

CSE 250: Data Structures

Course Reference

Chapter 0 Contents

1. Introduction	4
1.1. What to get out of this Data Structures class?	4
1.1.1. An intuition for data structures	4
1.1.2. Practice with formal proofs and recursion	4
1.2. What this class is not.	4
1.3. A word on Lies and Trickery	4
2. Math Refresher	6
3. Asymptotic Runtime Complexity	7
3.1. Some examples of asymptotic runtime complexity	7
3.2. Asymptotic Analysis in General	7
3.3. Runtime Growth Functions	8
3.4. Complexity Classes	8
3.4.1. Example	9
3.5. Formal Notation	10
3.5.1. Polynomials and Dominant Terms	11
3.5.2. Θ in mathematical formulas	12
3.6. Code Complexity	12
3.7. Complexity Bounds	13
3.8. Big- O and Big- Ω	14
3.9. Formalizing Big- O	16
3.9.1. Proving a function has a specific Big- O bound	17
3.9.2. Proving a function does not have a specific Big- O bound	18
3.10. Formalizing Big- Ω	18
3.11. Formalizing Big- Θ	18
3.12. Tight Bounds	19
3.13. Which Variable?	20
3.13.1. Related Variables	20
3.14. Summary	20
3.14.1. Formal Definitions	21
3.14.2. Interpreting Code	21
3.14.3. Simple Complexity Classes	22
3.14.4. Dominant Terms	22
3.14.5. Multiclass Asymptotics	22
4. The Sequence and List ADTs	24
4.1. What is an ADT?	24
4.2. The Sequence ADT	24
4.2.1. Sequences by Rule	25
4.2.2. Arbitrary Sequences	26
4.2.3. Arrays	27
4.2.3.1. Array Runtime Analysis	27
4.2.3.2. Side Note: Memory Allocation	28
4.2.3.3. Arrays In Java	28
4.2.4. Mutable Sequences	29
4.2.4.1. Mutable Array	29

4.2.5. Array Summary	30
4.3. The List ADT	30
4.3.1. A simple Array as a List	31
4.4. Linked Lists	32
4.4.1. Side Note: How Java uses Memory	33
4.4.2. A Linked List as a List	33
4.4.2.1. Linked List get	34
4.4.2.2. Linked List update	36
4.4.2.3. Linked List add	37
4.4.2.4. Aside: Invariants and Rule Preservation	39
4.4.2.5. Linked List add runtime	39
4.4.2.6. Linked List size (Take 1)	40
4.4.2.7. Aside: Loop Invariants	42
4.4.2.8. Linked List size (Take 1) runtime	42
4.4.2.9. Linked List size (Take 2)	43
4.5. Iterators	44
4.6. Doubly Linked Lists	44

Chapter 1. Introduction

Data Structures classes often have a mixed reception. One of the most common things I hear from students is that it's hard to see how the concepts they learn get deployed into practice: We don't solve specific problems. Although we will try to map many of the ideas we discuss to specific problems that you'll encounter as a computer scientist, it's also important to think about the content you learn in this class in a more general sense: as a set of tools for your toolbox.

This class will introduce you to the hammers, wrenches, and screwdrivers of computer science. We'll show you how to decide whether a phillips-head or a flathead screwdriver is better for your use case; when to use a socket wrench or a crescent wrench.

A bit less metaphorically, this class will present a set of organizational strategies (data structures) suitable for different use cases (abstract data types). We'll also introduce asymptotic complexity, a way to quickly summarize the performance characteristics of data structures (and a tool for thinking about algorithms).

1.1. What to get out of this Data Structures class?

1.1.1. An intuition for data structures

We'll talk about the material at a precise, formal level, but we honestly hope that you learn the material to a point where you internalize it. If, after this class, you never again consciously think about the fact that prepending to a linked list is $O(1)$, that's fine. What we care about is that you develop the intuition

- ... that you develop an instinct for which data structure is right for a given situation.
- ... that you get a little cringe in the back of your brain when you see in the documentation that a method you need to call repeatedly is $O(N)$.
- ... that you get a little cringe in the back of your brain when you see a doubly nested loop, or another piece of $O(N^2)$ or worse code.

1.1.2. Practice with formal proofs and recursion

By the time you take this class, you should have already taken a discrete math class and gotten your first exposure to proofs and recursive thinking. This class is intended to develop those same skills:

- We'll review an assortment of proofs regarding algorithm runtimes using specific data structures.
- We'll discuss specific strategies you can use while proving things.
- We'll review recursion and discuss approaches to identifying problems that can be solved through recursion and how to use recursion as a problem solving technique.

1.2. What this class is not.

In contrast to many data structures classes, which introduce C programming, memory management, and other related concepts, UB's 250 is intended as a concepts/theory-style class. To be clear, we will use code. Example code will be given in Java (or some cases Python), and we'll make extensive use of Java's class inheritance model. We will use code code to reinforce and motivate concepts that you will learn throughout the class. However, the primary focus of this class is on understanding asymptotic analysis and adding the standard suite of data structures to your toolbox.

1.3. A word on Lies and Trickery

"All models are wrong, but some are useful"

This is an introductory text. As such, you should expect that many of the things we say are simplified for the purposes of presentation. Some assertions (e.g., array accesses are constant-time) will be outright lies with respect to code running on modern computers. Nevertheless, these simplifications are here for a reason. It will be far easier to first grok the simpler model of code and data organization, before delving into the more nuanced details and special cases.

In general, we highlight some of the more blatant lies with footnotes that also hint at some of the nuanced details. In a few cases (e.g., constant time array accesses), we'll peel back the lies later on in the book.

That being said, these footnotes are primarily here for the pedants and excessively curious among you. You should still be able to understand the rest of the book even if you ignore every single footnote in the text.

Chapter 2. Math Refresher

Chapter 3. *Asymptotic Runtime Complexity*

Data Structures are the fundamentals of algorithms: How efficient an algorithm is depends on how the data is organized. Think about your workspace: If you know exactly where everything is, you can get to it much faster than if everything is piled up randomly. Data structures are the same: If we organize data in a way that meets the need of an algorithm, the algorithm will run much faster (remember the array vs linked list comparison from earlier)?

Since the point of knowing data structures is to make your algorithms faster, one of the things we'll need to talk about is how fast the data structure (and algorithm) is for specific tasks. Unfortunately, "how fast" is a bit of a nuanced comparison. I could time how long algorithms **A** and **B** take to run, but what makes a fair comparison depends on a lot of factors:

- How big is the data that the algorithm is running on?
 - **A** might be faster on small inputs, while **B** might be faster on big inputs.
- What computer is running the algorithm?
 - **A** might be much faster on one computer, **B** might be much faster on a network of computers.
 - **A** might be especially tailored to Intel x86 CPUs, while **B** might be tailored to the non-uniform memory latencies of AMD x86 CPUs.
- How optimized is the code?
 - Hand-coded optimizations can account for multiple orders of magnitude in algorithm performance.

In short, comparing two algorithms requires a lot of careful analysis and experimentation. This is important, but as computer scientists, it can also help to take a more abstract view. We would like to have a shorthand that we can use to quickly convey the 50,000-ft view of "how fast" the algorithm is going to be. That shorthand is asymptotic runtime complexity.

3.1. Some examples of asymptotic runtime complexity

Look at the documentation for data structures in your favorite language's standard library. You'll see things like:

- The cost of appending to a Linked List is $O(1)$
- The cost of finding an element in a Linked List is $O(N)$
- The cost of appending to an Array List is $O(N)$, but amortized $O(1)$
- The cost of inserting into a Tree Set is $O(\log N)$
- The cost of inserting into a Hash Map is Expected $O(1)$, but worst case $O(N)$
- The cost of retrieving an element from a Cuckoo Hash is always $O(1)$

These are all examples of asymptotic runtimes, and they give you a quick at-a-glance idea of how well the data structure handles specific operations. Knowing these facts about the data structures involved can help you plan out the algorithms you're writing, and avoid picking a data structure that tanks the performance of your algorithm.

3.2. Asymptotic Analysis in General

Although our focus in this book is mainly on asymptotic **runtime** complexity, asymptotic analysis is a general tool that can be used to discuss all sorts of properties of code and algorithms. For example:

- How fast is an algorithm?
- How much space does an algorithm need?
- How much data does an algorithm read from disk?
- How much network bandwidth will an algorithm use?

- How many CPUs will a parallel algorithm need?

3.3. Runtime Growth Functions

Let's start by talking about **Runtime Growth Functions**. A runtime growth function looks like this:

$$T(N)$$

Here, T is the name of the function (We usually use T for runtime growth functions, and N is the *size* of the input).

You can think of this function as telling us “For an input of size N , this algorithm will take $T(N)$ seconds to run” This is a little bit of an inexact statement, since the actual number of seconds it takes depends on the type of computer, implementation details, and more. We'll eventually generalize, but for now, you can assume that we're talking about a specific implementation, on a specific computer (like e.g., your computer).

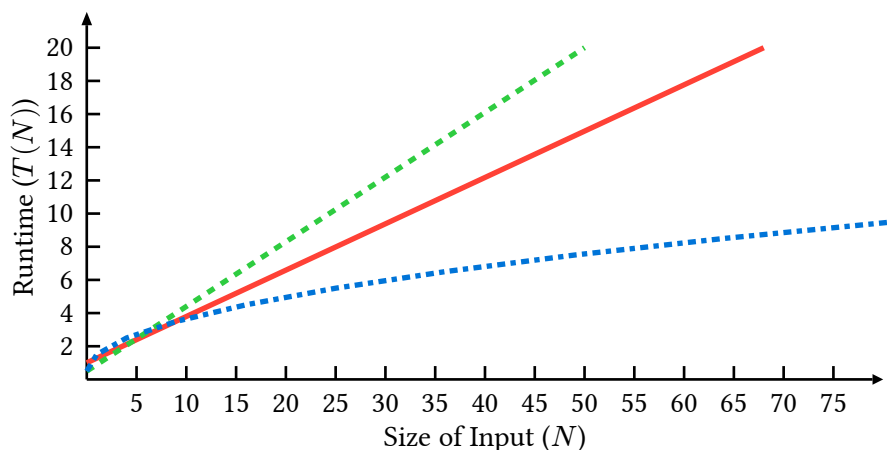
We call this a growth function because it generally has to follow a few rules:

- For all $N \geq 0$, it must be the case that $T(N) \geq 0$
 - The algorithm can't take negative time to run.
- For all $N \geq N'$, it must be the case that $T(N) \geq T(N')$
 - It shouldn't be the case that the algorithm runs faster on a bigger input¹.

