In other words, an insertion into the right subtree of node 6. This is case 4 and the rotation applied to it is known as **right-right single rotation**.



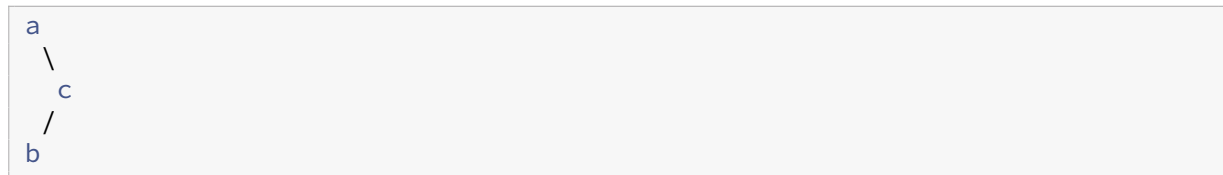
**Figure 30:** Fourth rotation after violation of AVL property

**Double rotation** The reason why single rotation doesn't work in cases 2 and 3 is because a subtree might be too deep, and a single rotation doesn't make it any less deep.

**NOTE:**  $k_1 < k_2 < k_3$

**Double rotation to fix case 2:****Figure 31:** Left-right double rotation to fix case 2

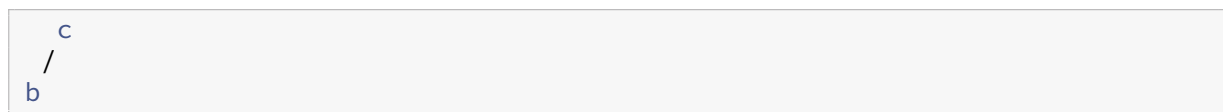
**Simplified example** Imagine we've the following situation:



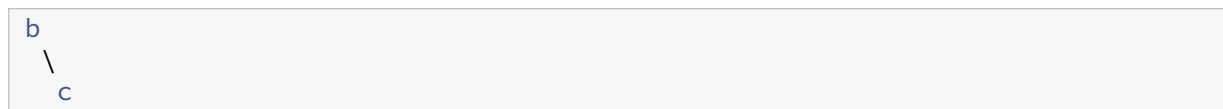
There's a height balance at node **a**, however performing a single rotation won't achieve the desired result.

The steps are the following:

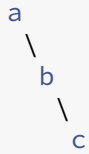
1. Perform a right rotation on the right subtree (we aren't rotating our current root).



becomes



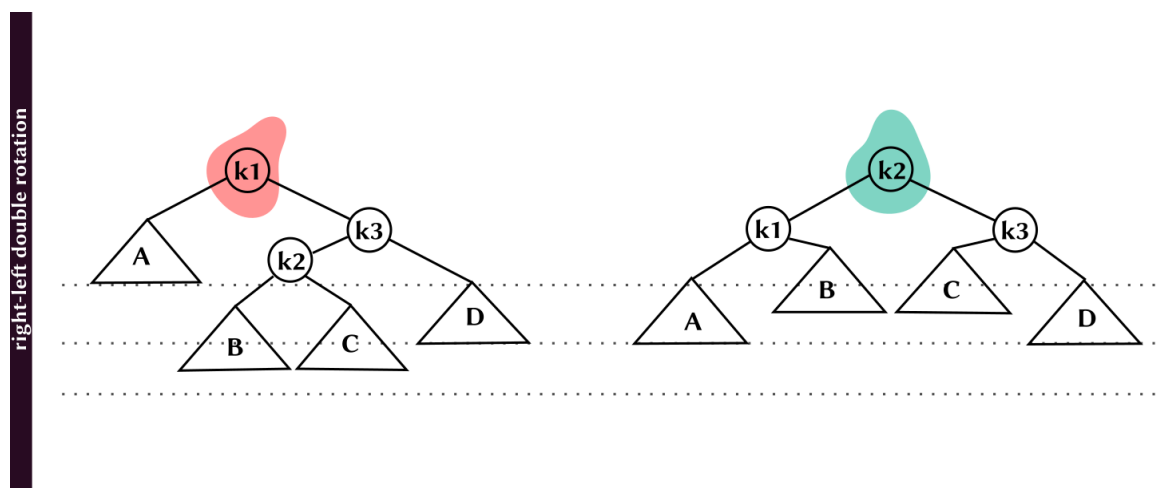
2. Performing the rotation on the right subtree has prepared the root to be rotated left. The tree now is:



3. Rotate the tree left. The result is:

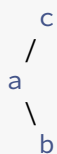


### Double rotation to fix case 3:



**Figure 32:** Right-left double rotation to fix case 3

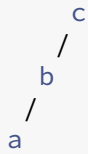
**Simplified example** Take the following tree:



There's a height balance at node **c**.

The steps are the following:

1. Perform a left rotation on the left subtree. Doing so leaves us with this situation:



2. The tree can now be balanced using a single right rotation. Perform the right rotation at **c**. The result is:



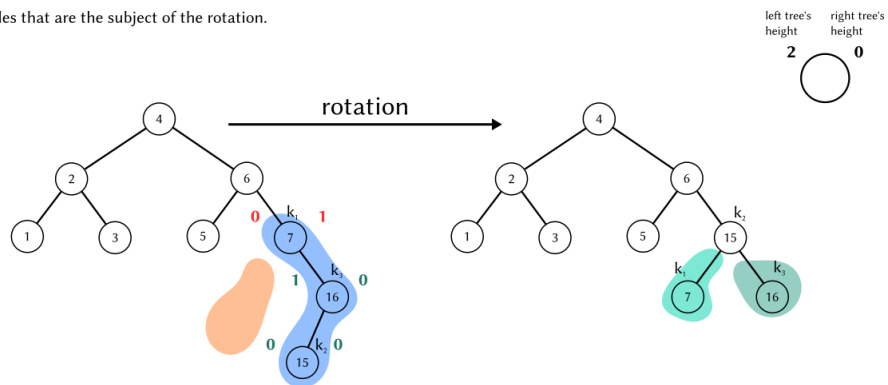
**Example** For this example, we'll continue with the AVL tree depicted in the "Fourth rotation after violation of AVL property" figure from the "Single rotation" section. We'll start by inserting 10 through 16 (in reverse order), followed by 8 and then 9.

Inserting 16 is easy since it doesn't violate the AVL balance property. However, inserting 15 causes a height imbalance at node 7, which is case 3 and solved by a right-left double rotation.

----- A dashed line joins the nodes that are the subject of the rotation.

