

# Computer Science II

Dr. Chris Bourke  
Department of Computer Science & Engineering  
University of Nebraska—Lincoln  
Lincoln, NE 68588, USA  
<http://chrisbourke.unl.edu>  
[cbourke@cse.unl.edu](mailto:cbourke@cse.unl.edu)

2017/06/01 13:58:39  
Version 0.1.0



This book is a *draft* covering Computer Science II topics as presented in CSCE 156 (Computer Science II) at the University of Nebraska—Lincoln.



This work is licensed under a [Creative Commons Attribution-ShareAlike 4.0 International License](https://creativecommons.org/licenses/by-sa/4.0/)



# Contents

<b>1</b>	<b>Introduction</b>	<b>1</b>
<b>2</b>	<b>Object Oriented Programming</b>	<b>3</b>
2.1	Introduction . . . . .	3
2.2	Objects . . . . .	3
2.3	The Big Four . . . . .	3
2.3.1	Abstraction . . . . .	3
2.3.2	Encapsulation . . . . .	3
2.3.3	Inheritance . . . . .	3
2.3.4	Polymorphism . . . . .	3
<b>3</b>	<b>Relational Databases</b>	<b>5</b>
<b>4</b>	<b>List-Based Data Structures</b>	<b>7</b>
4.1	Array-Based Lists . . . . .	8
4.1.1	Designing a Java Implementation . . . . .	8
4.2	Linked Lists . . . . .	16
4.2.1	Designing a Java Implementation . . . . .	22
4.2.2	Variations . . . . .	26
4.3	Stacks & Queues . . . . .	28
4.3.1	Stacks . . . . .	29
4.3.2	Queues . . . . .	33
4.3.3	Variations . . . . .	37
<b>5</b>	<b>Algorithm Analysis</b>	<b>39</b>
5.1	Introduction . . . . .	39
5.1.1	Example: Computing a Sum . . . . .	45
5.1.2	Example: Computing a Mode . . . . .	47
5.2	Pseudocode . . . . .	50
5.3	Analysis . . . . .	52
5.4	Asymptotics . . . . .	57
5.4.1	Big-O Analysis . . . . .	57
5.4.2	Other Notations . . . . .	59
5.4.3	Observations . . . . .	61
5.4.4	Limit Method . . . . .	64

5.5	Examples	66
5.5.1	Linear Search	66
5.5.2	Set Operation: Symmetric Difference	68
5.5.3	Euclid's GCD Algorithm	69
5.5.4	Selection Sort	71
5.6	Other Considerations	72
5.6.1	Importance of Input Size	72
5.6.2	Control Structures are Not Elementary Operations	75
5.6.3	Average Case Analysis	76
5.6.4	Amortized Analysis	77
5.7	Analysis of Recursive Algorithms	78
5.7.1	The Master Theorem	79
<b>6</b>	<b>Trees</b>	<b>83</b>
6.1	Introduction	83
6.2	Definitions & Terminology	83
6.3	Implementation	92
6.4	Tree Traversal	93
6.4.1	Preorder Traversal	94
6.4.2	Inorder Traversal	96
6.4.3	Postorder Traversal	99
6.4.4	Tree Walk Traversal	103
6.4.5	Breadth-First Search Traversal	104
6.5	Binary Search Trees	107
6.5.1	Retrieval	108
6.5.2	Insertion	110
6.5.3	Removal	110
6.5.4	In Practice	114
6.6	Heaps	115
6.6.1	Operations	116
6.6.2	Implementations	121
6.6.3	Variations	125
6.6.4	Applications	126
6.7	Exercises	129
	<b>Glossary</b>	<b>131</b>
	<b>Acronyms</b>	<b>133</b>
	<b>Index</b>	<b>136</b>
	<b>References</b>	<b>137</b>

# List of Algorithms

1	Insert-At-Head Linked List Operation	17
2	Insert Between Two Nodes Operation	18
3	Index-Based Retrieval Operation	20
4	Key-Based Delete Operation	22
5	Computing the Mean	51
6	Computing the Mode	52
7	Trivial Sorting (Bad Pseudocode)	52
8	Trivially Finding the Minimal Element	53
9	Finding the Minimal Element	53
10	Linear Search	67
11	Symmetric Difference of Two Sets	68
12	Euclid's GCD Algorithm	69
13	Selection Sort	71
14	Sieve of Eratosthenes	73
15	Fibonacci( $n$ )	79
16	Binary Search – Recursive	81
17	Merge Sort	81
18	Stack-based Preorder Tree Traversal	95
19	preOrderTraversal( $u$ ): Recursive Preorder Tree Traversal	96
20	Stack-based Inorder Tree Traversal	97
21	inOrderTraversal( $u$ ): Recursive Inorder Tree Traversal	99
22	Stack-based Postorder Tree Traversal	100
23	postOrderTraversal( $u$ ): Recursive Postorder Tree Traversal	103
24	Tree Walk based Tree Traversal	105
25	Queue-based BFS Tree Traversal	106
26	Search algorithm for a binary search tree	110
27	Finding the maximum key value in a node's left subtree.	113
28	Heapify	117
29	Find Next Open Spot - Numerical Technique	125
30	Heap Sort	127





# List of Code Samples

4.1	Parameterized Array-Based List in Java . . . . .	14
4.2	A linked list node Java implementation. Getter and setter methods have been omitted for readability. A convenience method to determine if a node has a next element is included. This implementation uses <code>null</code> as its terminating value. . . . .	24
5.1	Summing a collection of integers . . . . .	42
5.2	Summation Algorithm 1 . . . . .	45
5.3	Summation Algorithm 2 . . . . .	45
5.4	Summation Algorithm 3 . . . . .	45
5.5	Mode Finding Algorithm 1 . . . . .	47
5.6	Mode Finding Algorithm 2 . . . . .	48
5.7	Mode Finding Algorithm 3 . . . . .	49
5.8	Naive Exponentiation . . . . .	74
5.9	Computing an Average . . . . .	75



# List of Figures

4.1	A simple linked list containing 3 nodes. . . . .	17
4.2	Insert-at-head Operation in a Linked List. We wish to insert a new element, 42 at the head of the list. . . . .	18
4.3	Inserting Between Two Nodes in a Linked List. Here, we wish to insert a new element 42 between the given two nodes containing 8 and 6. . . . .	19
4.4	Delete Operation in a Linked List . . . . .	21
4.5	Key-Based Find and Remove Operation. We wish to remove the first node we find containing 42. . . . .	23
4.6	A Doubly Linked List Example . . . . .	26
4.7	A Circularly Linked List Example . . . . .	27
4.8	An Unrolled Linked List Example . . . . .	28
4.9	A stack holding integer elements. Push and pop operations are depicted as happening at the “top” of the stack. In actuality, a stack stored in a computer’s memory is not really oriented but this visualization is consistent with a physical stack growing “upwards.” . . . . .	30
4.10	An example of a queue. Elements are enqueued at the end of the queue and dequeued from the front of the queue. . . . .	34
4.11	Array-Based Queue . . . . .	36
4.12	Producer Consumer Pattern . . . . .	37
5.1	Quadratic Regression of Index-Based Linked List Performance . . . . .	44
5.2	Plot of two functions. . . . .	56
5.3	Expected number of comparisons for various success probabilities $p$ . . . . .	77
6.1	An undirected graph with labeled vertices. . . . .	84
6.2	Several examples of trees. It doesn’t matter how the tree is depicted or organized, only that it is acyclic. The final example, 6.2(d) represents a disconnected tree, called a <i>forest</i> . . . . .	86
6.3	Several possible rooted orientations for the tree from Figure 6.2(a). . . . .	87
6.4	A Tree Node’s Relations. . . . .	89
6.5	A Binary Tree . . . . .	89
6.6	A complete tree of depth $d = 3$ which has $1 + 2 + 4 + 8 = 15$ nodes. . . . .	90
6.7	A summation of nodes at each level of a complete binary tree up to depth $d$ . . . . .	91
6.8	A small binary tree. . . . .	94
6.9	A walkthrough of a preorder traversal on the tree from Figure 6.5. . . . .	95
6.10	A walkthrough of an inorder traversal on the tree from Figure 6.5. . . . .	98

*List of Figures*

6.11	A walkthrough of a postorder traversal on the tree from Figure 6.5. . . .	102
6.12	A tree walk on the tree from Figure 6.8. . . . .	104
6.13	The three general cases of when to process a node in a tree walk. . . . .	104
6.14	A Breadth First Search Example . . . . .	106
6.15	A Binary Search Tree . . . . .	108
6.16	Various Search Examples on a Binary Search Tree . . . . .	109
6.17	Binary Search Tree Insertion Operation . . . . .	111
6.18	Binary Search Tree Deletion Operation Examples . . . . .	113
6.19	A degenerate binary search tree. . . . .	114
6.20	A min-heap . . . . .	115
6.21	A Max-heap . . . . .	116
6.22	An Invalid Max-heap . . . . .	117
6.23	Insertion and Heapification . . . . .	118
6.24	Another Invalid Heap . . . . .	119
6.25	Removal of the root element (getMax) and Heapification . . . . .	120
6.26	Heap Node's Index Relations . . . . .	122
6.27	An array implementation of the heap from Figure 6.21 along with the generalized parent, left, and right child relations. . . . .	122
6.28	Tree-based Heap Analysis . . . . .	124

# 1 Introduction

To Come.



# **2 Object Oriented Programming**

## **2.1 Introduction**

To Come.

## **2.2 Objects**

## **2.3 The Big Four**

### **2.3.1 Abstraction**

### **2.3.2 Encapsulation**

### **2.3.3 Inheritance**

### **2.3.4 Polymorphism**





# 3 Relational Databases

To Come.



## 4 List-Based Data Structures

Most programming languages provide some sort of support for storing collections of similar elements. The most common way is to store elements in an *array*. That is, elements are stored together in a contiguous chunk of memory and individual elements are accessed using an *index*. An index represents an *offset* with respect to the first element which is usually stored at index 0 (referred to as *zero-indexing*).

There are several disadvantages to using arrays, however. In particular, once allocated, the capacity of an array is fixed. It cannot grow to accommodate new elements and it cannot shrink if we end up removing elements. Moreover, we may not need the full capacity of the array at any given point in a program, leading to wasted space. Though libraries may provide convenience functions, in general all of the “bookkeeping” in an array is up to us. If we remove an element in the middle of the array, our data may no longer be contiguous. If we add an element in the array, we have to make sure to find an available spot. Essentially, all of the organization of the array falls to the user.

A much better solution is to use a dynamic data structure called a *List*. A List is an [Abstract Data Type \(ADT\)](#) that stores elements in an *ordered* manner. That is, there is a notion of a “first” element, a “second” element, etc. This is *not* necessarily the same thing as being *sorted*. A list containing the elements 10, 30, 5 is not sorted, but it is ordered (10 is the first element, 30 is the second, and 5 is the third and final element). In contrast to an array, the list automatically organizes the elements in some underlying structure and provides an *interface* to the user that provides some set of core functionality, including:

- A way to add elements to the list
- A way to retrieve elements from the list
- A way to remove elements from the list

in some manner. We’ll examine the specifics later on, but the key aspect to a list is that, in contrast to an array, it dynamically expands and contracts automatically as the user adds/removes elements.

How a list supports this core functionality may vary among different implementations. For example, the list’s interface may allow you to add an element to the beginning of the list, or to the end of the list, or to add the new element at a particular index; or any combination of these options. The retrieval of elements could be supported by providing an index-based retrieval method or an iterator pattern that would allow a user

to conveniently iterate over every element in the list.

In addition, a list may provide secondary functionality as a convenience to users, making the implementation more flexible. For example, it may be useful for a user to tell how many elements are in the list; whether or not it is empty or full (if it is designed to have a constrained capacity). A list might also provide batch methods to allow a user to add a collection of elements to the list rather than just one at a time.

Most languages will provide a list implementation (or several) as part of their standard library. In general, the best practice is to use the built-in implementations unless there is a very good reason to “roll your own” and create your own implementation. However, understanding how various implementations of lists work and what properties they provide is very important. Different implementations provide advantages and disadvantages and so using the correct one for your particular application may mean the difference between an efficient algorithm and an inefficient or even infeasible one.

### 4.1 Array-Based Lists

Our first implementation is an obvious extension of basic arrays. We’ll still use a basic array to store data, but we’ll build a data structure around it to implement a full list. The details of how the list works will be encapsulated inside the list and users will interact with the list through publicly available methods.

The basic idea is that the list will own an array (via composition) with a certain *capacity*. When users add elements to the list, they will be stored in the array. When the array becomes full, the list will automatically reallocate a new, larger array (giving the list a larger capacity), copy over all the elements in the old array and then switch to this new array. We can also implement the opposite functionality and shrink the array if we desire.

#### 4.1.1 Designing a Java Implementation

To illustrate this design, consider the basic the following code sample in Java.

```
1 public class IntegerArrayList {
2
3     private int arr[];
4     private int size;
5
6 }
```

This array-based list is designed to store integers in the `arr` array. The second member variable, `size` will be used to track the number of elements stored in `arr`. Note that this is not the same thing as the size of the array (that is, `arr.length`). Elements

may or may not be stored in each array position. To distinguish between the size of the array-based list and the size of the internal array, we'll refer to them as the *size* and *capacity*

To initialize an empty array list, we'll create a default constructor and instantiate the array with an initial *capacity* of 10 elements with an initial size of 0.

```

1 public IntegerArrayList() {
2     this.arr = new int[10];
3     this.size = 0;
4 }

```

## Adding Elements

Now let's design a method to provide a way to add elements. As a first attempt, let's allow users to add elements to the *end* of the list. That is, if the list currently contains the elements 8, 6, 10 and the user adds the element 42 the list will then contain the elements 8, 6, 10, 42.

Since the `size` variable keeps track of number of elements in the list, we can use it to determine the index at which the element should be inserted. Moreover, once we insert the element, we'll need to be sure to increment the `size` variable since we are increasing the number of elements in the list. Before we do any of this, however, we need to check to ensure that the underlying array has enough room to hold the new element and if it doesn't, we need to increase the capacity by creating a new array and copying over the old elements. This is all illustrated in the following code snippet.

```

1 public void addAtEnd(int x) {
2
3     //if the array is at capacity, resize it
4     if(this.size == this.arr.length) {
5         //create a new array with a larger capacity
6         int newArr = new int[this.arr.length + 10];
7         //copy over all the old elements
8         for(int i=0; i<this.arr.length; i++) {
9             newArr[i] = this.arr[i];
10        }
11        //use the new array
12        this.arr = newArr;
13    }
14    this.arr[this.size] = x;
15    this.size++;
16 }

```

The astute Java programmer will note that lines 5–12 can be improved by utilizing

methods provided by Java's `Arrays` class. Adhering to the Don't Repeat Yourself (DRY) principle, a better version would be as follows.

```
1 public void addAtEnd(int x) {
2
3     if(this.size == this.arr.length) {
4         this.arr = Arrays.copyOf(this.arr, this.arr.length + 10);
5     }
6     this.arr[size] = x;
7     this.size++;
8 }
```

As another variation, we could allow users to add elements at an arbitrary index. That is, if the array contained the elements `8, 6, 10`, we could allow the user to insert the element `42` at any index 0 through 3. The list would automatically shift elements down to accommodate the new element. For example:

- Adding at index 0 would result in `42, 8, 6, 10`
- Adding at index 1 would result in `8, 42, 6, 10`
- Adding at index 2 would result in `8, 6, 42, 10`
- Adding at index 3 would result in `8, 6, 10, 42`

Note that though there is no element (initially) at index 3, we still allow the user to “insert” at that index to allow the user to insert the element at the end. However, any other index should be considered invalid as it would either be invalid (negative) or it would mean that the data is no longer contiguous. For example, adding `42` at index 5 may result in `8, 6, 10, null, null, 42`. Thus, we do some basic index checking and throw an exception for invalid indices.

```

1 public void insertAtIndex(int x, int index) {
2
3     if(index < 0 || index > this.size) {
4         throw new IndexOutOfBoundsException("invalid index: " + index);
5     }
6
7     if(this.size == this.arr.length) {
8         this.arr = Arrays.copyOf(this.arr, this.arr.length + 10);
9     }
10
11     //start at the end; shift elements to the right to accommodate x
12     for(int i=this.size-1; i>=index; i--) {
13         this.arr[i+1] = this.arr[i];
14     }
15     this.arr[index] = element;
16     this.size++;
17 }

```

At this point we note that these two methods have a lot of code in common. In fact, one could be implemented in terms of the other. Specifically, `insertAtIndex` is the more general of the two and so `addToEnd` should use it:

```

1 public void addAtEnd(int x) {
2
3     this.insertAtIndex(x, this.size);
4 }

```

This is a big improvement as it reduces the complexity of our code and thus the complexity of testing, maintenance, etc.

## Retrieving Elements

In a similar manner, we can allow users to retrieve elements using an index-based retrieval method. We will again take care to do index checking, but otherwise returning the element is straightforward.

```

1 public void getElement(int index) {
2
3     if(index < 0 || index >= this.size) {
4         throw new IndexOutOfBoundsException("invalid index: " + index);
5     }
6     return this.arr[index];
7 }

```

## Removing Elements

When a user removes an element, we want to take care that all the elements remain contiguous. If the array contains the elements 8, 6, 10 and the user removes 8. we want to ensure that the resulting list is 6, 10 and not null, 6, 10. Likewise, we'll want to make sure to *decrement* the size when we remove an element. We'll also return the value that is being removed. The user is free to ignore it if they don't want it, but this makes the list interface a bit more flexible as it means that the user doesn't have to make two method calls to retrieve-then-delete the elements. This is a common [idiom](#) with many collection data structures.

```

1 public int removeAtIndex(int index) {
2
3     if(index < 0 || index > this.size) {
4         throw new IndexOutOfBoundsException("invalid index: " + index);
5     }
6
7     //start at index and shift elements to the left
8     for(int i=index; i<size-1; i++) {
9         this.arr[i] = this.arr[i+1];
10    }
11    this.size--;
12 }

```

Note that we omitted automatically “shrinking” the underlying array if its unused capacity became too big. We leave that and other variations on the core functionality as an exercise.

## Secondary Functionality

In addition to the core functionality of a list, you could extend our design to provide more convenience methods to make the list implementation even more flexible. For example, you may implement the following methods, the details of which are left as an exercise.

- `public boolean isEmpty()` – this method would return `true` if no elements were stored in the list, `false` otherwise.
- `public int size()` – more generally, a user may want to know how many elements are in the list especially if they wanted to avoid an `IndexOutOfBoundsException`.
- `public void addAtBeginning(int element)` – the companion to the `addAtEnd(int)` method
- `public int replaceElementAt(int element, int index)` – a variation on the remove method that replaces rather than remove the element, returning the replaced



element.

- `public void addAll(int arr[], int index)` – a *batch* method that allows a user to add entire array to the list with one method call. Similarly, you could allow the user to add elements stored in another list instance, `public void addAll(IntegerArrayList list, int index)`. Sometimes this operation is referred to as “splicing.”
- `public void clear()` – another batch method that allows a user to remove all elements at once.

Many other possible variations exist. In addition to the interface, one could vary how the underlying array expands or shrinks to accommodate additions and deletions. In the example above, we increased the size by a constant size of 10. Variations may include expanding the array by a certain percentage or doubling it in size each time. Each strategy has its own advantages and disadvantages.

## A Better Implementation

The preceding list design was still very limited in that it only allowed the user to store integers. If we wanted to design a list to hold floating point numbers, strings, or a user defined type, we would need to create an implementation for every possible type that we wanted a list for. Obviously this is not a good approach. Each implementation would only differ in its name and the type of elements it stored in its array. This is a quintessential example of when to use [parameterized polymorphism](#). Instead of designing a list that holds a particular type of element, we can parameterize it to hold *any* type of element. Thus, only one array-based list implementation is needed.

In the example we designed before for integers, we never actually examined the content of the array inside the class. We never used the fact that the underlying array held integers and the only time we referred to the `int` type was in the method signatures which can all be parameterized to accept and return the same type.

Code Sample 4.1 contains a parameterized version of the array-based list we implemented before.

```

1 public class ArrayList<T> {
2
3     private T[] arr;
4     private int size;
5
6     public ArrayList() {
7         this.arr = (T[]) new Object[10];
8         this.size = 0;
9     }
10
11    public T getElement(int index) {
12
13        if(index < 0 || index >= this.size) {
14            throw new IndexOutOfBoundsException("invalid index: " + index);
15        }
16        return this.arr[index];
17    }
18
19    public void removeAtIndex(int index) {
20
21        if(index < 0 || index >= size) {
22            throw new IndexOutOfBoundsException("invalid index: " + index);
23        }
24
25        for(int i=index; i<size-1; i++) {
26            this.arr[i] = this.arr[i+1];
27        }
28        this.size--;
29    }
30
31    public void insertAtIndex(T x, int index) {
32        if(index < 0 || index > size) {
33            throw new IndexOutOfBoundsException("invalid index: " + index);
34        }
35
36        if(this.size == arr.length) {
37            this.arr = Arrays.copyOf(this.arr, this.arr.length + 10);
38        }
39
40        for(int i=this.size-1; i>=index; i--) {
41            this.arr[i+1] = this.arr[i];
42        }
43        this.arr[index] = element;
44        this.size++;
45    }
46
47    public void addAtEnd(T x) {
48        this.insertAtIndex(x, this.size);
49    }
50 }

```

Code Sample 4.1: Parameterized Array-Based List in Java

A few things to note about this parameterized implementation. First, with the original implementation since we were storing primitive `int` elements, `null` was not an issue. Now that we are using parameterized types, a user would be able to store `null` elements in our list. We could make the design decision to allow this or disallow it (by adding null pointer checks on the add/insert methods).

Another issue, particular to Java, is the instantiation of the array on line 7. We *cannot* invoke the `new` keyword on an array with an indeterminate type. This is because different types require a different number of bytes (integers take 4 bytes, `double`s take 8 bytes). The number of bytes may even vary between Java Virtual Machine (JVM)s (32-bit vs. 64-bit). Without knowing how many bytes each element takes, it would be impossible for the JVM to allocate the right amount of memory. Thus, we are forced to use a *raw type* (the `Object` type) and do an explicit type cast. This is not much of an issue because the parameterizations guarantee that the user would only ever be able to add elements of type `T`.

Another useful feature that we could add would be an `iterator` pattern. An iterator allows you to iterate over each element in a collection and process them. With a traditional array, a simple for loop can be used to iterate over elements. An iterator pattern relieves the user of the need to write such boilerplate code and instead use a *foreach* loop.<sup>1</sup>

In Java, an iterator pattern is achieved by implementing the `Iterable<T>` interface (which is also parameterized). The interface requires the implantation of a public method that returns an `Iterator<T>` which has several methods that need to be implemented. The following is an example that could be included in our `ArrayList` implementation.

```

1 public Iterator<T> iterator() {
2     return new Iterator<T>() {
3         private int currentIndex = 0;
4         @Override
5         public boolean hasNext() {
6             return (this.currentIndex < size);
7         }
8
9         @Override
10        public T next() {
11            this.currentIndex++;
12            return arr[currentIndex-1];
13        }
14
15    };
16 }

```

<sup>1</sup>Note that this is not necessarily mere *syntactic sugar*. As we will see later, some collections are unordered and would *require* the use of such a pattern. Yet still, some list implementations, in particular linked lists, an index-based get method is actually very *inefficient*. An iterator pattern allows you to encapsulate the most efficient logic for iterating over a particular list implementation.

Essentially, the iterator (an anonymous class declaration/definition) has an internal index, `currentIndex` that is initialized to 0 (the first element). Each time `next()` is called, the “current” element is returned, but the method also sets itself up for the next iteration by incrementing `currentIndex`. The main advantage of implementing an iterator is that you can then use a `foreach` loop (which Java calls an “enhanced for-loop”). For example:

```
1 ArrayList<Integer> list = new ArrayList<Integer>();
2 list.addAtEnd(8);
3 list.addAtEnd(6);
4 list.addAtEnd(10);
5
6 //prints "8 6 10"
7 for(Integer x : list) {
8     System.out.print(x + " ");
9 }
```

## 4.2 Linked Lists

The array-based list implementation offers a lot of advantages and improvements over a primitive array. However, it comes at some cost. In particular, when we need to expand the underlying array, we have to create a new array and copy over every last element. If there were 1 million elements in the underlying array and we wanted to add one more, we would essentially be performing 1 million copy operations. In general, if there are  $n$  elements, we would be performing  $n$  copy operations (or a copy operation proportional to  $n$ ). That is, some operations will induce a *linear* amount of work with respect to how many elements are already in the list. Even with different strategies for increasing the size of the underlying array (increasing the size by a percentage or doubling it), we still have some built-in overhead cost associated with the basic list operations.

For many applications this cost is well-worth it. A copy operation can be performed quite efficiently and the advantages that a list data structure provide to development may outweigh the efficiency issues (and in any case, the same issues would be present even if we used a primitive array). In some applications, however, an alternative implementation may be more desirable. In particular, a *linked list* implementation avoids the expansion/shrinking of an underlying array as it uses a series of linked *nodes* to store elements rather than a primitive array.

A simple linked list is depicted in Figure 4.1. Each element is stored in a node. Each node contains both an element (in this example, an integer) and a reference to the next node in the list. In this way, each node is *linked* together to form a chain. The start of the list is usually referred to as the *head* of the list. Similarly, the end of the list is usually referred to as the *tail* (in this case, the node containing 10). A special symbol or value is usually used to denote the end of the list so that that tail node does not

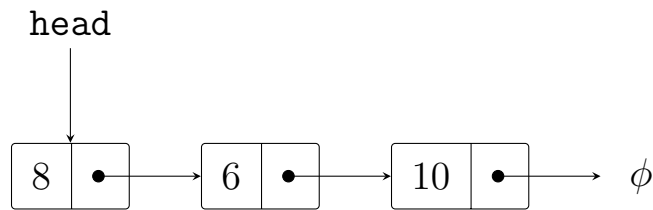


Figure 4.1: A simple linked list containing 3 nodes.

point to another node. In this case, we've used the value  $\phi$  to indicate the end of the list. In practice, a `null` value could be used or a special *sentinel* node can be created to indicate the end of a list.

To understand how a linked list works, we'll design several algorithms to support the core functionality of a list data structure. In order to refer to nodes and their elements, we'll use the following notation. Suppose that  $u$  is a node, then its value will be referred to as  $u.value$  and the next node it is linked to will be referred to as  $u.next$ .

### Adding Elements

With a linked list, there is no underlying array of a fixed size. To add an element, we simply need to create a new node and link it somewhere in the chain. Since there is no fixed array space to fill up, we no longer have to worry about expanding/shrinking it to accommodate new elements.

To start, an empty list will be represented by having the head reference refer to the special end-of-list symbol. That is,  $head \rightarrow \phi$ . To add a new element to an empty list, we simply create a new node and have the head reference point to it. In fact, we can more generally support an *insert at head* operation with the same idea. To insert at the head of the list, we create a new node containing the inserted value, make it point to the "old head" node of the list and then update the head reference to the new node. This operation is depicted in Algorithm 1.

INPUT : A linked list  $L$  with head,  $L.head$  and a new element to insert  $x$  at the head

- 1  $u \leftarrow$  a new node
- 2  $u.value \leftarrow x$
- 3  $u.next \leftarrow L.head$
- 4  $L.head \leftarrow u$

#### Algorithm 1: Insert-At-Head Linked List Operation

Just as with the array-based list, we could also support a more general insert-at-index

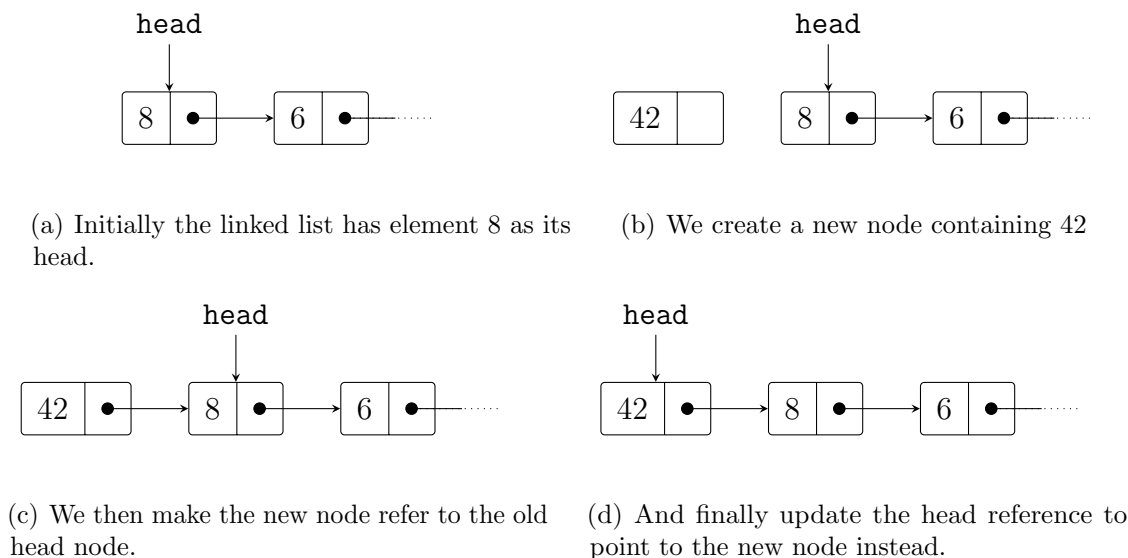


Figure 4.2: Insert-at-head Operation in a Linked List. We wish to insert a new element, 42 at the head of the list.

method that would allow the user to insert a new node at any position in the list. The first step, of course, would be to find the two nodes between which you wanted to insert the new node. We will save the details for this procedure as it represents a more general retrieval method. For now, suppose we have two nodes,  $a, b$  and we wish to insert a new node,  $u$  between them. To do this, we simply need to make  $a$  refer to the new node and to make the new node refer to  $b$ . However, we must be careful with the order in which we do this so as not to lose  $a$ 's reference to  $b$ . The general procedure is depicted in Algorithm 2.

INPUT : Two linked nodes,  $a, b$  in a linked list and a new value,  $x$  to insert between them

- 1  $u \leftarrow$  a new node
- 2  $u.value \leftarrow x$
- 3  $u.next \leftarrow a.next$
- 4  $a.next \leftarrow u$

**Algorithm 2:** Insert Between Two Nodes Operation

One special **corner case** occurs when we want to insert at the end of a list. In this scenario,  $b$  would not exist and so  $a.next$  would end up referring to  $\phi$ . This ends up working out with the algorithm presented. We never actually made any explicit reference to  $b$  in Algorithm 2. When we assigned  $u.next$  to refer to  $a.next$ , we took care of both the possibility that it was an actual node and the possibility that it was  $\phi$ .

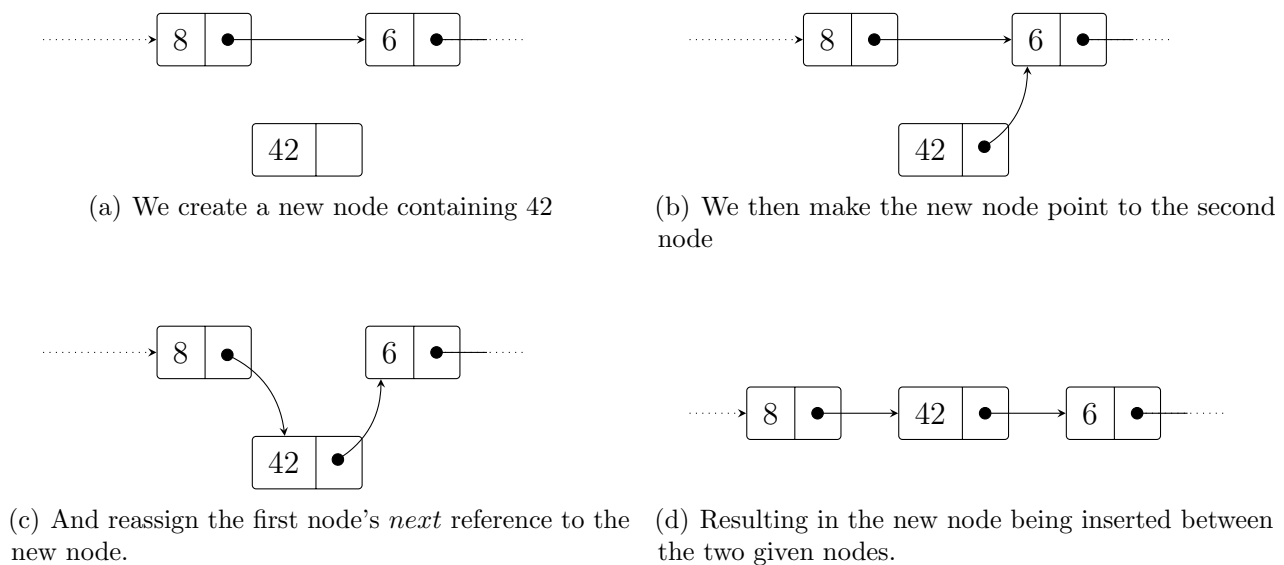


Figure 4.3: Inserting Between Two Nodes in a Linked List. Here, we wish to insert a new element 42 between the given two nodes containing 8 and 6.

Another corner case occurs if we wish to insert at the head of the list using this algorithm. In that scenario,  $a$  would not refer to an actual node while  $b$  would refer to the *head* element. In this case, lines 3–4 would be invalid as  $a.next$  would be an invalid reference. For this corner case, we would need to either fall back to our first insert-at-head (Algorithm 1) operation or we would need to handle it separately.

## Retrieving Elements

As previously noted, a linked list avoids the cost of expanding/shrinking of an underlying array. However, this does come at a cost. With an array-based list, we had “free” [random access](#) to elements. That is, if we wanted the element stored at index  $i$  it is a simple matter to compute a memory offset and “jump” to the proper memory location. With a linked list, we do *not* have the advantages of random access.