3.4. Complexity Classes

Before we define the idea of asymptotic runtimes precisely, let's start with an intuitive idea. We're going to take algorithms (including algorithms that perform specific operations on a data structure), and group them into what we call **Complexity Classes**².

Graph 1 shows three different runtime growth functions: Green (dashed), Red (solid), and Blue (dash-dotted). For an input of size $N = 2$, the green dashed function appears to be the best, while the blue dash-dotted function appears the worst. By the time we get to $N = 10$, the roles have reversed, and the blue dash-dotted function is the best.



Graph 1: Different types of growth functions

¹In practice, this is not actually the case. We'll see a few examples of functions whose runtime can sometimes be faster on a bigger input. Still, for now, it's a useful simplification.

²To be pedantic, what we'll be describing is called “simple complexity classes”, but throughout this book, we'll refer to them as just complexity classes.

So let's talk about these lines and what we can say about them. First, in this book, we're going to ignore what happens for "small" inputs. This isn't always the right thing to do, but we're taking the 50,000 ft view of algorithm performance. From this perspective, the blue dot-dashed line is the "best".

But why is it better? If we look closely, both the green dashed and the red solid line are straight lines. The blue dot-dashed line starts going up faster than both the other two lines, but bends downward. In short, the blue dot-dashed line draws a function of the form $a \log(N) + b$, while the other two lines draw functions of the form $a \cdot N + b$. For "big enough" values of N , any function of the form $a \log(N) + b$ will always be smaller than any function of the form $a \cdot N + b$. On the other hand, the value of any two functions of the form $a \cdot N + b$ will always "look" the same. No matter how far we zoom out, those functions will always be a constant factor different.

Our 50,000 foot view of the runtime of a function (in terms of N , the size of its input) will be to look at the "shape" of the curve as we plot it against N .

3.4.1. Example

Try the following code in python:

```
from random import randrange
from datetime import datetime
N = 10000
TRIALS = 1000

#### BEGIN INITIALIZE data
data = []
for x in range(N):
    data += [x]
data = list(data)
#### END INITIALIZE data

contained = 0
start_time = datetime.now()
for x in range(TRIALS):
    if randrange(N) in data:
        contained += 1
end_time = datetime.now()

time = (end_time - start_time).total_seconds() / TRIALS

print(f"For N = {N}, that took {time} seconds per lookup")
```

This code creates a list of N elements, and then does $TRIALS$ checks to see if a randomly selected value is somewhere in the list. This is a toy example, but see what happens as you increase the value of N . In most versions of python, you'll find that every time you multiply N by a factor of, for example 10, the total time taken per lookup grows by the same amount.

Now try something else. Modify the code so that the data variable is initialized as:

```
#### BEGIN INITIALIZE data
data = []
for x in range(N):
    data += [x]
```

```
data = set(data)
#### END INITIALIZE data
```

You'll find that now, as you increase N , the time taken **per lookup** grows at a much smaller rate. Depending on the implementation of python you're using, this will either grow as $\log N$ or only a tiny bit. The `set` data structure is much faster at checking whether an element is present than the `list` data structure.

Complexity classes are a language that we can use to capture this intuition. We might say that `set`'s implementation of the `in` operator belongs to the **logarithmic** complexity class, while `list`'s implementation of the operator belongs to the **linear** complexity class. Just saying this one fact about the two implementations makes it clear that, in general, `set`'s version of `in` is much better than `list`'s.

3.5. Formal Notation

Sometimes it's convenient to have a shorthand for writing down that a runtime belongs in a complexity class. We write:

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

... to mean that the mathematical function $g(N)$ belongs to the same **asymptotic complexity class** as $f(N)$. You may also see this written as an equality. For example

$$T(N) = \Theta(N)$$

... means that the runtime function $T(N)$ belongs to the **linear** complexity class. Continuing the example above, we would use our new shorthand to describe the two implementations of Python's `in` operator as:

- $T_{\text{set}} \in \Theta(\log N)$
- $T_{\text{list}} \in \Theta(N)$

Formalism: A little more formally, $\Theta(f(N))$ is the **set** of all mathematical functions $g(N)$ that belong to the same complexity class as $f(N)$. So, writing $g(N) \in \Theta(f(N))$ is saying that $g(N)$ is in (\in) the set of all mathematical functions in the same complexity class as $f(N)$ ³.

Here are some of the more common complexity classes that we'll encounter throughout the book:

- **Constant:** $\Theta(1)$
- **Logarithmic:** $\Theta(\log N)$
- **Linear:** $\Theta(N)$
- **Loglinear:** $\Theta(N \log N)$
- **Quadratic:** $\Theta(N^2)$
- **Cubic:** $\Theta(N^3)$
- **Exponential** $\Theta(2^N)$

³We are sweeping something under the rug here: We haven't precisely defined what it means for two functions to be in the same complexity class yet. We'll get to that shortly, after we introduce the concept of complexity bounds.

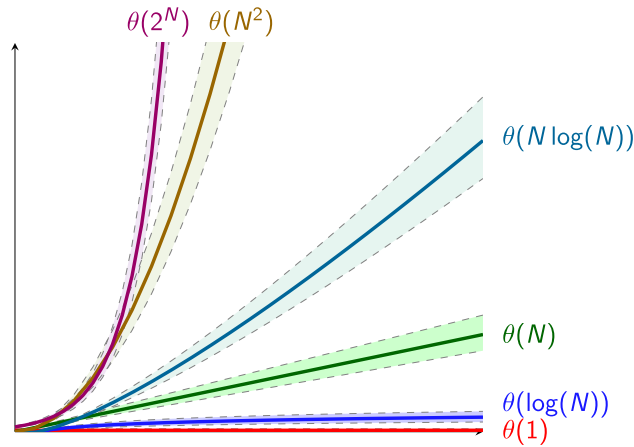
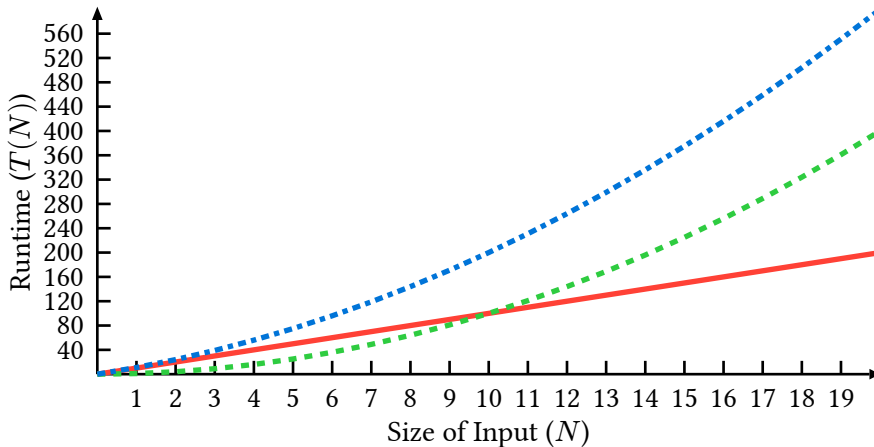


Figure 1: $\Theta(f(N))$ is the set of all mathematical functions including $f(N)$ and everything that has the same “shape”, represented in the chart above as a highlighted region around each line.

Complexity classes are given in order. The later they appear in the list, the faster they grow. Any function that has a **linear** shape, will always be smaller (for big enough values of N) than a function with a **loglinear** shape.

3.5.1. Polynomials and Dominant Terms

What complexity class does $10N + N^2$ fall into? Let’s plot it and see:



Graph 2: Comparing N (solid red), N^2 (dashed green), and $10N + N^2$ (dash-dotted blue)

Graph 2 compares these three functions. Observe that the dash-dotted blue line starts off very similar to the solid red line. However, as N grows, its shape soon starts resembling the dashed green N^2 more than the solid red $10N$.

This is a general pattern. In any polynomial (a sum of mathematical expressions), for really big values of N , the complexity class of the “biggest” term starts to win out once we get to really big values of N .

In general, for any sum of mathematical functions:

$$g(N) = f_1(N) + f_2(N) + \dots + f_k(N)$$

The complexity class of $g(N)$ is the greatest complexity class of any $f_{i(N)}$

For example:

- $10N + N^2 \in \Theta(N^2)$
- $2^N + 4N \in \Theta(2^N)$
- $1000 \cdot N \log(N) + 5N \in \Theta(N \log(N))$

3.5.2. Θ in mathematical formulas

Sometimes we'll write $\Theta(g(N))$ in regular mathematical formulas. For example, we could write:

$$\Theta(N) + 5N$$

You should interpret this as meaning any function that has the form:

$$f(N) + 5N$$

... where $f(N) \in \Theta(N)$.

3.6. Code Complexity

Let's see a few examples of how we can figure out the runtime complexity class of a piece of code.

```
def userFullName(users: List[User], id: int) -> str:
    user = users[id]
    fullName = user.firstName + " " + user.lastName
    return fullName
```

The `userFullName` function takes a list of users, and retrieves the `id`th element of the list and generates a full name from the user's first and last names. For now, we'll assume that looking up any element of any array (`users[id]`), string concatenation (`user.firstName + " " + user.lastName`), assignment (`user = ...`, and `fullName`), and returns are all constant-time operations⁴.

Under these assumptions, the first, second, and third lines can each be evaluated in constant time $\Theta(1)$. The total runtime of the function is the time required to run each line, one at a time, or:

$$T_{\text{userFullName}}(N) = \Theta(1) + \Theta(1) + \Theta(1)$$

Recall above, that $\Theta(1)$ in the middle of an arithmetic expression can be interpreted as $f(N)$ where $f(N) \in \Theta(1)$ (it is a constant). That is, $f(N) = c$. So, the above expression can be rewritten as⁵:

$$T_{\text{userFullName}}(N) = c_1 + c_2 + c_3$$

Adding three constant values together (even without knowing what they are, exactly) always gets us another constant value. So, we can say that $T_{\text{userFullName}}(N) \in \Theta(1)$.

```
def updateUsers(users: List[User]) -> None:
    x = 1
    for user in users:
        user.id = x
        x += 1
```

⁴Array lookups being constant-time is a huge simplification, called the RAM model, that we'll roll back at the end of the book.

⁵There's another simplification here. Technically, $f(N)$ is always within a bounded factor of a constant c_1 , and likewise for $g(N)$, but we'll clarify this when we get to complexity bounds below.