The insertion of 15 into the tree, causes a height imbalance at at node 7, which then needs to be rebalanced. Here onwards this node is referred as  $\alpha$ .

The insertion occurs at the left subtree of the right child of  $\alpha$ . In other words, an insertion into the left subtree of node 16. This is case 3 and the rotation applied to it is known as **right-left double rotation**.



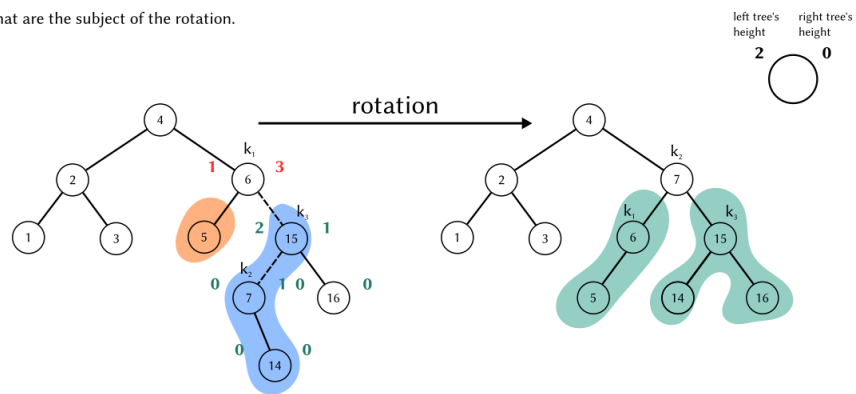
**Figure 33:** Fifth rotation after violation of AVL property

Next we insert 14, which also requires the same double rotation as before.

----- A dashed line joins the nodes that are the subject of the rotation.

The insertion of 14 into the tree, causes a height imbalance at node 6, which then needs to be rebalanced. Here onwards this node is referred as  $\alpha$ .

The insertion occurs at the left subtree of the right child of  $\alpha$ . In other words, an insertion into the left subtree of node 15. This is case 3 and the rotation applied to it is known as **right-left double rotation**.

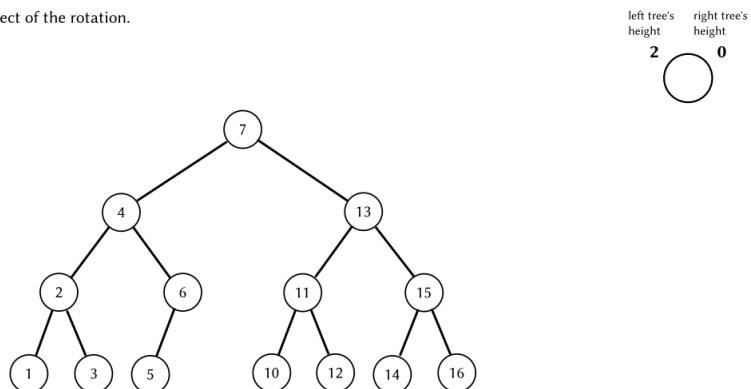


**Figure 34:** Sixth rotation after violation of AVL property

The insertions of 13, 12, 11, and 10 all cause height imbalance that are fixed using single rotations.

----- A dashed line joins the nodes that are the subject of the rotation.

This is the resulting tree after the insertions from 13 to 10. All of them cause some height imbalance that is fixed by applying either a left-left or a right-right single rotation.



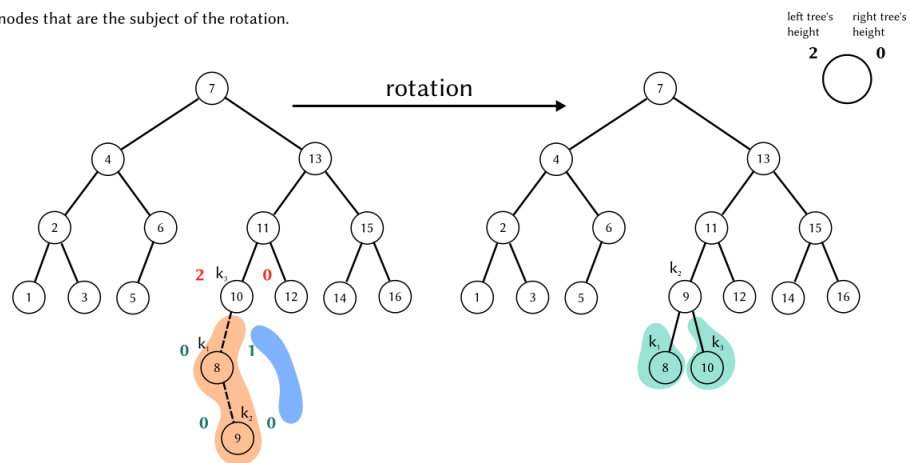
**Figure 35:** Tree after the seventh, eighth, ninth and tenth rotations after violation of AVL property

We can insert 8 without rotation, creating an almost perfectly balanced tree. The insertion of 9 causes a height imbalance at the node containing 10 which is fixed by a left-right double rotation.

----- A dashed line joins the nodes that are the subject of the rotation.

The insertion of 9 into the tree, causes a height imbalance at node 10, which then needs to be rebalanced. Here onwards this node is referred as  $\alpha$ .

The insertion occurs at the right subtree of the left child of  $\alpha$ . In other words, an insertion into the right subtree of node 8. This is case 2 and the rotation applied to it is known as **left-right double rotation**.



**Figure 36:** Tree after the twelfth rotation

### 6.7.2 When do you use rotations? Why?

Deciding when you need a tree rotation is usually easy, but determining which type of rotation you need requires a little thought.

A tree rotation is necessary when you have inserted or deleted a node which leaves the tree in an **unbalanced state**. An unbalanced state is defined as a state in which any subtree has a balance factor of greater than 1, or less than -1. That is, any tree with a difference between the heights of its two subtrees greater than 1, is considered unbalanced.

## 7 Lecture 7: AVL trees (cont.), splay trees, and B-trees

### 7.1 AVL tree implementation

implementation: [implementations/AVLTree](#)

### 7.1.1 Insertion

Input: X is the item to insert and T is the node that roots the subtree.

```
insert( X, T ):
    if T = null:
        T ← new AVLNode(X, null, null)
    else if X < T.item:
        insert(X, T.left)
    else if T.item < X:
        insert(X, T.right)

    balance(T)
```

### 7.1.2 Balance

Brief: Balance the subtree T by performing either a single or double rotation.