Instead, to retrieve an element, we must sequentially search through the list, starting from the head, until we find the element that we are trying to retrieve. Again, many variations exist, but we’ll illustrate the basic functionality by describing the same index-based retrieval method as before.

The key to this algorithm is to simply keep track of the “current” node. Each iteration we traverse to the next node in the chain and iterate a counter. When we have traversed  $i$  times, we stop as that is the node we are looking for.

In contrast to the “free” index-based retrieval method with an array-based list, the

```

INPUT   : A linked list  $L$  with head,  $L.head$  and in index  $i$ ,  $0 \leq i < n$  where  $n$  is
           the number of elements in  $L$ 
OUTPUT  : The  $i$ -th node in  $L$ 
1   $currentNode \leftarrow L.head$ 
2   $currentIndex \leftarrow 0$ 
3  WHILE  $currentIndex < i$  DO
4  |    $currentNode \leftarrow currentNode.next$ 
5  |    $currentIndex \leftarrow (currentIndex + 1)$ 
6  END
7  output  $currentNode$ 

```

**Algorithm 3:** Index-Based Retrieval Operation

operation of finding the  $i$ -th node in a linked list is much more expensive. We have to perform  $i$  operations to get the  $i$ -th node. In the worst case, when we are trying to retrieve the last (tail) element, we would end up performing  $n$  traversal operations. If the linked list held a lot of elements, say 1 million, then this could be quite expensive. In this manner, we see some clear trade-offs in the two implementations.

### Removing Elements

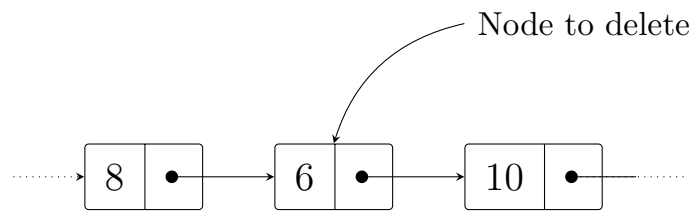
Like the insertion operation, the removal of an element begins by retrieving the node that contains it. Suppose we have found a node,  $u$  and wish to remove it. It actually suffices to simply circumvent the node, by making  $u$ 's predecessor point to  $u$ 's successor. This is illustrated in Figure 4.4.

Since we are changing a reference in  $u$ 's predecessor, however, we must make appropriate changes to the retrieval operation we developed before. In particular, we must keep track of two nodes: the current node as before but also its predecessor, or more generally, a *previous* node. Alternatively, if we were performing an index-based removal operation, we could easily find the predecessor by adjusting our index. If we wanted to delete the  $i$ -th node, we could retrieve the  $(i - 1)$ -th node and delete the next one in the list.

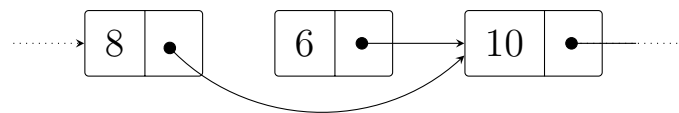
Again, we may need to deal with corner cases such as if we are deleting the head of the list or the tail of the list. As with the insertion, the case in which we delete the tail is essentially the same as the general case. If the successor of a tail node is  $\phi$ , thus when we make  $u$ 's predecessor refer to  $u$ 's successor, we are simply making it refer to  $\phi$ .

The more complicated part is when we delete the head of the list. By definition, there is no predecessor to the head, so handling it as the general case will result in an invalid reference. Instead, if we wanted to delete the head, we must actually change the list's head itself. This is easy enough, we simply need to make  $head$  refer to  $head \rightarrow next$ .





(a) We wish to delete the node containing 6.



(b) We make its predecessor node point to its successor.



(c) Though the node containing 6 still refers to the next node, from the perspective of the list, it has been removed.

Figure 4.4: Delete Operation in a Linked List

An example of a removal operation is presented in Algorithm 4. In contrast to previous examples, this operation is a *key-based* removal operation variation. That is, we search for the first instance of a node whose value matches a given key element and delete it. As another corner case for this variation, we also must deal with the situation in which no node matches the given key. In this case, we've decided to leave it as a no-operation (“noop”). This is achieved in lines 10–12 which will not delete the last node if the key does not match.

```

INPUT   : A linked list  $L$  with head,  $L.head$  and a key element  $k$ 
OUTPUT :  $L$  but with the first node  $u$  such that  $u.value = k$  removed if one exists
1 IF  $L.head.value = k$  THEN
2   |  $L.head. \leftarrow L.head.next$ 
3 ELSE
4   |  $previousNode \leftarrow \phi$ 
5   |  $currentNode \leftarrow L.head$ 
6   | WHILE  $currentNode.value \neq k$  and  $currentNode.next \neq \phi$  DO
7     |  $previousNode \leftarrow currentNode$ 
8     |  $currentNode \leftarrow currentNode.next$ 
9   | END
10  | IF  $currentNode.value = k$  THEN
11  |   |  $previousNode.next \leftarrow currentNode.next$ 
12  |   | END
13 END
14 output  $L$ 

```

**Algorithm 4:** Key-Based Delete Operation

### 4.2.1 Designing a Java Implementation

We now adapt these operations and algorithms to design a Java implementation of a linked list. Note that the standard collections library has both an array-based list (`java.util.ArrayList<E>`) as well as a linked list implementation (`java.util.LinkedList<E>`) that should be used in general.

First, we need a node class in order to hold elements as well as a reference to another node. Code Sample 4.2 gives a basic implementation with several convenience methods. The class is parameterized so that nodes can hold any type.

Given this `Node` class, we can now define the basic state of our linked list implementation:

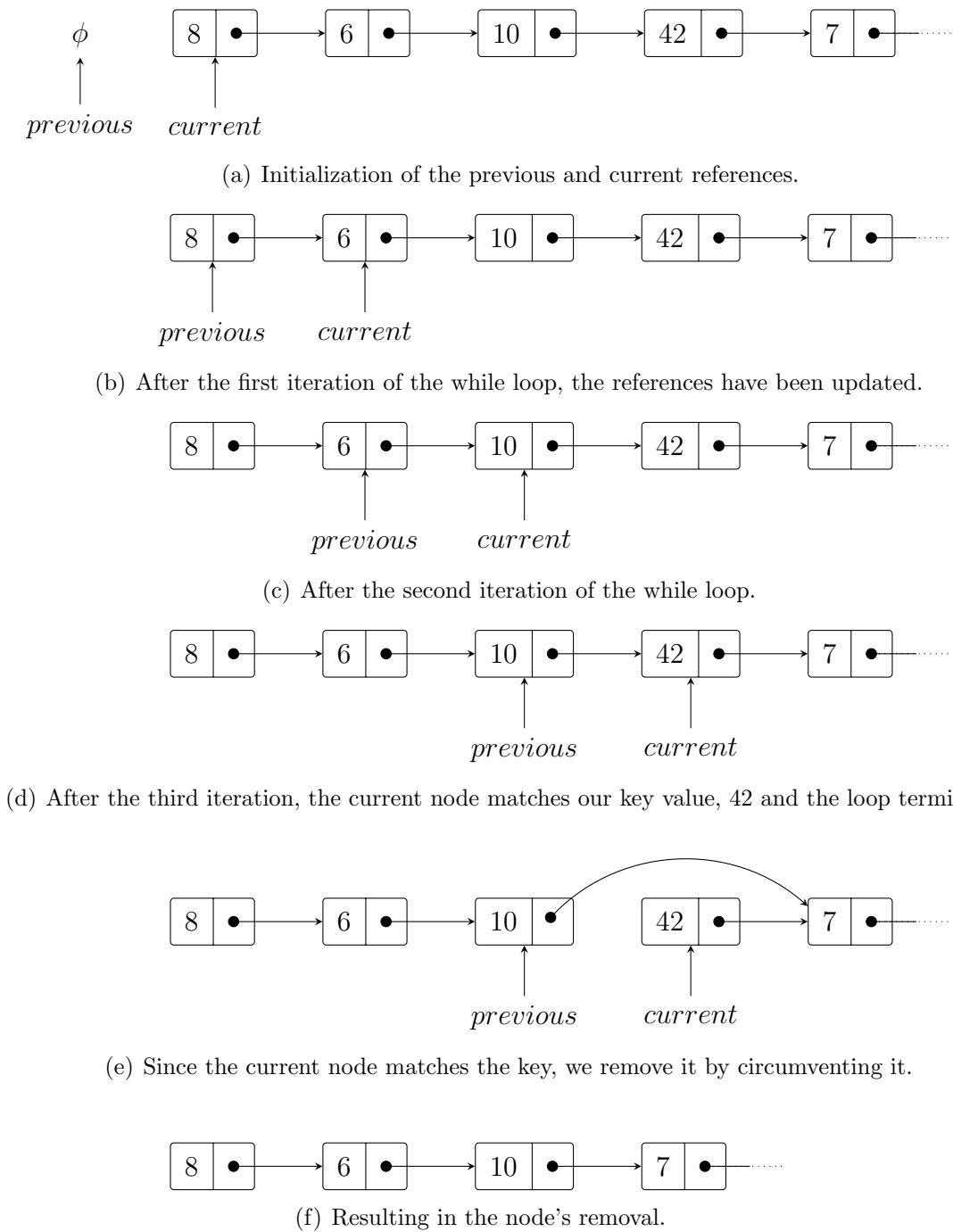


Figure 4.5: Key-Based Find and Remove Operation. We wish to remove the first node we find containing 42.

```

1 public class Node<T> {
2
3     private final T item;
4     private Node<T> next;
5
6     public Node(T item) {
7         this.item = item;
8         next = null;
9     }
10
11     //getters and setters omitted
12
13     public boolean hasNext() {
14         return (this.next == null);
15     }
16
17 }

```

Code Sample 4.2: A linked list node Java implementation. Getter and setter methods have been omitted for readability. A convenience method to determine if a node has a next element is included. This implementation uses `null` as its terminating value.

```

1 public class LinkedList<T> {
2
3     private Node<T> head;
4     private int size;
5
6     public LinkedList() {
7         this.head = null;
8         this.size = 0;
9     }
10
11     ...
12
13 }

```

We keep track of the size of the list and increment/decrement it on each add/delete operation so that we do not have to recompute it. To keep things simple, we will implement a two general purpose node retrieval methods, one index-based and one key based. These can be used as basic steps in other, more specific operations such as insertion, retrieval, and deletion methods. Note that both of these methods are `private` as they are intended for “internal” use by the class and not for external use. If we had made these methods `public` we would be exposing the internal structure of our linked

list to outside code. Such a design would be a typical example of a [leaky abstraction](#) and is, in general, considered bad practice and bad design.

```

1 private Node<T> getNodeAtIndex(int index) {
2
3     if(index < 0 || index >= this.size) {
4         throw new IndexOutOfBoundsException("invalid index: " + index);
5     }
6     Node<T> curr = this.head;
7     for(int i=0; i<index; i++) {
8         curr = curr.getNext();
9     }
10    return curr;
11
12 }
13
14 private void getNodeWithValue(T key) {
15
16     Node<T> curr = head;
17     while(curr != null && !curr.getItem().equals(key)) {
18         curr = curr.getNext();
19     }
20    return curr;
21 }

```

As previously mentioned, these two general purpose methods can be used or adapted to implement other, more specific operations. For example, index-based retrieval and removal methods become very easy to implement using these methods as subroutines.

```

1 public T getElement(int index) {
2     return getNodeAtIndex(index).getItem();
3 }
4
5 public T removeAtIndex(int index) {
6
7     if(index < 0 || index >= this.size) {
8         throw new IndexOutOfBoundsException("invalid index: " + index);
9     }
10    Node<T> previous = getNodeAtIndex(index-1);
11    Node<T> current = previous.getNext();
12    T removedItem = current.getItem();
13    previous.setNext(current.getNext());
14    return removedItem;
15 }

```

Adapting these methods and using them to implement other operations is left as an

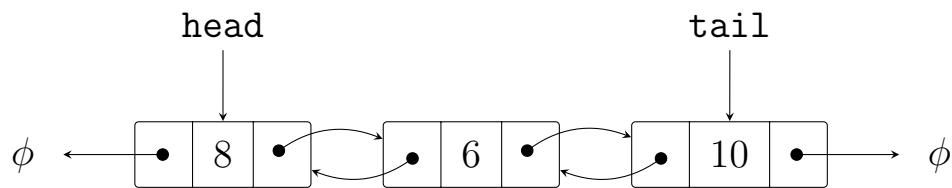


Figure 4.6: A Doubly Linked List Example

exercise.

## 4.2.2 Variations

In addition to the basic linked list data structure design, there are several variations that have different properties that may be useful in various applications. For example, one simple variation would be to not only keep track of the head element, but also a tail element. This would allow you to efficiently add elements to either end of the list without having to traverse the entire list. Other variations are more substantial and we will not look at several of them.

### Doubly Linked Lists

As presented, a linked list is “one-way.” That is, we can only traverse forward in the list from the head toward the tail. This is because our tree nodes only kept track of a *next* element, referring to the node’s predecessor. As an alternative, our list nodes could keep track of both the *next* element as well as a *previous* element so that it has access to both a node’s predecessor as well as its successor. This allows us to traverse the list two ways, forwards and backwards. This variation is referred to as a *doubly linked list*.

An example of a doubly linked list is depicted in Figure 4.6. In this example, we’ve also established a reference to the *tail* element to enable us to start at the end of the list and traverse backwards.

A doubly linked list may provide several advantages over a singly linked list. For example, when we designed our key-based delete operation we had to be sure to keep track of a previous element in order to manipulate its reference to the next node. In a doubly linked list, however, since we can always traverse to the previous node this is not necessary.

A doubly linked list has the potential to simplify a lot of algorithms, however it also means that we have to take greater care when we manipulate the next and previous references. For example, to insert a new node, we would have to modify up to four references rather than just two. Greater care must also be taken with respect to corner cases.

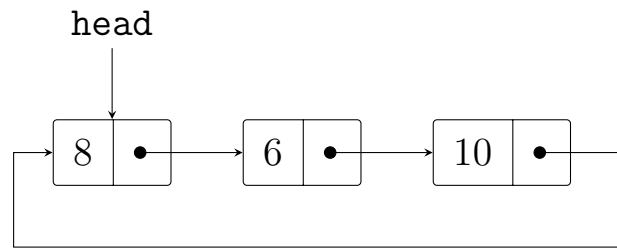


Figure 4.7: A Circularly Linked List Example

### Circularly Linked Lists

Another variation is a *circularly linked list*. Instead of ending the list using a sentinel value,  $\phi$ , the “last” node points back to the first node, closing the list into a large loop. An example is given in Figure 4.7. With such a structure, it is less clear that there is a head or tail to the list. Certain operations such as index-based operations may no longer make sense with such an implementation. However, we could still designate an arbitrary node in the list as the “head” as in the figure.

Alternatively, we could instead simply have a reference to a *current* node. At any point during the data structure’s life cycle the current node may reference any node in the list. The core operations may iterate through the list and end up at a different node each time. Care would have to be taken to ensure that operations such as searching do not result in an infinite loop, however. An unsuccessful key-based search operation would need to know when to terminate (when it has gone through the entire circle once and returned to where it started). It is easy enough to keep track of the size of the list and to ensure that no circular operations exceed this size.

Circularly linked lists are useful for applications where elements must be processed over and over each in turn. For example, an operating system may give time slices to running applications. Instead of a well-defined beginning and end, a continuous poll loop is run. When the “last” application has exhausted its time slot, the operating system returns to the “first” in a continuous loop.

### Unrolled Linked Lists

As presented, each node in a linked list holds a single element. However, we could design a node to hold any number of elements, in particular an array of  $m$  elements. Nodes would still be linked together, but each node would be a mini array-based list as well. An example is presented in Figure 4.8.

This hybrid approach is intended to achieve better performance with respect to memory. The size of each node may be designed to be limited to fit in a particular *cache line*, the amount of data typically transferred from main memory to a processor’s cache. The goal is to improve cache performance while reducing the overhead of a typical linked list. In a

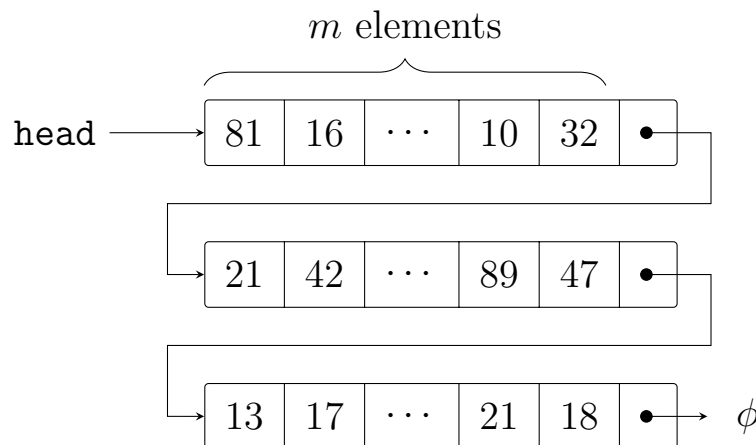


Figure 4.8: An Unrolled Linked List Example

typical linked list, every node has a reference (or two in the case of doubly linked lists) which can double the amount of memory required to store elements. By storing more elements in each node, the overall number of nodes is reduced and thus the number of references is reduced. At the same time, an unrolled linked list does not require large chunks of contiguous storage like an array-based list but provides some of the advantages (such as random access within a node). Overall, the goal of an unrolled linked list is to reduce the number of cache misses required to retrieve a particular element.

### 4.3 Stacks & Queues

Collection data structures such as lists are typically *unstructured*. A list, whether array-based, a linked list, or some other variation simply hold elements in an ordered manner. That is, there is a notion of a first element, second element, etc. However that ordering is not necessarily structured, it is simply the order in which the elements were added to the collection. In contrast, you can *impose* a structured ordering by sorting a list (thus, “sorted” is not the same thing as “ordered”). Sorting a list, however, does not give the data structure itself any more structure. The interface would still allow a user to insert or rearrange the elements so that they are no longer sorted. Sorting a list only change’s the collection’s *state*, not its behavior.

We now turn our attention to a different kind of collection data structure whose structure is defined by its behavior rather than its state. In particular, we will look at two data structures, stack and queues, that represent *restricted access data structures*. These data structures still store elements, but the general core functionality of a list (the ability to add, retrieve, and remove arbitrary elements) will be restricted in a very particular manner so as to give the data structure behavioral properties. This restriction is built into the data structure as part of the object’s interface. That is, users may only interact



with the collection in very specific ways. In this manner, structure is imposed through the collection's behavior.

### 4.3.1 Stacks

A *stack* is a data structure that stores elements in a last-in, first-out (or **Last-In First-Out (LIFO)**) manner. That is, the last element to be inserted into a stack is the first element that will come out of the stack. A stack data structure can be described as a stack of dishes. When dealing with such a stack, we can add a dish to it, but only at the *top* of the stack lest we risk causing the entire stack of dish to fall and break. Likewise, when removing a dish, we remove the top-most dish rather than pulling a dish from the middle or bottom of the stack for the same reason. Thus, the *last* dish that we added to the stack will be the *first* dish that we take off the stack. It may also be helpful to visualize such a stack of dishes as in a cafeteria where a spring-loaded cart holds the stack. When we add a dish, the entire stack moves down into the cart and when we remove one, the next one “pops” up to the top of the stack.

You are probably already familiar with the concept of stacks in the context of a program's *call stack*. As a program invokes functions, a new stack frame is created that contains all the “local” information (parameters, local variables, etc.) and is placed on top of the call stack. When a function is done executing and returns control back to the calling function the stack frame is removed from the top of the call stack, restoring the stack frame right below it. In this manner, information can be saved and restored from function call to function call efficiently. This idea goes all the way back to the very earliest computers and programming languages (specifically, the Information Processing Language in 1956).

#### Core Functionality

In order to have achieve the LIFO behavior, access to a stack's elements are restricted through its interface by only allowing two core operations:

- *push* adds a new element to the stop of the stack and
- *pop* removes the element at the top of the stack

The element removed at the top of the stack may also be “returned” in the context of a method call.

Similar to lists, stacks also have several corner cases that need to be considered. First, when implementing a stack you must consider what will happen when a user performs a pop operation on an empty stack. Obviously there is nothing to remove from the stack, but how might we handle this situation? In general, such an invalid operation is referred to as a stack *underflow*. In programming languages that support error handling via exceptions, we could choose to throw an exception. If we were to make this design decisions then we should, in addition, provide a way for the user to check that such an

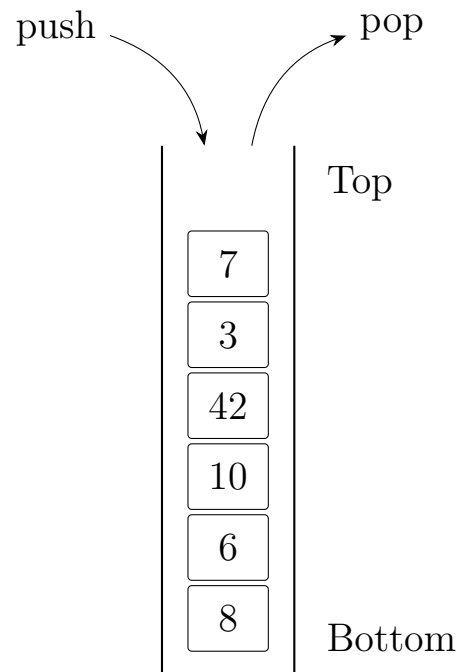


Figure 4.9: A stack holding integer elements. Push and pop operations are depicted as happening at the “top” of the stack. In actuality, a stack stored in a computer’s memory is not really oriented but this visualization is consistent with a physical stack growing “upwards.”

operation may result in an exception by providing a way for them to check if the stack is empty or not (see Secondary Functionality below). Alternatively, we could instead return a “flag” value such as `null` to indicate an empty stack. Though this design decision has consequences as well: we would either need to disallow the user from pushing a `null` value onto the stack (and decide how again to handle that) or we would need to provide a way for the user to distinguish the situation where `null` was actually popped off the stack or was returned because the stack was empty.

In addition, we could design our stack to be either *bounded* or *unbounded*. An unbounded stack means that there would be no practical restrictions on how large the stack could grow. A program’s use of our stack would only be limited by the amount of system memory available. A user could continue to push as many elements onto the stack as they like and a problem would only arise when the program or system itself runs out of memory.

A bounded stack means that we could design our stack to have a fixed capacity or limit of (say)  $n$  elements. If we went with such a design we would again have to deal with the corner case of a user pushing an element to the stack when it is full referred to as a *stack overflow*. Solutions similar to the popping from an empty stack could be used here as well. We could throw an exception (and give the user the ability to check if the stack is full or not) or we could make such an operation a no-op: we would not push the element to the stack, leaving it as it was before and then report the no operation to the user. Typically a boolean value is used, *true* to indicate that the operation was valid and had some side effect on the stack (the element was added) or *false* to indicate that the operation resulted in no side effects.

## Secondary Functionality

To make a stack more versatile it is common to include secondary functionality such as the following.

- A *peek* method that allows a user to access the element at the top of the stack without removing it. The same effect can be achieved with a pop-then-push operation, but it may be more convenient to allow such access directly. This is useful if a particular algorithm needs to make a decision based on what will be popped off the stack next.
- A means to iterate over all the elements in the stack or to allow read-only access arbitrary elements. We would not want to allow arbitrary write access as that would violate the LIFO behavior of a stack and defeat the purpose of using this particular data structure.
- A way to determine how many elements are on the stack and, related whether or not the stack is empty and, if it is bounded, whether or not it is full or its remaining capacity. Such methods would allow the user to use a more defensive-style

programming approach and not make invalid operations.

- A way to empty or “clear” the stack of all its elements with a single method call.
- General find or contains methods that would allow the user to determine if a particular element was already in the stack and more generally, if it is, how far down in the stack it is (its “index”).

### Implementations

A straightforward and efficient implementation for a stack is to simply use a list data structure “under the hood” and to restrict access to it through the stack’s interface. The push and pop operations can then be achieved in terms of the list’s add and remove operations, taking care that both work from the same “end” of the list. That is, if the push operation adds an element at the beginning of the list, then the pop operation must remove (and return) the element at the beginning of the list as well.

A linked list is ideal for bounded and unbounded stacks as adding and removing from the head of the list are both very efficient operations, requiring only the creation of a new node and the shuffling of a couple of references. There is also no expensive copy-and-expand operation over the life of the stack. For unbounded stacks, the capacity can be constrained by simply checking the size of the underlying list and handling the stack overflow appropriately. If the underlying linked list also keeps track of the tail element, adding and removing from the tail would also be an option.

An array-based list may also be a good choice for bounded stacks if the list is initialized to have a capacity equal to the capacity of the stack so as to avoid any copy-and-expand operations. An array-based list is less than ideal for unbounded stacks as expensive copy-and-expand operations may be common. However, care must be taken to ensure that the push and pop operations are efficient. If we designate the “first” element (the element at index 0) as the top of the stack, we would constantly be shifting elements with each and every push and pop operation which can be quite expensive. Instead, it would be more efficient to keep track of the top element and add elements to the “end” of the array.

In detail, we would keep track of a current index, say *top* that is initialized to  $-1$  indicating an empty stack. Then as elements are pushed, we would add them to the list at index ( $top + 1$ ) and increment *top*. As elements are popped, we return the element at index *top* and decrement the index. In this manner, each push and pop operation requires a constant number of operations rather than shifting up to  $n$  elements.

### Applications

As previously mentioned, stacks are used extensively in computer architecture. They are used to keep track of local variables and parameters as functions are called using a program

stack. In-memory stacks may also be used to simulate a series of function/method calls in order to avoid using (or misusing) the program stack. A prime example of such a use case is avoiding recursion. Instead of a sequence of function calls that may result in a stack overflow of the program stack (which is generally limited) an in-memory stack data structure, which is generally able to accommodate many more elements, can be used.

Stacks are also the core data structures used in many fundamental algorithms. Stacks are used extensively in algorithms related to parsing and processing of data such as in the Shunting Yard Algorithm used by compilers to evaluate an [Abstract Syntax Tree \(AST\)](#). Stacks are also used in many graph algorithms such as [Depth First Search \(DFS\)](#) and in a preorder processing of binary trees (see [Chapter 6](#)). Essentially any application in which something needs to be “tracked” or remembered in a particular LIFO ordering, a stack is an ideal data structure to use.

### 4.3.2 Queues

A similar data structure is a *queue* which provides a [First-In First-Out \(FIFO\)](#) ordering of elements. As its name suggests, a queue can be thought of as a line. Elements enter the queue at one end (the “end” of the line) and are removed from the other end (the “front” of the line). Thus, the first element to enter a queue is the first element to be removed from the queue. Each element in the “line” is served in the order in which they entered the line. Similar to a stack, a queue’s structure is defined by its interface.

#### Core Functionality

The two core operations in a queue are:

- *enqueue* which adds an element to the queue at its end and
- *dequeue* which removes the element at the front of the queue

An example of a queue and its operations is depicted in [Figure 4.10](#). Some programming languages and data structure implementations may use different terminology to describe these two operations. Some use the same terms as a stack (“push” would correspond to enqueue and “pop” would correspond to a dequeue operation). However, using the same terminology for fundamentally different data structures is confusing.<sup>2</sup> Some implementations use the terms “offer” (we are offering an element to the queue if it is able to handle it) and “poll” (we are asking or “polling” the queue to see if it has any elements to remove).

---

<sup>2</sup> `<opinion>and wrong</opinion>`

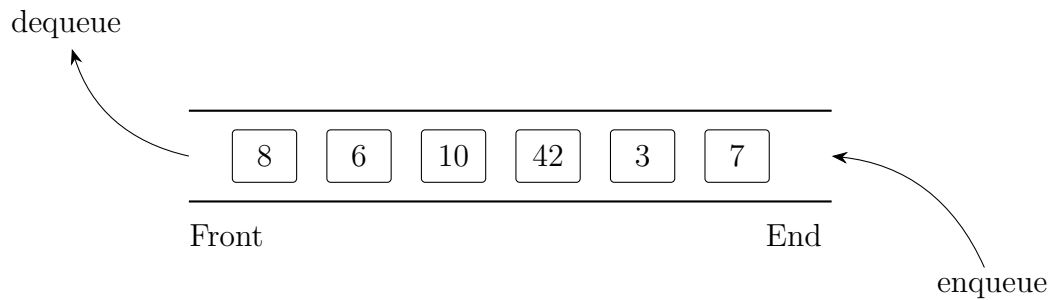


Figure 4.10: An example of a queue. Elements are enqueued at the end of the queue and dequeued from the front of the queue.

### Secondary Functionality

Secondary functionality is pretty much the same as with a stack. We could choose to make our queue bounded or unbounded and deal with corner cases (enqueueing to a full queue, dequeuing from an empty queue) similarly. In fact, we would probably want to design our queue with the same behavior for consistency. We could include methods to determine the size of the queue, its remaining capacity (if bounded), a way to determine if (and where) a certain element may be in the queue, etc.

As with stacks, in general we would not want to allow a user to arbitrarily insert elements into a queue. Doing so would be allowing “line jumpers” to jump ahead of other elements, violating **FIFO**. In some applications this does make sense and we discuss these variations in Section 4.3.3. However, we *could* allow arbitrary removal of certain elements. Strictly speaking, this would not violate **FIFO** as the remaining elements would still be processed in the order in which they were enqueued. This could model situations where those waiting in line got impatient and left (a process or request timed-out for example).

### Implementations

The obvious implementation for a queue is a linked list. Since we have to work from both ends, however, our linked list implementation will need to keep track of both the head element and the tail element so that adding a new node at either end has the same constant cost (creating a new node, shuffling some references). If our linked list only keeps track of the head, then to enqueue an element, we would need to traverse the entire list all the way to the end in order to add an element to the tail.

A linked list that offers constant-time add and remove methods to both the head and the tail can be oriented either way. We could add to the head and remove from the tail or we could add to the tail and remove from the head. As long as we are consistently adding to one end and removing from the other, there really is no difference. Our design decision may have consequences, however, on the secondary functionality. If we design

the array with an iterator, it makes the most sense to start at the front of the queue and iterate toward the end. If our linked list is a doubly linked list, then we can easily iterate in either direction. However, if our list is a singly link linked then we need to make sure that the front of our queue corresponds to the head of the list.

An array-based list can also be used to implement a queue. As with stacks, it is most appropriate to use array-based lists when you want a bounded queue so as to avoid expensive expand-and-copy operations. However, you need to be a bit clever in how you enqueue and dequeue items to the array in order to ensure efficient operations.

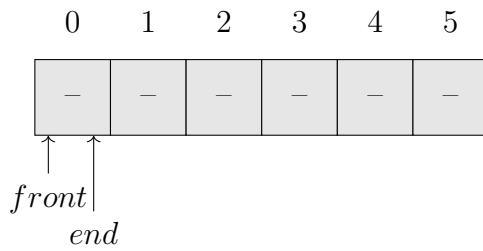
A naive approach would be to enqueue elements at one end of the array and dequeue them from the front (index 0). However, this would mean we need to shift elements down on every dequeue operation which is potentially very inefficient. A clever workaround would be to keep track of the *front* and *end* of the queue using two index variables. Initially, an empty queue would have both of these variables initialized to 0. As we add elements, the *end* index gets incremented (we add left-to-right). As elements are removed, the *front* index variable gets incremented. At some point, these variables will read the right-end of the array at which point we simply reset them back to zero. The queue is empty when  $front = end$  and it is full when  $(end - front) \bmod n = n - 1$  where  $n$  is the size of the array. The various states of such a queue are depicted in Figure 4.11. Using an array-based list may save a substantial amount of memory. However, the added complexity of implementation (and thus increased opportunities for bugs and errors) may not justify the savings.

## Applications

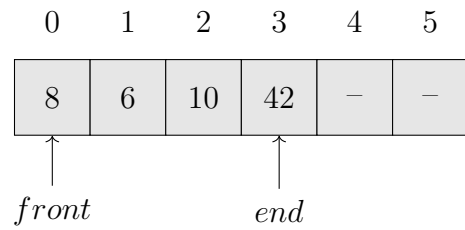
A queue is ideal for any application in which elements must be stored in order to be handled or processed in a particular order. For example, a queue is a natural data structure to implement buffers in which data is received but cannot be processed immediately (it may not be possible or it would be inefficient to process the data in small chunks). A queue ensures that the data remains in the order in which it was received.

Another typical example is when requests are received and must be processed. This is typical in a webserver for example where requests for resources or webpages are stored in a queue and processed in the order in which they are received. In an operating system, threads (or processes) may be stored in a job queue and then woken/executed.

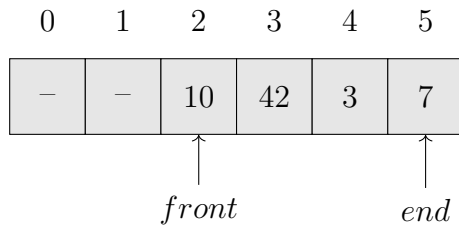
More generally, queues can facilitate communication in a Producer Consumer Pattern (see Figure 4.12). In this scenario we have independent producers and consumers. Producers may produce requests, tasks, works, or a resource that needs to be processed. For example, a producer may be a web browser requesting a particular web page, or it may be a thread making a request for a resource, etc. Consumers handle or service each of these requests. Each consumer and each producer acts independently and asynchronously. To facilitate thread-safe communication between the two groups, each request is enqueued to a *blocking queue*. As producers enqueue requests, they are stored in the order they



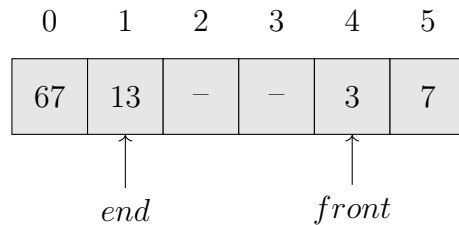
(a) An initially empty queue. Both index variables refer to index 0.