The `updateUsers` function takes a list of users and assigns each user a unique id. For now, we'll assume that the assignment operations (`x = 1` and `user.id`), and the increment operation (`x += 1`) all take constant ($\Theta(1)$) time. So, we can model the total time taken by the function as:

$$T_{\text{updateUsers}}(N) = O(1) + \sum_{\text{user}} (O(1) + O(1))$$

Simplifying as above, we get

$$T_{\text{updateUsers}}(N) = c_1 + \sum_{\text{user}} (c_2 + c_3)$$

Recalling the rule for summation of a constant, using N as the total number of users, and then the rule for sums of terms, we get:

$$T_{\text{updateUsers}}(N) = c_1 + N \cdot (c_2 + c_3) = \Theta(N)$$

3.7. Complexity Bounds

Not every mathematical function fits neatly into a single complexity class. Let's go to our python code example above. The `in` operator tests to see whether a particular value is present in our data. If `data` is a list, then the implementation checks every position in the list, in order. Internally, Python implements the expression `target in data` with something like:

```
def __in__(data, target):
    N = len(data)
    for i in range(N):
        if data[i] == target:
            return True
    return False
```

In the best case, the value we're looking for happens to be at the first position `data[0]`, and the code returns after a single check. In the worst case, the value we're looking for is at the last position `data[N-1]` or is not in the list at all, and we have to check every one of the N positions. Put another way, the **best case** behavior of the function is constant (exactly one check), while the **worst case** behavior is linear (N checks). We can write the resulting runtime using a case distinction:

$$T_{\text{in}}(N) = \begin{cases} a \cdot 1 + b & \text{if } \text{data}[0] = \text{target} \\ a \cdot 2 + b & \text{if } \text{data}[1] = \text{target} \\ \dots & \dots \\ a \cdot (N-1) + b & \text{if } \text{data}[N-2] = \text{target} \\ a \cdot N + b & \text{if } \text{data}[N-1] = \text{target} \end{cases}$$

We don't know the runtime exactly, as it is based on the computer and version of python we are using. However, we can model it, in general in terms of some upfront cost b (e.g., for computing $N = \text{len}(\text{data})$), and some additional cost a for every time we go through the loop (e.g., for computing `data[i] == target`). Since we don't know where the target is, exactly,

Let's do a quick experiment. The code below is like our example above, but measures the time for one lookup at a time. Each point it prints out is the runtime of a single lookup as the list gets bigger and bigger.

```
from random import randrange
from datetime import datetime
```

```
N = 100000
TRIALS = 400
```

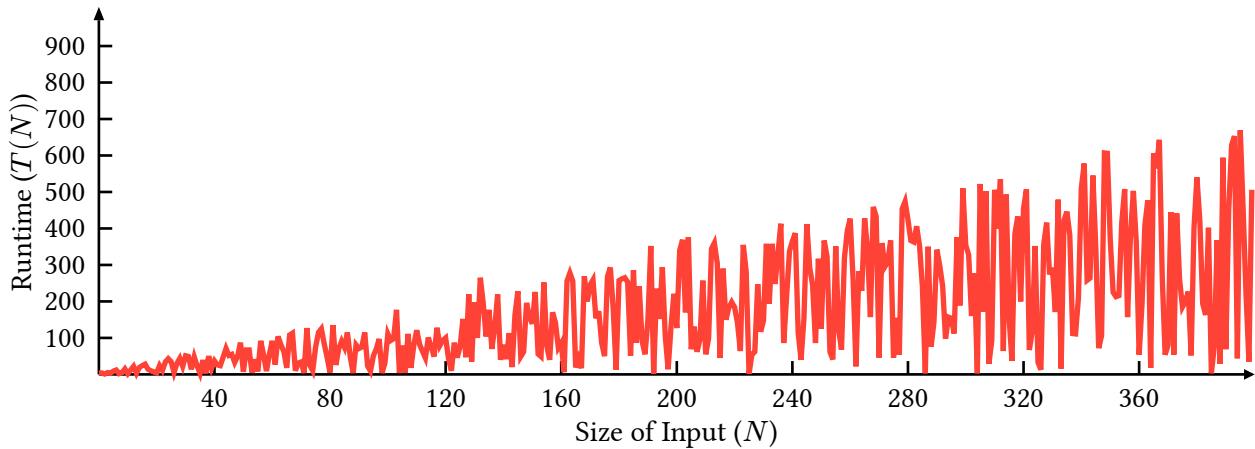
```

STEP = int(N/TRIALS)

data = list()

for i in range(TRIALS):
    # Increase the list size by STEP
    for j in range(STEP):
        data += [i * STEP + j]
    start = datetime.now()
    # Measure how long it takes to look up a random element
    if randrange(i * STEP + STEP) in data:
        pass
    end = datetime.now()
    # Print out the total time in microseconds
    microseconds = (end - start).total_seconds() * 1000000
    print(f"{i}, {microseconds}")

```

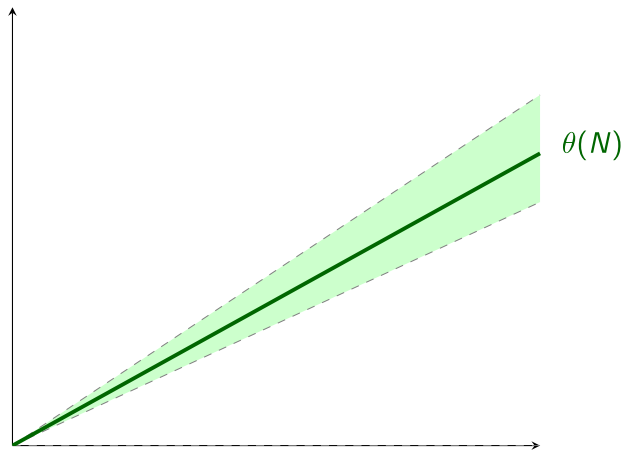


Graph 3: Scaling the List lookup.

Graph 3 shows the output of one run of the code above. You can see that it looks a lot like a triangle. The **worst case** (top of the triangle) looks a lot like the **linear** complexity class ($\Theta(N)$, or an angled line), but the **best case** (bottom of the triangle) looks a lot more like a flat line, or the **constant** complexity class ($\Theta(1)$, or a flat line). The runtime is *at least* constant, and *at most* linear: We can **bound** the runtime of the function between two complexity classes.

3.8. Big- O and Big- Ω

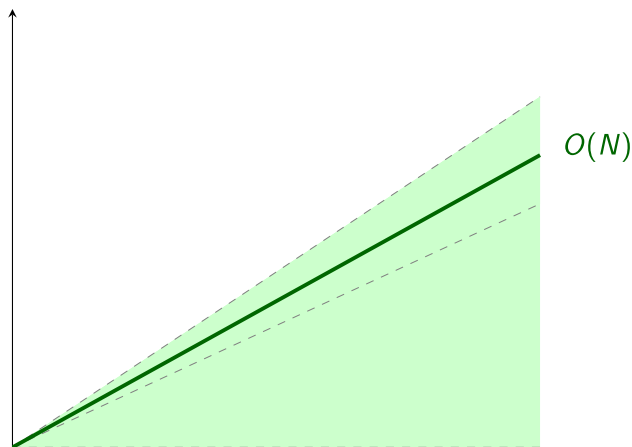
We capture this intuition of bounded runtime by introducing two new concepts: Worst-case (upper, or Big- O) and Best-case (lower, or Big- Ω) bounds. To see these in practice, let's take the linear complexity class as an example:



We write $O(N)$ to mean the set of all mathematical functions that are **no worse than** $\Theta(N)$. This includes:

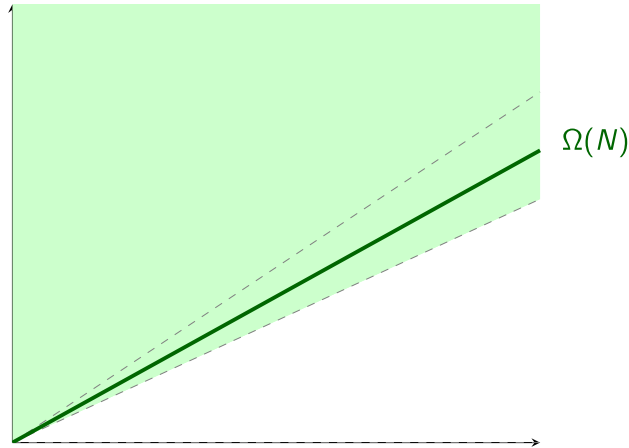
- all mathematical functions in $\Theta(N)$ itself
- all mathematical functions in lesser (slower-growing) complexity classes (e.g., $\Theta(1)$)
- any mathematical function that never grows faster than $O(N)$ (e.g., the runtime of each individual lookup in our Python example above)

The figure below illustrates $\Theta(N)$ (the dotted region) and all lesser complexity classes. Note the similarity to Graph 3.



Similarly, we write $\Omega(N)$ to mean the set of all mathematical functions that are **no better than** $\Theta(N)$. This includes:

- all functions in $\Theta(N)$ itself
- all greater (faster-growing) complexity classes (e.g., $\Theta(N^2)$), and anything in between.



To summarize, we write:

- $f(N) \in O(g(N))$ to say that $f(N)$ is in $\Theta(g(N))$ or a lesser complexity class.
- $f(N) \in \Omega(g(N))$ to say that $f(N)$ is in $\Theta(g(N))$ or a greater complexity class.

3.9. Formalizing Big-O

Before we formalize our bounds, let's first figure out what we want out of that formalism.

Let's start with the basic premise we outlined above: For a function $f(N)$ to be in the set $O(g(N))$, we want there to be some function in $\Theta(g(N))$ that is always bigger than $f(N)$.

The first problem we run into with this formalism is that we haven't really defined what exactly $\Theta(g(N))$ is yet, so we need to pin down something first. Let's start with the same assumption we made earlier: we can scale $g(N)$ by any constant value without changing its complexity class. Formally:

Formally: $\forall c : c \cdot g(N) \in \Theta(g(N))$

That is, for any constant c (\forall means 'for all'), the product $c \cdot g(N)$ is in $\Theta(g(N))$ (remember that \in means is in). This isn't meant to be all-inclusive: There are many more functions in $\Theta(g(N))$, but this gives us a pretty good range of functions that, at least intuitively, belong in $g(N)$'s complexity class.