Input: T is the tree node to be rebalanced.

```
balance( T ):
    if T = null:
        return
    if T.left.height - T.right.height > 1:
        if T.left.left.height ≥ T.left.right.height:
            rotate-with-left-child(T)
        else:
            double-with-left-child(T)
    else if T.right.height - T.left.height > 1:
        if T.right.right.height ≥ T.right.left.height:
            rotate-with-right-child(T)
        else:
            double-with-right-child(T)

    T.height = max(T.left.height, T.right.height) + 1
```



### 7.1.3 Left-left single rotation (case 1)

Brief: Rotate binary tree with left child. This is a single rotation **for case 1**.

Input: K2 is the tree node that needs to be rebalanced.

```
rotate-with-left-child( K2 ):
    K1      = K2.left
    K2.left = K1.right
    K1.right = K2
    k2.height = max(K1.left.height, K2.right.height) + 1
    K1.height = max(K1.left.height, K2.height) + 1
    K2      = K1
```

### 7.1.4 Left-right double rotation (case 2)

Brief: Rotate binary tree with left child. This is a double rotation **for case 2**.

Input: K3 is the tree node that needs to be rebalanced.

```
double-with-left-child( K3 ):
    rotate-with-right-child(K3.left)
    rotate-with-left-child(K3)
```

### 7.1.5 Deletion

## 7.2 Amortized cost

Consider a sequence of  $M$  operations (insert/delete/find) and suppose the total cost is  $O(M * f(N))$ , irrespective of the actual sequence of operations; the average cost is  $O(f(N))$  for each operation. This is called **amortized running time**. One caveat of this is that individual operations in the sequence can be expensive.

## 7.3 Splay trees

A **splay tree** is a tree with amortized time of  $O(\log N)$ . It guarantees that any  $M$  consecutive tree operations starting from an empty tree take at most  $O(M \log N)$  time (i.e., the tree has an  $O(\log N)$  amortized cost per operation).

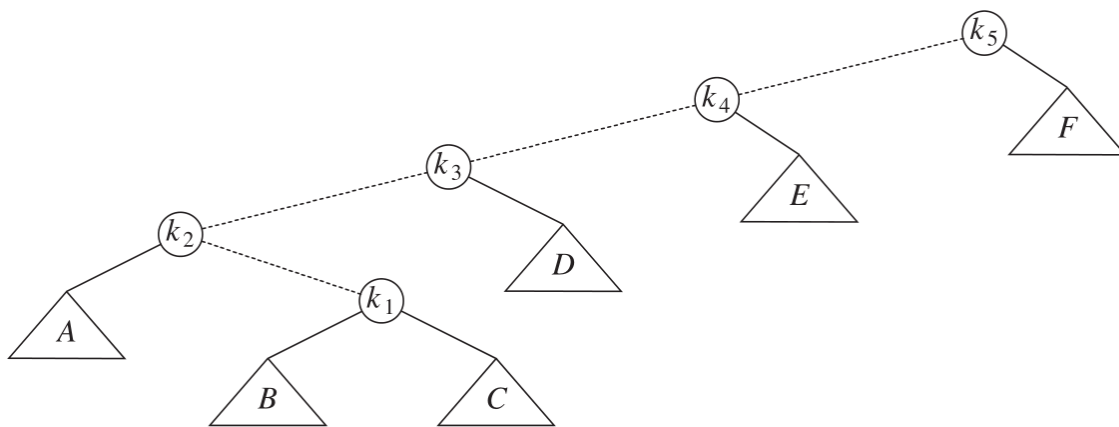
The trick is to rebalance the tree after the `find()` operation. This is done by bringing the item returned by `find()` to the tree's root while applying AVL rotations on the way to the root. The practical utility of

this method is that in many applications when a node is accessed, it is likely to be accessed again in the near future.

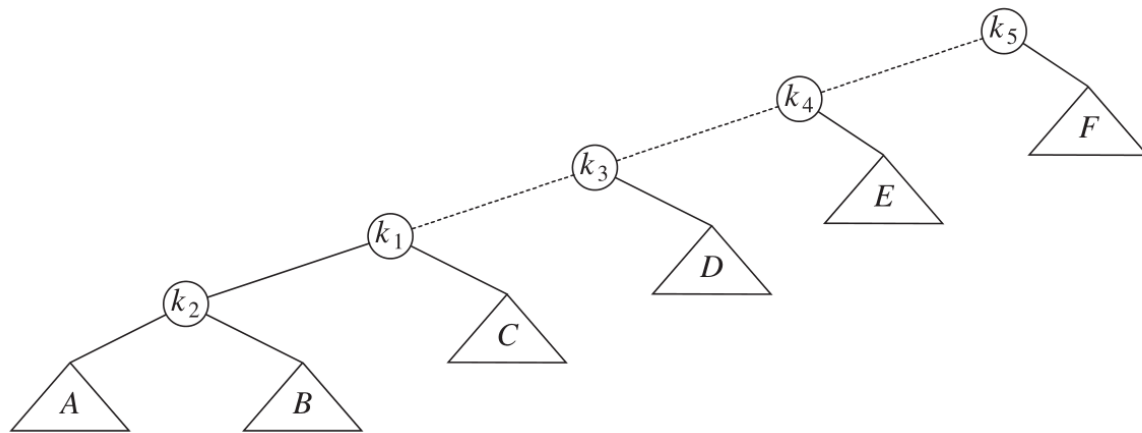
Furthermore, splay trees don't require the maintenance of height or balance information, thus saving space and simplifying the code to some extent.

### 7.3.1 A simple idea (that doesn't work)

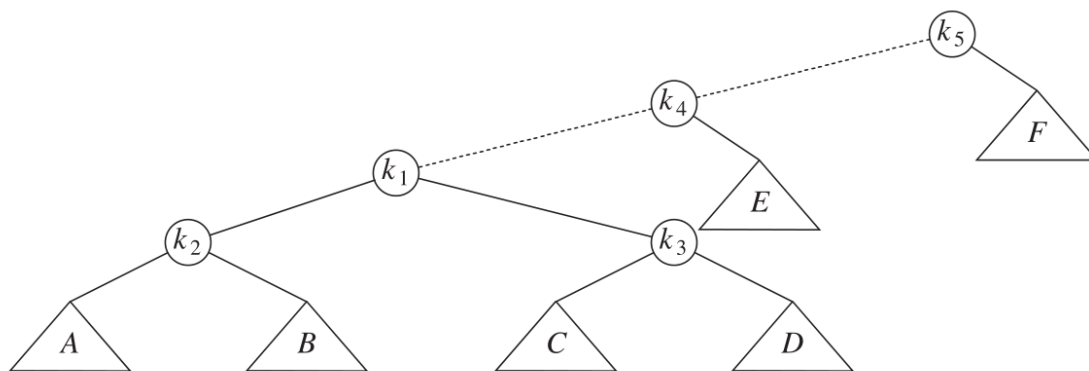
One way to perform the restructuring described above is by performing single rotations from the bottom up. This means that we rotate every node on the access path with its parent. For example, consider the following figures that show what happens after an access (a **find**) on  $k_1$  in the a tree. The rotations have the effect of pushing  $k_1$  all the way to the root making future accesses on  $k_1$  easy but unfortunately, it pushed another node ( $k_3$ ) almost as deep as  $k_1$  used to be.