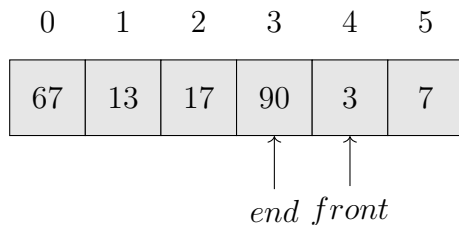
(b) As elements are enqueued the *end* index variable is incremented.



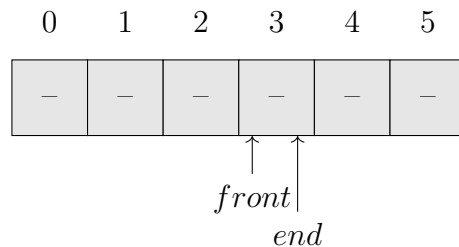
(c) More elements may be enqueued as well as dequeued, moving both index variables.



(d) At some point the index variables wrap around to the beginning of the array and the queue may “straddle” the array.



(e) The queue may also become full, at which point the two index variables are beside each other.



(f) The queue may again become empty and the two index variables refer to the same value, but not necessarily index 0.

Figure 4.11: Various states of an array-based queue. As the index variables, *front* and *end* become  $n$  or larger, they wrap around to the beginning of the array using a modulus operation.



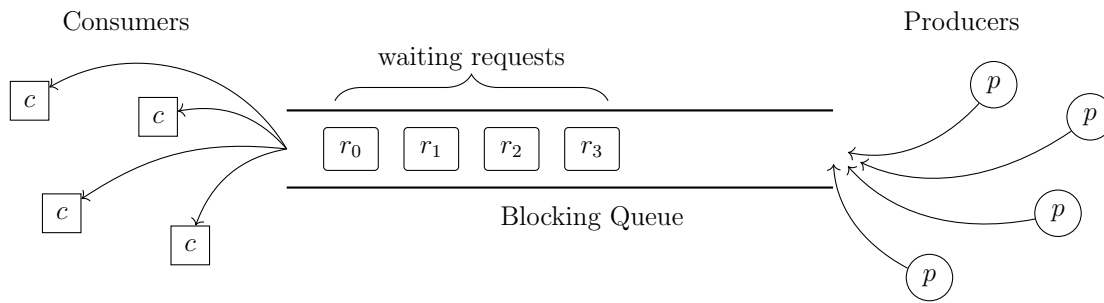


Figure 4.12: Producer Consumer Pattern. Requests (or tasks or resources) are enqueued into a blocking (or otherwise thread safe queue) by producers. Independently (asynchronously) consumers handle requests. Once done, consumers *poll* the queue for another request. If the queue is empty it *blocks* consumers until a new request becomes available.

are received. Then, as each consumer becomes available to service a request, it *polls* the queue for the next request.

A naive approach would be to do *busy polling* where a consumer keeps polling the queue over and over for another request even if it is empty. Using a thread-safe blocking queue means that we don't do busy polling. Instead, if no request is available for a consumer, the queue *blocks* the consumer (typically, the consumer is a thread and this puts the thread to sleep) until a request becomes available so that it is not continuously polling for more work.

### 4.3.3 Variations

In addition to the fundamental stack and queue data structures, there are numerous useful variations. Some simple variations involve allowing the user to perform more than just the two core operations on each.

For example, we could design a double-ended stack which, in addition to the push and pop operations at the top, we could allow a pop operation at the bottom. This provides a variation on the usual bounded stack where the “oldest” elements at the bottom drop out when the stack becomes full instead of rejecting the “newest” items at the top of the queue. A prime example of this is how a typical “undo” operation is facilitated in many applications. Take for example a word processor in which the user performs many actions in sequence (typing, highlighting, copy-paste, delete, etc.). A word processor typically allows the user to undo previous operations but in the reverse sequence that they were performed. However, at the same time we don't want to keep track of *every* change as it would start to take more and more memory and may impact performance. Using a double-ended stack means that the oldest action performed is removed at the bottom, allowing the last  $n$  operations to be tracked.

Relatedly, we could design a queue in which we are allowed to enqueue elements at *either* end, but only remove from one. This would allow the aforementioned “line jumpers” to jump ahead to the front of the line. This establishes a sort-of “fast lane” for high priority elements. As an example, consider system processes in an operating system waiting for their time slice to execute. User processes may be preempted by system-level processes as they have a high priority. More generally, we can establish more than two levels of priority (high/low) and build what is known as a *priority queue* (see below).

### Dequeues

A logical extension is to allow insert and remove operations at both ends. Such a generalized data structure is known as a *deque* (pronounced “deck”, also called a double-ended queue). In fact, the primary advantage to creating this generalized data structure is that it can be used as a stack, queue, etc. all with only a single implementation.

In fact in Java (as of version 6), it is recommended to use its `java.util.Deque<E>` interface rather than the its stack and queue implementations. The `Deque<E>` interface is extremely versatile, offering several different versions of each of the insert and remove methods to both ends (head and tail). One set of operations will throw exceptions for corner cases (inserting into a full deque and removing from an empty deque) and another set will return special values (`null` or `false`). Java also provides several different implementations including an array-based deque (`ArrayDeque<E>`) and its usual `LinkedList<E>` implementation as well as several concurrent and thread-safe versions.

### Priority Queues

Another useful variation is a priority queue in which elements are stored not in a FIFO manner, but with respect to some priority. When elements are dequeued, the highest priority element is removed from the queue first. When elements are enqueued, conceptually they are placed in the queue according to their priority. This may mean that the new element jumps ahead to the front of the queue (if it has a priority higher than all other elements) or it may mean it ends up at the end of the queue (if it has the lowest priority) or somewhere in between. In general, any scheme can be used to define priority but it is typical to use integer values.

A naive implementation of a priority queue would implement the enqueue operation by making comparisons and inserting the new element at the appropriate spot, requiring up to  $n$  comparisons/operations in a queue with  $n$  elements. There are much better implementations that we’ll look at later on (in particular, a heap implementation, see Chapter 6). Again, most programming languages will have a built-in implementation. Java provides an efficient heap-based `PriorityQueue<E>` implementation. Priority is defined using either a natural ordering or a custom `Comparator` object.

# 5 Algorithm Analysis

## 5.1 Introduction

An *algorithm* is a procedure or description of a procedure for solving a problem. An algorithm is a step-by-step specification of operations to be performed in order to compute an output, process data, or perform a function. An algorithm must always be *correct* (it must always produce a valid output) and it must be *finite* (it must terminate after a finite number of steps).

Algorithms are not code. Programs and code in a particular language are *implementations* of algorithms. The word, “algorithm” itself is derived from the latinization of Abū ‘Abdalāh Muhammad ibn Mūsā al-Khwārizmī, a Persian mathematician (c. 780 – 850). The concept of algorithms predates modern computers by several thousands of years. Euclid’s algorithm for computing the greatest common denominator (see Section 5.5.3) is 2,300 years old.

Often, to be useful an algorithm must also be *feasible*: given its input, it must execute in a reasonable amount of time using a reasonable amount of resources. Depending on the application requirements our tolerance may be on the order of a few milliseconds to several days. An algorithm that takes years or centuries to execute is certainly not considered feasible.

**Deterministic** An algorithm is deterministic if, when given a particular input, will always go through the exact same computational process and produce the same output. Most of the algorithms you’ve used up to this point are deterministic.

**Randomized** An algorithm is randomized is an algorithm that involves some form of random input. The random source can be used to make decisions such as random selections or to generate random state in a program as candidate solutions. There are many types of randomized algorithms including Monte-Carlo algorithms (that may have some error with low probability), Las Vegas algorithms (whose results are always correct, but may fail with a certain probability to produce any results), etc.

**Optimization** Many algorithms seek not only to find a solution to a problem, but to find the *best*, optimal solution. Many of these type of algorithms are *heuristics*: rather than finding the actual best solution (which may be infeasible), they can approximate a solution (*Approximation* algorithms). Other algorithms simulate

biological processes (Genetic algorithms, Ant Colony algorithms, etc.) to search for an optimal solution.

**Parallel** Most modern processors are multicore, meaning that they have more than one processor on a chip. Many servers have dozens of processors that work together. Multiple processors can be utilized by designing parallel algorithms that can split work across multiple processes or threads which can be executed in parallel to each other, improving overall performance.

**Distributed** Computation can also be distributed among completely separate devices that may be located half way across the globe. Massive distributed computation networks have been built for research such as simulating protein folding ([Folding@Home](#)).

An algorithm is a more abstract, generalization of what you might be used to in a typical programming language. In an actual program, you may have functions/methods, subroutines or procedures, etc. Each one of these pieces of code could be considered an algorithm in and of itself. The combination of these smaller pieces create more complex algorithms, etc. A program is essentially a concrete implementation of a more general, theoretical algorithm.

When a program executes, it expends some amount of resources. For example:

**Time** The most obvious resource an algorithm takes is time: how long the algorithm takes to finish its computation (measured in seconds, minutes, etc.). Alternatively, time can be measured in how many CPU cycles or floating-point operations a particular piece of hardware takes to execute the algorithm.

**Memory** The second major resource in a computer is memory. An algorithm requires memory to store the input, output, and possibly extra memory during its execution. How much memory an algorithm uses in its execution may be even more of an important consideration than time in certain environments or systems where memory is extremely limited such as embedded systems.

**Power** The amount of power a device consumes is an important consideration when you have limited capacity such as a battery in a mobile device. From a consumer's perspective, a slower phone that offered twice the batter life may be preferable. In certain applications such as wireless sensor networks or autonomous systems power may be more of a concern than either time or memory.

**Bandwidth** In computer networks, efficiency is measured by how much data you can transmit from one computer to another, called *throughput*. Throughput is generally limited by a network's bandwidth: how much a network connection can transmit under ideal circumstances (no data loss, no retransmission, etc.)

**Circuitry** When designing hardware, resources are typically measured in the number of gates or wires are required to implement the hardware. Fewer gates and wires means you can fit more chips on a silicon die which results in cheaper hardware. Fewer wires and gates also means faster processing.

**Idleness** Even when a computer isn't computing anything, it can still be "costing" you something. Consider purchasing hardware that runs a web server for a small user base. There is a substantial investment in the hardware which requires maintenance and eventually must be replaced. However, since the user base is small, most of the time it sits idle, consuming power. A better solution may be to use the same hardware to serve multiple virtual machines (VMs). Now several small web servers can be served with the same hardware, increasing our utilization of the hardware. In scenarios like this, the lack of work being performed is the resource.

**Load** Somewhat the opposite of idleness, sometimes an application or service may have occasional periods of high demand. The ability of a system to service such high *loads* may be considered a resource, even if the capacity to handle them goes unused most of the time.

These are all very important engineering and business considerations when designing systems, code, and algorithms. However, we'll want to consider the complexity of algorithms in a more abstract manner.

Suppose we have two different programs (or algorithms) *A* and *B*. Both of those algorithms are correct, but *A* uses fewer of the above resources than *B*. Clearly, algorithm *A* is the better, more *efficient* solution. However, how can we better quantify this efficiency?

## List Operations

To give a concrete example, consider the list data structures from Chapter ???. The list could be implemented as an array-based list (where the class owns a static array that is resized/copied when full) or a linked list (with nodes containing elements and linking to the next node in the list). Some operations are "cheap" on one type of list while other operations may be more "expensive."

Consider the problem of inserting a new element into the list at the beginning (at index 0). For a linked list this involves creating a new node and shuffling a couple of references. The number of operations in this case is not contingent on the *size* of the list. In contrast, for an array-based list, if the list contains  $n$  elements, each element will need to be *shifted* over one position in the array in order to make room for the element to be inserted. The number of shifts is proportional to the number of elements in the array,  $n$ . Clearly for this operation, a linked list is better (more efficient).

Now consider a different operation: given an index  $i$ , retrieve the  $i$ -th element in the list. For an array-based list we have the advantage of having **random access** to the array. When we index an element, `arr[i]`, it only takes one memory address computation to "jump" to the memory location containing the  $i$ -th element. In contrast, a linked list would require us to start at the head, and traverse the list until we reach the  $i$ -th node. This requires  $i$  traversal operations. In the worst case, retrieving the last element, the  $n$ -th element, would require  $n$  such operations. A summary of these operations can be

List Type	Insert at start	Index-based Retrieve
Array-based List	$n$	1
Linked List	2	$i \approx n$

Table 5.1: Summary of the Complexity of List Operations

found in Table 5.1.

We will now demonstrate the consequences of this difference in performance for an index-based retrieval operation. Consider the Java code in Code Sample 5.1. This method is using a naive index-based retrieval (line 5) to access each element.

```

1 public static int sum(List<Integer> items) {
2
3     int total = 0;
4     for(int i=0; i<items.size(); i++) {
5         total += items.get(i);
6     }
7     return total;
8 }
```

Code Sample 5.1: Summing a collection of integers

Suppose that the `List` passed to this method is an `ArrayList` in which each `get(i)` method call take a constant number of operations and thus roughly a constant amount of time. Since we perform this operation once for each element in the list, we can assume that the amount of time that the entire algorithm will take will be proportional to  $n$ , the number of elements in the list. That is, the time  $t$  that the algorithm will take to execute will be some linear function of  $n$ , say

$$t_{lin}(n) = an + b$$

Here, the constants  $a$  and  $b$  are placeholders for some values that are dependent on the particular machine that we run it on. These two values may vary depending on the speed of the machine, the available memory, etc. and will necessarily change if we run the algorithm on a different machine with different specifications. However, the overall complexity of the algorithm, the fact that it takes a linear amount of work to execute, will not change.

Now let's contrast the scenario where a `LinkedList` is passed to this method instead. On each iteration, each `get(i)` method will take  $i$  node traversals to retrieve the  $i$ -th element. On the first iteration, only a single node traversal will be required. On the second, 2 traversals are required, etc. Adding all of these operations together gives us the following

$$1 + 2 + 3 + 4 + \dots + (n - 1) + n$$

That is, the sum of natural numbers 1 up to  $n$ . We can use the well-known Gauss's Formula to get a closed form of this summation.

**Theorem 1** (Gauss's Formula).

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

Thus, the total number of operations will be proportional to some *quadratic* function,

$$t_{quad}(n) + an^2 + bn + c$$

Again, we don't know what the constants,  $a, b, c$  are as they may vary depending on the languages, system and other factors. However, we can compute them experimentally.

The code in Code Sample was run on a laptop and its performance was timed. The experiment generated random lists of various sizes of lists, starting at  $n = 50,000$  up to 1 million by increments of 50,000. The time to execute was recorded in seconds. The experiment was repeated for both types of lists multiple times and an average was taken to reduce the influence of other factors (such as other processes running on the machine at the same time).

To find the coefficients in both functions, a linear and quadratic regression was performed, giving us the following functions

$$\begin{aligned} t_{lin}(n) &= 5.138e-5n + 0.004 \\ t_{quad}(n) &= 6.410e-4n^2 - 0.112n + 9.782 \end{aligned}$$

Both had very high correlation coefficients (the quadratic regression was nearly perfect at 0.994) which gives strong empirical evidence for our theoretical analysis. The quadratic regression along with the experimental data points is given in Figure 5.1. The linear data is not graphed as it would remain essentially flat on this scale.

These experimental results allow us to estimate and project how this code will perform on larger and larger lists. For example, we can plug  $n = 10$  million into  $t_{lin}(n)$  and find that it will still take less than 1 second for the method to execute. Likewise, we can find out how long it would take for the linked list scenario. For a list of the same size,  $n = 10$  million, the algorithm would take 17.498 *hours*! More values are depicted in Table 5.2.

The contrast between the two algorithms gets extreme even with moderately large lists of integers. With 10 billion numbers, the quadratic approach is not even feasible. Arguably, even at 10 million, a run time of over 17 hours is not acceptable, especially when an alternative algorithm can perform the same calculation in less than 1 second.

Being able to identify the complexity of algorithms and the data structures they use in order to avoid such inefficient solutions is an essential skill in Computers Science.<sup>1</sup> In the following examples, we'll begin to be a little bit more formal about this type of analysis.

---

<sup>1</sup>Using an index-based iteration is a common mistake in Java code. The proper solution would have been to use an enhanced for-loop or iterator pattern so that the method would perform the same on either array-based or linked lists.

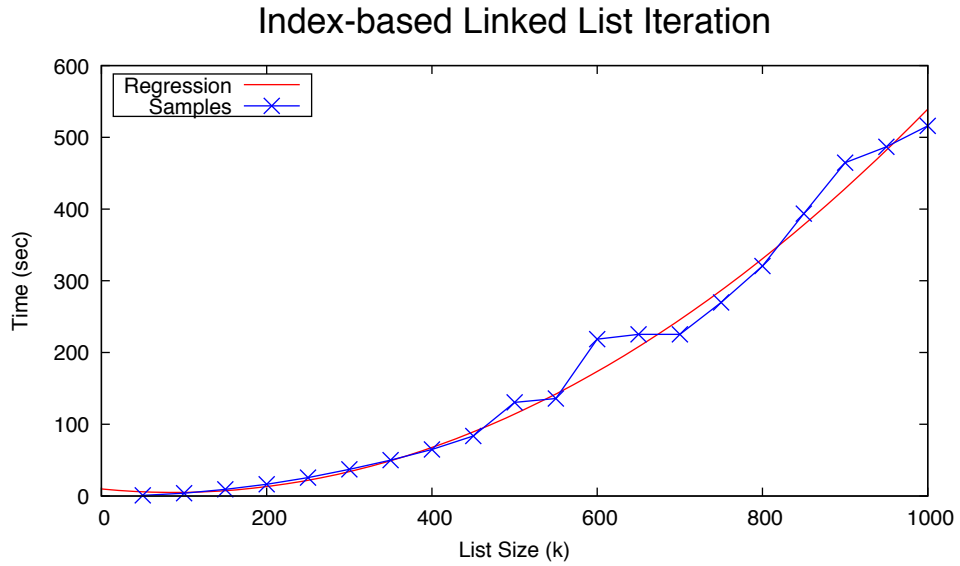


Figure 5.1: Quadratic Regression of Index-Based Linked List Performance.

Table 5.2: Various Projected Runtimes for Code Sample 5.1.

List Size	Execution Time	
	Linear	Quadratic
10 million	1 second	17.498 hours
100 million	10 seconds	74.06 days
1 billion	2 minutes	20.32 years
10 billion	8.56 minutes	2031.20 years



### 5.1.1 Example: Computing a Sum

The following is a toy example, but its easy to understand and straightforward. Consider the following problem: given an integer  $n \geq 0$ , we want to compute the arithmetic series,

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

As a naive approach, consider the algorithm in Code Sample 5.2. In this algorithm, we iterate over each possible number  $i$  in the series. For each number  $i$ , we count 1 through  $i$  and add one to a result variable.

```

1 int result = 0;
2 for(int i=1; i<=n; i++) {
3     for(int j=1; j<=i; j++) {
4         result = result + 1;
5     }
6 }
```

Code Sample 5.2: Summation Algorithm 1

As an improvement, consider the algorithm in Code Sample 5.3. Instead of just adding one on each iteration of the inner loop, we omit the loop entirely and simply just add the index variable  $i$  to the result.

```

1 int result = 0;
2 for(int i=1; i<=n; i++) {
3     result = result + i;
4 }
```

Code Sample 5.3: Summation Algorithm 2

Can we do even better? Yes. Recall Theorem 1 that gives us a direct closed form solution for this summation:

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

Code Sample 5.4 uses this formula to directly compute the sum without any loops.

```

1 int result = n * (n + 1) / 2;
```

Code Sample 5.4: Summation Algorithm 3

All three of these algorithms were run on a laptop computer for various values of  $n$  from 10 up to 1,000,000. Table 5.3 contains the resulting run times (in milliseconds) for each of these three algorithms on the various input sizes.

## 5 Algorithm Analysis

Algorithm	Number of Additions	Input Size					
		10	100	1,000	10,000	100,000	1,000,000
1	$\approx n^2$	0.003ms	0.088ms	1.562ms	2.097ms	102.846ms	9466.489ms
2	$n$	0.002ms	0.003ms	0.020ms	0.213ms	0.872ms	1.120ms
3	1	0.002ms	0.001ms	0.001ms	0.001ms	0.001ms	0.000ms

Table 5.3: Empirical Performance of the Three Summation Algorithms

With small input sizes, there is almost no difference between the three algorithms. However, that would be a naive way of analyzing them. We are more interested in how each algorithm performs as the input size,  $n$  increases. In this case, as  $n$  gets larger, the differences become very stark. The first algorithm has two nested for loops. On average, the inner loop will run about  $\frac{n}{2}$  times while the outer loop runs  $n$  times. Since the loops are nested, the inner loop executes about  $\frac{n}{2}$  times *for each* iteration of the outer loop. Thus, the total number of iterations, and consequently the total number of additions is about

$$n \times \frac{n}{2} \approx n^2$$

The second algorithm saves the inner for loop and thus only makes  $n$  additions. The final algorithm only performs a constant number of operations.

Observe how the running time grows as the input size grows. For Algorithm 1, increasing  $n$  from 100,000 to 1,000,000 (10 times as large) results in a running time that is about 100 times as slow. This is because it is performing  $n^2$  operations. To see this, consider the following. Let  $t(n)$  be the time that Algorithm 1 takes for an input size of  $n$ . From before we know that

$$t(n) \approx n^2$$

Observe what happens when we increase the input size from  $n$  to  $10n$ :

$$t(10n) \approx (10n)^2 = 100n^2$$

which is 100 times as large as  $t(n)$ . The running time of Algorithm 1 will grow quadratically with respect to the input size  $n$ .

Similarly, Algorithm 2 grows linearly,

$$t(n) \approx n$$

Thus, a 10 fold increase in the input,

$$t(10n) \approx 10n$$

leads to a 10 fold increase in the running time. Algorithm 3's runtime does not depend on the input size, and so its runtime does not grow as the input size grows. It essentially remains flat—constant.

Of course, the numbers in Table 5.3 don't follow this trend exactly, but they are pretty close. The actual experiment involves a lot more variables than just the algorithms: the laptop may have been performing other operations, the compiler and language may have optimizations that change the algorithms, etc. Empirical results only provide general evidence as to the runtime of an algorithm. If we moved the code to a different, faster machine or used a different language, etc. we would get different numbers. However, the general trends in the rate of growth *would* hold. Those rates of growth will be what we want to analyze.

### 5.1.2 Example: Computing a Mode

As another example, consider the problem of computing the *mode* of a collection of numbers. The mode is the most common element in a set of data.<sup>2</sup>

```

1 public static int mode01(int arr[]) {
2
3     int maxCount = 0;
4     int modeIndex = 0;
5     for(int i=0; i<arr.length; i++) {
6         int count = 0;
7         int candidate = arr[i];
8         for(int j=0; j<arr.length; j++) {
9             if(arr[j] == candidate) {
10                count++;
11            }
12        }
13        if(count > maxCount) {
14            modeIndex = i;
15            maxCount = count;
16        }
17    }
18    return arr[modeIndex];
19 }

```

Code Sample 5.5: Mode Finding Algorithm 1

Consider the strategy as illustrated in Code Sample 5.5. For each element in the array, we iterate through all the other elements and count how many times it appears (its *multiplicity*). If we find a number that appears more times than the candidate mode we've found so far, we update our variables and continue. As with the previous algorithm,

<sup>2</sup>In general there may be more than one mode, for example in the set {10, 20, 10, 20, 50}, 10 and 20 are both modes. The problem will simply focus on finding *a* mode, not all modes.

the nested nature of our loops leads to an algorithm that performs about  $n^2$  operations (in this case, the comparison on line 9).

Now consider the following variation in Code Sample 5.6. In this algorithm, the first thing we do is *sort* the array. This means that all equal elements will be contiguous. We can exploit this to do less work. Rather than going through the list a second time for each possible mode, we can count up contiguous runs of the same element. This means that we need only examine each element exactly once, giving us  $n$  comparison operations (line 8).

```

1  public static int mode02(int arr[]) {
2      Arrays.sort(arr);
3      int i=0;
4      int modeIndex = 0;
5      int maxCount = 0;
6      while(i < arr.length-1) {
7          int count=0;
8          while(i < arr.length-1 && arr[i] == arr[i+1]) {
9              count++;
10             i++;
11         }
12         if(count > maxCount) {
13             modeIndex = i;
14             maxCount = count;
15         }
16         i++;
17     }
18     return arr [modeIndex];
19 }

```

Code Sample 5.6: Mode Finding Algorithm 2

We can't, however, ignore the fact that to exploit the ordering, we needed to first “invest” some work upfront by sorting the array. Using a typical sorting algorithm, we would expect that it would take about  $n \log(n)$  comparisons. Since the sorting phase and mode finding phase were separate, the total number of comparisons is about

$$n \log(n) + n$$

The highest order term here is the  $n \log(n)$  term for sorting. However, this is still lower than the  $n^2$  algorithm. In this case, the investment to sort the array pays off! To compare with our previous analysis, what happens when we increase the input size 10 fold? For simplicity, let's only consider the highest order term:

$$t(n) = n \log(n)$$

Then

$$t(10n) = 10n \log(10n) = 10n \log(n) + 10n \log(10)$$

Ignoring the lower order term, the increase in running time is essentially linear! We cannot discount the additive term in general, but it is so close to linear that terms like  $n \log(n)$  are sometimes referred to as *quasilinear*.

```

1  public static int mode03(int arr[]) {
2      Map<Integer, Integer> counts = new HashMap<Integer, Integer>();
3      for(int i=0; i<arr.length; i++) {
4          Integer count = counts.get(arr[i]);
5          if(count == null) {
6              count = 0;
7          }
8          count++;
9          counts.put(arr[i], count);
10     }
11     int maxCount = 0;
12     int mode = 0;
13     for(Entry<Integer, Integer> e : counts.entrySet()) {
14         if(e.getValue() > maxCount) {
15             maxCount = e.getValue();
16             mode = e.getKey();
17         }
18     }
19     return mode;
20 }

```

Code Sample 5.7: Mode Finding Algorithm 3

Yet another solution, presented in Code Sample 5.7, utilizes a map data structure to compute the mode. A map is a data structure that allows you to store key-value pairs. In this case, we map elements in the array to a counter that represents the element's multiplicity. The algorithm works by iterating over the array and entering/updating the elements and counters.

There is some cost associated with inserting and retrieving elements from the map, but this particular implementation offers *amortized* constant running time for these operations. That is, some particular entries/retrievals may be more expensive (say linear), but when averaged over the life of the algorithm/data structure, each operation only takes a constant amount of time.

Once built, we need only go through the elements in the map (at most  $n$ ) and find the one with the largest counter. This algorithm, too, offers essentially linear runtime for all inputs. Similar experimental results can be found in Table 5.4.

Algorithm	Number of Additions	Input Size					
		10	100	1,000	10,000	100,000	1,000,000
1	$\approx n^2$	0.007ms	0.155ms	11.982ms	45.619ms	3565.570ms	468086.566ms
2	$n$	0.143ms	0.521ms	2.304ms	19.588ms	40.038ms	735.351ms
3	$n$	0.040ms	0.135ms	0.703ms	10.386ms	21.593ms	121.273ms

Table 5.4: Empirical Performance of the Three Mode Finding Algorithms

The difference in performance is even more dramatic than in the previous example. For an input size of 1,000,000 elements, the  $n^2$  algorithm took nearly *8 minutes*! This is certainly unacceptable performance for most applications. If we were to extend the experiment to  $n = 10,000,000$ , we would expect the running time to increase to about *13 hours*! For perspective, input sizes in the millions are *small* by today’s standards. Algorithms whose runtime is quadratic are *not* considered feasible for today’s applications.

## 5.2 Pseudocode

We will want to analyze algorithms in an abstract, general way independent of any particular hardware, framework, or programming language. In order to do this, we need a way to specify algorithms that is also independent of any particular language. For that purpose, we will use *pseudocode*.

Pseudocode (“fake” code) is similar to some programming languages that you’re familiar with, but does not have any particular syntax rules. Instead, it is a higher-level description of a process. You may use familiar control structures such as loops and conditionals, but you can also utilize natural language descriptions of operations.

There are no established rules for pseudocode, but in general, good pseudocode:

- Clearly labels the algorithm
- Identifies the input and output at the top of the algorithm
- Does not involve any language or framework-specific syntax—no semicolons, declaration of variables or their types, etc.
- Makes liberal use of mathematical notation and natural language for clarity

Good pseudocode abstracts the algorithm by giving enough details necessary to understand the algorithm and subsequently implement it in an actual programming language. Let’s look at some examples.

Algorithm 5 describes a way to compute the average of a collection of numbers. Observe:

- The input does not have a specific *type* (such as `int` or `double`), it uses set notation which also indicates how large the collection is.

INPUT : A collection of numbers,  $A = \{a_1, \dots, a_n\}$

OUTPUT: The *mean*,  $\mu$  of the values in  $A$

```

1  $sum \leftarrow 0$ 
2 FOREACH  $a_i \in A$  DO
3   |  $sum \leftarrow sum + a_i$ 
4 END
5  $\mu \leftarrow \frac{sum}{n}$ 
6 output  $\mu$ 

```

### Algorithm 5: Computing the Mean

- There is no language-specific syntax such as semicolons, variable declarations, etc.
- The loop construct doesn't specify the details of incrementing a variable, instead using a "foreach" statement with some set notation<sup>3</sup>
- Code blocks are not denoted by curly brackets, but are clearly delineated by using indentation and vertical lines.
- Assignment and compound assignment operators do not use the usual syntax from C-style languages, instead using a left-oriented arrow to indicate a value is assigned to a variable.<sup>4</sup>

Consider another example of computing the mode, similar to the second approach in a previous example.

Some more observations about Algorithm 6:

- The use of natural language to specify that the collection should be sorted and how
- The usage of  $-\infty$  as a placeholder so that any other value would be greater than it
- The use of natural language to specify that an iteration takes place over contiguous elements (line 3) or that a sub-operation such as a count/summation (line 4) is performed

In contrast, bad pseudocode would be have the opposite elements. Writing a full program or code snippet in Java for example. Bad pseudocode may be unclear or it may overly simplify the process to the point that the description is trivial. For example, suppose we wanted to specify a sorting algorithm, and we did so using the pseudocode in Algorithm 7. This trivializes the process. There are many possible sorting algorithms (insertion sort, quick sort, etc.) but this algorithm doesn't specify *any* details for how to go about sorting it.

<sup>3</sup> To review,  $a_i \in A$  is a predicate meaning the element  $a_i$  is *in* the set  $A$ .

<sup>4</sup>Not all languages use the familiar single equals sign `=` for the assignment operator. The statistical programming language R uses the left-arrow operator, `<-` and Maple uses `:=` for example.

```

INPUT   : A collection of numbers,  $A = \{a_1, \dots, a_n\}$ 
OUTPUT : A mode of  $A$ 
1 Sort the elements in  $A$  in non-decreasing order
2  $multiplicity \leftarrow -\infty$ 
3 FOREACH run of contiguous equal elements  $a$  DO
4    $m \leftarrow$  count up the number of times  $a$  appears
5   IF  $m > multiplicity$  THEN
6      $mode \leftarrow a$ 
7      $multiplicity \leftarrow m$ 
8   END
9 END
10 output  $m$ 

```

**Algorithm 6:** Computing the Mode

On the other hand, in Algorithm 6, we *did* essentially do this. In that case it was perfectly fine: sorting was a side operation that could be achieved by a separate algorithm. The point of the algorithm was not to specify how to sort, but instead how sorting could be used to solve another problem, finding the mode.

```

INPUT   : A collection of numbers,  $A = \{a_1, \dots, a_n\}$ 
OUTPUT :  $A'$ , sorted in non-decreasing order
1  $A' \leftarrow$  Sort the elements in  $A$  in non-decreasing order
2 output  $A'$ 

```

**Algorithm 7:** Trivial Sorting (Bad Pseudocode)

Another example would be if we need to find a minimal element in a collection. Trivial pseudocode may be like that found in Algorithm 8. No details are presented on *how* to find the element. However, if finding the minimal element were an operation used in a larger algorithm (such as selection sort), then this terseness is perfectly fine. If the primary purpose of the algorithm is to find the minimal element, then details *must* be presented as in Algorithm 9.

## 5.3 Analysis

Given two competing algorithms, we could empirically analyze them like we did in previous examples. However, it may be infeasible to implement both just to determine



```

INPUT  : A collection of numbers,  $A = \{a_1, \dots, a_n\}$ 
OUTPUT: The minimal element of  $A$ 
1  $m \leftarrow$  minimal element of  $A$ 
2 output  $m$ 

```

**Algorithm 8:** Trivially Finding the Minimal Element

```

INPUT  : A collection of numbers,  $A = \{a_1, \dots, a_n\}$ 
OUTPUT: The minimal element of  $A$ 
1  $m \leftarrow \infty$ 
2 FOREACH  $a_i \in A$  DO
3   | IF  $a_i < m$  THEN
4   |   |  $m \leftarrow a_i$ 
5   | END
6 END
7 output  $m$ 

```

**Algorithm 9:** Finding the Minimal Element

which is better. Moreover, by analyzing them from a more abstract, theoretical approach, we have a better more mathematically-based *proof* of the relative complexity of two algorithms.

Given an algorithm, we can analyze it by following this step-by-step process.

1. Identify the input
2. Identify the input size,  $n$
3. Identify the *elementary operation*
4. Analyze how many times the elementary operation is executed with respect to the input size  $n$
5. Characterize the algorithm's complexity by providing an asymptotic (Big-O, or Theta) analysis

### Identifying the Input

This step is pretty straightforward. If the algorithm is described with good pseudocode, then the input will already be identified. Common types of inputs are single numbers, collections of elements (lists, arrays, sets, etc.), data structures such as graphs, matrices, etc.

However, there may be some algorithms that have *multiple* inputs: two numbers or a collection and a key, etc. In such cases, it simplifies the process if you can, without loss of generality, restrict attention to a single input value, usually the one that has the most relevance to the elementary operation you choose.

### Identifying the Input Size

Once the input has been identified, we need to identify its size. We'll eventually want to characterize the algorithm as a function  $f(n)$ : given an input size, how many resources does it take. Thus, it is important to identify the number corresponding to the domain of this function.