Now we have a basis for formalizing Big-O: We can say that $f(N) \in O(g(N))$ if there is **some** multiple of $g(N)$ that is always bigger than $f(N)$. Formally:

$\exists c > 0, \forall N : f(N) \leq c \cdot g(N)$

That is, there exists (\exists means there exists) some positive constant c , such that for each value of N , the value of $f(N)$ is smaller than the corresponding value of $c \cdot g(N)$.

Let's look at some examples:

Example: $f(N) = 2N$ vs $g(N) = N$

Can we find a c and show that for this c , for all values of N , the Big-O inequality holds for f and g ?

- $f(N) \leq c \cdot g(N)$
- $2N \leq c \cdot N$
- $2 \leq c$

We start with the basic inequality, substitute in the values of $f(N)$ and $g(N)$, and then divide both sides by N . So, the inequality is always true for any value of $c \geq 2$.

Example: $f(N) = 100N^2$ vs $g(N) = N^2$

Can we find a c and show that for this c , for all values of N , the Big- O inequality holds for f and g ?

- $f(N) \leq c \cdot g(N)$
- $100N^2 \leq c \cdot N^2$
- $100 \leq c$

We start with the basic inequality, substitute in the values of $f(N)$ and $g(N)$, and then divide both sides by N . So, the inequality is always true for any value of $c \geq 100$.

Example: $f(N) = N$ vs $g(N) = N^2$

Can we find a c and show that for this c , for all values of N , the Big- O inequality holds for f and g ?

- $f(N) \leq c \cdot g(N)$
- $N \leq c \cdot N^2$
- $1 \leq c \cdot N$

Uh-oh! For $N = 0$, there is no possible value of c that we can plug into that inequality to make it true ($0 \cdot c$ is never bigger than 1 for any c).

Attempt 2: So what went wrong? Well, we mainly care about how $f(N)$ behaves for really big values of N . In fact, for the example, for any $N \geq 1$ (and $c \geq 1$), the inequality is satisfied! It's just that pesky $N = 0$!

So, we need to add one more thing to the formalism: the idea that we only care about “big” values of N . Of course, that leaves the question of how big is “big”? Now, we could pick specific cutoff values, but any specific cutoff we picked would be entirely arbitrary. So, instead, we just make the cutoff definition part of the proof: When proving that $f(N) \in O(g(N))$, we just need to show that **some** cutoff exists, beyond which $f(N) \leq c \cdot g(N)$. **The formal definition of Big- O is:**

$$f(N) \in O(g(N)) \Leftrightarrow \exists c > 0, N_0 \geq 0 : \forall N \geq N_0 : f(N) \leq c \cdot g(N)$$

This is the same as our first attempt, with only one thing added: N_0 . In other words, $f(N) \in O(g(N))$ if we can pick some cutoff value ($\exists N_0 \geq 0$) so that for every bigger value of N ($N \geq N_0$), $f(N)$ is smaller than $c \cdot g(N)$.

3.9.1. Proving a function has a specific Big- O bound

To show that a mathematical function is **in** $O(g(N))$, we need to find a c and an N_0 for which we can prove the Big- O inequality. A generally useful strategy is:

1. Write out the Big- O inequality
2. “plug in” the values of $f(N)$ and $g(N)$
3. “Solve for” c , putting it on one side of the inequality, with everything else on the other side.
4. Try a safe default of $N_0 = 1$.
5. Use the $A \leq B$ and $B \leq C$ imply $A \leq C$ trick (transitivity of inequality) to replace any term involving N with N_0

⁶Recall that we're only allowing non-negative input sizes (i.e., $N \geq 0$), so negative values of N aren't a problem.

5. Use the resulting inequality to find a lower bound on c

Continuing the above example:

- $f(N) \leq c \cdot g(N)$
- $N \leq c \cdot N^2$
- $\frac{1}{N} \leq c$

For that last step, we have $N_0 \leq N$, or $1 \leq N$, so dividing both sides by N , we get $\frac{1}{N} \leq 1$. So, if we pick $1 \leq c$, then $\frac{1}{N} \leq 1 \leq c$, and $\frac{1}{N} \leq c$.

3.9.2. Proving a function does not have a specific Big-O bound

To show that a mathematical function is **not in** $O(g(N))$, we need to prove that there can be **no** c or N_0 for which we can prove the Big-O inequality. A generally useful strategy is:

1. Write out the Big-O inequality
2. “plug in” the values of $f(N)$ and $g(N)$
3. “Solve for” c , putting it on one side of the inequality, with everything else on the other side.
4. Simplify the equation on the opposite side and show that it is strictly growing. Generally, this means that the right-hand-side is in a complexity class at least N .

Flipping the above example:

- $f(N) \leq c \cdot g(N)$
- $N^2 \leq c \cdot N$
- $N \leq c$

N is strictly growing: for bigger values of N , it gets bigger. There is no constant that can upper bound the mathematical function N .

3.10. Formalizing Big-Ω

Now that we’ve formalized Big-O (the upper bound), we can formalize Big-Ω (the lower bound) in exactly the same way:

$$f(N) \in O(g(N)) \Leftrightarrow \exists c > 0, N_0 \geq 0 : \forall N \geq N_0 : f(N) \geq c \cdot g(N)$$

The only difference is the direction of the inequality: To prove that a function exists in Big-Ω, we need to show that $f(N)$ is bigger than some constant multiple of $g(N)$.

3.11. Formalizing Big-Θ

Although we started with an intuition for Big-Θ, we haven’t yet formalized it. To understand why, let’s take a look at the following runtime:

$$T(N) = 10N + \text{rand}(10)$$

Here $\text{rand}(10)$ means a randomly generated number between 0 and 10 for each call. If the function were **just** $10N$, we’d be fine in using our intuitive definition of $\Theta(N)$ being all multiples of N . However, this function still “behaves like” $g(N) = N$... just with a little random noise added in. For big values of N (e.g., 10^{10}), the random noise is barely perceptible. Although we can’t say that $T(N)$ is **equal to** some multiple $c \cdot N$, we can say that it is **close to** that multiple (in fact, it’s always between $10N$ and $10N + 10$). In other words, we can bound it from both above and below!

Let's try proving this with the tricks we developed for Big- O and Big- Ω :

- $T(N) \leq c_{\text{upper}} \cdot N$
- $10N + \text{rand}(10) \leq c_{\text{upper}} \cdot N$ (plug in $T(N)$)
- $10 + \frac{\text{rand}(10)}{N} \leq c_{\text{upper}}$ (divide by N)

Looking at this formula, we can make a few quick observations. First, by definition $\text{rand}(10)$ is never bigger than 10. Second, if $N_0 = 1$, $\frac{1}{N}$ can never be bigger than 1. In other words, $\frac{\text{rand}(10)}{N}$ can not be bigger than 10. Let's prove that to ourselves.

Taking the default $N_0 = 1$ we get:

- $1 \leq N$ (plug in N_0)
- $10 \leq 10N$ (multiply by 10)
- $\text{rand}(10) \leq 10 \leq 10N$ (transitivity with $\text{rand}(10) \leq 10$)
- $\frac{\text{rand}(10)}{N} \leq 10$ (divide by N)
- $10 + \frac{\text{rand}(10)}{N} \leq 10 + 10$ (add 10)

So, if we pick $c_{\text{upper}} \geq 20$, we can show (again, by transitivity):

$$10 + \frac{\text{rand}(10)}{N} \leq c_{\text{upper}}$$

Which gets us $T(N) \leq c_{\text{upper}} \cdot N$ for all $N > N_0$.

Now let's try proving a lower (Big- Ω) bound:

- $T(N) \geq c_{\text{lower}} \cdot N$
- $10N + \text{rand}(10) \geq c_{\text{lower}} \cdot N$ (plug in $T(N)$)
- $10 + \frac{\text{rand}(10)}{N} \geq c_{\text{lower}}$ (divide by N)
- $10 \geq c_{\text{lower}}$ (By transitivity: $10 \geq 10 + \frac{\text{rand}(10)}{N}$)

This inequality holds for any $10 \geq c_{\text{lower}} > 0$ (recall that c has to be strictly bigger than zero).

So, we've shown that $T(N) \in O(N)$ **and** $T(N) \in \Omega(N)$. The new thing is that we've shown that the upper and lower bounds **are the same**. That is, we've shown that $T(N) \in O(g(N))$ and $T(N) \in \Omega(g(N))$ **for the same mathematical function** g . If we can prove that an upper and lower bound for some mathematical function $f(N)$ that is the same mathematical function $g(N)$, we say that $f(N)$ and $g(N)$ are in the same complexity class. Formally, $f(N) \in \Theta(g(N))$ if and only if $f(N) \in O(g(N))$ **and** $f(N) \in \Omega(g(N))$.

3.12. Tight Bounds

In the example above, we said that $\text{rand}(10) \leq 10$. We could have just as easily said that $\text{rand}(10) \leq 100$. The latter inequality is just as true, but somehow less satisfying; yes, the random number will always be less than 100, but we can come up with a "tighter" bound (i.e., 10).

Similarly Big- O and Big- Ω are bounds. We can say that $N \in O(N^2)$ (i.e., N is no worse than N^2). On the other hand, this bound is just as unsatisfying as $\text{rand}(10) \leq 100$, we can do better.

If it is not possible to obtain a better Big- O or Big- Ω bound, we say that the bound is **tight**. For example:

- $10N^2 \in O(N^2)$ and $10N^2 \in \Omega(N^2)$ are tight bounds.
- $10N^2 \in O(2^N)$ is correct, but **not** a tight bound.

- $10N^2 \in \Omega(N)$ is correct, but **not** a tight bound.

Note that since we define Big- Θ as the intersection of Big- O and Big- Ω , all Big- Θ bounds are, by definition tight. As a result, we sometimes call Big- Θ bounds “tight bounds”.

3.13. Which Variable?

We define asymptotic complexity bounds in terms of **some** variable, usually the size of the collection N . However, it’s also possible to use other bounds. For example, consider the function, which computes factorial:

```
public int factorial(int idx)
{
    if(idx <= 0){ return 0; }
    int result = 1;
    for(i = 1; i <= idx; i++) { result *= i; }
    return result
}
```

The runtime of this loop depends on the input parameter `idx`, performing one math operation for each integer between 1 and `idx`. So, we could give the runtime as $\Theta(\text{idx})$.

When the choice of variable is implicit, it’s customary to just use N , but this is not always the case. For example, when we talk about sequences and lists, the size of the sequence/list is most frequently the variable of interest. However, there might be other parameters of interest:

- If we’re searching the list for a specific element, what position to we find the element at?
- If we’re looking through a linked list for a specific element, what index are we looking for?
- If we have two or more lists (e.g., in an Edge List data structure), each list may have a different size.