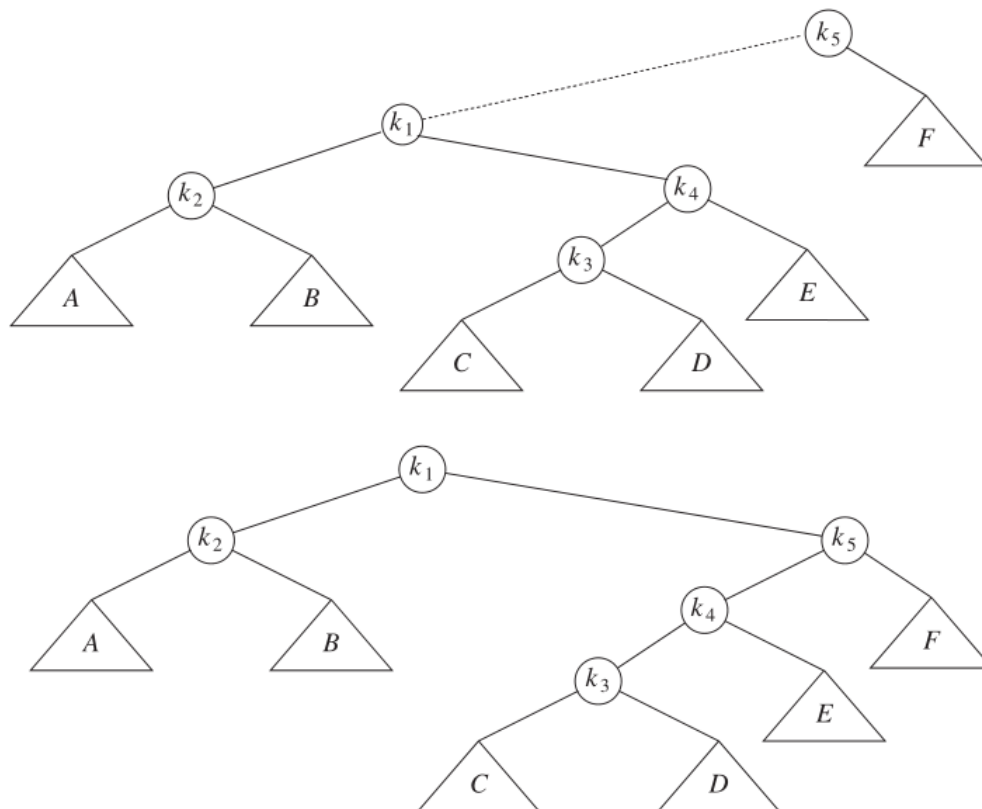
**Figure 37:** Before rotation: pushing  $k_1$  to the root



**Figure 38:** 1st rotation: pushing  $k_1$  to the root



**Figure 39:** 2nd rotation: pushing  $k_1$  to the root



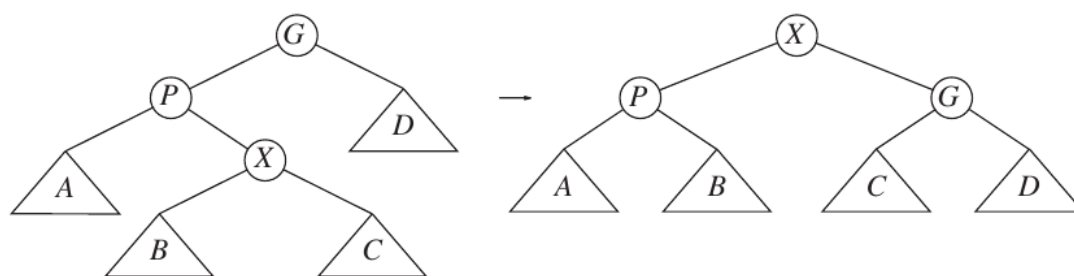
**Figure 40:** 3rd and 4th rotations: pushing  $k_1$  to the root

### 7.3.2 Splaying

The splaying strategy is similar to the rotation from the bottom up, except that we're a little more selective about how rotations are performed.

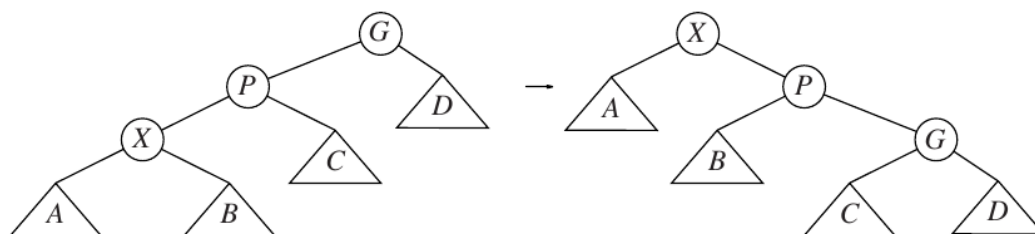
Let  $X$  be a (non-root) node on the access path at which we are rotating. If the parent of  $X$  is the root of the tree, we merely rotate  $X$  and the root. This is the last rotation along the access path. Otherwise,  $X$  has both a parent ( $P$ ) and a grandparent ( $G$ ), and there are two cases, plus symmetries, to consider:

- The first case is the **zig-zag** case. Here  $X$  is a right child and  $P$  is a left child (or vice versa). If this is the case, we perform a double rotation, exactly like an AVL double rotation.



**Figure 41:** Zig-Zag

- The second case is a **zig-zig** case:  $X$  and  $P$  are both left children (or, in the symmetric case, both right children).



**Figure 42:** Zig-Zig

### 7.3.3 Fundamental property of splay trees

- When access paths are long, they lead to longer-than normal search time.
- When accesses are cheap, the rotations are not as good and can be bad.
- The extreme case is the initial tree formed by the insertions.

### 7.3.4 Summary

8 The analysis of splay trees is difficult, because it must take into account the ever-changing structure of the tree.

- Splay trees are much simpler to program than most balanced search trees, since there are fewer cases to consider and no balance information to maintain.
- Some empirical evidence suggests that this translates into faster code in practice, although the case for this is far from complete.