This step is also pretty straightforward, but may be dependent on the type of input or even its representation. Examples:

- For collections (sets, lists, arrays), the most natural is to use the number of elements in the collection (cardinality, size, etc.). The size of individual elements is not as important as number of elements since the size of the collection is likely to grow more than individual elements do.
- An  $n \times m$  matrix input could be measured by one or both  $nm$  of its dimensions.
- For graphs, you could count either the number of vertices or the number of edges in the graph (or both!). How the graph is represented may also affect its input size (an adjacency matrix vs. an adjacency list).
- If the input is a number  $x$ , the input size is typically the number of bits required to represent  $x$ . That is,

$$n \approx \log_2(x)$$

To see why, recall that if you have  $n$  bits, the maximum unsigned integer you can represent is  $2^n - 1$ . Inverting this expression gives us  $\lceil \log_2(x + 1) \rceil$ .

Some algorithms may have multiple inputs. For example, a collection and a number (for searching) or two integers as in Euclid's algorithm. The general approach to analyzing such algorithms to simplify things by only considering *one* input. If one of the inputs is larger, such as a collection vs. a single element, the larger one is used in the analysis. Even if it is not clear which one is larger, it may be possible to assume, without loss of generality, that one is larger than the other (and if not, the inputs may be switched). The input size can then be limited to one variable to simplify the analysis.

### Identifying the Elementary Operation

We also need to identify what part of the algorithm does the actual work (where the most resources will be expended). Again, we want to keep the analysis simple, so we generally only identify one *elementary operation*. There may be several reasonable candidates

for the elementary operation, but in general it should be the most common or most expensive operation performed in the algorithm. For example:

- When performing numeric computations, arithmetic operations such as additions, divisions, etc.
- When sorting or searching, comparisons are the most natural elementary operations. Swaps may also be a reasonable choice depending on how you want to analyze the algorithm.
- When traversing a data structure such as a linked list, tree, or graph a node traversal (visiting or processing a node) may be considered the elementary operation.

In general, operations that are necessary to control structures (such as loops, assignment operators, etc.) are not considered good candidates for the elementary operation. An extended discussion of this can be found in Section 5.6.2.

## Analysis

Once the elementary operation has been identified, the algorithm must be analyzed to count the number of times it is executed with respect to the input size. That is, we analyze the algorithm to find a function  $f(n)$  where  $n$  is the input size and  $f(n)$  gives the number of times the elementary operation is executed.

The analysis may involve deriving and solving a summation. For example, if the elementary operation is performed within a for loop and the loop runs a number of times that depends on the input size  $n$ .

If there are multiple loops in which the elementary operation is performed, it may be necessary to setup multiple summations. If two loops are separate and independent (one executes *after* the other), then the *sum rule* applies. The total number of operations is the sum of the operations of each loop.

If two loops are nested, then the *product rule* applies. The inner loop will execute fully *for each* iteration of the outer loop. Thus, the number of operations are multiplied with each other.

Sometimes the analysis will not be so clear cut. For example, a while loop may execute until some condition is satisfied that does not directly depend on the input size but also on the nature of the input. In such cases, we can simplify our analysis by considering the *worst-case* scenario. In the while loop, what is the *maximum* possible number of iterations for any input?

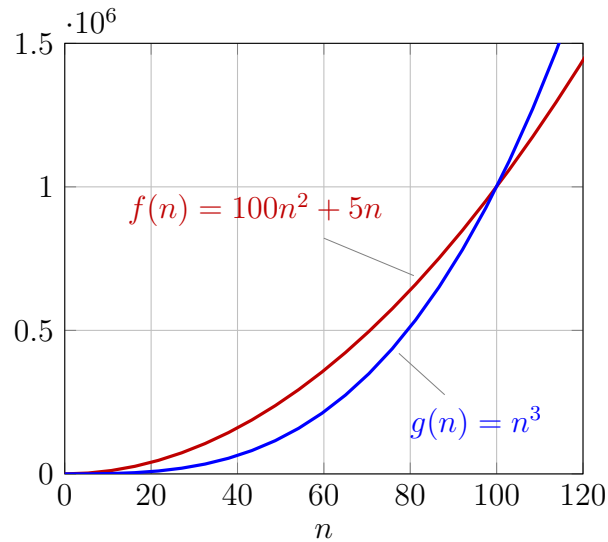


Figure 5.2: Plot of two functions.

### Asymptotic Characterization

As computers get faster and faster and resources become cheaper, they can process more and more information in the same amount of time. However, the characterization of an algorithm should be *invariant* with respect to the underlying hardware. If we run an algorithm on a machine that is twice as fast, that doesn't mean that the algorithm has improved. It still takes the same *number* of operations to execute. Faster hardware simply means that the time it takes to execute those operations is half as much as it was before.

To put it in another perspective, performing Euclid's algorithm to find the GCD of two integers took the same number of steps 2,300 years ago when he performed them on paper as it does today when they are executed on a digital computer. A computer is obviously faster than Euclid would have been, but both Euclid and the computer are performing the same number of steps when executing the same algorithm.

For this reason, we characterize the number of operations performed by an algorithm using *asymptotic analysis*. Improving the hardware by a factor of two only affects the "hidden constant" sitting outside of the function produced by the analysis in the previous step. We want our characterization to be invariant of those constants.

Moreover, we are really more interested in how our algorithm performs for larger and larger input sizes. To illustrate, suppose that we have two algorithms, one that performs

$$f(n) = 100n^2 + 5n$$

operations and one that performs

$$g(n) = n^3$$

operations. These functions are graphed in Figure 5.2. For inputs of size less than 100, the first algorithm performs *worse* than the second (the graph is higher indicating “more” resources). However, for inputs of size greater than 100, the first algorithm is better. For small inputs, the second algorithm may be better, but small inputs are not the norm for any “real” problems.<sup>5</sup> In any case, on modern computers, we would expect small inputs to execute fast anyway as they did in our empirical experiments in Section 5.1.1 and 5.1.2. There was essentially no discernible difference in the three algorithms for sufficiently small inputs.

We can rigorously quantify this by providing an asymptotic characterization of these functions. An asymptotic characterization essentially characterizes the *rate of growth* of a function or the *relative* rate of growth of functions. In this case,  $n^3$  grows much faster than  $100n^2 + 5n$  as  $n$  grows (tends toward infinity). We formally define these concepts in the next section.

## 5.4 Asymptotics

### 5.4.1 Big-O Analysis

We want to capture the notion that one function grows faster than (or at least as fast as) another. Categorizing functions according to their growth rate has been done for a long time in mathematics using *big-O notation*.<sup>6</sup>

**Definition 1.** Let  $f$  and  $g$  be two functions,  $f, g : \mathbb{N} \rightarrow \mathbb{R}^+$ . We say that

$$f(n) \in O(g(n))$$

read as “ $f$  is big-O of  $g$ ,” if there exist constants  $c \in \mathbb{R}^+$  and  $n_0 \in \mathbb{N}$  such that for every integer  $n \geq n_0$ ,

$$f(n) \leq cg(n)$$

First, let’s make some observations about this definition.

- The “O” originally stood for “order of”, Donald Knuth referred to it as the capital greek letter omicron, but since it is indistinguishable from the Latin letter “O” it makes little difference.

---

<sup>5</sup>There are problems where we can apply a “hybrid” approach: we can check for the input size and choose one algorithm for small inputs and another for larger inputs. This is typically done in hybrid sorting algorithms such as when merge sort is performed for “large” inputs but switches over to insertion sort for smaller arrays.

<sup>6</sup>The original notation and definition are attributed to Paul Bachmann in 1894 [2]. Definitions and notation have been refined and introduced/reintroduced over the years. Their use in algorithm analysis was first suggested by Donald Knuth in 1976 [5].

- Some definitions are more general about the nature of the functions  $f, g$ . However, since we're looking at these functions as characterizing the resources that an algorithm takes to execute, we've restricted the domain and codomain of the functions. The domain is restricted to non-negative integers since there is little sense in negative or fractional input sizes. The codomain is restricted to nonnegative reals as it doesn't make sense that an algorithm would potentially consume a negative amount of resources.
- We've used the set notation  $f(n) \in O(g(n))$  because, strictly speaking,  $O(g(n))$  is a *class* of functions: the set of all functions that are asymptotically bounded by  $g(n)$ . Thus the set notation is the most appropriate. However, you will find many sources and papers using notation similar to

$$f(n) = O(g(n))$$

This is a slight abuse of notation, but common nonetheless.

The intuition behind the definition of big-O is that  $f$  is *asymptotically* less than or equal to  $g$ . That is, the rate of growth of  $g$  is *at least as fast* as the growth rate of  $f$ . Big-O provides a means to express that one function is an *asymptotic upper bound* to another function.

The definition essentially states that  $f(n) \in O(g(n))$  if, after some point (for all  $n \geq n_0$ ), the value of the function  $g(n)$  will always be larger than  $f(n)$ . The constant  $c$  possibly serves to “stretch” or “compress” the function, but has no effect on the growth rate of the function.

### Example

Let's revisit the example from before where  $f(n) = 100n^2 + 5n$  and  $g(n) = n^3$ . We want to show that  $f(n) \in O(g(n))$ . By the definition, we need to show that there exists a  $c$  and  $n_0$  such that

$$f(n) \leq cg(n)$$

As we observed in the graph in Figure 5.2, the functions “crossed over” somewhere around  $n = 100$ . Let's be more precise about that. The two functions cross over when they are equal, so we setup an equality,

$$100n^2 + 5n = n^3$$

Collecting terms and factoring out an  $n$  (that is, the functions have one crossover point at  $n = 0$ ), we have

$$n^2 - 100n - 5 = 0$$

The values of  $n$  satisfying this inequality can be found by applying the quadratic formula, and so

$$n = \frac{100 \pm \sqrt{10000 + 20}}{2}$$

Which is  $-0.049975\dots$  and  $100.0499\dots$ . The first root is negative and so irrelevant. The second is our cross over point. The next largest integer is 101. Thus, for  $c = 1$  and  $n_0 = 101$ , the inequality is satisfied.

In this example, it was easy to find the intersection because we could employ the quadratic equation to find roots. This is much more difficult with higher degree polynomials. Throw in some logarithmic functions, exponential functions, etc. and this approach can be difficult.

Revisit the definition of big-O: the inequality doesn't have to be tight or precise. In the previous example we essentially fixed  $c$  and tried to find  $n_0$  such that the inequality held. Alternatively, we could fix  $n_0$  to be small and then find the  $c$  (essentially compressing the function) such that the inequality holds. Observe:

$$\begin{aligned} 100n^2 + 5n &\leq 100n^2 + 5n^2 && \text{since } n \leq n^2 \text{ for all } n \geq 0 \\ &= 105n^2 \\ &\leq 105n^3 && \text{since } n^2 \leq n^3 \text{ for all } n \geq 0 \\ &= 105g(n) \end{aligned}$$

By adding positive values, we make the equation larger until it looks like what we want, in this case  $g(n) = n^3$ . By the end we've got our constants: for  $c = 105$  and  $n_0 = 0$ , the inequality holds. There is nothing special about this  $c$ ,  $c = 1000000$  would work too. The point is we need only find at least one  $c, n_0$  pair that the inequality holds (there are an infinite number of possibilities).

## 5.4.2 Other Notations

Big-O provides an asymptotic upper bound characterization of two functions. There are several other notations that provide similar characterizations.

### Big-Omega

**Definition 2.** Let  $f$  and  $g$  be two functions,  $f, g : \mathbb{N} \rightarrow \mathbb{R}^+$ . We say that

$$f(n) \in \Omega(g(n))$$

read as " $f$  is big-Omega of  $g$ ," if there exist constants  $c \in \mathbb{R}^+$  and  $n_0 \in \mathbb{N}$  such that for every integer  $n \geq n_0$ ,

$$f(n) \geq cg(n)$$

Big-Omega provides an asymptotic lower bound on a function. The only difference is the inequality has been reversed. Intuitively  $f$  has a growth rate that is bounded below by  $g$ .

**Big-Theta**

Yet another characterization can be used to show that two functions have the *same* order of growth.

**Definition 3.** Let  $f$  and  $g$  be two functions  $f, g : \mathbb{N} \rightarrow \mathbb{R}^+$ . We say that

$$f(n) \in \Theta(g(n))$$

read as “ $f$  is Big-Theta of  $g$ ,” if there exist constants  $c_1, c_2 \in \mathbb{R}^+$  and  $n_0 \in \mathbb{N}$  such that for every integer  $n \geq n_0$ ,

$$c_1g(n) \leq f(n) \leq c_2g(n)$$

Big- $\Theta$  essentially provides an asymptotic equivalence between two functions. The function  $f$  is bounded above *and* below by  $g$ . As such, both functions have the same rate of growth.

**Soft-O Notation**

Logarithmic factors contribute very little to a function’s rate of growth especially compared to larger order terms. For example, we called  $n \log(n)$  *quasilinear* since it was nearly linear. Soft-O notation allows us to simplify terms by removing logarithmic factors.

**Definition 4.** Let  $f, g$  be functions such that  $f(n) \in O(g(n) \cdot \log^k(n))$ . Then we say that  $f(n)$  is *soft-O* of  $g(n)$  and write

$$f(n) \in \tilde{O}(g(n))$$

For example,

$$n \log(n) \in \tilde{O}(n)$$

**Little Asymptotics**

Related to big- $O$  and big- $\Omega$  are their corresponding “little” asymptotic notations, little- $o$  and little- $\omega$ .

**Definition 5.** Let  $f$  and  $g$  be two functions  $f, g : \mathbb{N} \rightarrow \mathbb{R}^+$ . We say that

$$f(n) \in o(g(n))$$

read as “ $f$  is little- $o$  of  $g$ ,” if

$$\lim_{n \rightarrow \infty} \frac{f(n)}{g(n)} = 0$$



The little- $o$  is sometimes defined as for every  $\epsilon > 0$  there exists a constant  $N$  such that

$$|f(n)| \leq \epsilon |g(n)| \quad \forall n \geq N$$

but given the restriction that  $g(n)$  is positive, the two definitions are essentially equivalent.

Little- $o$  is a much stronger characterization of the relation of two functions. If  $f(n) \in o(g(n))$  then not only is  $g$  an asymptotic upper bound on  $f$ , but they are *not* asymptotically equivalent. Intuitively, this is similar to the difference between saying that  $a \leq b$  and  $a < b$ . The second is a stronger statement as it implies the first, but the first does not imply the second. Analogous to this example, little- $o$  provides a “strict” asymptotic upper bound. The growth rate of  $g$  is *strictly* greater than the growth rate of  $f$ .

Similarly, a little- $\omega$  notation can be used to provide a *strict* lower bound characterization.

**Definition 6.** Let  $f$  and  $g$  be two functions  $f, g : \mathbb{N} \rightarrow \mathbb{R}^+$ . We say that

$$f(n) \in \omega(g(n))$$

read as “ $f$  is little-omega of  $g$ ,” if

$$\lim_{n \rightarrow \infty} \frac{f(n)}{g(n)} = \infty$$

### 5.4.3 Observations

As you might have surmised, big- $O$  and big- $\Omega$  are duals of each other, thus we have the following.

**Lemma 1.** Let  $f, g$  be functions. Then

$$f(n) \in O(g(n)) \iff g(n) \in \Omega(f(n))$$

Because big- $\Theta$  provides an asymptotic equivalence, both functions are big- $O$  *and* big- $\Theta$  of each other.

**Lemma 2.** Let  $f, g$  be functions. Then

$$f(n) \in \Theta(g(n)) \iff f(n) \in O(g(n)) \text{ and } f(n) \in \Omega(g(n))$$

Equivalently,

$$f(n) \in \Theta(g(n)) \iff g(n) \in O(f(n)) \text{ and } g(n) \in \Omega(f(n))$$

With respect to the relationship between little- $o$  and little- $\omega$  to big- $O$  and big- $\Omega$ , as previously mentioned, little asymptotics provide a stronger characterization of the growth rate of functions. We have the following as a consequence.

**Lemma 3.** Let  $f, g$  be functions. Then

$$f(n) \in o(g(n)) \Rightarrow f(n) \in O(g(n))$$

and

$$f(n) \in \omega(g(n)) \Rightarrow f(n) \in \Omega(g(n))$$

Of course, the converses of these statements do not hold.

### Common Identities

As a direct consequence of the definition, constant coefficients in a function can be ignored.

**Lemma 4.** For any constant  $c$ ,

$$c \cdot f(n) \in O(f(n))$$

In particular, for  $c = 1$ , we have that

$$f(n) \in O(f(n))$$

and so any function is an upper bound on itself.

In addition, when considering the sum of two functions,  $f_1(n), f_2(n)$ , it suffices to consider the one with a larger rate of growth.

**Lemma 5.** Let  $f_1(n), f_2(n)$  be functions such that  $f_1(n) \in O(f_2(n))$ . Then

$$f_1(n) + f_2(n) \in O(f_2(n))$$

In particular, when analyzing algorithms with independent operations (say, loops), we only need to consider the operation with a higher complexity. For example, when we presorted an array to compute the mode, the presort phase was  $O(n \log(n))$  and the mode finding phase was  $O(n)$ . Thus the total complexity was

$$n \log(n) + n \in O(n \log(n))$$

When dealing with a polynomial of degree  $k$ ,

$$c_k n^k + c_{k-1} n^{k-1} + c_{k-2} n^{k-2} + \cdots + c_1 n + c_0$$

The previous results can be combined to conclude the following lemma.

**Lemma 6.** Let  $p(n)$  be a polynomial of degree  $k$ ,

$$p(n) = c_k n^k + c_{k-1} n^{k-1} + c_{k-2} n^{k-2} + \cdots + c_1 n + c_0$$

then

$$p(n) \in \Theta(n^k)$$

Class Name	Asymptotic Characterization	Algorithm Examples
Constant	$O(1)$	Evaluating a formula
Logarithmic	$O(\log(n))$	Binary Search
Polylogarithmic	$O(\log^k(n))$	
Linear	$O(n)$	Linear Search
Quasilinear	$O(n \log(n))$	Mergesort
Quadratic	$O(n^2)$	Insertion Sort
Cubic	$O(n^3)$	
Polynomial	$O(n^k)$ for any $k > 0$	
Exponential	$O(2^n)$	Computing a powerset
Super-Exponential	$O(2^{f(n)})$ for $f(n) \in \Omega(n)$ For example, $n!$	Computing permutations

Table 5.5: Common Algorithmic Efficiency Classes

## Logarithms

When working with logarithmic functions, it suffices to consider a single base. As Computer Scientists, we always work in base-2 (binary). Thus when we write  $\log(n)$ , we implicitly mean  $\log_2(n)$  (base-2). It doesn't really matter though because all logarithms are the same to within a constant as a consequence of the change of base formula:

$$\log_b(n) = \frac{\log_a(n)}{\log_a(b)}$$

That means that for any valid bases  $a, b$ ,

$$\log_b(n) \in \Theta(\log_a(n))$$

Another way of looking at it is that an algorithm's complexity is the same regardless of whether or not it is performed by hand in base-10 numbers or on a computer in binary.

Other logarithmic identities that you may find useful remembering include the following:

$$\log(n^k) = k \log(n)$$

$$\log(n_1 n_2) = \log(n_1) + \log(n_2)$$

## Classes of Functions

Table 5.5 summarizes some of the complexity functions that are common when doing algorithm analysis. Note that these classes of functions form a hierarchy. For example, linear and quasilinear functions are also  $O(n^k)$  and so are polynomial.

### 5.4.4 Limit Method

The method used in previous examples directly used the definition to find constants  $c, n_0$  that satisfied an inequality to show that one function was big- $O$  of another. This can get quite tedious when there are many terms involved. A much more elegant proof technique borrows concepts from calculus.

Let  $f(n), g(n)$  be functions. Suppose we examine the limit, as  $n \rightarrow \infty$  of the ratio of these two functions.

$$\lim_{n \rightarrow \infty} \frac{f(n)}{g(n)}$$

One of three things could happen with this limit.

The limit could converge to 0. If this happens, then by Definition 5 we have that  $f(n) \in o(g(n))$  and so by Lemma 3 we know that  $f(n) \in O(g(n))$ . This makes sense: if the limit converges to zero that means that  $g(n)$  is growing much faster than  $f(n)$  and so  $f$  is big- $O$  of  $g$ .

The limit could diverge to infinity. If this happens, then by Definition 6 we have that  $f(n) \in \omega(g(n))$  and so again by Lemma 3 we have  $f(n) \in \Omega(g(n))$ . This also makes sense: if the limit diverges,  $f(n)$  is growing much faster than  $g(n)$  and so  $f(n)$  is big- $\Omega$  of  $g$ .

Finally, the limit could converge to some positive constant (recall that both functions are restricted to positive codomains). This means that both functions have essentially the same order of growth. That is,  $f(n) \in \Theta(g(n))$ . As a consequence, we have the following Theorem.

**Theorem 2** (Limit Method). Let  $f(n)$  and  $g(n)$  be functions. Then if

$$\lim_{n \rightarrow \infty} \frac{f(n)}{g(n)} = \begin{cases} 0 & \text{then } f(n) \in O(g(n)) \\ c > 0 & \text{then } f(n) \in \Theta(g(n)) \\ \infty & \text{then } f(n) \in \Omega(g(n)) \end{cases}$$

### Examples

Let's reuse the example from before where  $f(n) = 100n^2 + 5n$  and  $g(n) = n^3$ . Setting up our limit,

$$\begin{aligned}\lim_{n \rightarrow \infty} \frac{f(n)}{g(n)} &= \lim_{n \rightarrow \infty} \frac{100n^2 + 5n}{n^3} \\ &= \lim_{n \rightarrow \infty} \frac{100n + 5}{n^2} \\ &= \lim_{n \rightarrow \infty} \frac{100n}{n^2} + \lim_{n \rightarrow \infty} \frac{5}{n^2} \\ &= \lim_{n \rightarrow \infty} \frac{100}{n} + 0 \\ &= 0\end{aligned}$$

And so by Theorem 2, we conclude that

$$f(n) \in O(g(n))$$

Consider the following example: let  $f(n) = \log_2 n$  and  $g(n) = \log_3(n^2)$ . Setting up our limit we have

$$\begin{aligned}\lim_{n \rightarrow \infty} \frac{f(n)}{g(n)} &= \frac{\log_2 n}{\log_3 n^2} \\ &= \frac{\log_2 n}{\frac{2 \log_2 n}{\log_2 3}} \\ &= \frac{\log_2 3}{2} \\ &= .7924 \dots > 0\end{aligned}$$

And so we conclude that  $\log_2(n) \in \Theta(\log_3(n^2))$ .

As another example, let  $f(n) = \log(n)$  and  $g(n) = n$ . Setting up the limit gives us

$$\lim_{n \rightarrow \infty} \frac{\log(n)}{n}$$

The rate of growth might seem obvious here, but we still need to be mathematically rigorous. Both the denominator and numerator are monotone increasing functions. To solve this problem, we can apply l'Hôpital's Rule:

**Theorem 3** (l'Hôpital's Rule). Let  $f$  and  $g$  be functions. If the limit of the quotient  $\frac{f(n)}{g(n)}$  exists, it is equal to the limit of the derivative of the denominator and the numerator. That is,

$$\lim_{n \rightarrow \infty} \frac{f(n)}{g(n)} = \lim_{n \rightarrow \infty} \frac{f'(n)}{g'(n)}$$

Applying this to our limit, the denominator drops out, but what about the numerator? Recall that  $\log(n)$  is the logarithm base-2. The derivative of the natural logarithm is well known,  $\ln'(n) = \frac{1}{n}$ . We can use the change of base formula to transform  $\log(n) = \frac{\ln(n)}{\ln(2)}$  and then take the derivative. That is,

$$\log'(n) = \frac{1}{\ln(2)n}$$

Thus,

$$\begin{aligned} \lim_{n \rightarrow \infty} \frac{\log(n)}{n} &= \lim_{n \rightarrow \infty} \frac{\log'(n)}{n'} \\ &= \lim_{n \rightarrow \infty} \frac{1}{\ln(2)n} \\ &= 0 \end{aligned}$$

Concluding that  $\log(n) \in O(n)$ .

### Pitfalls

l'Hôpital's Rule is not always the most appropriate tool to use. Consider the following example: let  $f(n) = 2^n$  and  $g(n) = 3^n$ . Setting up our limit and applying l'Hôpital's Rule we have

$$\begin{aligned} \lim_{n \rightarrow \infty} \frac{2^n}{3^n} &= \lim_{n \rightarrow \infty} \frac{(2^n)'}{(3^n)'} \\ &= \lim_{n \rightarrow \infty} \frac{(\ln 2)2^n}{(\ln 3)3^n} \end{aligned}$$

which doesn't get us anywhere. In general, we should look for algebraic simplifications *first*. Doing so we would have realized that

$$\lim_{n \rightarrow \infty} \frac{2^n}{3^n} = \lim_{n \rightarrow \infty} \left(\frac{2}{3}\right)^n$$

Since  $\frac{2}{3} < 1$ , the limit of its exponent converges to zero and we have that  $2^n \in O(3^n)$ .

## 5.5 Examples

### 5.5.1 Linear Search

As a simple example, consider the problem of searching a collection for a particular element. The straightforward solution is known as Linear Search and is featured as

## Algorithm 10

```

INPUT   : A collection  $A = \{a_1, \dots, a_n\}$ , a key  $k$ 
OUTPUT  : The first  $i$  such that  $a_i = k$ ,  $\phi$  otherwise
1 FOR  $i = 1, \dots, n$  DO
2   | IF  $a_i = k$  THEN
3   |   | output  $i$ 
4   | END
5 END
6 output  $\phi$ 

```

**Algorithm 10:** Linear Search

Let's follow the prescribed outline above to analyze Linear Search.

1. Input: this is clearly indicated in the pseudocode of the algorithm. The input is the collection  $A$ .
2. Input Size: the most natural measure of the size of a collection is its cardinality; in this case,  $n$
3. Elementary Operation: the most common operation is the comparison in line 2 (assignments and iterations necessary for the control flow of the algorithm are not good candidates).
4. How many times is the elementary operation executed with respect to the input size,  $n$ ? The situation here actually depends not only on  $n$  but also the contents of the array.
  - Suppose we get lucky and find  $k$  in the first element,  $a_1$ : we've only made one comparison.
  - Suppose we are unlucky and find it as the last element (or don't find it at all). In this case we've made  $n$  comparisons
  - We could also look at the *average* number of comparisons (see Section 5.6.3)

In general, algorithm analysis considers at the worst case scenario unless otherwise stated. In this case, there are  $C(n) = n$  comparisons.

5. This is clearly linear, which is why the algorithm is called Linear Search,

$$\Theta(n)$$

### 5.5.2 Set Operation: Symmetric Difference

Recall that the *symmetric difference* of two sets,  $A \oplus B$  consists of all elements in  $A$  or  $B$  but not both. Algorithm 11 computes the symmetric difference.

```

INPUT  : Two sets,  $A = \{a_1, \dots, a_n\}$ ,  $B = \{b_1, \dots, b_m\}$ 
OUTPUT : The symmetric difference,  $A \oplus B$ 
1   $C \leftarrow \emptyset$ 
2  FOREACH  $a_i \in A$  DO
3  |   IF  $a_i \notin B$  THEN
4  |   |    $C \leftarrow C \cup \{a_i\}$ 
5  |   END
6  END
7  FOREACH  $b_j \in B$  DO
8  |   IF  $b_j \notin A$  THEN
9  |   |    $C \leftarrow C \cup \{b_j\}$ 
10 |  END
11 END
12 output  $C$ 

```

**Algorithm 11:** Symmetric Difference of Two Sets

Again, following the step-by-step process for analyzing this algorithm,

1. Input: In this case, there are two sets as part of the input,  $A, B$ .
2. Input Size: As specified, each set has cardinality  $n, m$  respectively. We could analyze the algorithm with respect to both input sizes, namely the input size could be  $n + m$ . For simplicity, to work with a single variable, we could also define  $N = n + m$ .

Alternatively, we could make the following observation: without loss of generality, we can assume that  $n \geq m$  (if not, switch the sets). If one input parameter is bounded by the other, then

$$n + m \leq 2n \in O(n)$$

That is, we could simplify the analysis by only considering  $n$  as the input size. There will be no difference in the final asymptotic characterization as the constants will be ignored.

3. Elementary Operation: In this algorithm, the most common operation is the set membership query ( $\notin$ ). Strictly speaking, this operation may not be trivial depending on the type of data structure used to represent the set (it may entail a



series of  $O(n)$  comparisons for example). However, as our pseudocode is concerned, it is sufficient to consider it as our elementary operation.

- How many times is the elementary operation executed with respect to the input size,  $n$ ? In the first for-loop (lines 2–6) the membership query is performed  $n$  times. In the second loop (lines 7–11), it is again performed  $m$  times. Since each of these loops is independent of each other, we would *add* these operations together to get

$$n + m$$

total membership query operations.

- Whether or not we consider  $N = n + m$  or  $n$  to be our input size, the algorithm is clearly linear with respect to the input size. Thus it is a  $\Theta(n)$ -time algorithm.

### 5.5.3 Euclid's GCD Algorithm

The greatest common divisor (or GCD) of two integers  $a, b$  is the largest positive integer that divides both  $a$  and  $b$ . Finding a GCD has many useful applications and the problem has one of the oldest known algorithmic solutions: Euclid's Algorithm (due to the Greek mathematician Euclid c. 300 BCE).

```

INPUT   : Two integers,  $a, b$ 
OUTPUT  : The greatest common divisor,  $\text{gcd}(a, b)$ 
1 WHILE  $b \neq 0$  DO
2   |  $t \leftarrow b$ 
3   |  $b \leftarrow a \bmod b$ 
4   |  $a \leftarrow t$ 
5 END
6 Output  $a$ 

```

#### Algorithm 12: Euclid's GCD Algorithm

The algorithm relies on the following observation: any number that divides  $a, b$  must also divide the remainder of  $a$  since we can write  $b = a \cdot k + r$  where  $r$  is the remainder. This suggests the following strategy: iteratively divide  $a$  by  $b$  and retain the remainder  $r$ , then consider the GCD of  $b$  and  $r$ . Progress is made by observing that  $b$  and  $r$  are necessarily smaller than  $a, b$ . Repeating this process until we have a remainder of zero gives us the GCD because once we have that one evenly divides the other, the larger must be the GCD. Pseudocode for Euclid's Algorithm is provided in Algorithm 12.

The analysis of Euclid's algorithm is seemingly straightforward. It is easy to identify the division in line 3 as the elementary operation. But how many times is it executed with respect to the input size? What is the input size?

When considering algorithms that primarily execute numerical operations the input is usually a number (or in this case a pair of numbers). How big is the input of a number? The input size of 12,142 is not 12,142. The number 12,142 has a compact representation when we write it: it requires 5 digits to express it (in base 10). That is, the input size of a number is the number of symbols required to represent its magnitude. Considering a number's input size to be equal to the number would be like considering its representation in unary where a single symbol is used and repeated for as many times as is equal to the number (like a prisoner marking off the days of his sentence).

Computers don't "speak" in base-10, they speak in binary, base-2. Therefore, the input size of a numerical input is the number of bits required to represent the number. This is easily expressed using the base-2 logarithm function:

$$\lceil \log_2(n) \rceil$$

But in the end it doesn't really matter if we think of computers as speaking in base-10, base-2, or any other integer base greater than or equal to 2 because as we've observed that all logarithms are equivalent to within a constant factor using the change of base formula. In fact this again demonstrates again that algorithms are an abstraction independent of any particular platform: that the same algorithm will have the same (asymptotic) performance whether it is performed on paper in base-10 or in a computer using binary!

Back to the analysis of Euclid's Algorithm: how many times does the while loop get executed? We can first observe that each iteration reduces the value of  $b$ , but by how much? The exact number depends on the input: some inputs would only require a single division, other inputs reduce  $b$  by a different amount on each iteration. The important thing to realize is that we want a general characterization of this algorithm: it suffices to consider the worst case. That is, at maximum, how many iterations are performed? The number of iterations is maximized when the reduction in the value of  $b$  is minimized at each iteration. We further observe that  $b$  is reduced by at least half at each iteration. Thus, the number of iterations is maximized if we reduce  $b$  by at most half on each iteration. So how many iterations  $i$  are required to reduce  $n$  down to 1 (ignoring the last iteration when it is reduced to zero for the moment)? This can be expressed by the equation:

$$n \left(\frac{1}{2}\right)^i = 1$$

Solving for  $i$  (taking the log on either side), we get that

$$i = \log_2(n)$$

Recall that the size of the input is  $\log_2(n)$ . Thus, Euclid's Algorithm is *linear* with respect to the input size.

### 5.5.4 Selection Sort

Recall that Selection Sort is a sorting algorithm that sorts a collection of elements by first finding the smallest element and placing it at the beginning of the collection. It continues by finding the smallest among the remaining  $n - 1$  and placing it second in the collection. It repeats until the “first”  $n - 1$  elements are sorted, which by definition means that the last element is where it needs to be.

The Pseudocode is presented as Algorithm 13.

```

INPUT  : A collection  $A = \{a_1, \dots, a_n\}$ 
OUTPUT : A sorted in non-decreasing order
1 FOR  $i = 1, \dots, n - 1$  DO
2    $min \leftarrow a_i$ 
3   FOR  $j = (i + 1), \dots, n$  DO
4     IF  $min < a_j$  THEN
5        $min \leftarrow a_j$ 
6     END
7     swap  $min, a_i$ 
8   END
9 END
10 output  $A$ 

```

**Algorithm 13:** Selection Sort

Let's follow the prescribed outline above to analyze Selection Sort.

1. Input: this is clearly indicated in the pseudocode of the algorithm. The input is the collection  $A$ .
2. Input Size: the most natural measure of the size of a collection is its cardinality; in this case,  $n$
3. Elementary Operation: the most common operation is the comparison in line 4 (assignments and iterations necessary for the control flow of the algorithm are not good candidates). Alternatively, we could have considered swaps on line 7 which would lead to a different characterization of the algorithm.
4. How many times is the elementary operation executed with respect to the input size,  $n$ ?
  - Line 4 does one comparison each time it is executed
  - Line 4 itself is executed multiple times for each iteration of the for loop in line 3 (for  $j$  running from  $i + 1$  up to  $n$  inclusive).