In these cases, and others like them, it’s important to be clear about which variable you’re talking about.

3.13.1. Related Variables

When using multiple variables, we can often bound one variable in terms of another. Examples include:

- If we’re looking through a linked list for the element at a specific index, the index must be somewhere in the range $[0, N)$, where N is the size of the list. As a result, we can always replace $O(\text{index})$ with $O(N)$ and $\Omega(\text{index})$ with $\Omega(1)$, since `index` is bounded from above by a linear function of N and from below by a constant.
- The number of edges in a graph can not be more than the square of the number of vertices. As a result, we can always replace $O(\text{edges})$ with $O(\text{vertices}^2)$ and $\Omega(\text{edges})$ with $\Omega(1)$.

Note: Even though $O(\text{index})$ in the first example may be a tighter bound than $O(N)$, the $O(N)$ bound is still tight **in terms of N** : We can not obtain a tighter bound that is a function only of N .

3.14. Summary

We defined three ways of describing runtimes (or any other mathematical function):

- Big- O : The worst-case complexity:
 - $T(N) \in O(g(N))$ means that the runtime $T(N)$ scales **no worse than** the complexity class of $g(N)$

- Big- Ω : The best-case complexity
 - $T(N) \in \Omega(g(N))$ means that the runtime $T(N)$ scales **no better than** the complexity class of $g(N)$
- Big- Θ : The tight complexity
 - $T(N) \in \Theta(g(N))$ means that the runtime $T(N)$ scales **exactly as** the complexity class of $g(N)$

We'll introduce amortized and expected runtime bounds later on in the book; Since these bounds are given without qualifiers, and so are sometimes called the **Unqualified** runtimes.

3.14.1. Formal Definitions

For any two functions $f(N)$ and $g(N)$ we say that:

- $f(N) \in O(g(N))$ if and only if $\exists c > 0, N_0 \geq 0 : \forall N > N_0 : f(N) \leq c \cdot g(N)$
- $f(N) \in \Omega(g(N))$ if and only if $\exists c > 0, N_0 \geq 0 : \forall N > N_0 : f(N) \geq c \cdot g(N)$
- $f(N) \in \Theta(g(N))$ if and only if $f(N) \in O(g(N))$ and $f(N) \in \Omega(g(N))$

Note that a simple $\Theta(g(N))$ may not exist for a given $f(N)$, specifically when the tight Big- O and Big- Ω bounds are different.

3.14.2. Interpreting Code

In general⁷, we will assume that most simple operations: basic arithmetic, array accesses, variable access, string operations, and most other things that aren't function calls will all be $\Theta(1)$.

Other operations are combined as follows...

Sequences of instructions

```
{
  op1;
  op2;
  op3;
}
...
```

Sum up the runtimes. $T(N) = T_{\text{op1}}(N) + T_{\text{op2}}(N) + T_{\text{op3}}(N) + \dots$

Loops

```
for(i = min; i < max; i++)
{
  block;
}
```

Sum up the runtimes for each iteration. Make sure to consider the effect of the loop variable on the runtime of the inner block. $T(N) = \sum_{i=\text{min}}^{\text{max}} T_{\text{block}}(N, i)$

As a simple shorthand, if (i) the number of iterations is predictable (e.g., if the loop iterates N times) and (ii) the complexity of the loop body is independent of which iteration the loop is on (i.e., i does not appear in the loop body), you can just multiply the complexity of the loop by the number of iterations.

⁷All of these are lies. The cost of basic arithmetic is often $O(\log N)$, array access runtimes are affected by caching (we'll address that later in the book), and string operations are proportional to the length of the string. However, these are all useful simplifications for now.

Conditionals

```
if(condition){  
    block1;  
} else {  
    block2;  
}
```

The total runtime is the cost of either `block1` or `block2`, depending on the outcome of `condition`. Make sure to add the cost of evaluating `condition`.

$$T(N) = T_{\text{condition}(N)} + \begin{cases} T_{\text{block1}}(N) & \text{if condition is true} \\ T_{\text{block2}}(N) & \text{otherwise} \end{cases}$$

The use of a `cases` block is especially important here, since if $T_{\text{block1}}(N)$ and $T_{\text{block2}}(N)$ belong to different asymptotic complexity classes, the overall block of code belongs to multiple classes (and thus does not have a simple Θ bound).

3.14.3. Simple Complexity Classes

We will refer to the following specific complexity classes:

- **Constant:** $\Theta(1)$
- **Logarithmic:** $\Theta(\log N)$
- **Linear:** $\Theta(N)$
- **Loglinear:** $\Theta(N \log N)$
- **Quadratic:** $\Theta(N^2)$
- **Cubic:** $\Theta(N^3)$
- **Exponential:** $\Theta(2^N)$

These complexity classes are listed in order.

3.14.4. Dominant Terms

In general, any function that is a sum of simpler functions will be dominated by one of its terms. That is, for a polynomial:

$$f(N) = f_1(N) + f_2(N) + \dots + f_k(N)$$

The asymptotic complexity of $f(N)$ (i.e., its Big- O and Big- Ω bounds, and its Big- Θ bound, if it exists) will be the **greatest** complexity of any individual term $f_{i(N)}$ ⁸.

Remember: If the dominant term in a polynomial belongs to a single simple complexity class, then the entire polynomial belongs to this complexity class, and the Big- O , Big- Ω , and Big- Θ bounds are all the same.

3.14.5. Multiclass Asymptotics

A mathematical function may belong to multiple simple classes, depending on an unpredictable input or the state of a data structure. Generally, multiclass functions arise in one of two situations. First, the branches of a conditional may have different complexity classes:

$$T(N) = \begin{cases} T_1(N) & \text{if a thing is true} \\ T_2(N) & \text{otherwise} \end{cases}$$

⁸Note that this is only true when k is fixed. If the number of polynomial terms depends on N , we need to consider the full summation.

If $T_1(N)$ and $T_2(N)$ belong to different complexity classes, then $T(N)$ as a whole belongs to **either** class. In this case, we can only bound the runtime $T(N)$. Specifically, if $\Theta(T_1(N)) > \Theta(T_2(N))$, then:

- $T(N) \in \Omega(T_2(N))$
- $T(N) \in O(T_1(N))$.

Second, the number of iterations of a loop may depend on an input that is bounded by multiple complexity classes. For example, if $\text{idx} \in [1, N]$ (idx is somewhere between 1 and N , inclusive), then the following code does not belong to a single complexity class:

```
for(i = 0; i < idx; i++){ do_a_thing(); }
```

In this case, we can bound the runtime based on idx . Assuming $\text{do_a_thing}()$ is $\Theta(1)$, then $T(N) \in \Theta(\text{idx})$. However, since we can't bound idx in terms of N , then we can only provide weaker bounds with respect to N :

- $T(N) \in \Omega(1)$
- $T(N) \in O(N)$

Remember that if we can not obtain identical, tight upper and lower bounds in terms of a given input variable, there is no simple Θ -bound in terms of that variable.

Chapter 4. *The Sequence and List ADTs*

Now that we have the right language to talk about algorithm runtimes (asymptotic complexity), we can start talking about actual data structures that these algorithms can use. Since the choice of data structure can have a huge impact on the complexity of an algorithm, it's often useful to group specific data structures together by the roles that they can fulfill in an algorithm. We refer to such a grouping as an **abstract data type** or ADT.

4.1. What is an ADT?

An ADT is, informally, a contract that states what we can expect from a data structure. Take, for example, a Java interface like the following:

```
public interface Narf<T>
{
    public void foo(T poit, int zort);
    public T bar(int zort);
    public int baz();
}
```

This interface states that any class that implements Moof must provide `foo`, `bar`, and `baz` methods, with arguments and return values as listed above. This is not especially helpful, since it doesn't give us any idea of what these methods are supposed to do. Contrast Moof with the following interface:

```
public interface MutableSequence<T>
{
    public void update(int index, T element);
    public T get(int index);
    public int size();
}
```

Sequence is semantically identical to Narf, but far more helpful. Just by reading the names of these methods, you get some idea of what each method does, and how it interacts with the state represented by a Sequence. You can build a mental model of what a Sequence is from these names.

An Abstract Data Type is:

1. A mental model of some sort of data (a data type)
2. One or more operations for accessing or modifying that state (an interface)
3. Any other rules for the state (constraints)

For most ADTs discussed in this book, the ADT models sort of collection of elements. The difference between them is how you're allowed to interact with the data.

A data structure is a **specific** strategy for organizing the data modeled by an ADT. We say that the data structure implements (or conforms to) an ADT if:

1. The data structure stores the same type of data as the ADT
2. Each of the operations required by the ADT can be implemented on the data structure
3. We can prove that the operation implementation is guaranteed to respect the rules for the ADT's state.

4.2. The Sequence ADT

A very common example of an ADT is a sequence. Examples include:

- The English alphabet ('A', 'B', ..., 'Z')
- The Fibonacci Sequence (1, 1, 2, 3, 5, 8, ...)
- An arbitrary collection of numbers (42, 19, 86, 23, 19)

What are some commonalities in these examples?

- Every sequence is a collection of elements of some type T (integer, character, etc...)
- The same element may appear multiple times.
- Every sequence has a size (that may be infinite). We often write N to represent this size.
- Every element (or occurrence of an element) is assigned to an index: $0 \leq \text{index} < N$.
- Every index in the range $0 \leq \text{index} < N$ has exactly one element assigned to it.

What kind of operations can we do on a sequence⁹?

```
public interface Sequence<T>
{
    /** Retrieve the element at a specific index or throw an IndexOutOfBoundsException if
    index < 0 or index >= size() */
    public T get(int index);
    /** Obtain the size of the sequence */
    public int size();
}
```

To recap, we have a mental model, and a series of rules:

- A sequence contains N elements.
- Every index i from $0 \leq i < N$ identifies exactly one element (no more, no less).

The `get` method returns the element identified by `index`, and the `size` method returns N .

4.2.1. Sequences by Rule

For some of the example sequences listed above, we can implement this interface directly. For example, for the English alphabet:

```
public class Alphabet implements Sequence<Character>
{
    /** Retrieve the element at a specific index */
    public Character get(int index)
    {
        if(index == 0){ return 'A'; }
        if(index == 1){ return 'B'; }
        /* ... */
        if(index == 25){ return 'Z'; }
        throw IndexOutOfBoundsException("No character at index "+index)
    }
    /** Obtain the size of the sequence */
    public int size()
    { return 26; }
}
```

Is this a sequence?