- There are several variations of splay trees that can perform even better in practice.

## 7.4 B-trees

A **B-tree** is a generalization of a binary tree, which is efficient in the *external memory model*.

So far, we have assumed that we can store an entire data structure in the main memory of a computer. Now imagine the scenario of a tree that is immensely huge and can't fit into memory. Now the data structure must reside on the hard, which changes the rules of the game, because the Big-Oh model is no longer meaningful. The problem is that a Big-Oh analysis assumes that all operations are equal. However, this is not true, especially when disk I/O is involved.

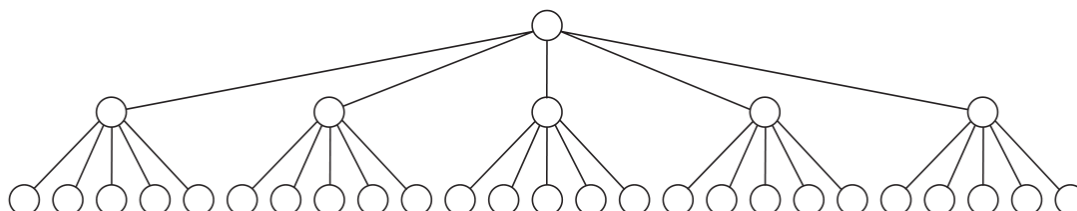
Here is how the typical search tree performs on disk: Suppose we want to access the driving records for citizens in the state of Florida. We assume that we have 10,000,000 items, that each key is 32 bytes (representing a name), and that a record is 256 bytes. We assume this does not fit in main memory and that we are 1 of 20 users on a system (so we have 1/20 of the resources). Thus, in 1 sec we can execute many millions of instructions or perform six disk accesses.

The first idea might be to use a unbalanced binary search tree but it turns into a disaster. In the worst case, it has linear depth and thus could require 10,000,000 disk accesses. On average, a successful search would require  $1.38 \log N$  disk accesses, and since  $\log 10000000 \approx 24$ , an average search would require 32 disk accesses, or 5 sec.

In a typical randomly constructed tree, we would expect that a few nodes are three times deeper; these would require about 100 disk accesses, or 16 sec.

An AVL tree is somewhat better. The worst case of  $1.44 \log N$  is unlikely to occur, and the typical case is very close to  $\log N$ . Thus an AVL tree would use about 25 disk accesses on average, requiring 4 sec.

We want to reduce the number of disk accesses to a very small constant, such as three or four. Basically, we need smaller trees (i.e., trees with more branching and thus less height). An  **$M$ -ary search tree** allows  $M$ -way branching. As branching increases, the depth decreases. Whereas a complete binary tree has height that is roughly  $\log_2 N$ , a complete  $M$ -ary tree has height that is roughly  $\log_M N$ .

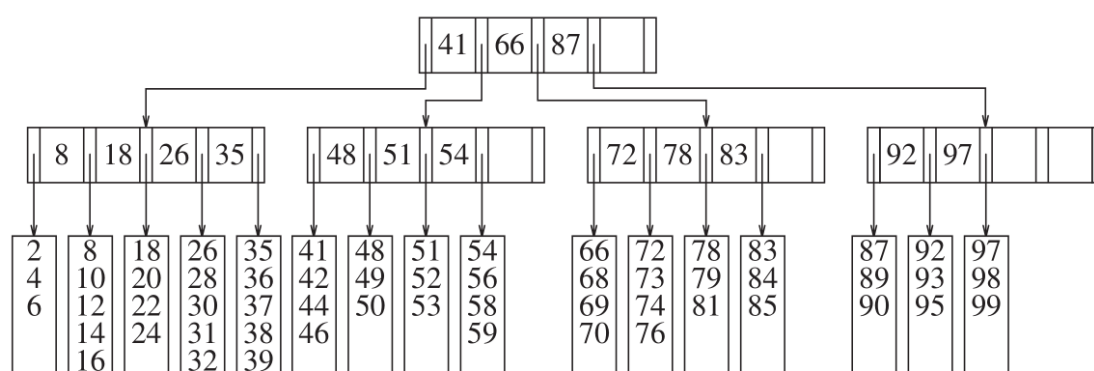


**Figure 43:** 5-ary tree with 31 nodes

A **B-tree** of order  $M$  is an  $M$ -ary tree with the following properties:

1. The data items are stored at leaves.
2. The nonleaf nodes store up to  $M - 1$  keys to guide the searching; key  $i$  represents the smallest key in subtree  $i + 1$ .
3. The root is either a leaf or has between two and  $M$  children.
4. All nonleaf nodes (except the root) have between  $M/2$  and  $M$  children.
5. All leaves are at the same depth and have between  $L/2$  and  $L$  data items, for some  $L$ .

$M$  and  $L$  are determined based on disk block (one access should load a whole node).



**Figure 44:** B-tree of order 5

The “B-tree of order 5” figure depicts a B-tree of order 5:

- Each represents a disk block.
- All nonleaf nodes have between 3 and 5 children.
- $L = M = 5$ , however this not need to be the case.
- Since  $L = 5$ , each leaf has between 3 and 5 data items.

In the Florida driving records example:

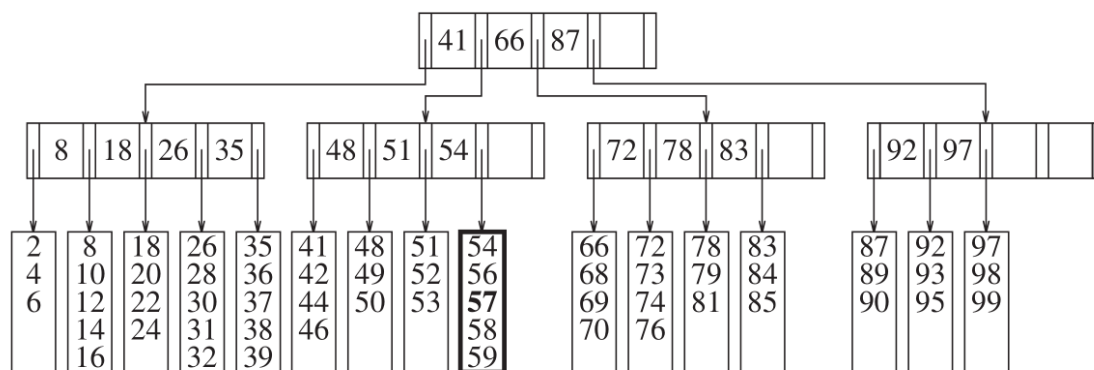
- A single block is supposed to hold 8,192 bytes.
- Each key uses 32 bytes. In a B-tree of order  $M$ , we would’ve  $M - 1$  keys, for a total of  $32 \times (M - 1) = 32M - 32$  bytes, plus  $M$  branches.
- Each branch is essentially a number of another disk block, so we can assume a branch is 4 bytes. Thus the branches use  $4M$  bytes.
- The total memory requirement for a nonleaf node is thus  $32M - 32 + 4M = 36M - 32$ .
- Since  $8192 \leq 36M - 32$ , the largest value of  $M$  for which this is no more than 8192 is 228. Thus we choose  $M = 228$ .