- line 3 (and subsequent blocks of code) are executed multiple times for each iteration of the for loop in line 1 (for  $i$  running from 1 up to  $n - 1$ )

This gives us the following summation:

$$\underbrace{\sum_{i=1}^{n-1} \underbrace{\sum_{j=i+1}^n 1}_{\text{line 4}}}_{\text{line 3}}_{\text{line 1}}$$

Solving this summation gives us:

$$\begin{aligned} \sum_{i=1}^{n-1} \sum_{j=i+1}^n 1 &= \sum_{i=1}^{n-1} n - i \\ &= \sum_{i=1}^{n-1} n - \sum_{i=1}^{n-1} i \\ &= n(n-1) - \frac{n(n-1)}{2} \\ &= \frac{n(n-1)}{2} \end{aligned}$$

Thus for a collection of size  $n$ , Selection Sort makes

$$\frac{n(n-1)}{2}$$

comparisons.

5. Provide an asymptotic characterization: the function determined is clearly  $\Theta(n^2)$

## 5.6 Other Considerations

### 5.6.1 Importance of Input Size

The second step in our algorithm analysis outline is to identify the input size. Usually this is pretty straightforward. If the input is an array or collection of elements, a reasonable input size is the cardinality (the number of elements in the collection). This is usually the case when one is analyzing a sorting algorithm operating on a list of elements.

Though seemingly simple, sometimes identifying the appropriate input size depends on the nature of the input. For example, if the input is a data structure such as a graph,

the input size could be either the number of vertices *or* number of edges. The most appropriate measure then depends on the algorithm and the details of its analysis. It may even depend on how the input is represented. Some graph algorithms have different efficiency measures if they are represented as adjacency lists or adjacency matrices.

Yet another subtle difficulty is when the input is a single numerical value,  $n$ . In such instances, the input size is *not* also  $n$ , but instead the number of symbols that would be needed to represent  $n$ . That is, the number of digits in  $n$  or the number of bits required to represent  $n$ . We'll illustrate this case with a few examples.

### Sieve of Eratosthenes

A common beginner mistake is made when analyzing the Sieve of Eratosthenes (named for Eratosthenes of Cyrene, 276 BCE – 195 BCE). The Sieve is an ancient method for prime number factorization. A brute-force algorithm, it simply tries every integer up to a point to see if it is a factor of a given number.

```

INPUT  : An integer  $n$ 
OUTPUT : Whether  $n$  is prime or composite
1 FOR  $i = 2, \dots, n$  DO
2   | IF  $i$  divides  $n$  THEN
3   |   | Output composite
4   |   END
5 END
6 Output prime

```

#### Algorithm 14: Sieve of Eratosthenes

The for-loop only needs to check integers up to  $\sqrt{n}$  because any factor greater than  $\sqrt{n}$  would necessarily have a corresponding factor  $\sqrt{n}$ . A naive approach would observe that the for-loop gets executed  $\sqrt{n} - 1 \in O(\sqrt{n})$  times which would lead to the (incorrect) impression that the Sieve is a polynomial-time (in fact sub-linear!) running time algorithm. Amazing that we had a primality testing algorithm over 2,000 years ago! In fact, primality testing was a problem that was not known to have a deterministic polynomial time running algorithm until 2001 (the AKS Algorithm [1]).

The careful observer would realize that though  $n$  is the input, the actual input size is again  $\log(n)$ , the number of bits required to represent  $n$ . Let  $N = \log(n)$  be a placeholder for our actual input size (and so  $n = 2^N$ ). Then the running time of the Sieve is actually

$$O(\sqrt{n}) = O(\sqrt{2^N})$$

which is exponential with respect to the input size  $N$ .

This distinction is subtle but crucial: the difference between a polynomial-time algorithm and an exponential algorithm is huge even for modestly sized inputs. What may take a few milliseconds using a polynomial time algorithm may take billions and billions of years with an exponential time algorithm as we'll see with our next example.

### Computing an Exponent

As a final example, consider the problem of computing a modular exponent. That is, given integers  $a$ ,  $n$ , and  $m$ , we want to compute

$$a^n \bmod m$$

A naive (but common!) solution might be similar to the Java code snippet in Code Sample 5.8.

Whether one chooses to treat multiplication or integer division as the elementary operation, the for-loop executes exactly  $n$  times. For “small” values of  $n$  this may not present a problem. However, for even moderately large values of  $n$ , say  $n \approx 2^{256}$ , the performance of this code will be terrible.

```

1  int a = 45, m = 67;
2  int result = 1;
3  for(int i=1; i<=n; i++) {
4      result = (result * a % m);
5  }
```

Code Sample 5.8: Naive Exponentiation

To illustrate, suppose that we run this code on a 14.561 petaFLOP (14 quadrillion floating point operations per second) super computer cluster (this throughput was achieved by the Folding@Home distributed computing project in 2013). Even with this power, to make  $2^{256}$  floating point operations would take

$$\frac{2^{256}}{14.561 \times 10^{15} \cdot 60 \cdot 60 \cdot 24 \cdot 365.25} \approx 2.5199 \times 10^{53}$$

or 252 sexdecilliion *years* to compute!

For context, this sort of operation is performed by millions of computers around the world every second of the day. A 256-bit number is not really all that “large”. The problem again lies in the failure to recognize the difference between an input,  $n$ , and the input’s *size*. We will explore the better solution to this problem when we examine Repeated Squaring in Section ??.

## 5.6.2 Control Structures are Not Elementary Operations

When considering which operation to select as the elementary operation, we usually do *not* count operations that are necessary to the control structure of the algorithm. For example, assignment operations, or operations that are necessary to execute a loop (incrementing an index variable, a comparison to check the termination condition).

To see why, consider method in Code Sample 5.9. This method computes a simple average of an array of `double` variables. Consider some of the minute operations that are performed in this method:

- An assignment of a value to a variable (lines 2, 3, and 4)
- An increment of the index variable `i` (line 3)
- A comparison (line 3) to determine if the loop should terminate or continue
- The addition (line 4) and division (line 6)

```

1 public static double average(double arr[]) {
2     double sum = 0.0;
3     for(int i=0; i<arr.length; i++) {
4         sum = sum + arr[i];
5     }
6     return sum / arr.length;
7 }

```

Code Sample 5.9: Computing an Average

A proper analysis would use the addition in line 4 as the elementary operation, leading to a  $\Theta(n)$  algorithm. However, for the sake of argument, let's perform a detailed analysis with respect to each of these operations. Assume that there are  $n$  values in the array, thus the for loop executes  $n$  times. This gives us

- $n + 2$  total assignment operations
- $n$  increment operations
- $n$  comparisons
- $n$  additions and
- 1 division

Now, suppose that each of these operations take time  $t_1, t_2, t_3, t_4$ , and  $t_5$  milliseconds each (or whatever time scale you like). In total, we have a running time of

$$t_1(n + 2) + t_2n + t_3n + t_4n + t_5 = (t_1 + t_2 + t_3 + t_4)n + (2t_1 + t_5) = cn + d$$

Which doesn't change the asymptotic complexity of the algorithm: considering the

additional operations necessary for the control structures only changed the constant sitting out front (as well as some additive terms).

The amount of resources (in this case time) that are expended for the assignment, increment and comparison for the loop control structure are proportional to the true elementary operation (addition). Thus, it is sufficient to simply consider the most common or most expensive operation in an algorithm. The extra resources for the control structures end up only contributing constants which are ultimately ignored when an asymptotic analysis is performed.

### 5.6.3 Average Case Analysis

The behavior of some algorithms may depend on the nature of the input rather than simply the input size. From this perspective we can analyze an algorithm with respect to its best-, average-, and worst-case running time.

In general, we prefer to consider the worst-case running time when comparing algorithms.

#### Example: searching an array

Consider the problem of searching an array for a particular element (the array is unsorted, contains  $n$  elements). You could get lucky (best-case) and find it immediately in the first index, requiring only a single comparison. On the other hand, you could be unlucky and find the element in the last position or not at all. In either case  $n$  comparisons would be required (worst-case).

What about the average case? A naive approach would be to average the worst and best case to get an average of  $\frac{n+1}{2}$  comparisons. A more rigorous approach would be to define a probability of a successful search  $p$  and the probability of an unsuccessful search,  $1 - p$ .

For a successful search, we further define a uniform probability distribution on finding the element in each of the  $n$  indices. That is, we will find the element at index  $i$  with probability  $\frac{p}{n}$ . Finding the element at index  $i$  requires  $i$  comparisons. Thus the total number of expected comparisons in a successful search is

$$\sum_{i=1}^n i \cdot \frac{p}{n} = \frac{p(n+1)}{2}$$

For an unsuccessful search, we would require  $n$  comparisons, thus the number of expected comparisons would be

$$n(1 - p)$$

Since these are mutually exclusive events, we sum the probabilities:

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



$p$	$C$
0	$n$
$\frac{1}{4}$	$\frac{7}{8}n + \frac{1}{8}$
$\frac{1}{2}$	$\frac{3}{4}n + \frac{1}{4}$
$\frac{3}{4}$	$\frac{5}{8}n + \frac{3}{8}$
1	$\frac{n+1}{2}$

Table 5.6: Expected number of comparisons  $C$  for various values of the probability of a successful search  $p$ .

We cannot remove the  $p$  terms, but we can make some observations for various values (see Table 5.6.3). When  $p = 1$  for example, we have the same conclusion as the naive approach. As  $p$  decreases, the expected number of comparisons grows to  $n$ .

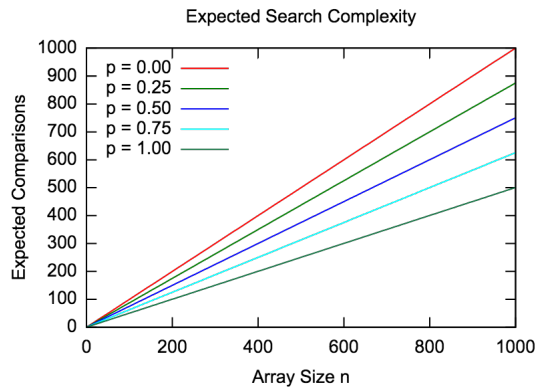


Figure 5.3: Expected number of comparisons for various success probabilities  $p$ .

#### 5.6.4 Amortized Analysis

Sometimes it is useful to analyze an algorithm not based on its worst-case running time, but on its *expected* running time. That is, on average, how many resources does an algorithm perform? This is a more practical approach to analysis. Worst-case analysis assumes that all inputs will be difficult, but in practice, difficult inputs may be rare. The average input may require fewer resources.

Amortized algorithm analysis is similar to the average-case analysis we performed before, but focuses on how the *cost* of operations may change over the course of an algorithm. This is similar to a loan from a bank. Suppose the loan is taken out for \$1,000 at a 5% interest rate. You don't actually end up paying \$50 the first year. You pay slightly less than that since you are (presumably) making monthly payments, reducing the balance and thus reducing the amount of interest accrued each month.

We will not go into great detail here, but as an example, consider Heap Sort. This algorithm uses a *Heap* data structure. It essentially works by inserting elements into the heap and then removing them one by one. Due to the nature of a heap data structure, the elements will come out in the desired order.

An amortized analysis of Heap Sort may work as follows. First, a heap requires about  $\log(n)$  comparisons to insert an element (as well as remove an element, though we'll focus on insertion), where  $n$  is the size of the heap (the number of elements currently in the heap). As the heap grows (and eventually shrinks), the cost of inserts will change. With an initially empty heap, there will only be 0 comparisons. When the heap has 1 element, there will be about  $\log(1)$  comparisons, then  $\log(2)$  comparisons and so on until we insert the last element. This gives us

$$C(n) = \sum_{i=1}^{n-1} \log(i)$$

total comparisons.

## 5.7 Analysis of Recursive Algorithms

Recursive algorithms can be analyzed with the same basic 5 step process. However, because they are recursive, the analysis (step 4) may involve setting up a recurrence relation in order to characterize how many times the elementary operation is executed.

Though cliché, as a simple example consider a naive recursive algorithm to compute the  $n$ -th Fibonacci number. Recall that the Fibonacci sequence is defined as

$$F_n = F_{n-1} + F_{n-2}$$

with initial conditions  $F_0 = F_1 = 1$ . That is, the  $n$ -th Fibonacci number is defined as the sum of the two previous numbers in the sequence. This defines the sequence

$$1, 1, 2, 3, 5, 8, 13, 21, 34, 55, \dots$$

The typical recursive solution is presented as Algorithm 15

The elementary operation is clearly the addition in line 4. However, how do we analyze the number of additions performed by a call to  $\text{Fibonacci}(n)$ ? To do this, we setup

```

INPUT   : An integer  $n$ 
OUTPUT  : The  $n$ -th Fibonacci Number,  $F_n$ 
1 IF  $n = 0$  or  $n = 1$  THEN
2   |   output 1
3 ELSE
4   |   output Fibonacci( $n - 1$ ) + Fibonacci( $n - 2$ )
5 END

```

**Algorithm 15:** Fibonacci( $n$ )

a recurrence relation. Let  $A(n)$  be a function that counts the number of additions performed by Fibonacci( $n$ ). If  $n$  is zero or one, the number of additions is zero (the base case of our recursion performs no additions). That is,  $A(0) = A(1) = 0$ . But what about  $n > 1$ ?

When  $n > 1$ , line 4 executes. Line 4 contains one addition. However, it also contains two recursive calls. How many additions are performed by a call to Fibonacci( $n - 1$ ) and Fibonacci( $n - 2$ )? We defined  $A(n)$  to be the number of additions on a call to Fibonacci( $n$ ), so we can reuse this function: the number of additions is  $A(n - 1)$  and  $A(n - 2)$  respectively. Thus, the total number of additions is

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

This is a recurrence relation,<sup>7</sup> in particular, it is a second-order linear non-homogeneous recurrence relation. This particular relation can be solved. That is,  $A(n)$  can be expressed as a non-recursive closed-form solution.

The techniques required for solving these type of recurrence relations are beyond the scope of this text. However, for many common recursive algorithms, we can use a simple tool to characterize their running time, the “Master Theorem.”

### 5.7.1 The Master Theorem

Suppose that we have a recursive algorithm that takes an input of size  $n$ . The recursion may work by dividing the problem into  $a$  subproblems each of size  $\frac{n}{b}$  (where  $a, b$  are constants). The algorithm may also perform some amount of work before or after the recursion. Suppose that we can characterize this amount of work by a polynomial function,  $f(n) \in \Theta(n^d)$ .

This kind of recursive algorithm is common among “divide-and-conquer” style algorithms that *divide* a problem into subproblems and *conquers* each subproblem. Depending on

<sup>7</sup>Also called a *difference equation*, which are sort of discrete analogs of differential equations.

the values of  $a, b, d$  we can categorize the runtime of this algorithm using the Master Theorem.

**Theorem 4** (Master Theorem). Let  $T(n)$  be a monotonically increasing function that satisfies

$$\begin{aligned} T(n) &= aT\left(\frac{n}{b}\right) + f(n) \\ T(1) &= c \end{aligned}$$

where  $a \geq 1, b \geq 2, c > 0$ . If  $f(n) \in \Theta(n^d)$  where  $d \geq 0$ , then

$$T(n) = \begin{cases} \Theta(n^d) & \text{if } a < b^d \\ \Theta(n^d \log n) & \text{if } a = b^d \\ \Theta(n^{\log_b a}) & \text{if } a > b^d \end{cases}$$

“Master” is a bit of a misnomer. The Master Theorem can only be applied to recurrence relations of the particular form described. It can’t, for example, be applied to the previous recurrence relation that we derived for the Fibonacci algorithm. However, it can be used for several common recursive algorithms.

### Example: Binary Search

As an example, consider the Binary Search algorithm, presented as Algorithm 16. Recall that binary search takes a *sorted* array (random access is required) and searches for an element by checking the middle element  $m$ . If the element being searched for is larger than  $m$ , a recursive search on the upper half of the list is performed, otherwise a recursive search on the lower half is performed.

Let’s analyze this algorithm with respect to the number of comparisons it performs. To simplify, though we technically make two comparisons in the pseudocode (lines 5 and 7), let’s count it as a single comparison. In practice a comparator pattern would be used and logic would branch based on the result. However, there would still only be one invocation of the comparator function/object. Let  $C(n)$  be a function that equals the number of comparisons made by Algorithm 16 while searching an array of size  $n$  *in the worst case* (that is, we never find the element or it ends up being the last element we check).

As mentioned, we make one comparison (line 5, 7) and then make one recursive call. The recursion roughly cuts the size of the array in half, resulting in an array of size  $\frac{n}{2}$ . Thus,

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

Applying the master theorem, we find that  $a = 1, b = 2$ . In this case,  $f(n) = 1$  which is bounded by a polynomial: a polynomial of degree  $d = 0$ . Since

$$1 = a = b^d = 2^0$$

by case 2 of the Master Theorem applies and we have

$$C(n) \in \Theta(\log(n))$$

```

INPUT   : A sorted collection of elements  $A = \{a_1, \dots, a_n\}$ , bounds  $1 \leq l, h \leq n$ ,
          and a key  $e_k$ 
OUTPUT  : An element  $a \in A$  such that  $a = e_k$  according to some criteria;  $\phi$  if no
          such element exists

1 IF  $l > h$  THEN
2   | output  $\phi$ 
3 END
4  $m \leftarrow \lfloor \frac{h+l}{2} \rfloor$ 
5 IF  $a_m = e_k$  THEN
6   | output  $a_m$ 
7 ELSE IF  $a_m < e_k$  THEN
8   | BINARYSEARCH( $A, m + 1, h, e$ )
9 ELSE
10  | BINARYSEARCH( $A, l, m - 1, e$ )
11 END

```

**Algorithm 16:** Binary Search – Recursive**Example: Merge Sort**

Recall that Merge Sort is an algorithm that works by recursively splitting an array of size  $n$  into two equal parts of size roughly  $\frac{n}{2}$ . The recursion continues until the array is trivially sorted (size 0 or 1). Following the recursion back up, each subarray half is sorted and they need to be *merged*. This basic divide and conquer approach is presented in Algorithm 17. We omit the details of the merge operation on line 4, but observe that it can be achieved with roughly  $n - 1$  comparisons in the worst case.

```

INPUT   : An array, sub-indices  $1 \leq l, r \leq n$ 
OUTPUT  : An array  $A'$  such that  $A[l, \dots, r]$  is sorted

1 IF  $l < r$  THEN
2   | MERGESORT( $A, l, \lfloor \frac{r+l}{2} \rfloor$ )
3   | MERGESORT( $A, \lceil \frac{r+l}{2} \rceil, r$ )
4   | Merge sorted lists  $A[l, \dots, \lfloor \frac{r+l}{2} \rfloor]$  and  $A[\lceil \frac{r+l}{2} \rceil, \dots, r]$ 
5 END

```

**Algorithm 17:** Merge Sort

Let  $C(n)$  be the number of comparisons made by Merge Sort. The main algorithm makes *two* recursive calls on subarrays of size  $\frac{n}{2}$ . There are also  $n + 1$  comparisons made after the recursion. Thus we have

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

Again, applying the Master Theorem we have that  $a = 2, b = 2$  and that  $f(n) = (n + 1) \in \Theta(n)$  and so  $d = 1$ . Thus,

$$2 = a = b^d = 2^1$$

and by case 2,

$$C(n) \in \Theta(n \log(n))$$

We will apply the Master Theorem to several other recursive algorithms later on.

# 6 Trees

## 6.1 Introduction

One of our fundamental goals is to design a data structure to store elements that offers efficient and arbitrary retrieval (search), insertion, and deletion operations. We've already seen several examples of list-based data structures that offer different properties. Array-based lists offer fast index-based retrieval and, if sorted even more efficient,  $O(\log n)$  key-based retrieval using binary search. However, arbitrary insertion and deletion may result in shifts leading to  $O(n)$  behavior. In contrast, linked lists offer efficient,  $O(1)$  insertion and deletion at the head/tail of the list, but inefficient,  $O(n)$  arbitrary search/insert/delete.

Stacks and queues are efficient in that their core functionality offers  $O(1)$  operations (push/pop, enqueue/dequeue) but they are restricted-access data structures and do not offer arbitrary and efficient operations.

We now turn our attention to another fundamental data structure based on *trees*. Trees are a very special type of graph. Graphs are useful because they can model many types of relations and processes in physics, biology, finance, and pretty much every other discipline. Moreover, we have over two centuries of foundational work on graph results and algorithms to work with.<sup>1</sup> General graphs, however, are not necessarily as structured as they need to be for an efficient data structure. Trees, however, are highly structured and have several useful properties that we can exploit to design an efficient data structure. As we will see, trees have the potential to offer efficient  $O(\log n)$  behavior for all three fundamental operations.

## 6.2 Definitions & Terminology

A *graph* is a collection of nodes such that pairs of nodes may be connected by edges. More formally,

**Definition 7** (Undirected Graph). A *graph* is a two-tuple,  $G = (V, E)$  where

- $V = \{v_1, \dots, v_n\}$  is a set of vertices

---

<sup>1</sup>Graph Theory usually credited to Euler who studied the *Seven Bridges of Königsberg* [4] in the 18th century.

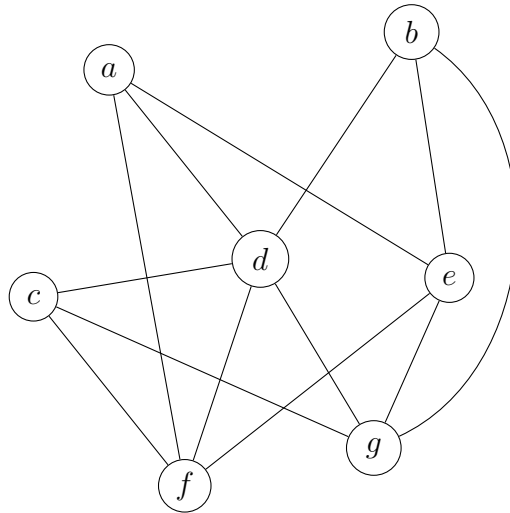


Figure 6.1: An undirected graph with labeled vertices.

- $E \subseteq V \times V$  is a set of edges connecting nodes

In general, the edges connecting two vertices may be directed or undirected. However, we'll only focus on undirected edges so that if two vertices  $u, v$  are connected, then the edge that connects them can be written as an unordered pair,  $e = (u, v) = (v, u)$ . An example of a graph can be found in Figure 6.1.

In the example, the nodes have been labeled with single characters. In general, we'll want to store elements of any type in a graph's vertices. In the context of data structures when vertices hold elements, they are usually referred to as *nodes* and we'll use this terminology from here on. Another point to highlight is that the size of a graph is usually measured in terms of the number of nodes in it. In the definition we have designated the variable  $n$  to denote this. In general, a graph may have up to

$$\binom{n}{2} = \frac{n(n-1)}{2}$$

edges connecting pairs of nodes.<sup>2</sup>

## Trees

General graphs may model very complex binary relations. Other than the graph structure itself (nodes possibly connected by edges) there is not much structure to exploit. For this reason, we will instead focus on a subclass of graphs called trees.

<sup>2</sup>This notation is from combinatorics,  $\binom{n}{k}$  is a binomial coefficient and can be read as “ $n$  choose  $k$ .” This operation counts the number of ways there are to choose an unordered subset of  $k$  elements from a set of  $n$  unique elements. In this case,  $\binom{n}{2}$  is counting the number of ways there are to choose a pair of vertices from a set of  $n$  vertices.



**Definition 8** (Trees). A *tree* is an acyclic graph.

A *cycle* in a graph is any sequence of pairwise connected nodes that begin and end at the same node. For example, the sequence of vertices  $aegcda$  in the graph in Figure 6.1 forms a cycle. Likewise, the sequences  $begdb$ ,  $bgdb$ , and  $cdgc$  also form cycles as well as many others. A graph in which there are no cycles is referred to as an *acyclic* graph and has a much simpler structure. This simplicity gives a tree several properties that can be exploited. An example of several trees can be found in Figure 6.2.

The example in Figure 6.2(d) is actually a *disconnected* tree in that there are several components whose nodes are not connected by any edges. Such graphs can be seen as a collection of trees and are usually called *forests*. Though forests are useful in modeling some problems, we will again restrict our attention to *connected trees* such that there is only *one* connected component.

### Structural Properties

By definition, trees are acyclic which already gives them a high degree of structure. There are several related properties that follow from this structure. First, since we are restricting attention to connected trees, it will always be the case that the number of edges in a tree will be equal to one less than the number of vertices. That is:

**Lemma 7.** In any tree,

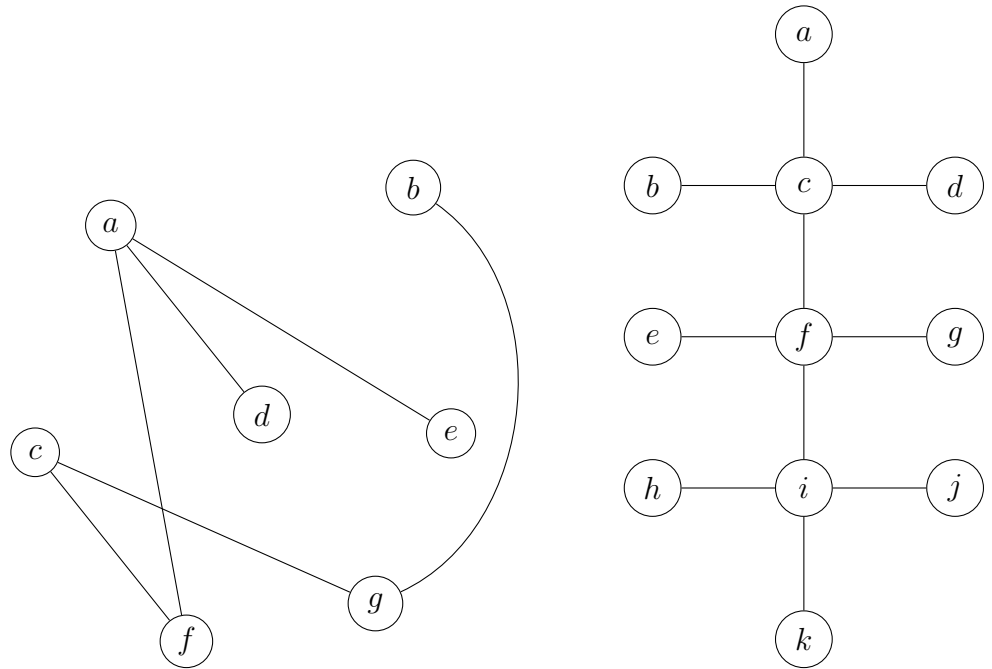
$$|E| = n - 1$$

where  $n = |V|$ .

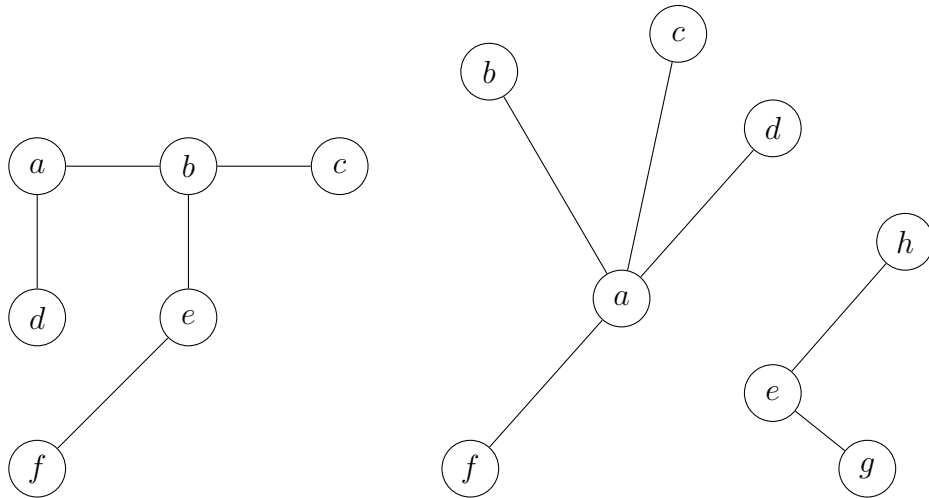
Recall that in general, graphs may have as many as  $O(n^2)$  vertices. However, trees will necessarily have fewer  $O(n)$  edges. This makes trees *sparse graphs*. Many graph problems that are hard or do not admit themselves to efficient solutions become easy or even trivial on sparse graphs and trees in particular.

Another consequence of a lack of cycles means that any two nodes in a tree will be connected by exactly one, unique *path*. A path is much like a cycle except that it need not begin/end at the same node. In Figure 6.2(a), the sequence of vertices  $afcgb$  form a path. This is, in fact, the only path connecting the vertices  $a$  and  $b$ . Furthermore, the *length* of a path is defined as the number of edges on it. The aforementioned path  $afcgb$  has length 4. Note that the length of a path is also equal to the number of nodes in the path minus one. In the same tree in Figure 6.2(a),  $afcgcfad$  is also technically a path, however, it traverses several edges more than once, backtracking from  $g$ . Such a path is referred to as a non-simple path. A path that does not traverse an edge more than once is a *simple* path and we will restrict our consideration to simple paths.

**Lemma 8.** Let  $T = (V, E)$  be a connected tree and let  $u, v \in V$  be two nodes in  $T$ . Then there exists exactly one single unique path connecting  $u, v$ .



(a) The previous graph from Figure 6.1 with enough edges removed to make it a tree. (b) Another tree drawn in a more organized manner.



(c) Yet another little tree, happier. (d) A disconnected tree with two components.

Figure 6.2: Several examples of trees. It doesn't matter how the tree is depicted or organized, only that it is acyclic. The final example, 6.2(d) represents a disconnected tree, called a *forest*.

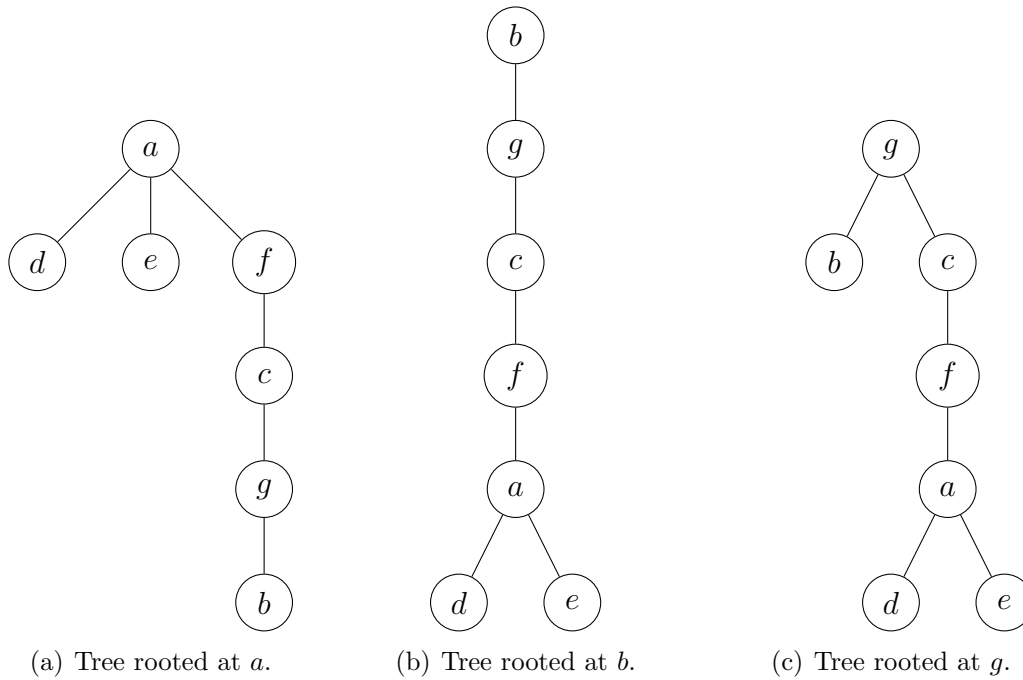


Figure 6.3: Several possible rooted orientations for the tree from Figure 6.2(a).

*Proof.* First we observe that there has to be at least one path between  $u, v$ . Otherwise, if there were not, then  $T$  would not be a connected tree.

Now suppose, by way of contradiction that there exist at least two paths  $p$  and  $p'$  connecting  $u, v$  such that  $p \neq p'$ . However, this means that we can now form a cycle. Starting at  $u$  and taking the first path  $p$  to  $v$  and then from  $v$  returning to  $u$  via  $p'$ . Since this forms a cycle, it contradicts the fact that we started with a tree. Thus, there cannot be two distinct paths between two vertices.  $\square$

These properties will prove useful in adapting trees as collection data structures.

## Orienting Trees

To give trees even more structure, we can *orient* them. The first way that we'll do this is by *rooting* them. Imagine that you could take a tree (say any tree in Figure 6.2) and lift it from the page<sup>3</sup> by grabbing a single node. Then shake the tree and let gravity dangle the remaining nodes downward. Then place the tree back onto the page. The node that you shook the tree from will be designated as the tree's *root* and all other nodes are oriented downward. An example can be found in Figure 6.3 where we have done this operation on the tree in Figure 6.2(a) by dangling from various nodes.

<sup>3</sup>Or screen if you will.