⁹Note: Although several of the ADTs and data structures we'll present throughout this book correspond to actual Java interfaces and classes, `Sequence` is **not** one of them. However, it serves as a useful simplification of the `List` interface that we'll encounter later on.

- size: There are $N = 26$ elements.
- get: Every element from $0 \leq i < 26$ is assigned exactly one value.

Similarly, we can implement the Fibonacci sequence according to the Fibonacci rule $\text{Fib}(x) = \text{Fib}(x - 1) + \text{Fib}(x - 2)$:

```
public class Fibonacci implements Sequence<Int>
{
    /** Retrieve the element at a specific index */
    public Int get(int index)
    {
        if(index < 0){
            throw IndexOutOfBoundsException("Invalid index: "+index)
        }
        if(index == 0){ return 0; }
        if(index == 1){ return 1; }
        return get(index-1) + get(index-2);
    }
    /** Obtain the size of the sequence */
    public int size()
    { return INFINITY; }
}
```

Is this a sequence?

- size: There are $N = \infty$ elements.
- get: Every element from $0 \leq i < \infty$ is assigned exactly one value derived from the value of preceding elements of the sequence as $\text{Fib}(i - 1) + \text{Fib}(i - 2)$ ¹⁰.

4.2.2. Arbitrary Sequences

Often, however, we want to represent an sequence of **arbitrary** elements. We won't have simple rule we can use to translate the index into a value. Instead, we need to store the elements somewhere if we want to be able to retrieve them. The easiest way to store a sequence, when we know the size of the sequence upfront is called an **array**.

Let's take a brief diversion into the guts of a computer. Most computers store data in a component called RAM¹¹. You can think of RAM as being a ginormous sequence of bits (0s and 1s). Folks who work on computers have, over the course of many years, come to agree on some common guidelines for how to represent other types of information in bit sequences. The details of most of these guidelines (e.g., Little- vs Big-endian, Floating-point encodings, Characters and strings, etc...) are usually covered in a computer organization or computer architecture class. However, there's a few details that are really helpful for us to review.

First, we group every 8 bits into a 'byte'. A byte can take on $2^8 = 256$ different possible values. It turns out that this is enough to represent most simple printable characters. Similarly, bytes can be grouped

¹⁰If you're paying attention, you might notice that we're waving our hands a bit with this proof. Specifically, how do we convince ourselves that there is **exactly** one value at index i , and not zero (can $\text{Fib}(i)$ return `IndexOutOfBoundsException` for $i \geq 0$), or more than 1 (will $\text{Fib}(i)$ always return the same value)? Proving either of these will require a technique called **induction** that we'll come back to later in the book.

¹¹What you're about to read is called the RAM model of computing. It is a blatant lie: The way computers store and access data involves multiple levels of caches, disks, and networked storage. However, the RAM model is quite useful as a starting point. We'll walk the lie back when we introduce the External Memory model towards the end of the book.

together into 2, 4, or 8 bytes, allowing us to represent more complex data like emoji, larger integers, negative numbers, or floating point numbers. Strings are a special case that we'll come back to in a moment.

For example, the ASCII character encoding¹² forms the basis for most character representations we use today, and assigns each of the letters of the English alphabet (upper and lower case, and a few other symbols) to a particular sequence of 8 bits (one byte).

Most modern computers are 'byte-addressable', meaning that every byte is given its own address. The first byte is at address 0000, followed by address 0001, then address 0002, and so forth. If we want to represent a value that takes up multiple bytes (e.g., a 4-byte integer), we typically point at the address of the first byte of the value.

4.2.3. Arrays

Let's say we have a series of values we want to store. For example, let's say we wanted to store the sequence 'H', 'e', 'l', 'l', 'o'. The ASCII code for 'H' is hexadecimal 0x48, or binary 01001000. We could store that value at one specific address in RAM. The value takes up one byte of space, so we could put the second value in the sequence 'e' (hexadecimal 0x65, or binary 10101001) in the byte right after it. We then similarly store the remaining three elements of the sequence, each in the byte after the previous one.

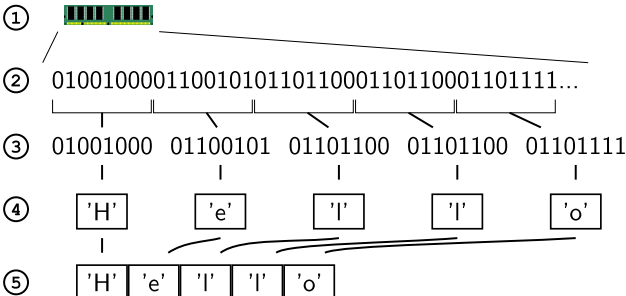


Figure 5: RAM (1) can be viewed as a sequence of bits (2). We can break these bits up into fixed size chunks (3). Each fixed size chunk can be used to identify some value, like for example a character (4). Thus a sequence of characters can be stored somewhere in RAM (5).

This arrangement, where we just concatenate elements side by side in memory is really useful, *if each element always uses up the same number of bytes*. Specifically, if we want to retrieve the *i*th element of the sequence, we need only two things:

1. The position of the 0th element in RAM (Let's call this *S*)
2. The number of bytes that each element takes up (Let's call this *E*).

The *i*th element's *E* bytes will always start at byte $S + i \cdot E$ (with 0 being the first element).

4.2.3.1. Array Runtime Analysis

In order to retrieve the *i*'th element of an array, we need one multiplication, and one addition. This is true regardless of how big the array is. If the array has a thousand, or a million, or a trillion elements, we can still obtain the address of the *i*th element with one multiplication and one addition. Given the

¹²The American Standard Code for Information Interchange forms the basis for most modern character 'encodings'. More recent encodings (e.g. UTF-8, UTF-16) extend the ASCII code, and sometimes require more than one byte per printable character.

address, we can then obtain the value in a single access to memory. In other words, the cost of retrieving an element from an array is constant (i.e., $\Theta(1)$)¹³.

4.2.3.2. Side Note: Memory Allocation

Most operating systems provide a way to allocate contiguous regions of memory (in C, this operation is called `malloc`, in Java it is called `new`). Eventually, when the data is no longer necessary (C's `free` operation, or based on Java's automatic 'garbage collector'), the allocated memory is returned to the general pool of 'free' memory. We will assume that both allocating and freeing memory are constant-time ($\Theta(1)$) operations¹⁴

It's important to note that once memory is allocated, it has a fixed size. It is not generally possible to simply resize a previously allocated chunk of memory, since that memory is located at a specific position in RAM, and it is possible that the space immediately following the array may already be in-use¹⁵.

4.2.3.3. Arrays In Java

Java provides a convenient shorthand for creating and accessing arrays using square brackets. To instantiate an array, use `new`:

```
int[] myArray = new int[100];
```

Note that the type of an integer array is `int[]`¹⁶, and the number of elements in the array must be specified upfront (100 in the example above). To access an element of an array use `array[index]`.

```
myArray[10] = 29;
int myValue = myArray[10];
assert(myValue == 29);
```

Java arrays are **bounds-checked**. That is, Java stores the size of the array as part of the array itself. The array actually starts 4 bytes after the official address of the array; these bytes are used to store the size of the array. Whenever you access an array element, Java checks to make sure that the index is greater than or equal to zero and less than the size. If not, it throws an `IndexOutOfBoundsException` otherwise. As a convenient benefit, bounds checking means we can get the array size from Java.

```
int size = myArray.size
assert(size == 100)
```

Since the size is stored at a known location, we can always retrieve it in constant ($O(1)$) time.

¹³If you're paying close attention, you might note that we still need to retrieve E bytes, and so technically the cost of an array access is $\Theta(E)$. However, we're mainly interested in how runtime scales with the size of the collection, as E is fixed upfront, and usually very small. So, more precisely, **with respect to N** , the cost of retrieving an element from an array is $\Theta(1)$. The even more attentive reader might be aware of things like caches, which as previously mentioned, we're going to ignore for now.

¹⁴Again, this is a lie. Depending on the operating system's implementation (e.g., if it allocates pages lazily or not), allocating memory may require linear time in the size of the allocation, or non-constant (e.g., logarithmic or linear) time in the number of preceding allocations. Nevertheless, for the rest of the book, we'll treat it as being constant

¹⁵C provides a `realloc` operation that opportunistically *tries* to extend an allocated chunk of memory if space exists for it. However, if this is not possible, the entire allocation will be copied byte-for-byte to a new position in memory where space does exist. As we'll emphasize shortly, copying an array is a linear ($\Theta(N)$) operation.

¹⁶Java sort of allows you to store variable size objects in an array. For example, `String[]` is a valid array type. The trick to this is that Java allocates the actual string *somewhere else in RAM* called the heap. What it stores in the actual array is the address of the string (also called a pointer).

To prove to ourselves that the Array implements the Sequence ADT, we can implement its methods:

```
public class ArraySequence<T> implements Sequence<T>
{
    T[] data;

    public ArraySequence(T[] data){ this.data = data; }
    public T get(int index) { return data[index]; }
    public int size() { return data.size; }
}
```

The ArraySequence, initialized by an array, follows the rules for a Sequence:

- size: There are $N = \text{data.size}$ elements in the array
- get: The i th element of the sequence is `data[i]`.

4.2.4. Mutable Sequences

The Fibonacci sequence and the English alphabet are examples of **immutable** sequences. Immutable sequences are pre-defined and can not be changed; We can't arbitrarily decide that the 10th Fibonacci number should instead be 3. However, if the sequence is given explicitly (e.g., as an array), then it's physically possible to just modify the bytes in the array to a new value. To accommodate such edits, we can make our Sequence ADT a little more general by adding a way to modify its contents. We'll call the resulting ADT a MutableSequence¹⁷.

```
public interface MutableSequence<T> extends Sequence<T>
{
    public void update(int index, T element)
}
```

The extends keyword in java can be used to indicate that an interface takes in all of the methods of the extended interface, and potentially adds more of its own. In this case MutableSequence has all of the methods of Sequence, plus its own update method.

A Mutable Sequence introduces one new rule:

- After we call `update(i, value)` with a valid index (i.e., $0 \leq i < N$), every following call to `get(i)` for the same index will return the value passed to update (until the next time index i is updated).