- Since each data record is 256 bytes, we'd be able to fit 32 records in a block. Thus we'd choose  $L = 32$ .
- Each leaf has between 16 and 32 data records and each internal node (except the root) branches in at least 114 ways.
- Since there are 10,000,000 records, there are, at most, 625,000 leaves. Consequently, in the worst case, leaves would be on level 4 of the tree.
  - The worst-case number of accesses is given by approximately  $\log_{M/2} N$ .
  - The root and the next level could be *cached* in main memory, so that over the long run, disk accesses would be needed only for level tree and deeper.

### 7.4.1 Example

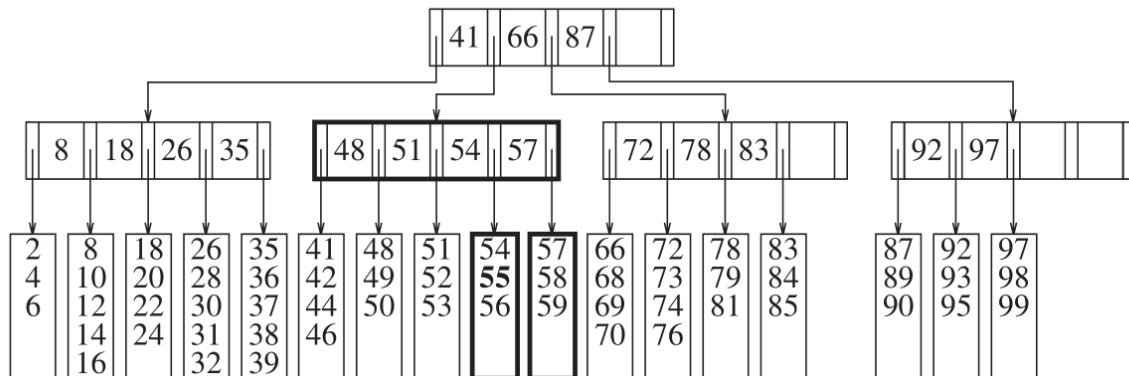
### 7.4.2 Insertion

- Put it in the appropriate leaf.
- If the leaf is full, break it in two, adding a child to the parent.
- If this puts the parent over the limit, split upwards recursively.
- If you need to split the root, add a new one with two children. This is the only way you add depth.

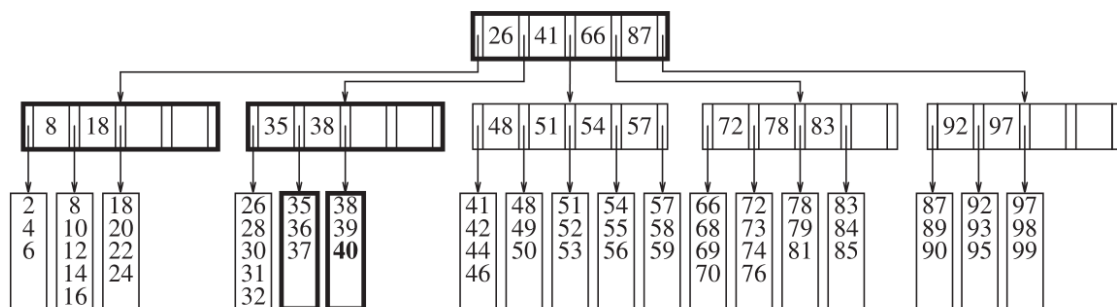


**Figure 45:** B-tree of order 5: Inserting 57





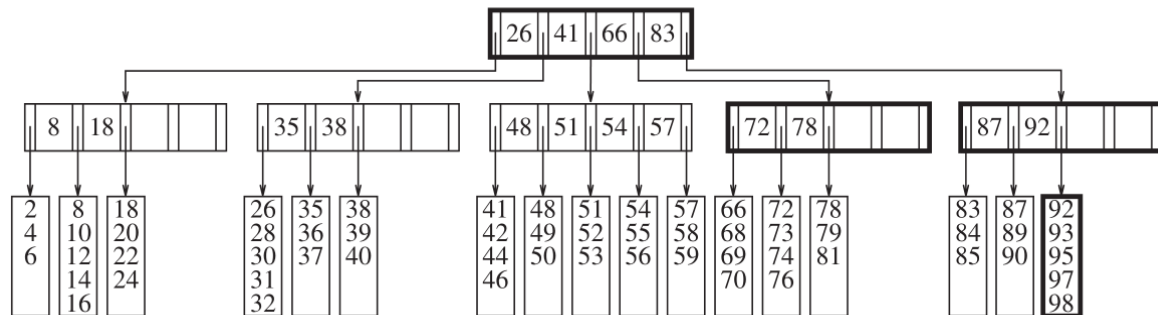
**Figure 46:** B-tree of order 5: Inserting 55



**Figure 47:** B-tree of order 5: Inserting 40

### 7.4.3 Deletion

- Delete from appropriate leaf.
- If the leaf is below its minimum, adopt from a neighbor if possible.
- If that's not possible, you can merge with the neighbor. This causes the parent to lose a branch and you continue upward recursively.



**Figure 48:** B-tree of order 5: Inserting 40

## 8 Lecture 08: Sets and maps

### 8.1 STL containers

The *vector* and *list* containers are inefficient for search and insert. Thus, the STL provides the *set* and *map* containers which guarantee logarithmic time for insertion, deletion and searching.

### 8.2 The set container

A *set* is an ordered container that doesn't allow duplicates:

- It stores objects of type `Key` in sorted order.
- The `Value` is the `Key` itself, without any additional data.

#### 8.2.1 Insertion

The `insert` method returns an `iterator` which represents the position of the newly inserted item or the existing item that caused the insertion to fail (i.e., no duplicates).