The alert reader and professional arborists<sup>4</sup> alike will notice that we've drawn our trees with the root at the top and its "branches" growing downward, the exact opposite of how most real, organic trees grow. Trees in Computer Science "grow" downward because that is how we read. Since trees will ultimately be used to store data it is easier for a human to read them top-to-bottom, left-to-right. Of course, when represented in a computer, there is no such orientation; it is only a convention that we humans use when drawing and describing them. This convention is quite old; in fact Arthur Cayley, the original mathematician who first studied trees drew them like this in his original paper [3].<sup>5</sup>

The vertical orientation of a trees lends itself to some obvious terminology, borrowed from the concept of a family tree. Let  $u$  be a tree node. The node immediately above it is called its *parent*. Any node(s) immediately below it are called  $u$ 's *children*. Likewise, all nodes on the path from  $u$  all the way back up to the root are known as  $u$ 's ancestors;  $u$ 's children and all nodes connected below them are called  $u$ 's descendants. Suppose we were to remove all nodes in a tree except  $u$  and  $u$ 's descendant(s). The new tree would form a *subtree* rooted at  $u$ . Other tree-related terminology will be used. For example, the root of the tree will be the only node without a parent. If a node has no children, it will be referred to as a *leaf* node.

## Binary Trees

We can place more useful restrictions on our tree. An oriented tree such that every node has *at most* two children is a *binary tree*. This means that every node in a binary tree may have 0 children (a leaf), 1 child, or 2 children. We can further horizontally orient the children in a binary tree. Since the number of children is restricted to at most 2, we can refer to them as a left child and right child. All the identified relations in a binary tree node are depicted in Figure 6.4. Since a node may only have 1 child, using this orientation means that it may have a left child and missing its right child, or it may have a right child and missing its left child.

A larger binary tree is depicted in Figure 6.5. In this tree, the nodes  $k, l, n, o, p, q, r$  are all leaf nodes. Many nodes have both a left and right child ( $b, g, j$  as examples) while some have a left child, but no right child ( $m, i$  for example) and some have a right child, but no left child ( $e, h$  for example).

Given this orientation, it is natural to define the *depth* of a node. Let  $u$  be a node in a binary tree with root  $r$ . The depth of  $u$  is the length of the unique path from  $r$  to  $u$ . The depth of the root,  $r$  itself is 0 by convention. The depth of a tree  $T$  itself is defined as the maximal depth (i.e. "deepest") of any node in the tree. Thus, if a tree consists of a single root node, its depth is zero.<sup>6</sup> It is possible to have an empty tree (where  $|V| = 0$ ) and so by convention and for consistency, the depth of an empty tree is usually

<sup>4</sup>Or dendrologists.

<sup>5</sup>Specifically in part LVIII. *On the Analytical Forms called Trees.*—Part II in which he essentially defines binary trees.

<sup>6</sup>The author proposes to call this a *seed*; let's hope it catches on.

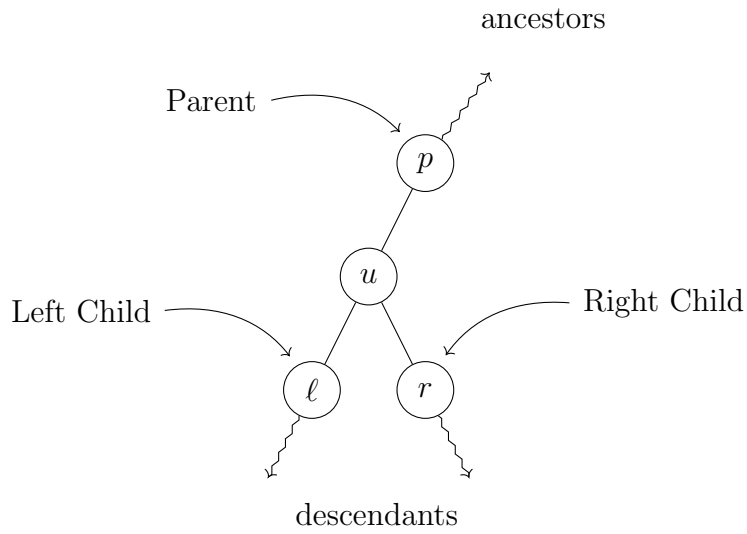


Figure 6.4: The various relations of a tree node. For a given node  $u$ , its parent, left and right child are labelled. In addition, nodes above  $u$  can collectively be referred to as ancestors (including its parent) and nodes in the subtree at  $u$  can collectively be referred to as descendants (including its immediate children).

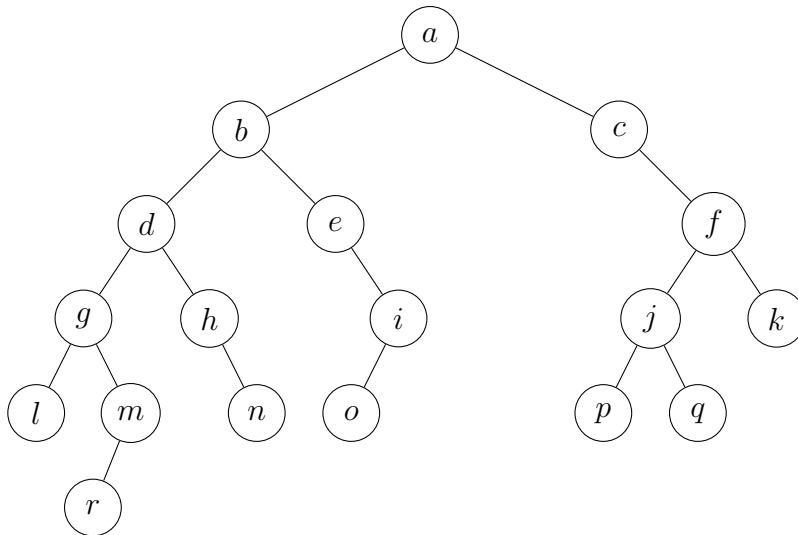
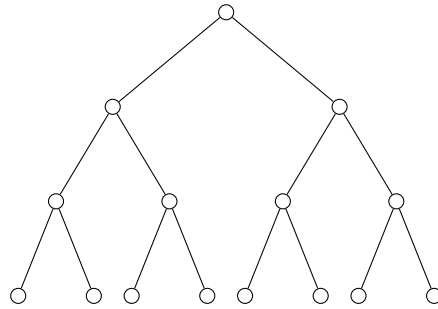


Figure 6.5: A Binary Tree

Table 6.1: Depths of nodes in Figure 6.5

Node(s)	Depth
$a$	0
$b, c$	1
$d, e, f$	2
$g, h, i, j, k$	3
$l, m, n, o, p, q$	4
$r$	5

Figure 6.6: A complete tree of depth  $d = 3$  which has  $1 + 2 + 4 + 8 = 15$  nodes.

–1. The depth of every node in the tree in Figure 6.5 is given in Table 6.1. Since  $r$  is the deepest node, the depth of the tree itself is 5.

When depicting trees, we draw all nodes of equal depth on the same horizontal line. All nodes on this line are said to be at the same *level* in the tree. Levels are numbered as the depth so the root is at level 0, the next depth nodes are at level 1, etc. Structurally we could draw the tree however we want, but placing nodes of a different depth at different levels makes the tree less readable.

Now, consider a *complete* (also called *full* or *perfect*) binary tree of depth  $d$ . That is, a binary tree in which all nodes are present at all levels; from level 0 up to level  $d$ , the depth of the tree. How many nodes total are there in such a tree? At level 0, there would only be the root. At level 1, there would be 2 children from the root. Since all nodes are present, at level 2, both children from the previous 2 nodes would be present giving a total of 4, etc. That is, at each level, the number of nodes *on that level* doubles from the previous level. A concrete example for  $d = 3$  is given in Figure 6.6 which has a total of 15 nodes.

We can generalize this by adding up powers of 2. At level 0, there are  $2^0$  nodes, level 1, there are  $2^1$  nodes, level 2 there are  $2^2$  nodes, etc. Summing all of these nodes up to

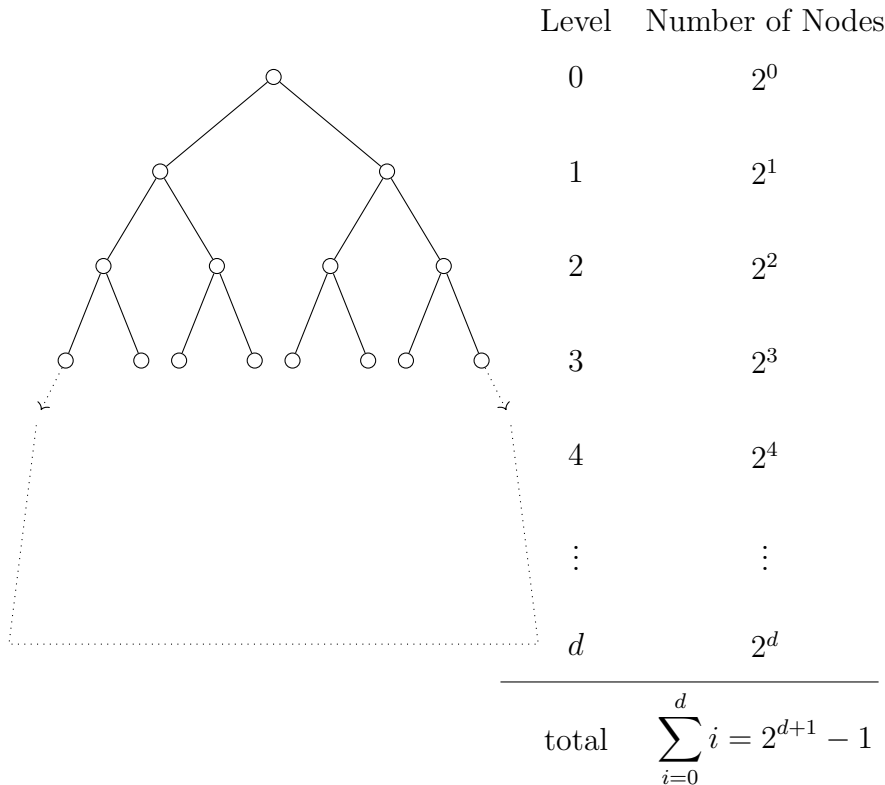


Figure 6.7: A summation of nodes at each level of a complete binary tree up to depth  $d$ .

depth  $d$  gives us the following geometric series which has a well-known solution.<sup>7</sup>

$$n = 2^0 + 2^1 + 2^2 + \dots + 2^{d-1} + 2^d = \sum_{k=0}^d 2^k = 2^{d+1} - 1$$

This is visualized in Figure 6.7.

Taken from another perspective, if we have a complete binary tree with  $n$  nodes, its depth is

$$d = \log(n + 1) - 1$$

That is, the depth is logarithmic,  $d \in O(\log n)$ . This observation gives us our motivation for using binary trees as a collection data structure. If we can develop a tree-based data structure such that insertion, retrieval, and removal operations are all proportional to

<sup>7</sup>As a Computer Scientist, you should also find it familiar—it represents the maximum number representable by  $d$  bits.

the depth  $d$  then we have the *potential* for a very efficient, generalized collection data structure. Of course it is not necessarily a simple matter as we will see shortly.

## 6.3 Implementation

A binary tree implementation is fairly straightforward and similar to that of a linked list. Each tree node will hold references to its two children, its parent (which is optional, but simplifies a lot of operations) and the key element stored in the tree node. In pseudocode we'll use the same dot operator syntax as with a linked list, so that each of these components can be referenced as follows. Let  $u$  be a tree node, then

- $u.key$  is its key,
- $u.parent$  is its parent,
- $u.leftChild$  is its left child, and
- $u.rightChild$  is its right child.

The tree itself only needs a reference to the root element, represented using  $T.root$ . Note that we do *not* keep a reference to every single tree node. Keeping track of and updating every tree node would bring us right back to a linked list or array-based list to store all the nodes, defeating the purpose of exploiting the tree structure. We could, however, keep track of how many nodes are in the tree as this could easily be updated with each insert/remove operation just as with a linked list.

As a concrete example, a pair of Java classes may be written to implement these two objects. In particular, the tree node class may look something like the following (getters, setters and other standard methods are omitted).

```

1 public class BinaryTreeNode<T> {
2
3     private T key;
4     private TreeNode<T> parent;
5     private TreeNode<T> leftChild;
6     private TreeNode<T> rightChild;
7
8     ...
9 }
```

While a tree itself would look something like the following.



```

1 public class BinaryTree<T> {
2
3     private TreeNode<T> root;
4     private int size;
5
6     ...
7 }

```

## 6.4 Tree Traversal

Before we develop the basic operations that will give us a collection data structure we need to ensure that we have a way to process all the elements in a general binary tree. With a linked list, we could easily iterate over the elements stored in it by starting at the head and iterating over each element until the tail. With a binary tree the order in which we iterate over elements is not all that clear. It is natural to start at the root node, but where do we go from there? To the left child or the right? If we go to the left, then we must remember to eventually come back and iterate over the right child.

More generally given a node  $u$  and its children,  $u.leftChild$ ,  $u.rightChild$  in which order do we enumerate them? In fact, there are several tree traversal strategies each with its own advantages and applications. Each one traverses the three different nodes in a different order but all are based on a more general graph traversal algorithm called **DFS** which attempts to explore a graph as deeply as possible before backtracking to explore the other areas. In general we will prefer to descend left in the tree before we backtrack and explore the right subtree. There is nothing particularly special about this preference (other than reading left-to-right is more natural). We *could* prefer to go right before going left. However, this would be equivalent to going left while also “reversing” the tree (reversing the left/right child at each node).

To understand the traversal strategies better, we will present the traversal choices locally. At any particular node,  $u$ , we’ll refer to this node as the root (think of the subtree rooted at  $u$ ) and its left and right child respectively. Given that we’ll always go left-before-right, this gives three possible orders in which to enumerate the three nodes:

- Preorder: root-left-right
- Inorder: left-root-right
- Postorder: left-right-root

Of course, any given node may be missing one or both of its children in which case we do not enumerate it, but otherwise, we will preserve these orderings.

When we describe a traversal strategy, it is not restricted to merely enumerating all the elements in a tree. The purpose of using a binary tree as a data structure is to store data in each of its nodes. When visiting a node, in general we can perform any operation

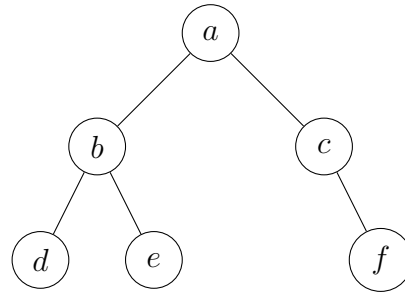


Figure 6.8: A small binary tree.

or process on the data stored in it rather than simply enumerating the element. The details of how to process a node would be specific to the type of data being stored and the overall goal of the algorithm being executed.

### 6.4.1 Preorder Traversal

A preorder traversal enumerates nodes in a root-left-right manner. This means that we process a node immediately when we first visit it, then after processing it, we next process its left child. Only after we have processed its left child and *all of its descendants*, do we return to process the right child and its descendants. However, we need a way to “remember” that the right child and its subtree still need to be traversed. As a small example, consider the tree in Figure 6.8. Starting at the root, we enumerate  $a$ . We next visit its left subtree rooted at  $b$ . When we visit  $b$ , we immediately enumerate it and again traverse left and enumerate its left child,  $d$ . Since  $d$  has no children, we backtrack to  $b$ . Having already enumerated  $b$  and its left child  $d$ , we next traverse to its right child,  $e$ . We enumerate  $e$  and since it has no children, we backtrack all the way back up to the root  $a$  and continue to its right child,  $c$ . We immediately enumerate  $c$ , but since it has no left child to visit, we traverse to its right child  $f$  as the final node. Altogether, the order of enumeration was

$$a, b, d, e, c, f$$

We can implement a preorder traversal algorithm using a stack to “remember” each right child that we need to visit and what order to visit them. Initially, we push the tree’s root onto the stack. We then go into a loop that executes until the stack is emptied and we have processed all the nodes in the tree. On each iteration, we pop the node on the top of the stack and process it. We then need to process its children before we continue with other elements on the stack. To set up the next iteration, we push these child nodes onto the stack so that they are the next ones to be processed. We need to take care to exploit the **LIFO** ordering of the stack properly, however. We want the left child to be processed before the right child, so we push the right child first, then the left child second. The full

push <i>a</i>	push <i>m</i>	(enter loop)	print <i>i</i>	push <i>q</i>
(enter loop)	push <i>l</i>	pop, <i>node = h</i>	(enter loop)	push <i>p</i>
pop, <i>node = a</i>	print <i>g</i>	push <i>n</i>	pop, <i>node = o</i>	print <i>j</i>
push <i>c</i>	(enter loop)	(no push)	(no push)	(enter loop)
push <i>b</i>	pop, <i>node = l</i>	print <i>h</i>	(no push)	pop, <i>node = p</i>
print <i>a</i>	(no push)	(enter loop)	print <i>o</i>	(no push)
(enter loop)	(no push)	pop, <i>node = n</i>	(enter loop)	(no push)
pop, <i>node = b</i>	print <i>l</i>	(no push)	pop, <i>node = c</i>	print <i>p</i>
push <i>e</i>	(enter loop)	(no push)	push <i>f</i>	(enter loop)
push <i>d</i>	pop, <i>node = m</i>	print <i>n</i>	(no push)	pop, <i>node = q</i>
print <i>b</i>	(no push)	(enter loop)	print <i>c</i>	(no push)
(enter loop)	push <i>r</i>	pop, <i>node = e</i>	(enter loop)	(no push)
pop, <i>node = d</i>	print <i>m</i>	push <i>i</i>	pop, <i>node = f</i>	print <i>q</i>
push <i>h</i>	(enter loop)	(no push)	push <i>k</i>	(enter loop)
push <i>g</i>	pop, <i>node = r</i>	print <i>e</i>	push <i>j</i>	pop, <i>node = k</i>
print <i>d</i>	(no push)	(enter loop)	print <i>f</i>	(no push)
(enter loop)	(no push)	pop, <i>node = i</i>	(enter loop)	(no push)
pop, <i>node = g</i>	print <i>r</i>	(no push)	pop, <i>node = j</i>	print <i>k</i>
		push <i>o</i>		

Figure 6.9: A walkthrough of a preorder traversal on the tree from Figure 6.5.

traversal is presented as Algorithm 18.

```

INPUT  : A binary tree, T
OUTPUT : A preorder traversal of the nodes in T
1 S ← empty stack
2 push T.root onto S
3 WHILE S is not empty DO
4   | u ← S.pop
5   | process u
6   | push u.rightChild onto S
7   | push u.leftChild onto S
8 END

```

**Algorithm 18:** Stack-based Preorder Tree Traversal. If a node does not exist, we implicitly do not push it onto the stack.

A full preorder traversal on the binary tree in Figure 6.5 would result in the following ordering of vertices.

*a, b, d, g, l, m, r, h, n, e, i, o, c, f, j, p, q, k*

A full walkthrough of the algorithm on this example is presented in Figure 6.9.

A common alternative way of presenting a preorder traversal is to use recursion. This alternative version is presented as Algorithm 19. We are still implicitly using a stack to keep track of where we've come from in the tree, but instead of using a stack data structure, we are exploiting the system call stack to do this.

```

INPUT  : A binary tree node  $u$ 
OUTPUT : A preorder traversal of the nodes in the subtree rooted at  $u$ 
1 IF  $u = null$  THEN
2   | return
3 ELSE
4   | process  $u$ 
5   | preOrderTraversal( $u.leftChild$ )
6   | preOrderTraversal( $u.rightChild$ )
7 END

```

**Algorithm 19:** preOrderTraversal( $u$ ): Recursive Preorder Tree Traversal

## Applications

A preorder traversal is a straightforward and arguably the simplest to implement tree traversal strategy. As such, it is the most common traversal strategy for many applications that do not require tree nodes to be traversed in any particular order. It can be used to build a tree, enumerate elements in a tree collection, copy a tree, etc.

TODO: more?

### 6.4.2 Inorder Traversal

An inorder traversal strategy visits nodes in a left-root-right order. This means that instead of immediately processing a node the first time we visit it, we wait until we have processed its left child and all of its descendants before processing it and moving on to the right subtree. From another we process a node when we backtrack to it (after having processed the left subtree).

Again, consider the small tree in Figure 6.8. We again start at the root  $a$ , but do not immediately enumerate it. Instead, we traverse to its left child,  $b$ . However, again we do not enumerate  $b$  until after we have enumerated its left subtree. We traverse to  $d$  and since there is no left child, we enumerate  $d$  and return to  $b$ . Now that we have enumerated  $b$ 's left subtree, we now enumerate  $b$  and traverse down to its right child  $e$  and enumerate it. We now backtrack all the way back up to  $a$ . Since we have enumerated its entire left subtree, we can now enumerate  $a$  and continue to its right subtree. Since  $c$  has no left child,  $c$  is enumerated next, and then  $f$ . Altogether, the order of enumeration

was

$$d, b, e, a, c, f$$

We can use the same basic idea of using a stack to keep track of tree nodes, however we want to delay processing the node until we've explored the left subtree. To do this, we need a way to tell if we are visiting the node for the first time or returning from exploring the left subtree (so that we may process it). To achieve this, we allow the node to be *null* as a flag to indicate that we are backtracking from the left subtree. Thus, if the current node is not null, we push it back onto the stack for later processing and explore the left subtree. If it is null, we pop the node on the top of the stack and process it, then push its right child to setup the next iteration. The full process is described in Algorithm 21.

```

INPUT  : A binary tree,  $T$ 
OUTPUT : An inorder traversal of the nodes in  $T$ 
1  $S \leftarrow$  empty stack
2  $u \leftarrow T.root$ 
3 WHILE  $S$  is not empty OR  $u \neq null$  DO
4   IF  $u \neq null$  THEN
5     push  $u$  onto  $S$ 
6      $u \leftarrow u.leftChild$ 
7   ELSE
8      $u \leftarrow S.pop$ 
9     process  $u$ 
10     $u \leftarrow u.rightChild$ 
11  END
12 END

```

**Algorithm 20:** Stack-based Inorder Tree Traversal

In lines 6 and 10, if no such child exists,  $u$  becomes null, indicating that we are backtracking on the next iteration.

A full inorder traversal on the binary tree in Figure 6.5 would result in the following order of vertices.

$$l, g, r, m, d, h, n, b, e, o, i, a, c, p, j, q, f, k$$

Again, a full walkthrough of the algorithm on this example is presented in Figure 6.10.

As with a preorder traversal, we can use recursion to implicitly use the call stack to keep track of the ordering. The only difference is when we process the given node. In the preorder case, we did it before the two recursive calls. In the inorder case, we process it

## 6 Trees

(enter loop, $u = a$ )	(enter loop, $u = \text{null}$ )	(enter loop, $u = \text{null}$ )	$u = \text{null}$ pop $i$ , update $u = i$	(enter loop, $u = \text{null}$ )
push $a$	pop $r$ , update	pop $b$ , update	process $i$	pop $j$ , update
update $u = b$	$u = r$	$u = b$	update $u = \text{null}$	$u = j$
(enter loop, $u = b$ )	process $r$	process $b$	(enter loop, $u = \text{null}$ )	process $j$
push $b$	update $u = \text{null}$	update $u = e$	pop $a$ , update	update $u = q$
update $u = d$	(enter loop, $u = \text{null}$ )	(enter loop, $u = e$ )	$u = a$	(enter loop, $u = q$ )
(enter loop, $u = d$ )	pop $m$ , update	push $e$	process $a$	push $q$
push $d$	$u = m$	update $u = \text{null}$	update $u = c$	update $u = \text{null}$
update $u = g$	process $m$	(enter loop, $u = \text{null}$ )	(enter loop, $u = c$ )	(enter loop, $u = \text{null}$ )
(enter loop, $u = g$ )	update $u = \text{null}$	pop $e$ , update	push $c$	pop $q$ , update
push $g$	(enter loop, $u = \text{null}$ )	$u = e$	update $u = \text{null}$	$u = q$
update $u = l$	pop $d$ , update	process $e$	(enter loop, $u = \text{null}$ )	process $q$
(enter loop, $u = l$ )	$u = d$	update $u = i$	pop $c$ , update	update $u = \text{null}$
push $l$	process $d$	(enter loop, $u = i$ )	$u = c$	(enter loop, $u = \text{null}$ )
update $u = \text{null}$	update $u = h$	push $i$	process $c$	pop $f$ , update
(enter loop, $u = \text{null}$ )	(enter loop, $u = h$ )	update $u = o$	update $u = f$	$u = f$
pop $l$ , update	push $h$	(enter loop, $u = o$ )	(enter loop, $u = f$ )	process $f$
$u = l$	update $u = \text{null}$	push $o$	push $f$	update $u = k$
process $l$	(enter loop, $u = \text{null}$ )	update $u = \text{null}$	update $u = j$	(enter loop, $u = k$ )
update $u = \text{null}$	pop $h$ , update	(enter loop, $u = \text{null}$ )	(enter loop, $u = j$ )	push $k$
(enter loop, $u = \text{null}$ )	$u = h$	pop $o$ , update	push $j$	update $u = \text{null}$
pop $g$ , update	process $h$	$u = o$	update $u = p$	(enter loop, $u = \text{null}$ )
$u = g$	update $u = n$	process $o$	(enter loop, $u = p$ )	pop $k$ , update
process $g$	(enter loop, $u = n$ )	update $u = \text{null}$	push $p$	$u = k$
update $u = m$	push $n$	(enter loop, $u = \text{null}$ )	update $u = \text{null}$	process $k$
(enter loop, $u = m$ )	update $u = \text{null}$	pop $i$ , update	(enter loop, $u = \text{null}$ )	update $u = \text{null}$
push $m$	(enter loop, $u = \text{null}$ )	$u = i$	(done)	
update $u = r$	pop $n$ , update	process $i$	pop $p$ , update	
(enter loop, $u = r$ )	$u = n$	update $u = \text{null}$	$u = p$	
push $r$	process $n$	(enter loop, $u = \text{null}$ )	process $p$	
update $u = \text{null}$	update $u = \text{null}$		update $u = \text{null}$	

Figure 6.10: A walkthrough of a inorder traversal on the tree from Figure 6.5.

between them. The recursive version is presented as Algorithm 21.

```

INPUT   : A binary tree node  $u$ 
OUTPUT  : An inorder traversal of the nodes in the subtree rooted at  $u$ 
1 IF  $u = null$  THEN
2   | return
3 ELSE
4   | inOrderTraversal( $u.leftChild$ )
5   | process  $u$ 
6   | inOrderTraversal( $u.rightChild$ )
7 END

```

**Algorithm 21:** inOrderTraversal( $u$ ): Recursive Inorder Tree Traversal

### Applications

The primary application of an inorder tree traversal is from where it derives its name. If we perform an inorder traversal on a binary search tree (see the next section), then we get an enumeration of elements that is sorted or “in order.”

TODO: more?

### 6.4.3 Postorder Traversal

The final depth-first-search based traversal strategy is a postorder traversal in which nodes are visited in a left-right-root manner. That is, we hold off on processing a node until both of its children and all of their respective descendants have been processed.

We turn once again to the small binary tree in Figure 6.8. We start at the root  $a$  and descend all the way to the left and process  $d$  since it has no children. Backtracking to  $b$  we again hold off on processing it until we’ve traversed its right child,  $e$ . Only after we’ve enumerated both  $d$  and  $e$  do we finally enumerate  $b$ . At this point,  $a$ ’s left subtree has been enumerated, but again we hold off on processing  $a$  until we are done with its right subtree.  $c$  has no left child, so we enumerate its right child,  $f$  first and only then do we process  $c$ . Finally, the last node to be processed is the root itself,  $a$ .

We can again utilize a stack data structure to perform a postorder traversal, but it gets a bit more complicated as we need to distinguish 3 different cases. When visiting a node, we need to know if it is the first time we’ve traversed it, the second time (returning from its left subtree) or the third time (returning from its right subtree, thus it needs to be processed). To do this, we keep track of not only the current node but also the previous node to know *where* we came from, either the parent, the left child, or the right child.

That way, we can tell if it is the first time we've visited the node (the previous node would be the parent), the second time (the previous node is the left child) or the third time (the previous node is the right child). The full postorder traversal is presented in Algorithm 22.

```

INPUT   : A binary tree,  $T$ 
OUTPUT  : A postorder traversal of the nodes in  $T$ 
1  $S \leftarrow$  empty stack
2  $prev \leftarrow null$ 
3 push  $T.root$  onto  $S$ 
4 WHILE  $S$  is not empty DO
5    $curr \leftarrow S.peek$ 
6   IF  $prev = null$  OR  $prev.leftChild = curr$  OR  $prev.rightChild = curr$  THEN
7     IF  $curr.leftChild \neq null$  THEN
8       | push  $curr.leftChild$  onto  $S$ 
9     ELSE IF  $curr.rightChild \neq null$  THEN
10      | push  $curr.rightChild$  onto  $S$ 
11      END
12    ELSE IF  $curr.leftChild = prev$  THEN
13      | IF  $curr.rightChild \neq null$  THEN
14        | push  $curr.rightChild$  onto  $S$ 
15        END
16    ELSE
17      | process  $curr$ 
18      |  $S.pop$ 
19    END
20     $prev \leftarrow curr$ 
21 END

```

**Algorithm 22:** Stack-based Postorder Tree Traversal

A full postorder traversal on the binary tree in Figure 6.5 would result in the following order of vertices.

$l, r, m, g, n, h, d, o, i, e, b, p, q, j, k, f, c, a$

Again, a full walkthrough of the algorithm on this example is presented in Figure 6.11.