4.2.4.1. Mutable Array

To prove to ourselves that the array implements the MutableSequence ADT, we can implement its methods:

```
public class MutableArraySequence<T> implements MutableSequence<T>
{
    T[] data;

    public ArraySequence(int size){ this.data = new T[size]; }
    public void update(int index, T element) { data[index] = value; }
    public T get(int index) { return data[index]; }
    public int size() { return data.size; }
}
```

¹⁷Note: Like Sequence, the MutableSequence is not a native part of Java. Its role is subsumed by List, which we'll discuss shortly.

As before, we can show that the array satisfies the rules on `get` and `size`, leaving the new rule:

- `update`: Calling `update` overwrites the array element at `data[index]` with `value`, which is the value returned by `get(index)`.

4.2.5. Array Summary

Observe that each of the three `MutableSequence` methods are implemented in a single, primitive operation. Thus:

- `get`: Retrieving an element of an Array (i.e., `array[index]`) is $\Theta(1)$
- `update`: Updating an element of an array (i.e., `array[index] = ...`) is $\Theta(1)$
- `size`: Retrieving the array's size (i.e., `array.size`) is $\Theta(1)$

Recall that `size` accesses a pre-computed value, stored with the array in Java.

4.3. The List ADT

Although we can change the individual elements of an array, once it's allocated, the size of the array is fixed. This is reflected in the `MutableSequence` ADT, which does not provide a way to change the sequence's size. Let's design our next ADT by considering how we might want to change an array's size:

- Inserting a new element at a specific position
- Removing an existing element at a specific position

It's also useful to treat inserting at/removing from the front and end of the sequence as special cases, since these are both particularly common operations.

We can summarize these operations in the `List` ADT¹⁸:

```
public interface List<T> extends MutableSequence<T>
{
    /** Append an element to the end of a list */
    public void add(T element);
    /** Insert an element at an arbitrary position */
    public void add(int index, T element);
    /** Remove an element at a specific index */
    public void remove(int index);

    // ... and more operations that are not relevant to us for now.
}
```

`List` brings with it a new set of rules:

- After calling `add(index, element)` with a valid $0 \leq \text{index} \leq N$ (note that N is an allowable index):
 - Every element previously at an index $i \geq \text{index}$ will be moved to index $i + 1$
 - The value `element` will be the new element at index `index`,
 - `size()` will increase by 1.
- After calling `remove(index)` with a valid $0 \leq \text{index} < N$:
 - Every element previously at an index $i > \text{index}$ will be moved to index $i - 1$
 - The value previously at index `index` will be removed

¹⁸See Java's `List` interface for the full list of operations supported on `Lists`. Most of these operations are so-called syntactic sugar on top of these basic operations, offering cleaner code, but no new functionality.

- ▶ `size()` will decrease by 1.

Note that calling `add(element)` is the same as calling `add(element, size())`.

4.3.1. A simple Array as a List

We can still use Arrays to implement the List ADT. However, recall that it's not (generally) possible to resize a chunk of memory once it's been allocated. Since we can't change the size of an Array once it's allocated¹⁹, we'll need to allocate an entirely new array to store the updated list. Once we allocate the new array, we'll need to copy over everything in our original array to the new array.

```
public class SimpleArrayAsList<T> extends MutableArraySequence implements List<T>
{
    // data, get, update, size() inherited from MutableArraySequence

    public void add(T element){ add(size(), element); }
    public void add(int index, T element)
    {
        // Skipped: Check that index is in-bounds.
        T[] newData = new data[size() + 1];
        for(i = 0; i < newData.size; i++){
            if(i < index){ newData[i] = data[i]; }
            else if(i == index){ newData[i] = element; }
            else { newData[i] = data[i-1]; }
        }
        data = newData;
    }
    public void remove(int index)
    {
        // Skipped: Check that index is in-bounds.
        T[] newData = new data[size() - 1];
        for(i = 0; i < newData.size; i++){
            if(i < index){ newData[i] = data[i]; }
            else { newData[i] = data[i+1]; }
        }
    }
}
```

Does this satisfy our rules?

- `add`: `newData` is one larger, elements at positions `index` and above are shifted right by one position, and `element` is inserted at position `index`.
- `remove`: `newData` is one smaller, and elements at positions `index` and above are shifted left by one position.

Let's look at the runtime of the `add` method:

- We'll assume that memory allocation is constant-time ($\Theta(1)$)²⁰.
- We already said that array updates and math operations are constant-time ($\Theta(1)$)

¹⁹Again, the C language has a method called `realloc` that can **sometimes** change the size of an array... if you're lucky and the allocator happens to have some free space right after the array. However, in this book we try to avoid relying purely on unconstrained luck.

²⁰Lies! Lies and trickery! Memory allocation may require zeroing pages, multiple calls into the kernel, page faults, and a whole mess of other nonsense that scale with the size of the memory allocated. Still, especially for our purposes here, it's usually safe to assume that the runtime of memory allocation is a bounded constant.

So, we can view the the add method as

```
public void add(int index, T element)
{
    /* Theta(1) */
    for(i = 0; i < newData.size; i++)
    {
        if(/* Theta(1) */){ /* Theta(1) */ }
        else if(/* Theta(1) */){ /* Theta(1) */ }
        else { /* Theta(1) */ }
    }
}
```

Since every branch of the `if` inside the loop is the same, we can simplify the loop body to just $\Theta(1)$.

Recalling how we use Θ bounds in an arithmetic expression, we can rewrite the runtime and simplify it as:

- $T_{\text{add}(N)} = \Theta(1) + \sum_{i=0}^{\text{newData.size}} \Theta(1)$
- $= \Theta(1) + \sum_{i=0}^{N+1} \Theta(1)$ (`newData` is one bigger than the original N)
- $= \Theta(1) + (N + 2) \cdot \Theta(1)$ (Plugging in the formula for summation of a constant)
- $= \Theta(1 + N + 2)$ (Merging Θ s)
- $= \Theta(N)$ (Dominant term)

If you analyze the runtime of the `remove` method similarly, you'll likewise get $\Theta(N)$.

4.4. Linked Lists

$\Theta(N)$ is not a particularly good runtime for simple “building block” operations like `add` and `remove`. Since these will be called in loops, the `Array` data structure is not ideal for situations where a `List` is required.

The main difficulty with the `Array` is that the entire list is stored in the same memory allocation. A single chunk of memory holds all elements in the list. So, instead of allocating one chunk for the entire list, we can go to the opposite extreme and give each element its own chunk.

Giving each element its own chunk of memory means that we can allocate (or free) space for new elements without having to copy existing elements around. However, it also means that the elements of the list are scattered throughout RAM. If we want to be able to find the i th element of the list (which we need to do to implement `get`), we need some way to keep track of where the elements are stored. One approach would be to keep a list of all N addresses somewhere, but this brings us back to our original problem: we need to be able to store a variable-size list of N elements.

Another approach is to use the chunks of memory as part of our lookup strategy: We can have each chunk of memory that we allocate store the address of (i.e., a pointer to) the **next** element in the list. That way, we only need to keep track of a single address: the position of the **first** element of the list (also called the list head). The resulting data structure is called a linked list. Figure 6 contrasts the linked list approach with an `Array`.

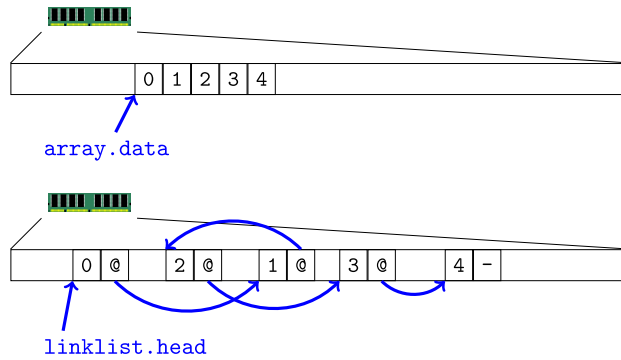


Figure 6: Instead of allocating a fixed unit of memory like an array, a linked list consists of chunks of memory scattered throughout RAM. Each element gets one chunk that has a pointer to the next (and/or previous) element in the sequence.

4.4.1. Side Note: How Java uses Memory

Java represents chunks of memory as classes, so we can implement a linked list by defining a new class that contains a value, and a reference to the next node on the list.

It's important here to take a step back and think about how Java represents chunks of memory. Specifically, primitive values (`int`, `long`, `float`, `double`, `byte`, and `char`) are stored **by value**. This means that if you have a variable with type `long`, that variable directly stores a number. If you assign the value to a new variable, the new variable stores a **copy** of the number.

```
long x = 12;
long y = x;
x = x + 5;
System.out.println(x); // prints 17
System.out.println(y); // prints 12
```

On the other hand, objects (anything that extends `Object`) and arrays are stored **by reference**. This means that what you're storing is the *address* of the actual object value (aka a pointer). When you assign the object to a new variable, what you're *really* doing is copying the address of the object or array; both variables will now point to the **same** object.

```
int[] x = new [1, 2, 3, 4, 5];
int[] y = x;
x[2] = 10;
System.out.println(y[1]); // prints 2
System.out.println(y[2]); // prints 10
System.out.println(y[3]); // prints 4
```

4.4.2. A Linked List as a List

We refer to the chunk of memory that we allocate for each element of a linked list as a `Node`. In java, we can allocate `Node` objects if we define them as a class, here storing the element represented by the `Node`, and a reference (pointer) to the next `Node`. We use java's `Optional` type to represent situations where a `Node` is not present²¹

²¹An `Optional` extends existing types to indicate that the value may be present or absent (empty). Although java already allows for object variables to take a `null` value, using `Optional` instead makes it easier for people writing code to know when it's necessary to explicitly check for a missing value (e.g., `value.isPresent()`). Since `Optional` was added to

A simple linked list, also known as a ‘singly’ linked list, is just a reference to a head node (which is empty if the list is empty). The last element of the list has an empty next.

```
public class SinglyLinkedList<T> implements List<T>
{
    class Node
    {
        T value;
        Optional<Node> next = Optional.empty();
        public this(T value) { this.value = value; }
    }
    /** The first element of the list, or None if empty */
    Optional<Node> head = Option.none();

    /* method implementations */
}
```

To help us convince ourselves that we’re actually creating a list, let’s state some rules for this structure we’re building:

1. If the list is empty then `head.isEmpty()`.
2. If the list is non-empty, then `head` refers to the 0th node.
3. The x th node stores the x th element of the list.
 - **Corrolary:** There are N nodes.
4. If the x th node is the last node (i.e., $x = N - 1$), then the node’s `next.isEmpty()`.
5. If the x th node is not the last node, then the node’s `next` points to the $x + 1$ th node.