The STL defines a class template called `pair` that is little more than a struct with members `first` and `second` to access the two items in the pair. There are two different `insert` routines:

- `pair<iterator,bool> insert( const Object & x );`
- `pair<iterator,bool> insert( iterator hint, const Object & x );`

`pair<T1, T2>` is a heterogenous pair; it holds one object of type `T1` and another of type `T2`. For example:

```
pair<bool, double> result; // a pair with key as boolean and value as
    double
result.first = true; // the method first access the pair's key
result.second = 2.5; // the method second accesses the pair's value
```

### 8.2.2 Insertion with hint

The two-parameter `insert` allows the specification of a hint, which represents the position where `x` should go. If the hint is accurate, the insertion is fast, often  $O(1)$ . If not, the insertion is done using the normal insertion algorithm and performs comparably with the one-parameter `insert`. For instance, the following code might be faster using the two-parameter `insert` rather than the one-parameter `insert`:

```
set<int> s;
for (int i = 0; i < 10; i++) {
    const auto result = s.insert(s.end(), i);
}
```

### 8.2.3 erase

There are several versions of `erase`:

- `int erase( const Object & x );` This version removes `x` (if found) and returns the number of items actually removed, which is obviously either 0 or 1.
- `iterator erase( iterator itr );` This one behaves the same as in `vector` and `list`, and it removes the object at the position given by the iterator, returns an iterator representing the element that followed `itr` immediately prior to the call to `erase`, and invalidates `itr`, which becomes stale.
- `iterator erase( iterator start, iterator end );` This one behaves the same as in a `vector` or `list`, removing all the items starting at `start`, up to but not including the item at `end`.

### 8.2.4 find

The set provides a `find` routine that returns an iterator representing the location of the item (or the endmarker if the search fails). This provides considerably more information, at no cost in running time. The signature of `find` is `iterator find( const Object & x ) const`;

By default, ordering uses the `less<Object>` function object, which itself is implemented by invoking `operator<` for the `Object`. An alternative ordering can be specified by instantiating the `set` template

with a function object type. For instance, we can create a `set` that stores string objects, ignoring case distinctions by using the `CaseInsensitiveCompare` function object:

```
class CaseInsensitiveCompare {
public:
    bool operator()( const string & lhs, const string & rhs ) const {
        return strcasecmp( lhs.c_str( ), rhs.c_str( ) ) < 0;
    }
};

set<string, CaseInsensitiveCompare> s;
s.insert("Hello");
s.insert("HeLLo");
std::cout << "The size is: " << s.size( ) << std::endl; // size is 1
```

### 8.3 The map container

A `map` is used to store a collection of ordered entries that consists of keys and their values:

- Keys must be unique, but several keys can map to the same values. Thus values need not be unique.
- The keys in the `map` are maintained in logically sorted order.

The `map` also supports `insert`, `find`, `erase`, `begin`, `end`, `size`, and `empty`. For `insert`, one must provide a `pair<KeyType, ValueType>` object. Although `find` requires only a key, the `iterator` it returns references a `pair`.

There's another operation that yields simple syntax. The array-indexing operator is overloaded for maps as follows:

```
ValueType& operator[] ( const KeyType & key );
```

The semantics of `operator[]` are as follows:

- If `key` is present in the `map`, a reference to the corresponding value is returned.
- If `key` is not present in the map, it is inserted with a default value into the map and then a reference to the inserted default value is returned. The default value is obtained by applying a zero-parameter constructor or is zero for the primitive types.

These semantics do not allow an accessor version of `operator[]`, so `operator[]` cannot be used on a `map` that is constant. For instance, if a map is passed by constant reference, inside the routine, `operator[]` is unusable.

### 8.3.1 Example

```
#include <iostream>
#include <string>
#include <map>

int main() {
    std::map<std::string, double> salaries;

    salaries["Pat"] = 78000.0;

    std::cout << salaries["Pat"] << std::endl;
    std::cout << salaries["Per"] << std::endl;

    std::map<std::string, double>::const_iterator itr;
    itr = salaries.find("Paul");

    if (itr == salaries.end()) {
        std::cout << "Not an employee" << std::endl;
    }
    else {
        std::cout << itr->second << std::endl;
    }

    return 0;
}
```

## 8.4 Implementation of set/map in C++

C++ requires that set and map support the basic `insert`, `erase`, and `find` operations in logarithmic worst-case time. The underlying implementation is a **balanced binary search tree**.

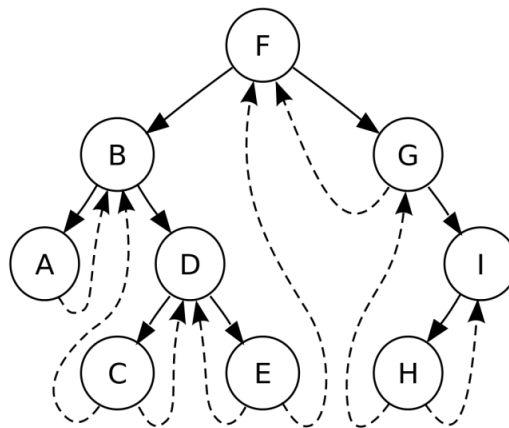
An important issue in implementing `set` and `map` is providing support for the `iterator` classes. Of course, internally, the `iterator` maintains a pointer to the “current” node in the iteration. The hard part is efficiently advancing to the next node. There are several possible solutions, but the smart solution is to use **threaded** trees which are used throughout different implementations in the STL.

This solution involves maintaining the extra links only for nodes that have `nullptr` left or right links by using extra `Boolean` variables to allow the routines to tell if a left link is being used as a standard binary search tree left link or a link to the next smaller node, and similarly for the right link.

### 8.4.1 Threaded binary trees

A binary tree is threaded by

- making all right child pointers (that would normally be `nullptr`) point to the in-order successor of the node (if it exists), and
- making all left child pointers (that would normally be `nullptr`) point to the in-order predecessor of the node.



**Figure 49:** Threaded tree, special threaded links are dashed arrows

In this type of tree, we have the pointers reference the next node in a in-order traversal. These are known as threads.

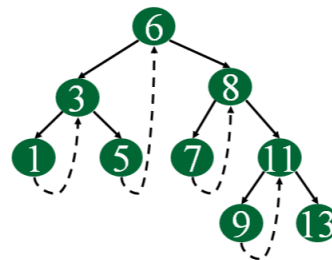
We need to know if a pointer is an actual link or a thread, so we keep a boolean for each pointer.

## Why do we need threaded binary trees?

- Binary trees have a lot of wasted space.
- Threaded binary trees make the tree traversal faster since we don't need a stack or use recursion for traversal.