$prev = null$	push (b)	$prev.leftChild = curr$	update $curr = (d)$
push a	update $prev = a$	$((b).leftChild = (b))$	check:
(enter loop)	(enter loop)	push (d)	$prev.leftChild = curr$
update $curr = (a)$	update $curr = (b)$	update $prev = b$	$((d).leftChild = (d))$
check: $prev = null$	check:	(enter loop)	push (g)
			update $prev = d$



(enter loop)	(enter loop)	(enter loop)	(enter loop)
update <i>curr</i> = ( <i>g</i> )	update <i>curr</i> = ( <i>r</i> )	update <i>curr</i> = ( <i>n</i> )	update <i>curr</i> = ( <i>i</i> )
check:	check:	check:	check:
<i>prev.leftChild</i> = <i>curr</i>	<i>prev.rightChild</i> = <i>curr</i>	<i>prev.rightChild</i> = <i>curr</i>	<i>prev.rightChild</i> = <i>curr</i>
(( <i>g</i> ). <i>leftChild</i> = ( <i>g</i> ))	( <i>null.rightChild</i> = ( <i>r</i> ))	(( <i>n</i> ). <i>rightChild</i> = ( <i>n</i> ))	(( <i>i</i> ). <i>rightChild</i> = ( <i>i</i> ))
push ( <i>l</i> )	process <i>r</i>	(noop)	push ( <i>o</i> )
update <i>prev</i> = <i>g</i>	update <i>prev</i> = <i>r</i>	update <i>prev</i> = <i>n</i>	update <i>prev</i> = <i>i</i>
(enter loop)	(enter loop)	(enter loop)	(enter loop)
update <i>curr</i> = ( <i>l</i> )	update <i>curr</i> = ( <i>m</i> )	update <i>curr</i> = ( <i>n</i> )	update <i>curr</i> = ( <i>o</i> )
check:	check:	check:	check:
<i>prev.leftChild</i> = <i>curr</i>	<i>prev.rightChild</i> = <i>curr</i>	<i>prev.rightChild</i> = <i>curr</i>	<i>prev.leftChild</i> = <i>curr</i>
(( <i>l</i> ). <i>leftChild</i> = ( <i>l</i> ))	( <i>null.rightChild</i> = ( <i>m</i> ))	( <i>null.rightChild</i> = ( <i>n</i> ))	(( <i>o</i> ). <i>leftChild</i> = ( <i>o</i> ))
(noop)	check:	process <i>n</i>	(noop)
update <i>prev</i> = <i>l</i>	<i>curr.leftChild</i> = <i>prev</i>	update <i>prev</i> = <i>n</i>	update <i>prev</i> = <i>o</i>
(enter loop)	(( <i>r</i> ). <i>leftChild</i> = ( <i>r</i> ))	(enter loop)	(enter loop)
update <i>curr</i> = ( <i>l</i> )	update <i>prev</i> = <i>m</i>	update <i>curr</i> = ( <i>h</i> )	update <i>curr</i> = ( <i>o</i> )
check:	(enter loop)	check:	check:
<i>prev.rightChild</i> = <i>curr</i>	update <i>curr</i> = ( <i>m</i> )	<i>prev.rightChild</i> = <i>curr</i>	<i>prev.rightChild</i> = <i>curr</i>
( <i>null.rightChild</i> = ( <i>l</i> ))	check:	( <i>null.rightChild</i> = ( <i>h</i> ))	( <i>null.rightChild</i> = ( <i>o</i> ))
process <i>l</i>	<i>prev.rightChild</i> = <i>curr</i>	process <i>h</i>	process <i>o</i>
update <i>prev</i> = <i>l</i>	( <i>null.rightChild</i> = ( <i>m</i> ))	update <i>prev</i> = <i>h</i>	update <i>prev</i> = <i>o</i>
(enter loop)	process <i>m</i>	(enter loop)	(enter loop)
update <i>curr</i> = ( <i>g</i> )	update <i>prev</i> = <i>m</i>	update <i>curr</i> = ( <i>d</i> )	update <i>curr</i> = ( <i>i</i> )
check:	(enter loop)	check:	check:
<i>prev.rightChild</i> = <i>curr</i>	update <i>curr</i> = ( <i>g</i> )	<i>prev.rightChild</i> = <i>curr</i>	<i>prev.rightChild</i> = <i>curr</i>
( <i>null.rightChild</i> = ( <i>g</i> ))	check:	(( <i>n</i> ). <i>rightChild</i> = ( <i>d</i> ))	( <i>null.rightChild</i> = ( <i>i</i> ))
check:	<i>prev.rightChild</i> = <i>curr</i>	process <i>d</i>	check:
<i>curr.leftChild</i> = <i>prev</i>	( <i>null.rightChild</i> = ( <i>g</i> ))	update <i>prev</i> = <i>d</i>	<i>curr.leftChild</i> = <i>prev</i>
(( <i>l</i> ). <i>leftChild</i> = ( <i>l</i> ))	process <i>g</i>	(enter loop)	(( <i>o</i> ). <i>leftChild</i> = ( <i>o</i> ))
push ( <i>m</i> )	update <i>prev</i> = <i>g</i>	update <i>curr</i> = ( <i>b</i> )	update <i>prev</i> = <i>i</i>
update <i>prev</i> = <i>g</i>	(enter loop)	check:	(enter loop)
(enter loop)	update <i>curr</i> = ( <i>d</i> )	<i>prev.rightChild</i> = <i>curr</i>	update <i>curr</i> = ( <i>i</i> )
update <i>curr</i> = ( <i>m</i> )	check:	(( <i>h</i> ). <i>rightChild</i> = ( <i>b</i> ))	check:
check:	<i>prev.rightChild</i> = <i>curr</i>	check:	<i>prev.rightChild</i> = <i>curr</i>
<i>prev.rightChild</i> = <i>curr</i>	(( <i>m</i> ). <i>rightChild</i> = ( <i>d</i> ))	<i>curr.leftChild</i> = <i>prev</i>	( <i>null.rightChild</i> = ( <i>i</i> ))
(( <i>m</i> ). <i>rightChild</i> = ( <i>m</i> ))	check:	(( <i>d</i> ). <i>leftChild</i> = ( <i>d</i> ))	process <i>i</i>
push ( <i>r</i> )	<i>curr.leftChild</i> = <i>prev</i>	push ( <i>e</i> )	update <i>prev</i> = <i>i</i>
update <i>prev</i> = <i>m</i>	(( <i>g</i> ). <i>leftChild</i> = ( <i>g</i> ))	update <i>prev</i> = <i>b</i>	(enter loop)
(enter loop)	push ( <i>h</i> )	(enter loop)	update <i>curr</i> = ( <i>e</i> )
update <i>curr</i> = ( <i>r</i> )	update <i>prev</i> = <i>d</i>	update <i>curr</i> = ( <i>e</i> )	check:
check:	(enter loop)	check:	<i>prev.rightChild</i> = <i>curr</i>
<i>prev.leftChild</i> = <i>curr</i>	update <i>curr</i> = ( <i>h</i> )	<i>prev.rightChild</i> = <i>curr</i>	( <i>null.rightChild</i> = ( <i>e</i> ))
(( <i>r</i> ). <i>leftChild</i> = ( <i>r</i> ))	check:	(( <i>e</i> ). <i>rightChild</i> = ( <i>e</i> ))	process <i>e</i>
(noop)	<i>prev.rightChild</i> = <i>curr</i>	push ( <i>i</i> )	update <i>prev</i> = <i>e</i>
update <i>prev</i> = <i>r</i>	(( <i>h</i> ). <i>rightChild</i> = ( <i>h</i> ))	update <i>prev</i> = <i>e</i>	(enter loop)
(enter loop)	push ( <i>n</i> )	(enter loop)	update <i>curr</i> = ( <i>b</i> )
update <i>curr</i> = ( <i>r</i> )	update <i>prev</i> = <i>h</i>	update <i>curr</i> = ( <i>b</i> )	

check:	<i>prev.leftChild = curr</i>	<i>((q).rightChild = (q))</i>	(noop)
<i>prev.rightChild = curr</i>	<i>((j).leftChild = (j))</i>	(noop)	update <i>prev = k</i>
<i>((i).rightChild = (b))</i>	push ( <i>p</i> )	update <i>prev = q</i>	(enter loop)
process <i>b</i>	update <i>prev = j</i>	(enter loop)	update <i>curr = (k)</i>
update <i>prev = b</i>	(enter loop)	update <i>curr = (q)</i>	check:
(enter loop)	update <i>curr = (p)</i>	check:	<i>prev.rightChild = curr</i>
update <i>curr = (a)</i>	check:	<i>prev.rightChild = curr</i>	<i>(null.rightChild = (k))</i>
check:	<i>prev.leftChild = curr</i>	<i>(null.rightChild = (q))</i>	process <i>k</i>
<i>prev.rightChild = curr</i>	<i>((p).leftChild = (p))</i>	process <i>q</i>	update <i>prev = k</i>
<i>((e).rightChild = (a))</i>	(noop)	update <i>prev = q</i>	(enter loop)
check:	update <i>prev = p</i>	(enter loop)	update <i>curr = (f)</i>
<i>curr.leftChild = prev</i>	(enter loop)	update <i>curr = (j)</i>	check:
<i>((b).leftChild = (b))</i>	(enter loop)	check:	<i>prev.rightChild = curr</i>
push ( <i>c</i> )	update <i>curr = (p)</i>	<i>prev.rightChild = curr</i>	<i>(null.rightChild = (f))</i>
update <i>prev = a</i>	check:	<i>(null.rightChild = (j))</i>	process <i>f</i>
(enter loop)	<i>prev.rightChild = curr</i>	process <i>j</i>	update <i>prev = f</i>
update <i>curr = (c)</i>	<i>(null.rightChild = (p))</i>	update <i>prev = j</i>	(enter loop)
check:	process <i>p</i>	(enter loop)	update <i>curr = (c)</i>
<i>prev.rightChild = curr</i>	update <i>prev = p</i>	update <i>curr = (f)</i>	check:
<i>((c).rightChild = (c))</i>	(enter loop)	check:	<i>prev.rightChild = curr</i>
push ( <i>f</i> )	update <i>curr = (j)</i>	<i>prev.rightChild = curr</i>	<i>((k).rightChild = (c))</i>
update <i>prev = c</i>	check:	<i>((q).rightChild = (f))</i>	process <i>c</i>
(enter loop)	<i>prev.rightChild = curr</i>	check:	update <i>prev = c</i>
update <i>curr = (f)</i>	<i>(null.rightChild = (j))</i>	<i>curr.leftChild = prev</i>	(enter loop)
check:	check:	<i>((j).leftChild = (j))</i>	update <i>curr = (a)</i>
<i>prev.rightChild = curr</i>	<i>curr.leftChild = prev</i>	push ( <i>k</i> )	check:
<i>((f).rightChild = (f))</i>	<i>((p).leftChild = (p))</i>	update <i>prev = f</i>	<i>prev.rightChild = curr</i>
push ( <i>j</i> )	push ( <i>q</i> )	(enter loop)	<i>((f).rightChild = (a))</i>
update <i>prev = f</i>	update <i>prev = j</i>	update <i>curr = (k)</i>	process <i>a</i>
(enter loop)	(enter loop)	check:	update <i>prev = a</i>
update <i>curr = (j)</i>	update <i>curr = (q)</i>	<i>prev.rightChild = curr</i>	
check:	check:	<i>((k).rightChild = (k))</i>	
	<i>prev.rightChild = curr</i>		

Figure 6.11: A walkthrough of a postorder traversal on the tree from Figure 6.5.

A recursive version would follow the same pattern as before. With a postorder traversal, we would process the node after *both* of the recursive calls. A recursive postorder traversal

is presented in Algorithm 23.

```

INPUT   : A binary tree node  $u$ 
OUTPUT  : A postorder traversal of the nodes in the subtree rooted at  $u$ 
1 IF  $u = null$  THEN
2   | return
3 ELSE
4   | postOrderTraversal( $u.leftChild$ )
5   | postOrderTraversal( $u.rightChild$ )
6   | process  $u$ 
7 END

```

**Algorithm 23:** postOrderTraversal( $u$ ): Recursive Postorder Tree Traversal

## Applications

TODO: more?

- Topological sorting
- Destroying a tree when manual memory management is necessary (roots are the last thing that get cleaned up)
- Reverse polish notation (operand-operand-operator, unambiguous, used in old HP calculators)
- PostScript (Page Description Language)

### 6.4.4 Tree Walk Traversal

It turns out that we can perform any of the three depth-first-search based tree traversals without using any extra space at all (that is, a stack is not necessary) by exploiting the oriented structure of a binary tree. As with the postorder traversal, we simply need to keep track of where we came from (a previous node) and where we are (a current node) to determine where to go next. By only keeping track of two variables, we can also determine when a node should be processed. This traversal is generally called a “tree walk” and it is illustrated on the small tree example we’ve been using in Figure 6.12. The traversal is essentially a “walk” around the perimeter of the tree.

The rules that we follow are quite simple and only require *local* information. In particular, there are only three cases that we need to distinguish. Suppose we are at a node  $u$  as our current node with parent  $p$ , and left/right child as  $\ell$ ,  $r$  respectively as in Figure 6.13. The rules are outlined in Table 6.2.

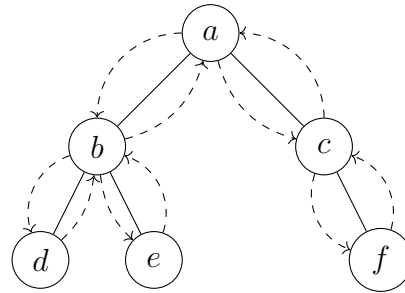


Figure 6.12: A tree walk on the tree from Figure 6.8.

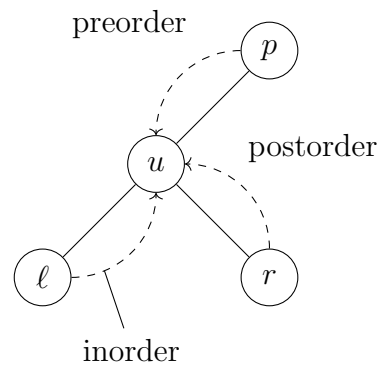


Figure 6.13: The three general cases of when to process a node in a tree walk.

Table 6.2: Rules for Tree Walk

If previous is...	Then traverse to...	And process $u$ if order is...
$p$	$l$	Preorder
$l$	$r$	Inorder
$r$	$p$	Postorder

The full tree walk is presented as Algorithm 24 and includes all three traversal strategies. The algorithm actually only keeps track of the previous node's *type* (whether it was a parent, left, or right child) rather than the node itself. It is a lot more complex than the three simple rules in Table 6.2 because we need to take care of many corner cases where the left and/or right children may not exist.

### 6.4.5 Breadth-First Search Traversal

An alternative to the depth-first-search based traversal strategies, we can explore a binary tree using what is known as a **Breadth First Search (BFS)** traversal. **BFS** is also a general graph traversal algorithm that explores a graph by visiting the closest vertices first before

```

INPUT  : A binary tree,  $T$ 
OUTPUT : A Tree Walk around  $T$ 
1  $curr \leftarrow T.root$ 
2  $prevType \leftarrow parent$ 
3 WHILE  $curr \neq null$  DO
4   IF  $prevType = parent$  THEN
5     //preorder: process  $curr$  here
6     IF  $curr.leftChild$  exists THEN
7       //Go to the left child:
8        $curr \leftarrow curr.leftChild$ 
9        $prevType \leftarrow parent$ 
10    ELSE
11      $curr \leftarrow curr$ 
12      $prevType \leftarrow left$ 
13  ELSE IF  $prevType = left$  THEN
14    //inorder: process  $curr$  here
15    IF  $curr.rightChild$  exists THEN
16     //Go to the right child:
17      $curr \leftarrow curr.rightChild$ 
18      $prevType \leftarrow parent$ 
19    ELSE
20      $curr \leftarrow curr$ 
21      $prevType \leftarrow right$ 
22  ELSE IF  $prevType = right$  THEN
23    //postorder: process  $curr$  here
24    IF  $curr.parent = null$  THEN
25     //root has no parent, we're done traversing
26      $curr \leftarrow curr.parent$ 
27    //are we at the parent's left or right child?
28    ELSE IF  $curr = curr.parent.leftChild$  THEN
29      $curr \leftarrow curr.parent$ 
30      $prevType \leftarrow left$ 
31    ELSE
32      $curr \leftarrow curr.parent$ 
33      $prevType \leftarrow right$ 
34  END

```

**Algorithm 24:** Tree Walk based Tree Traversal

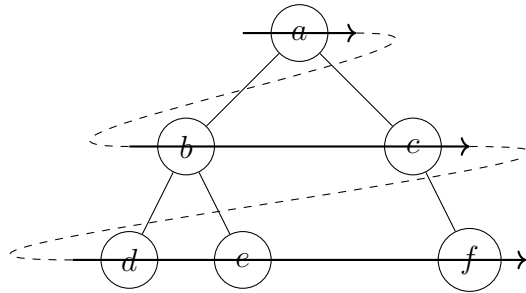


Figure 6.14: A Breadth First Search Example

venturing deeper into the graph. When applied to an oriented binary tree, **BFS** ends up exploring a graph in a top-to-bottom, left-to-right ordering. In this manner, all nodes at the same depth are enumerated before moving on to the next deepest level. The basic enumeration is depicted in Figure 6.14 which represents a BFS traversal on the same graph we've been using in previous examples. Unsurprisingly, the order of enumeration matches the labeling which was chosen to be top-bottom/left-right for readability.

The visualization of the ordering should give a hint as to the implementation of BFS. If we were to stretch out the ordering it would be clear that the nodes are enumerated in one straight long line, suggesting a queue should be used. Indeed, we start by enqueueing the root node. Then on each iteration, we dequeue the next node and enqueue its children to be processed later. Since a queue is **FIFO**, we enqueue in the left child first. This is presented as Algorithm 25.

```

INPUT  : A binary tree,  $T$ 
OUTPUT : A BFS traversal of the nodes in  $T$ 
1  $Q \leftarrow$  empty queue
2 enqueue  $T.root$  into  $Q$ 
3 WHILE  $Q$  is not empty DO
4    $u \leftarrow Q.dequeue$ 
5   process  $u$ 
6   enqueue  $u.leftChild$  into  $Q$ 
7   enqueue  $u.rightChild$  into  $Q$ 
8 END

```

**Algorithm 25:** Queue-based BFS Tree Traversal

Contrast the BFS algorithm with the original stack-based preorder algorithm (Algorithm 18). In fact, they are nearly identical. The only difference is that we use a queue instead of a stack (and its corresponding operations) and reverse the operations on the left/right child. This calls back to one of the major themes of this text on the importance of “smart” data structures. Simply by changing the underlying data structure used in an

algorithm, it results in a fundamentally different but complementarily useful property.

## 6.5 Binary Search Trees

Binary trees have a lot of structure but there is still not enough to make them efficient data structures. As we've presented them so far, we can store elements and retrieve them, but without any additional ordered structure, using any one of the general traversal strategies we've developed so far is still going to be  $O(n)$  in the worst case. We will now add some additional structure to our binary trees in order to attempt to make the general operations of insertion, retrieval and deletion more efficient. In particular, it will be our goal to make all three of these operations have complexity that is proportional to the depth of the tree,  $O(d)$ . To do that, we introduce [Binary Search Trees \(BSTs\)](#).

**Definition 9** (Binary Search Tree). A *binary search tree* is a binary tree such that every node  $u$  has an associated key,  $u.key$  which satisfies the *binary search tree property*:

1. Every node in the left subtree of  $u$  has a key value *less* than  $u.key$
2. Every node in the right subtree of  $u$  has a key value *greater* than  $u.key$

Implicitly, this definition does not allow for duplicate key values in a binary tree. In general, we could amend this definition to allow duplicate keys but we would need to be consistent on if we place duplicates in the left or right subtree. In any case, it is not that big of a deal. In general, we can “break ties” using some unique value to each object to induce a *total ordering* on all elements stored in a binary search tree.<sup>8</sup>

Furthermore, though we will use integer key values in all of our examples, in practice binary search trees need not be restricted to such values as they are intended to hold any type of object or data that we wish to store. Typically, an implementation will use a *hash value* of an object or data as a key value. A hash value is simply the result of applying a hash function that maps an object to an integer value. This is such a common operation it is built into many programming languages.

The binary search tree property is an *inductive* property. Since it holds for all nodes, it follows that every rooted subtree in a [BST](#) is also a binary search tree. A full example can be found in [Figure 6.15](#). Note that a binary search tree need not be full or complete as in the example. Many children (and thus entire subtrees) are not present. The binary search tree property merely guarantees that all keys in the left subtree are less and those in the right subtree are greater than the key stored in the root respectively.

---

<sup>8</sup>In practice, all things being equal, we could use the memory address of two objects stored in a tree to break ties.

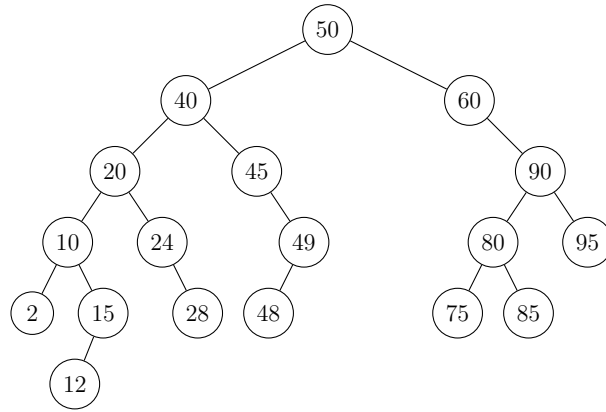


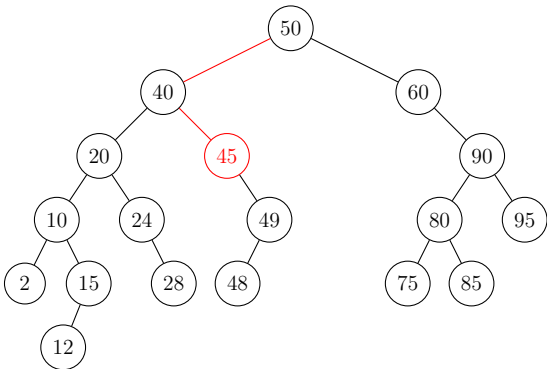
Figure 6.15: A Binary Search Tree

### 6.5.1 Retrieval

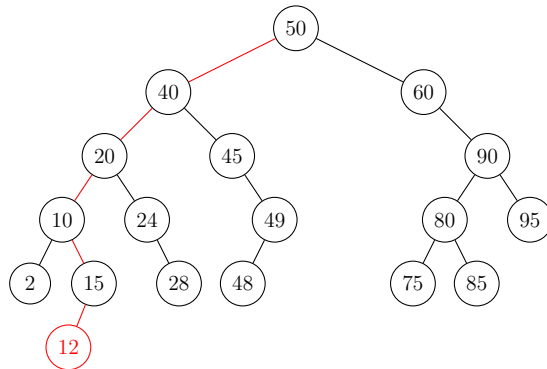
We can exploit the binary search tree property similar to how binary search exploits a sorted array. With binary search, we check the middle element and either find what we were search for or we have effectively cut the array in half as we know that we should either search in the lower half or the upper half. Similarly, if with a binary search tree, we start a search for a particular key  $k$  at the root. If the root's key value is greater than  $k$  ( $k < T.root.key$ ) then we know that the key lies in its left subtree while if the search key value is greater ( $T.root.key < k$ ) then we know that it must lie in the right subtree. In either case, we continue the search on the tree rooted at the left or right child respectively. If we ever find a node such that the key value matches we can stop the search and return the node's value. Of course, not every search will necessarily be successful. If we search for a key value that does not exist in the tree, we will eventually reach the end of the tree at some leaf node. If this happens, we can stop the process and return a flag indicating an unsuccessful search. Several examples are illustrated in Figure 6.16.

In some of the examples we got lucky and did not have to search too deeply in the tree. In fact, if we had searched for the key value  $k = 50$  only a single key comparison would have been required. Thus, in the best case only  $O(1)$  comparisons are necessary for a search/retrieval operation. In the worst case, however, we had to search to the deepest part of the tree as in the case where we successfully searched for  $k = 12$  and unsuccessfully searched for  $k = 13$ . In both cases we made 6 key comparisons before ending the search. In general, we may need to make  $d$  comparisons where  $d$  is the depth of the tree. This achieves what we set out to do as a retrieval operation is  $O(d)$ . The full

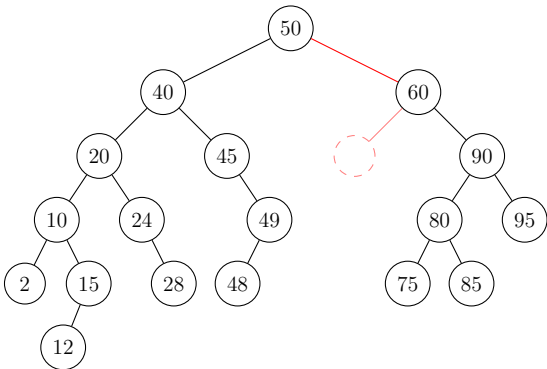




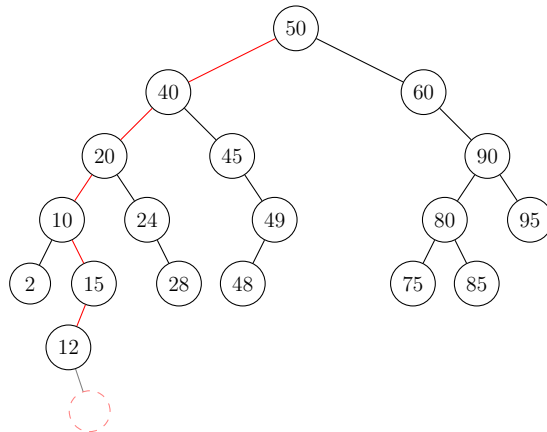
(a) Searching a BST Example 1. In this example, we search for the key value  $k = 45$ . The search starts at the root. Since  $45 < 50$  we traverse to the left child. The next comparison finds that  $40 < 45$  so we traverse right. Finally, the last comparison finds that the key value matches at which point the search process terminates.



(b) Searching a BST Example 2. In this example, we search for the key value  $k = 12$ . The search starts at the root. Since  $12 < 50$  we traverse to the left child. The search then traverses left twice more. When compared to 10, it traverses right and finally left, finding the key value at the deepest leaf in the tree. A total of 6 comparisons was necessary.



(c) Searching a BST Example 3. In this example, we search for the key value  $k = 55$  which will result in an unsuccessful search. Starting at the root, we traverse right, then left at which point we are no longer in the tree. With only 2 comparisons, we are able to conclude and return a flag value indicating an unsuccessful search.



(d) Searching a BST Example 4. Another unsuccessful search when we search for the key value  $k = 13$ . Similar to example 6.16(b). With the same 6 comparisons, we've traversed off the tree and return our flag value.

Figure 6.16: Various Search Examples on a Binary Search Tree

search process is presented as Algorithm 26.

```

INPUT  : A binary search tree,  $T$ , a key  $k$ 
OUTPUT : The tree node  $u \in T$  such that  $u.key = k$ ,  $\phi$  if no such node exists.
1  $u \leftarrow T.root$ 
2 WHILE  $u \neq \phi$  DO
3   IF  $u.key = k$  THEN
4     | output  $u$ 
5   ELSE IF  $u.key > k$  THEN
6     |  $u \leftarrow u.leftChild$ 
7   ELSE IF  $u.key < k$  THEN
8     |  $u \leftarrow u.rightChild$ 
9   END
10 END
11 output  $\phi$ 

```

**Algorithm 26:** Search algorithm for a binary search tree

### 6.5.2 Insertion

Inserting a new node into a binary search tree is a simple extension to searching for a node. In fact, it is exactly the same process. Since we will not allow duplicate key values, if we perform a search and find that a node with the given key value already exists, we can make the same design decisions we used with lists (throw an exception, return an error flag, or treat it as a no op). Otherwise, we perform the same search and when we reach the end of the tree, we insert the node at the spot where we would otherwise expect it to be.

This is illustrated in Figures 6.16(c) and 6.16(d) where a dashed “phantom” node is indicated where the key values would otherwise have been. When performing a search we will always be guaranteed to be inserting the new node as a leaf node. We’ll never insert a new node in the middle of the tree so no other considerations are necessary. Once we have performed the search, creating a new node and inserting it only requires a few reference shuffles. Thus, using the same analysis as before, in the worst case, insertion is also  $O(d)$ . Developing the pseudocode for the insertion operation is left as an exercise (see Exercise 6.12). A small example is demonstrated in Figure 6.17.

### 6.5.3 Removal

The removal of keys also follows a similar initial search operation. Given a key  $k$  to delete, we traverse the binary search tree for the node containing  $k$  which is an  $O(d)$

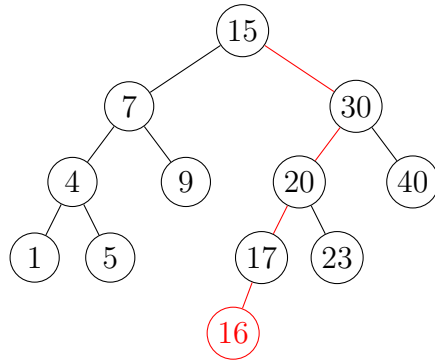


Figure 6.17: Insertion of a new node into a binary search tree. A search is performed for the key  $k = 16$ , when the open slot is found, the node is inserted.

operation. Suppose the node we wish to delete is  $u$ . There are several cases to take care of depending on how many children  $u$  has.

**Case 1:** Suppose that  $u$  is a leaf node and has no children. It is a simple matter to delete it by simply changing its parent node's reference. Though we do need to check if  $u$  is the left or right child. In particular,

- If  $u = u.parent.leftChild$  then set  $u.parent.leftChild \leftarrow \phi$ ; otherwise
- If  $u = u.parent.rightChild$  then set  $u.parent.rightChild \leftarrow \phi$ .

**Case 2:** Suppose that  $u$  has only one child (left or right) and is missing the other child. We don't want to "prune" the subtree of its child, but at the same time we want to preserve the binary search tree property. To do this we simply "promote" the only child up to  $u$ 's position. Specifically, suppose that  $u$ 's child is a left child. The steps to promote it are nearly exactly the same as used to delete (circumvent) a node in a linked list.

If  $u$  is a left child then:

- $u.parent.leftChild \leftarrow u.leftChild$
- $u.leftChild.parent \leftarrow u.parent$

Otherwise if  $u$  is a right child then:

- $u.parent.rightChild \leftarrow u.rightChild$
- $u.rightChild.parent \leftarrow u.parent$

The operations for the case that  $u$  only has a right child is analogous. An example of this operation is depicted in Figures 6.18(c) and fig:bstDelete004

**Case 3:** Suppose that  $u$  has both of its children. Our first instinct may be to promote one of its children up to its place similar to the previous case. In some instances this might work. For example, consider the tree in Figure 6.15. If we wanted to delete the

root, 50, we could promote its right child, 60 up to its place. However, that is only because 60 has no left child. If it had a left child (say 55), then it would not be possible as promoting 60 up would mean it would then have to have 3 children (40, 55, 90).

There are many potential solutions to this problem, but we want to ensure that the operation, in addition to preserving the binary search tree property, is simple, efficient, and causes a minimal change to the tree's structure. Instead of immediately deleting the node, let us instead remove the key, resulting in an "empty" node. We will want to backfill this empty node with a key value in the tree that preserves the BST property.

Again, we exploit the binary search tree property. Because the subtree rooted at  $u$  is also a BST, every key in its left subtree is less than its key and every key in its right subtree is greater. Thus, there are two candidates that we could "promote": the maximum value in the left subtree or the minimum value in the right subtree. Suppose we go with the first option, By replacing  $u$  with the the maximum value in the left subtree, we still ensure that all elements in the left subtree are less than it (because it was the maximum value) and all nodes in the right subtree are still greater than it (since it was less than  $u$  to begin with).

There are two issues that we need to deal with, however. One is finding the maximum (or minimum) value in a subtree and the other is how do we delete the node containing the maximum value from the tree so that it can take the place of  $u$ ? We'll deal with the second problem first. Once we've found the node containing the maximum, if it has zero or two children, we already know how to deal with that from the two previous cases. The only case we really have to deal with is if it has two children. However, we only have to think for a moment that this is, in fact, impossible! Suppose that a node,  $x$  contained the maximum key value and it had both children. By the binary search tree property, its right child,  $x.rightChild$  necessarily has a key value greater than  $x$ , thus  $x$  cannot contain the maximum key value! That is, necessarily the maximum must only have at most 1 child and we know how to deal with that situation.

How do we actually find the maximum value? Again, we exploit the binary search tree property. We first traverse to the left child. From there, we only have to keep traversing right until there is no longer a right child. Conversely, to find the minimum value in the right subtree, we traverse to the right child and then keep traversing left until there is no

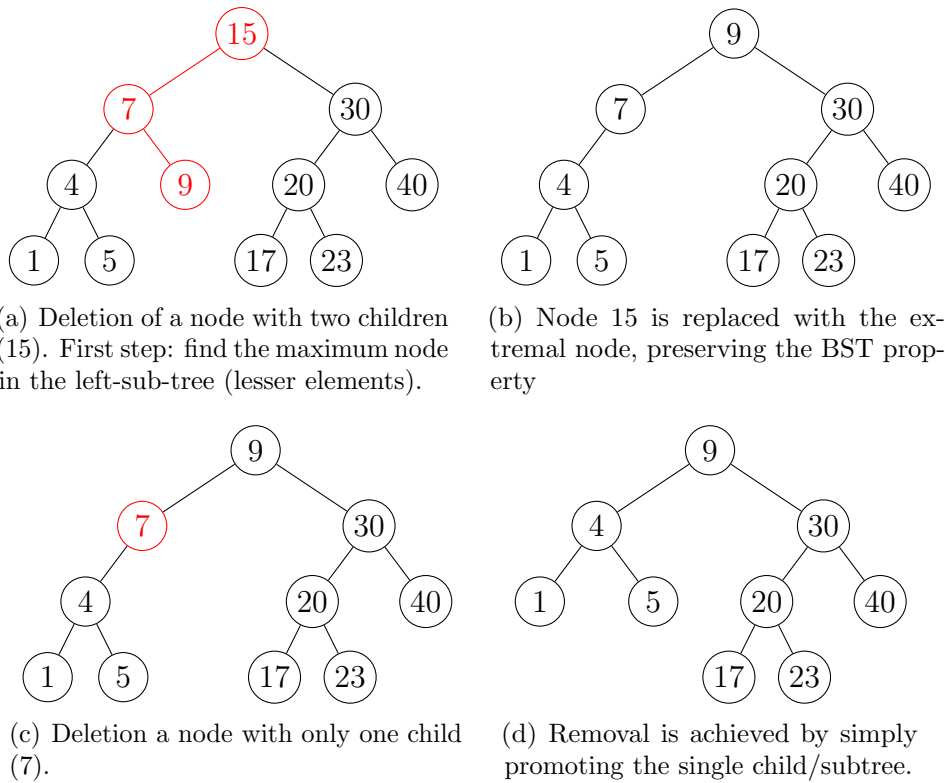


Figure 6.18: BST Deletion Operation Examples. Figures 6.18(a) and 6.18(b) depict the deletion of a node (15) with two children. Figures 6.18(c) and 6.18(d) depict the deletion of a node with only one child (7).

left child. The process for the former choice is presented as Algorithm 27.

INPUT : A node  $u$  in a binary search tree with two children.  
 OUTPUT : The node containing the maximum key value in  $u$ 's left subtree.

```

1  $curr \leftarrow u.leftChild$ 
2 WHILE  $curr.rightChild \neq \phi$  DO
3   |  $curr \leftarrow curr.rightChild$ 
4 END
5 output  $curr$ 

```

**Algorithm 27:** Finding the maximum key value in a node's left subtree.

Examples of cases 2 and 3 can be found in Figure 6.18. The full delete operation has involves several subroutines. Finding the node to be deleted (or determining it does not exist) is  $O(d)$ . Finding the minimum in the left subtree is also at most  $O(d)$ . Ultimately, swapping the keys and/or references is  $O(1)$ . Since these are all independent operations, we have the total complexity as  $O(d) + O(d) + O(1) = O(d)$ .

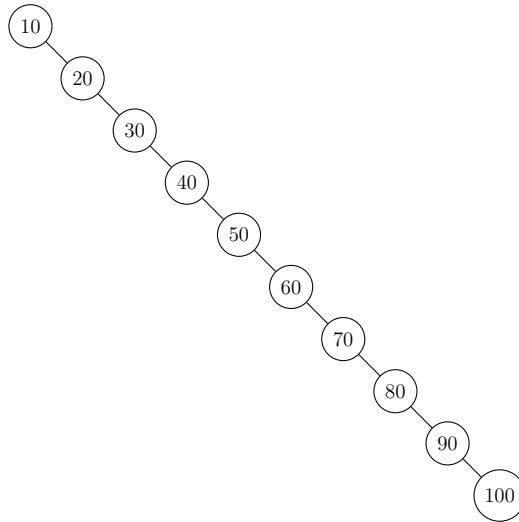


Figure 6.19: A degenerate binary search tree.

### 6.5.4 In Practice

We achieved our goal of developing a tree-based data structure that offered  $O(d)$  insert, retrieval, and update operations. However, we have still fallen short of our main goal of ensuring that these operations are all  $O(\log n)$ . The reason for this is that a binary search tree may become *skewed* or *degenerate*. Consider the binary search tree that would be formed if we inserted the elements 10, 20, 30,  $\dots$ , 100 in that order. The resulting tree would look that that in Figure 6.19. This tree is degenerate because its depth is linear with respect to the number of nodes in the tree;  $d = n - 1 \in O(n)$ . This should look familiar: a degenerate binary search tree is essentially a linked list!

In practice trees will not always be degenerate. However, there is no guarantee that the depth will be logarithmic. Yet more structure would be necessary to make this guarantee. There are plenty of examples of *balanced binary search trees* such that insert and remove operations reorder or *rebalance* a tree so that for every node the left/right subtrees are roughly equal in depth, ensuring that  $d \in O(\log n)$ . Each one has its own advantages and complexity in operations such as AVL trees, 2-3/B-trees, Red-Black trees, splay trees, and treaps. We will not examine such data structures here, but know that it is possible to guarantee that the depth, and thus the basic operations are all efficient,  $O(\log n)$ .

Many programming languages also offer standard implementations for binary search trees, in particular balanced BSTs. Java, for example, provides a `TreeSet` (as well as a `TreeMap`) which is a red-black tree implementation with guaranteed  $O(\log n)$  operations. Moreover, the default iterator pattern is an inorder traversal, providing a sorted ordering. It is usually used with one of its interfaces, a `SortedSet` which allows you to use a `Comparator` for ordering.