### 8.4.2 Types of threaded binary trees

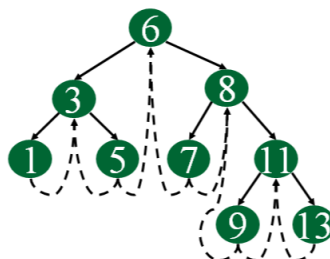
- **Single threaded:** Each node is threaded towards either the in-order predecessor or successor (left or right). This means all right null pointers will point to the in-order successor **or** all left null pointers will point to the in-order predecessor.



Single Threaded Binary Tree

**Figure 50:** Single threaded binary tree

- **Double threaded:** Each node is threaded towards both the in-order predecessor and successor (left and right). This means that all right null pointers will point to the in-order successor **and** all left null pointers will point to the in-order predecessor.



Double Threaded Binary Tree

**Figure 51:** Double threaded binary tree

### 8.4.3 An example

We would like to write a program to find all words that can be changed into at least 15 other words by a single one-character substitution. We assume that we have a dictionary consisting of approximately 89,000 different words of varying lengths. Most words are between 6 and 11 characters.

**First approach:** The most straightforward strategy is to use a [map](#) in which the keys are words and

the values are vectors containing the words that can be changed from the key with a one-character substitution.

```
/*
@brief The function prints the contents of the map only for the elements
for which the vector of strings has size greater than or equal to
    min_words.
@param adjacentWords input map from string to vector of strings.
@param min_words minimum number of words to consider printing.
*/
void print_high_changeables(
    const std::map<std::string>, std::vector<std::vector>> adjacent_words
    , int min_words = 15
) {
    for (auto& entry : adjacent_words) {
        const std::vector<std::string>& words = entry.second;
        if (words.size() >= min_ords) {
            std::cout << entry.first << "(" << words.size() << "): ";
            for (auto& str : words) {
                std::cout << " " << str << std::endl;
            }
        }
    }
}
```

The function above only shows how to print the the `map` that is eventually produced. How do we construct the `map` though?

One way is to test if two words are identical except for a one-character substitution. Then, we can brute-force our way testing all pairs of words and constructing the `map`. If we find a pair of words that differ in only one character, we can update the `map` by adding this new word (that differs by a single character) to the `map`.



```
/*
@brief Check if two words differ by a single character.
@return bool
*/
bool nearly_equal( const std::string& word1, const std::string& word2 ) {
    if (word1.length() != word2.length()) return false;
    int diffs = 0;
    for (int i = 0; i < word1.length(); i++) {
        if (word1[i] != word2[i]) {
            if (++diffs > 1) return false;
        }
    }
    return diffs == 1;
}

std::map<std::string, std::vector<std::string>> compute_adjacent_words(
    const std::vector<std::string>& words
) {
    std::map<std::string, std::vector<std::string>> adj_words;
    for (int i = 0; i < words.size(); i++) {
        for (int j = i; j < words.size(); j++) {
            adj_words[ words[i] ].push_back(words[j]);
            adj_words[ words[j] ].push_back(words[i]);
        }
    }
    return adj_words;
}
```

### Second approach:

The problem with the algorithm above is that it's slow. An obvious improvement is to avoid comparing words of different lengths. We can do this by grouping words by their length, and then running the previous algorithm on each of the separate groups. To do this, we can use a second `map` where the key is an integer representing a word length, and the value is a collection of all words of that length.

```

/*
@brief Compute a map in which the keys are words and values are
vectors of words that differ in only one character from the
corresponding key. Uses a quadratic algorithm, but speeds things
up a little by maintaining an additional map that groups words
by their length.
*/
std::map<std::string, std::vector<std::string>> compute_adjacent_words(
    const std::vector<std::string>& words
) {
    std::map<std::string, std::vector<std::string>> adj_words;
    std::map<int, std::vector<std::string>> words_by_length;

    // group the words by their length
    for (auto& word : words) {
        words_by_length[ word.length() ].push_back(word);
    }

    // work on each group separately
    for (auto& entry : words_by_length) {
        const std::vector<std::string>& groups_words = entry.second;

        for (int i = 0; i < groups_words.size(); i++) {
            for (int j = i + 1; j < groups_words.size(); j++) {
                if (nearly_equal(groups_words[i], groups_words[j])) {
                    adj_words[ groups_words[i] ].push_back( groups_words[j] );
                    adj_words[ groups_words[j] ].push_back( groups_words[i] );
                }
            }
        }
    }

    return adj_words;
}

```

### Third approach:

This third algorithm is more complex and uses additional `maps`. As before, words are grouped by lengths, and then each group is worked on separately.

```

std::map<std::string, std::vector<std::string>> compute_adjacent_words(
    const std::vector<std::string>& words
) {

    std::map<std::string, std::vector<std::string>> adj_words;
    std::map<int, std::vector<std::string>> words_by_length;

    // group words by their length
    for (auto& word : words) {
        word_by_length[ word.length() ].push_back(word);
    }

    // work on each group separately
    for (auto& entry : words_by_length) {
        const std::vector<std::string>& groups_words = entry.second;
        int group_num = entry.first;

        // work on each position in each group
        for (int i = 0; i < group_num; i++) {
            std::map<std::string, std::vector<std::string>> rep_to_word;

            // Remove one character in specified position, computing
            // representative.
            // Words with same representatives are adjacent; so populate a
            // map ...
            for (auto& word : groups_words) {
                std::string rep = word;
                rep.erase(i, 1);
                rep_to_word[ rep ].push_back(word);
            }

            // ... and then look for map values with more than one string
            for (auto& entry : rep_to_word) {
                const std::vector<std::string>& clique = entry.second;
                if (clique.size() >= 2) {
                    for (int p = 0; p < clique.size(); p++) {
                        for (int q = p + 1; q < clique.size(); q++) {
                            adj_words[ clique[p] ].push_back(clique[q]);
                            adj_words[ clique[q] ].push_back(clique[p]);
                        }
                    }
                }
            }
        }

        return adj_words;
    }
}

```

## 8.5 Summary to avoid performance issues

### 8.5.1 AVL trees

- The heights of left and right subtrees differ at most by 1; not too deep.
- Operations such as `insert` must restore the tree.

### 8.5.2 Splay trees

- Nodes can get arbitrarily deep but after every access the tree is adjusted either using *zig-zag* or *zig-zig*.
- The net effect is that any sequence of  $M$  operations takes  $O(M \log N)$  which is the same as a balanced tree.

### 8.5.3 B-trees

- Balanced  $M$ -way (as opposed to binary trees) which are well suited to the external memory model.