## 6.6 Heaps

We may have fallen short of presenting a guaranteed efficient general tree data structure, however we will present a related data structure that does provide efficiency guarantees. Like stack and queues, however, this data structure's efficiency will come at the cost of restricting access to its elements.

**Definition 10** (Heap). A *heap* is a binary tree of depth  $d$  that satisfies the following properties.

1. It is a *full* up to level  $d - 1$ . That is, every node is present at every level up to and including level  $d - 1$ .
2. At level  $d$  all nodes are full-to-the-left. That is, all nodes at the deepest level are all as far left as possible.
3. It satisfies the *heap property*; every node has a key value that is greater than *both* of its children.

This definition actually defines what is referred to as a *max-heap*. This is because as a consequence of the heap property, the maximum element is always guaranteed to be at the root of the heap. Alternatively, we could redefine a heap so that every node's children has key values that are greater than it. This would define a *min-heap*. In practice the distinction is unimportant and a comparator is usually used to define order. Thus, you can get a min-heap implementation using a max-heap with a comparator that reverses the usual ordering. A larger min-heap example can be found in Figure 6.20.

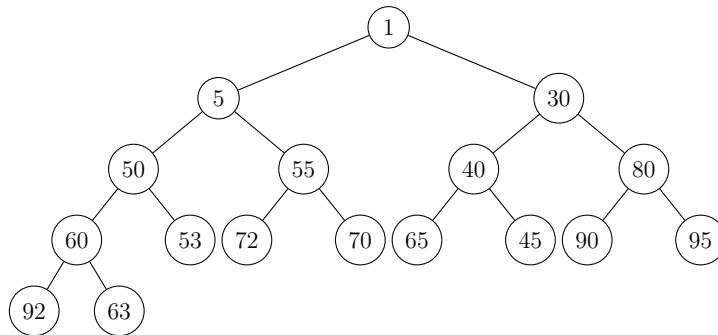


Figure 6.20: A min-heap

As depicted it is easy to see the fullness properties. Every node is present at every level except the last one. In the deepest level, the two remaining nodes are as far left as possible (the children of 60). Essentially there are no “gaps” in any of the levels. This fullness property guarantees that a heap's depth will always be logarithmic,  $O(\log n)$ .

However, it is easy to see that a heap does not provide too much additional structure on the ordering of elements other than the fact that the maximum element is at the top of

the heap. In fact, the structure of a heap resembles a real heap: it is an unorganized pile of stuff. If we were to throw something on top of a real heap, we may not have much control on where it ends up. The operations on a heap data structure are similar. We cannot perform efficient arbitrary searches on a heap because of this lack of structure. We can, however, support two core restricted access operations. We can add elements to the heap and we can remove the top most (getMax) element.

### 6.6.1 Operations

To develop the two core operations, we'll assume that we are working with a max-heap as in Definition 10. We'll use the smaller max-heap example in Figure 6.21 in our examples.

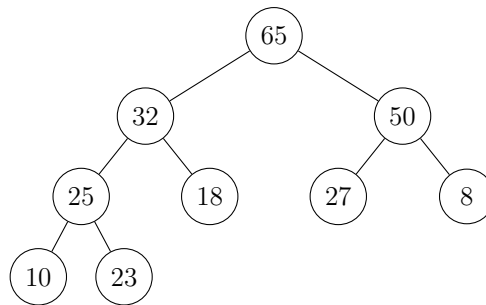


Figure 6.21: A Max-heap

#### Insert Key

The first operation is to insert a new key into the heap. Our first instinct may be to start at the root as we did with binary search trees. We could then move down left/right in the heap in some manner to insert a new key, but how? Suppose we were to insert a new key,  $k = 90$  into the max-heap in Figure 6.21. Clearly it would need to become the new root of the heap, displacing 65. We could continue to shove 65 down displacing either its left or right child. However, what criteria would we use? With only local information to go on, we may choose to exchange with the larger of the two children. Following this criteria we would then exchange 65 and 50, 50 and 27, resulting in something that looks like Figure 6.22.

The problem is that this is not a valid heap as it does not fulfill the fullness property. One might be tempted to modify this approach and instead always exchange the inserted node with the lesser of the two children. This strategy would work with this *particular* example, but we could easily come up with a counter example in which it fails (in fact, the same example, just swap the keys 32 and 50 and it will still fail).

Let's take a step back and redevelop an approach that first ensures that the fullness property is satisfied but the heap property need not necessarily be preserved (at first



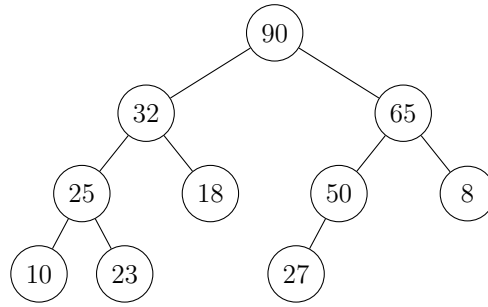


Figure 6.22: An Invalid Max-heap

at least). Preserving the fullness property means that we would necessarily insert the key at the deepest level at the left-most available spot. For the example in Figure 6.21, this would mean inserting as the left child of 18. Inserting 90 in this manner clearly violates the heap property, however, so we need to “fix” the heap. This process is known as *heapifying*. We exchange the inserted key with its parent until either 1) the heap property is satisfied or 2) we’ve reached the root node and the inserted key has become the new root of the heap. This process is fully illustrated in Figure 6.23.

The heapify process is presented as Algorithm 28 which assumes that the new key has already been inserted in the next available spot. The algorithm makes key comparisons, continually exchanging keys until the heap property is satisfied or it reaches the root. This pseudocode, as presented, assumes a tree-based implementation.

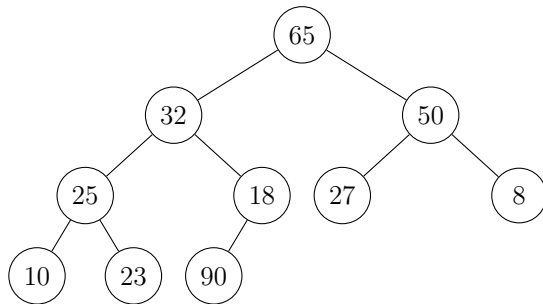
```

INPUT  : A heap  $H$  and an inserted node  $u$ 
1  $curr \leftarrow u$ 
2 WHILE  $curr.parent \neq \phi$  and  $curr.key > curr.parent.key$  DO
3   | swap  $curr.key, curr.parent.key$ 
4   |  $curr \leftarrow curr.parent$ 
5 END

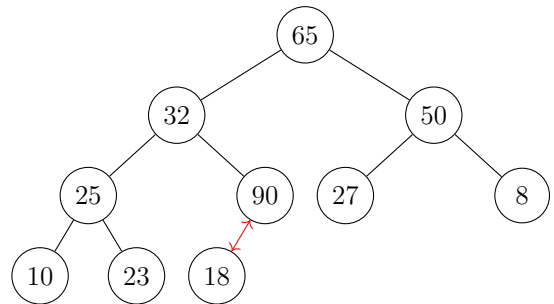
```

**Algorithm 28:** Heapify: fixing a heap data structure after the insertion of a new node,  $u$ .

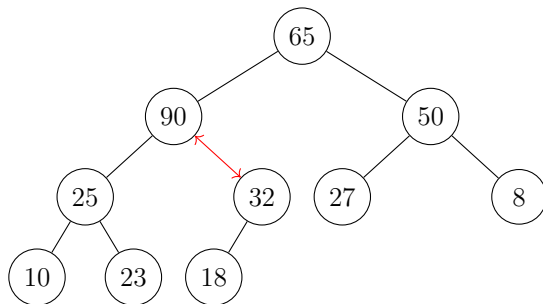
For the moment we will assume that we have an easy way to insert a node at the next available spot. The heapify algorithm itself makes at most  $d$  comparisons and swaps in the worst case (as was the case in the example in Figure 6.23). If we don’t need to heapify all the way up to the root node, then even fewer comparisons and swaps would be necessary, but ultimately the process is  $O(d)$ . However, since a heap guarantees the fullness property,  $d = O(\log n)$ , the process only requires a logarithmic number of comparisons/swaps in the worst case.



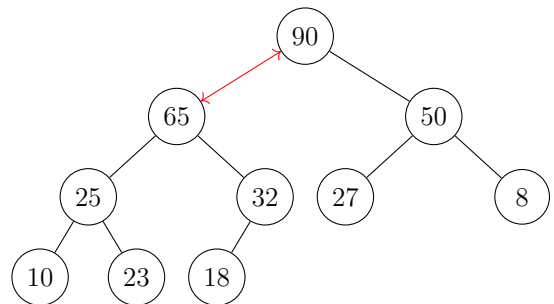
(a) The new key,  $k = 90$  is inserted at the next available spot.



(b) The first comparison finds that 90 and its parent 18 are out of order and exchanges them.



(c) The second comparison ends up exchanging 90 and 32.



(d) The third and final comparison ends up promoting 90 up to the new root.

Figure 6.23: Insertion and Heapification.

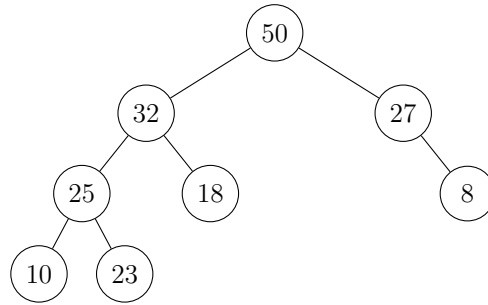


Figure 6.24: Another Invalid Heap

### Retrieve Max

The other core operation is to get the maximum element. Because of the heap property, this is guaranteed to be the root element. First, we remove the key/value from the root node (which we have immediate access to) and save it off to eventually return it. However, this leaves a gap in the tree that must be filled to preserve the fullness property.

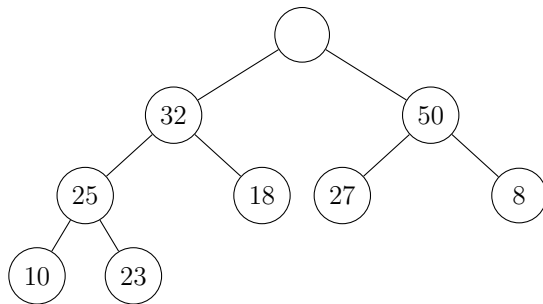
Once again we may be tempted promote one of the root’s children to fill its place. To preserve the heap property, we must necessary promote the larger of the two children. Of course we would have to continue down the heap as the promotion leaves another gap eventually promoting a leaf element up. This idea alone will not suffice however. Again, consider the heap from Figure 6.21. Were we to remove the root element 65 and proceed with the operation as described we would move up 50 then 27 resulting in a “heap” as in Figure 6.24. Obviously this does not fulfill the fullness property as there is a gap (27’s left child is missing).

As before, the solution is to prioritize the fullness property. Once we remove the root element, we need to backfill it with a node that will preserve the fullness property, specifically the “last” node (the rightmost node at the deepest level). Replacing the root element with the last node may not necessarily satisfy the heap property, but we can perform a top-down heapify process similar to what we did with the insert operation. This process is illustrated properly in Figure 6.25.

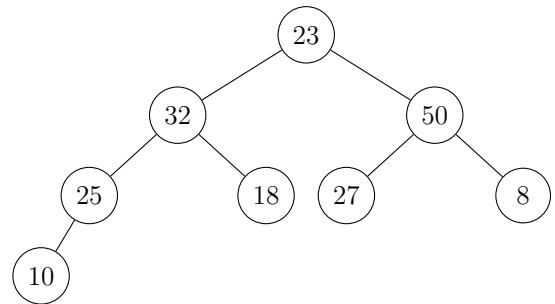
We leave the development of the algorithm and pseudocode as an exercise (see Exercise 6.13). Assuming, as before, that we have easy access to the last element, then swapping it and the the root is an  $O(1)$  operation. The heapification may again require up to  $d$  comparisons and swaps if the last element must be exchanged all the way down to the bottom of the heap. As before, since  $d = O(\log n)$  this operation is  $O(\log n)$ .

### Secondary Operations

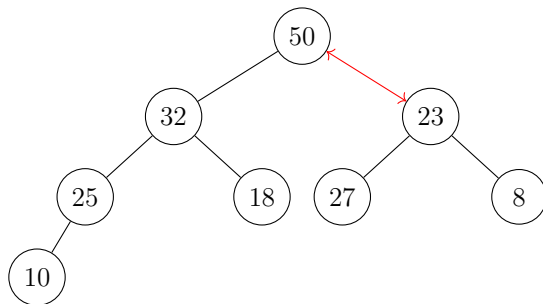
The heap properties don’t allow us to do general operations such as arbitrary search and remove efficiently. We can achieve them using any of our tree traversal strategies, but all



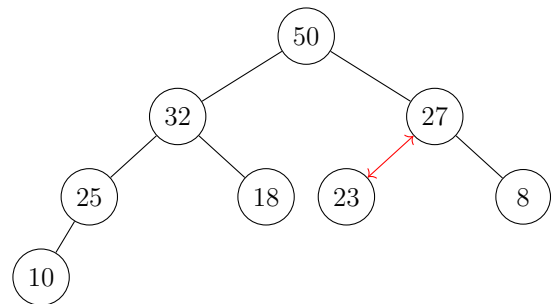
(a) The root element, 65 is saved off in a temporarily variable.



(b) The last element, 23 replaces the root node. However, the heap property is not satisfied.



(c) The node is exchanged with the *greater* of the two children to preserve the heap property.



(d) Another exchange must take place before the heap is fixed.

Figure 6.25: Removal of the root element (getMax) and Heapification

of these operations have the potential to be  $O(n)$ .

However, there are algorithms in which a heap is used that may require that keys in a heap be changed (either increase or decrease their key values). A primary example is when a heap is used to implement a priority queue and we wish to change the priority of an element that is already enqueued.

First, we assume that we have “free” access to the node whose key we want to change as finding an arbitrary node is going to be  $O(n)$  as previously noted. Once we have the node  $u$ , we can increase or decrease its key value. Suppose we increase it. Its key value was already greater than all of its descendants, so the heap property with respect to the subtree rooted at  $u$  is still satisfied. However, we may have increased the key value such it is now greater than its parent. We will need to once again, heapify and exchange keys with its parents and ancestors until the heap property is once again satisfied. This is exactly the same process as when we inserted a new node at the bottom.

Likewise, we can support a decrease key operation. Since the key value was already smaller than its parent, the heap property is unaffected with respect to  $u$ 's ancestors. In this case, however, the heap property could be violated with its children and/or descendants. We simply exchange the key with the larger of its two children if the heap property is violated and continue *downward* in the heap just as we did with the retrieve (and remove) the maximum element. Both of these operations are simply  $O(\log n)$  as before.

## 6.6.2 Implementations

So far we have presented heaps as binary trees. It is possible to implement heaps using the same tree nodes and structure as with binary trees, but it does present some challenges that we'll deal with later. Another, easier and simpler implementation, however is to come full circle back to array-based lists.

### Array-Based Implementation

Recall that the fullness property of a heap essentially means that all nodes are *contiguous*. It makes sense, then that we can store them in an array. In particular, we will store the root element at index 1 (we will leave index 0 unused for this presentation though you can easily shift all of the elements by one index if you'd rather not waste it).

Now, suppose a node  $u$  is stored at index  $i$ . Given this index, we need to locate  $u$ 's parent and left and right child. These will each be stored at:

- The left child is at index  $2i$
- The right child is at index  $2i + 1$
- The parent is at index  $\lfloor \frac{i}{2} \rfloor$

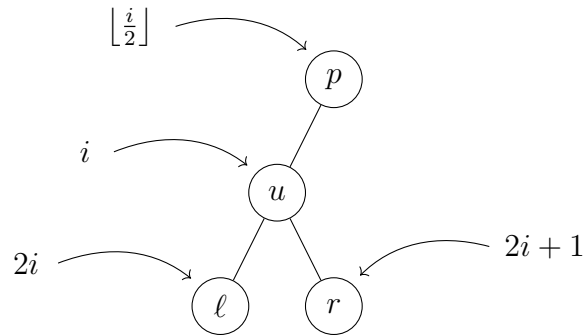


Figure 6.26: Heap Node's Index Relations

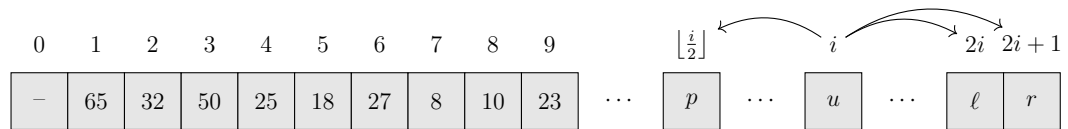


Figure 6.27: An array implementation of the heap from Figure 6.21 along with the generalized parent, left, and right child relations.

These relations are illustrated in Figure 6.26.

A full example is presented in Figure 6.27 which implements the same small max-heap example we've been using. It also demonstrates the backward and forward relationships to the parent and children nodes. Observe that the order in the array matches the order from the tree if we were to perform a breadth first search on it (recall Figure 6.14).

A clear advantage to this implementation is that it is extremely easy and we can even reuse an array-based list implementation. It is a simple matter to implement the heapify algorithms in terms of indices given the mappings above (see Exercise 6.14). One disadvantage is that it is costly to increase the size of the array when we need to add more elements. However, amortized over the life of the heap and given the "savings" in a simpler implementation, this may not be that big of a deal.

### Tree-Based Implementation

Though not common, you can implement a heap using the same binary tree structure using a node with references to a parent, left child, and right child. However, to ensure efficient operation for the heapify algorithms, we necessarily have to keep track of a parent element for every node.

One problem is that we do not have obvious access to the "last" element or next available

spot in the heap as we did with an array-based implementation. With an array, we could easily keep track of how many elements have been stored in the heap, say  $n$ . Then the last element is necessarily at index  $n$  and the next available spot is at index  $n + 1$ . Since we have random access, we can easily jump to either location to get the last element or to insert at the next available spot.

With a tree structure, however, we do not have such access. Though we keep track of the root, it is not straightforward to also keep track of the last element (or the next available spot). Searching for either using BFS or another tree traversal algorithm would kill our  $O(\log n)$  efficiency.

Fortunately, using a bit of mathematical analysis and exploiting the fullness property of a heap will allow us to find either the last element or the first available spot in  $O(\log n)$  time. We'll focus on finding the first available open spot as the same technique can be used to find the last element with minor modifications.

Without loss of generality, we'll assume that we've kept track of the number of nodes in the heap,  $n$  and thus we can compute the depth,

$$d = \lfloor \log n \rfloor$$

Due to the fullness property, simply knowing  $n$  and  $d$  give us a great deal of information. In particular, we know that there are

$$\sum_{k=0}^{d-1} 2^k = 2^d - 1$$

nodes in the first  $d$  levels (levels 0 up to  $d - 1$ ). Computing a simple difference,

$$m = n - (2^d - 1)$$

tells us how many nodes are in the last (deepest) level, level  $d$ . We also know that if the last level were full, it would have  $2^d$  nodes. This tells us whether or not the next available spot is in the left subtree or the right subtree by making a simple comparison:

- If  $m < \frac{2^d}{2} = 2^{d-1}$  then the left subtree is not “full” at level  $d$  and so it contains the next available spot.
- Otherwise if  $m \geq 2^{d-1}$  then the left subtree is full and the right subtree must contain the next available spot.

In each case we can traverse down to the left or right child respectively (which ever's tree contains the next available spot) and update both  $n$  and  $m$  (the number of elements at level  $d$ ) and repeat this process until we've found the next available spot. This analysis is visualized in Figure 6.28 and the full process is presented as Algorithm 29.

It is not difficult to see that the complexity of this algorithm is  $O(d) = O(\log n)$  since we iterate down the depth of the heap using simple arithmetic operations. Thus the problem

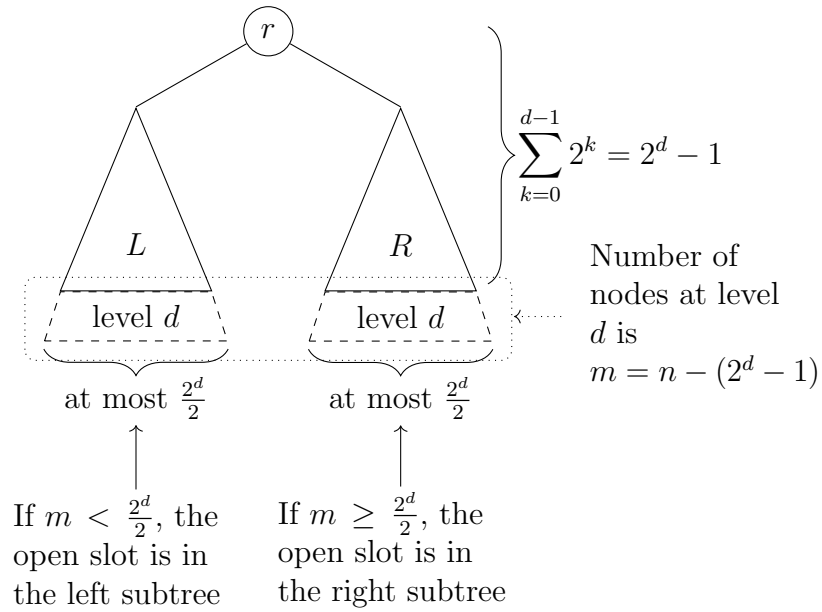


Figure 6.28: Tree-based Heap Analysis. Because of the fullness property, we can determine which subtree (left or right) the “open” spot in a heap’s tree is by keeping track of the number of nodes,  $n$ . This can be inductively extended to each subtree until the open spot is found.

of finding the last element or the next available spot does not significantly increase the complexity of the heap’s core operations.



```

INPUT   : A tree-based heap  $H$  with  $n$  nodes
OUTPUT  : The node whose child is the next available open spot in the heap
1  $curr \leftarrow T.head$ 
2  $d \leftarrow \lfloor \log n \rfloor$ 
3  $m \leftarrow n$ 
4 WHILE  $curr$  has both children DO
5   IF  $m = 2^{d+1} - 1$  THEN
6     //remaining tree is full, traverse all the way left
7     WHILE  $curr$  has both children DO
8       |  $curr \leftarrow curr.leftChild$ 
9     END
10  ELSE
11    //remaining tree is not full, determine if the next open
12    spot is in the left or right sub-tree
13    IF  $m \geq \frac{2^d}{2}$  THEN
14      //left sub-tree is full
15       $d \leftarrow (d - 1)$ 
16       $m \leftarrow (m - \frac{2^d}{2})$ 
17       $curr \leftarrow curr.rightChild$ 
18    ELSE
19      //left sub-tree is not full
20       $d \leftarrow (d - 1)$ 
21       $m \leftarrow m$ 
22       $curr \leftarrow curr.leftChild$ 
23    END
24  END
25 END
26 output  $curr$ 

```

**Algorithm 29:** Find Next Open Spot - Numerical Technique

### 6.6.3 Variations

As presented in Definition 10, the heaps we have been describing are sometimes referred to as *binary heaps* because they are based on a binary tree. There are other variations of heaps that uses different structures and offer different properties.

For example, Binomial Heaps are heaps that are a collection of binomial trees. A

binomial tree of order  $k$  is a tree such that its children are binomial trees of order  $k - 1, k - 2, \dots, 2, 1, 0$  (that is, it has  $k$  children). This is an inductive definition so that the base case of  $k = 0$  is a single node tree. A binomial tree of order  $k$  has  $2^k$  nodes and is of depth  $k$ . The core operations are a bit more complicated but have the same complexity of  $O(\log n)$ . The advantage of this implementation is that a *merge* operation of two heaps is also efficient ( $O(\log n)$ ).

A Fibonacci heap is a collection of trees that satisfy the heap property. However, the structure is a lot more flexible than binary heaps or binomial heaps. The main advantage is that some of the operations to keep track of elements are delayed until they are necessary. Thus, some operations may be quite expensive, but when looked at from an *amortized analysis*, the expected or average running time of the operations can be interpreted as constant,  $O(1)$ . This is similar to when we examined array-based lists. Sometimes (though not often) we may need to expand the underlying array. Though this is an expensive operation, since it is not too common, when you average the running time of the core operations over the lifetime of the data structure, it all “evens out.” and looks to be constant.

## 6.6.4 Applications

A heap is used in many algorithms as a data structure to efficiently hold elements to be processed. For example, several graph algorithms such as Prim’s Minimum Spanning Tree or Dijkstra’s Shortest Path algorithms use heaps as their fundamental building block. A heap is also used in Huffman’s Coding algorithm to efficiently compress a file without loss of information.

### Priority Queue

From the core operations (insert and getMax) it might be obvious that a heap is an efficient way to implement a priority queue. The enqueue operation is a simple insert and a dequeue operation is always guaranteed to result in the maximum (highest priority) element. Since both operations are  $O(\log n)$ , this is far more efficient than a list-based priority queue. Note that the Java Collections library provides a heap-based priority queue implementation in the `java.util.PriorityQueue<E>` class.

### Heap Sort

Suppose we have a collection of elements. Now suppose we threw them all into a heap and then, one-by-one, pulled them out using the getMax operation. As we pull them out, what order would they be in? Sorted of course! By using a sophisticated data structure like a heap, we can greatly simplify the code to achieve a more complex data operation,

sorting. This is referred to as Heap Sort and is presented as Algorithm 30.

```

INPUT  : A collection of elements  $A = \{a_1, \dots, a_n\}$ 
OUTPUT :  $A$ , sorted
1  $H \leftarrow$  empty min-heap
2 FOREACH  $a \in A$  DO
3   | insert  $a$  into  $H$ 
4 END
5  $i \leftarrow 1$ 
6 WHILE  $H$  is not empty DO
7   |  $a_i \leftarrow H.getMin$ 
8   |  $i \leftarrow (i + 1)$ 
9 END
10 output  $A$ 

```

**Algorithm 30:** Heap Sort

Examine at the code in this algorithm. It is extremely simple; we put stuff into a data structure, then we pull it out. It doesn't get much simpler than that. The real magic is in the data structure we used and exploited. This is a perfect illustration of one of the themes of this book, borrowed from Eric S. Raymond [6] that "Smart data structures and dumb code are a lot better than the other way around." There are many sophisticated sorting algorithms (Quick Sort, Merge Sort, etc.) that use clever code (recursion, partitioning, etc.) and simple arrays. In contrast, Heap Sort uses very dumb code: put stuff in, take stuff out and a very smart data structure. Beautiful.

Of course, this beauty is all for nothing if it is not also efficient. The complexity of the insertion and getMin operations in Heap Sort actually changes on each iteration of the algorithm because the data structure we're using is growing and shrinking. Let's first analyze the "put stuff in" phase (the foreach loop, lines 2–3). On the first iteration, no comparisons or swaps are made as  $H$  is initially empty. On the second iteration where we insert the second element,  $H$  has size 1 and so 1 comparison (and potentially 1 swap) is made. In general on the  $i$ -th iteration,  $H$  has  $i - 1$  elements in it, thus the insertion requires roughly  $\log i - 1$  (the depth of the heap) comparisons and/or swaps. This is all summarized in Table 6.3.

Thus, the total number of comparisons (or swaps) made by Heap Sort is the summation of the 4th column, or

$$\sum_{i=2}^{n-1} \log i = \log 2 + \log 3 + \dots + \log n - 1$$

Note that we start the index at  $i = 2$ , since the first iteration does not require any comparisons. Using logarithm identities, we can further simplify and bound this summation:

Table 6.3: Analysis of Heap Sort

Iteration	Size of $H$	Depth of $H$	Number of Comparisons
1	0	–	0
2	1	0	1
3	2	1	2
$\vdots$	$\vdots$	$\vdots$	$\vdots$
$i$	$i - 1$	$\log(i - 1)$	$\log(i - 1)$
$\vdots$	$\vdots$	$\vdots$	$\vdots$
$n$	$n - 1$	$\log(n - 1)$	$\log(n - 1)$

$$\begin{aligned}
\sum_{i=2}^{n-1} \log i &= \log(2) + \log(3) + \cdots + \log(n-1) \\
&= \log(2 \cdot 3 \cdot 4 \cdots (n-1)) \\
&= \log((n-1)!) \\
&\leq \log(n!) \\
&\leq \log(n^n) \\
&= n \log n
\end{aligned}$$

That is, the first phase is only  $O(n \log n)$ , matching the best case running time of Quick Sort and other fast sorting algorithms. The analysis for the “take stuff out” phase is similar, but in reverse. The size of the heap will be diminishing from  $n, n - 1, \dots, 1$  and will again contribute another factor of  $n \log n$ . Since these are independent phases, the entire algorithm is an efficient  $O(n \log n)$  sorting algorithm. Simple, efficient, elegant and beautiful.

## 6.7 Exercises

**Exercise 6.1.** Adapt and implement the Stack-Based Preorder Traversal (Algorithm 18) for the Java `BinaryTree<T>` class.

**Exercise 6.2.** Adapt and implement the Stack-Based Inorder Traversal (Algorithm 21) for the Java `BinaryTree<T>` class.

**Exercise 6.3.** Adapt and implement the Stack-Based Postorder Traversal (Algorithm 22) for the Java `BinaryTree<T>` class.

**Exercise 6.4.** Adapt and implement the Tree Walk Algorithm (Algorithm 24) for the Java `BinaryTree<T>` class.

**Exercise 6.5.** Let  $T$  be a binary tree of depth  $d$ . What is the maximum number of leaves that  $T$  could have with respect to  $d$ ?

**Exercise 6.6.** Develop an algorithm that inserts a new node into a binary tree at the shallowest available spot.

**Exercise 6.7.** Develop an algorithm that, given a binary tree  $T$ , counts the total number of nodes in the tree (assume it is not kept track of as part of the data structure).

**Exercise 6.8.** Develop an algorithm that, given a binary tree  $T$ , counts the number of leaves in it.

**Exercise 6.9.** Develop an algorithm that given a node  $u$  in a binary tree  $T$  determines its depth.

**Exercise 6.10.** Develop an algorithm that given a binary tree  $T$ , computes its depth.

**Exercise 6.11.** Let  $T$  be a binary tree with  $n$  nodes. What is the maximum number of leaves that  $T$  could have with respect to  $n$ ? Provide an example for  $n = 15$

**Exercise 6.12.** Develop an algorithm (write pseudocode) to insert a given key  $k$  into a binary search tree  $T$ .

**Exercise 6.13.** Develop an algorithm (write pseudocode) fix a heap after its root element has been removed.

**Exercise 6.14.** Rewrite Algorithm 28 (heapify) to work on an array instead of a tree structure.



# Glossary

**algorithm** a process or method that consists of a specified step-by-step set of operations.

**corner case** a scenario that occurs outside of typical operating parameters or situations; an exceptional case or situation that may need to be dealt with in a unique or different manner. [18](#)

**idiom** in the context of software an idiom is a common code or design pattern. [12](#)

**iterator** a pattern that allows a user to more easily and conveniently iterate over the elements in a collection. [15](#)

**leaky abstraction** a design that exposes details and limitations of an implementation that should otherwise be hidden through encapsulation.. [25](#)

**parameterized polymorphism** a means in which you can make a piece of code (a method, class, or variable) generic so that its type can vary depending on the context in which the code is used. [13](#)

**queue** a collection data structure that holds elements in a first-in first-out manner.

**random access** a mechanism by which elements in an array can be accessed by simply computing a memory offset of each element relative to the beginning of the array. [19](#), [41](#)

**stack** a collection data structure that holds elements in a last-in first-out manner. [29](#)

**syntactic sugar** syntax in a language or program that is not absolutely necessary (that is, the same thing can be achieved using other syntax), but may be shorter, more convenient, or easier to read/write. In general, such syntax makes the language “sweeter” for the humans reading and writing it. [15](#)





# Acronyms

**ACID** Atomicity, Concurrency, Isolation, Durability.

**ADT** Abstract Data Type. [7](#)

**AST** Abstract Syntax Tree. [33](#)

**BFS** Breadth First Search. [104](#), [106](#)

**CRUD** Create-Retrieve-Update-Destroy.

**DAG** Directed Acyclic Graph.

**DFS** Depth First Search. [33](#), [93](#)

**DRY** Don't Repeat Yourself. [10](#)

**FIFO** First-In First-Out. [33](#), [34](#), [106](#)

**JVM** Java Virtual Machine. [15](#)

**LIFO** Last-In First-Out. [29](#), [94](#)

**ORM** Object-Relational Mapping.

**RDMS** Relational Database Management System.

**SQL** Structured Query Language.



# Index

- array-based lists, 8
- balanced binary search trees, 114
- binary search tree, 107
  - definition, 107
- binary search trees
  - balanced binary search trees, 114
- binary tree, 88
- binomial heaps, 124
- consumer producer pattern, 35
- deque, 38
- fibonacci heaps, 125
- gauss's formula, 43
- heap sort, 125
- heaps, 115
  - binomial heaps, 124
  - definition, 115
  - fibonacci heaps, 125
  - heap sort, 125
- linked lists, 16
- lists, 7
  - array lists, 8
  - linked lists, 16
- object oriented programming, 3
- priority queue, 38
- queue
  - priority queue, 38
- queues, 33
- stacks, 29
  - peek, 31
- tree, 83
  - binary, 88
  - definition, 85



# Bibliography

- [1] Manindra Agrawal, Neeraj Kayal, and Nitin Saxena. PRIMES is in P. *Ann. of Math*, 2:781–793, 2002.
- [2] Paul Bachmann. *Die Analytische Zahlentheorie*. 1894.
- [3] Arthur Cayley. On the theory of the analytical forms called trees. *Philosophical Magazine Series 4*, 13(85):172–176, 1857.
- [4] Leonhard Euler. Solutio problematis ad geometriam situs pertinentis. *Commentarii academiae scientiarum Petropolitanae*, 8:128–140, 1741.
- [5] Donald E. Knuth. Big omicron and big omega and big theta. *SIGACT News*, 8(2):18–24, April 1976.
- [6] Eric S. Raymond. *The Cathedral and the Bazaar*. O'Reilly & Associates, Inc., Sebastopol, CA, USA, 1st edition, 1999.