Rules about how a data structure should behave are often called **invariants**. When writing your own code and/or data structures, a great way to avoid bugs is to write down the invariants that you expect your code to follow, and then proving that your code follows the rules (“preserves the invariant”).

4.4.2.1. Linked List get

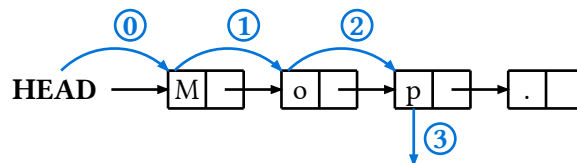


Figure 7: Retrieving the character at index 2 from a linked list representing the list ['M', 'o', 'p', '.']

To retrieve a specific element i of the list, we need to find the i th node. Since all we have to start is the 0th node (the head; step ①), we can start there. If that’s the node we’re looking for, great, we’re done; Otherwise, the only other place we can go is the 1st node (step ②). We repeat this process, moving from node to node until we get to the i th node (step ③), and return the corresponding value (step ④).

```
public T get(int index)
{
    if(index < 0){ throw new IndexOutOfBoundsException() }
    Optional<Node> currentNode = head;
```

java, it’s been considered bad practice to use `null`. As a further side note, `Optional` can actually be viewed as a `MutableSequence`. Implementing this is left as an exercise for the reader.

```

for(i = 0; i < index; i++){
    if(currentNode.isNone){ throw new IndexOutOfBoundsException() }
    currentNode = currentNode.get().next;
}
if(currentNode.isNone){ throw new IndexOutOfBoundsException() }
return currentNode.get().value;
}

```

Note that, if the rules we set for ourselves are being followed, we can prove to ourselves that this implementation is correct. Recall that we said `get(index)` is correct if it returns the `index`th element of the list.

1. We start with `currentNode` assigned to `head`; By the rules we set for ourselves, this is the 0th node.
2. The i th node has a pointer to the $(i + 1)$ th node, which we follow i times.
3. After i steps, `currentNode` has been reassigned to the $0 + \underbrace{1 + \dots + 1}_{i \text{ times}} = 0 + \sum_i 1 = i$ th node.
4. It follows from step 4 and the fact that the loop iterates `index` times that after the loop, `currentNode` is the `index`th node.
5. It follows from step 5 and the rule that the i th node (if it exists) contains the i th element (if it exists), that the value in `currentNode` which `get` returns is the `index`th node
6. Since a correct `get` is one that returns the `index`th node, we can conclude that `get` is correct.

As usual, we can figure out the runtime of a snippet of code, starting by replacing every primitive operation with $\Theta(1)$

```

public T get(int index)
{
    if( /* Theta(1) */ ){ throw /* Theta(1) */ }
    /* Theta(1) */
    for(i = 0; i < index; i++){
        if( /* Theta(1) */ ){ throw /* Theta(1) */ }
        /* Theta(1) */
    }
    if( /* Theta(1) */ ){ throw /* Theta(1) */ }
    return /* Theta(1) */
}

```

Before the `for` loop, we have only constant-time operations. Similarly, we have only constant-time operations inside the body of the loop. Both can be simplified to $\Theta(1)$ Let's start by figuring out the runtime of what's inside the `for` loop.

$$\sum_{i=0}^{\text{index}-1} \Theta(1) = ((\text{index} - 1) - 0 + 1) \cdot \Theta(1) = \text{index} \cdot \Theta(1) = \Theta(\text{index})$$

Since this is only the body of the `for` loop, we can add back the constant time operations before and after the loop to get the total runtime of `get`:

$$T_{\text{get}(\text{index})} = \Theta(1) + \Theta(\text{index}) + \Theta(1) = \Theta(\text{index}).$$

The $\Theta(\text{index})$ runtime is a little unusual: With arrays, most costs were expressed in terms of the size of the array itself. That is, our asymptotic runtime complexity was given in terms of N , while here, it's

given in terms of index. However, we know that $0 \leq \text{index} < N$, and we can use that fact to create an analogous bound on the runtime of `get`²².

Since $\text{index} < N$, `get` is at worst linear in the size of the array, and so we get the upper bound:

$$T_{\text{get}(N)} = O(N)$$

Similarly, since $\text{index} \geq 0$, so `get` could be constant time (e.g., `get(0)` requires only constant time), and we get the lower bound: $T_{\text{get}(N)} = \Omega(1)$

4.4.2.2. Linked List update

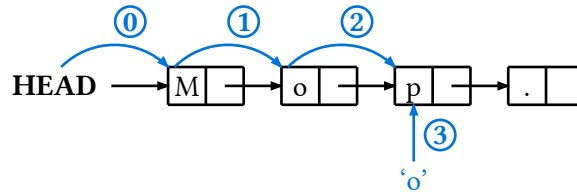


Figure 8: Updating the character at index 2 to 'o', on the linked list representing the list ['M', 'o', 'p', '.']. After the update, the linked list represents the list ['M', 'o', 'o', '.']

To update the i th element of a linked list, just like `get`, we need to first find the i th element. As a result, the code for `update` is almost exactly the same as `get`, except for the very last line.

```
public void update(int index, T element)
{
    if(index < 0){ throw new IndexOutOfBoundsException() }
    Option<Node> currentNode = head;
    for(i = 0; i < index; i++){
        if(currentNode.isNone){ throw new IndexOutOfBoundsException() }
        currentNode = currentNode.get().next;
    }
    if(currentNode.isNone){ throw new IndexOutOfBoundsException() }
    currentNode.get().value = element;
}
```

Using our Linked List invariants, we can similarly prove that our implementation of `update` follows the rules that we set out for `List`'s `update` operation. Remember that `update(index, element)` is supposed to change the linked list so that the next time we call `get(index)`, we get `element` instead of what was there before.

1. We start with `currentNode` assigned to `head`; By the rules we set for ourselves, this is the 0th node.
2. The i th node has a pointer to the $(i + 1)$ th node, which we follow i times.
3. After i steps, `currentNode` has been reassigned to the $0 + \underbrace{1 + \dots + 1}_{i \text{ times}} = 0 + \sum_i 1 = i$ th node.
4. It follows from step 4 and the fact that the loop iterates `index` times that after the loop, `currentNode` is the `index`th node.
5. It follows from step 5 and the rule that the i th node (if it exists) contains the i th element (if it exists), that the `value` in `currentNode` which `update` changes is the `index`th node

²²The proof here is given only for *valid* inputs. For values of `index` less than zero, the function aborts immediately, and we can prove to ourselves (based on the rule that the last element in the list has an empty next node) that we will never go through the loop more than N times. Strictly speaking, the runtime should be, $\Theta(\min(N, \text{index}))$. However, this added complexity makes no immediate impact on our discussion, and so we omit it.

6. Since any subsequent calls to get will return the indexth node, we can conclude that update is correct.

Side Note: Observe how the proof for update is almost exactly the same as the proof for get.

As usual we can figure out the runtime by replacing every primitive operation with $\Theta(1)$

```
public void update(int index, T element)
{
    if( /* Theta(1) */ ){ throw /* Theta(1) */ }
    /* Theta(1) */
    for(i = 0; i < index; i++){
        if( /* Theta(1) */ ){ throw /* Theta(1) */ }
        /* Theta(1) */
    }
    if( /* Theta(1) */ ){ throw /* Theta(1) */ }
    /* Theta(1) */
}
```

Once we replace all primitive operations with $\Theta(1)$, update and get are exactly the same. Taking the same steps gets us to a runtime of $\Theta(\text{index})$, or $O(N)$.

4.4.2.3. Linked List add

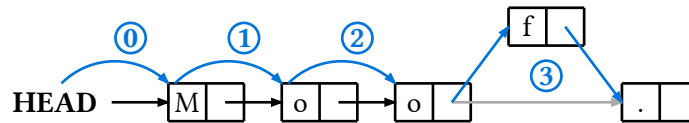


Figure 9: Inserting the character 'f' at index 3, on the linked list representing the list ['M', 'o', 'o', '.']. After the add, the linked list represents the list ['M', 'o', 'o', 'f', '.'] (the call of the noble dogcow).

Unlike the array, where we manually need to copy over every element of the array to a new position, in a linked list, the position of each node is given relative to the previous node. If we insert a new node by redirecting the next value of the previous node, we automatically shift every subsequent node by one position to the right.

```
public void add(int index, T element)
{
    if(index < 0){ throw new IndexOutOfBoundsException() }
    Node newNode = new Node(element);
    if(head.isEmpty){ /* Case 1: Initially empty list */
        head = Optional.of(newNode);
    } else if(index == 0){ /* Case 2: Insertion at head */
        newNode.next = head;
        head = Optional.of(newNode);
    } else { /* Case 3: Insertion elsewhere */
        Option<Node> currentNode = head;
        for(i = 0; i < index-1; i++){
            if(currentNode.isNone){ throw new IndexOutOfBoundsException() }
            currentNode = currentNode.get().next;
        }
        if(currentNode.isNone){ throw new IndexOutOfBoundsException() }
        newNode.next = currentNode.get().next;
    }
}
```

```

        currentNode.get().next = Optional.of(newNode);
    }
}

```

Let's prove to ourselves that this code is correct, and start by recalling our rules for how `add` is supposed to behave:

1. The size of the List grows by 1.
2. The value `element` will be the new element at index `index`.
3. Every element previously at an index $i \geq \text{index}$ will be moved to index $i + 1$.

Let's set aside the size rule for the moment and focus on the other two rules. We have our original list (L_{old} of size N_{old}), and the list we get after calling `add` (L_{new} of size $N_{\text{new}} = N_{\text{old}} + 1$). Collectively, the latter two rules give us four cases for what will happen if we call $L_{\text{new}}.\text{get}(i)$:

1. If $0 \leq i < \text{index}$: $L_{\text{new}}.\text{get}(i) = L_{\text{old}}.\text{get}(i)$
2. If $i = \text{index}$: $L_{\text{new}}.\text{get}(i) = \text{element}$
3. If $N_{\text{new}} > i > \text{index}$: $L_{\text{new}}.\text{get}(i) = L_{\text{old}}.\text{get}(i - 1)$
4. If $i \geq N_{\text{new}}$: $L_{\text{new}}.\text{get}(i)$ throws an exception

We'll want to show that after calling `add`, any subsequent call to `get(i)` will return the correct value. We already proved that `get` is correct if the linked list follows the rules we set forth for ourselves, so we need to show that the new linked list is correct for the updated list. In other words, we want to show that:

1. All nodes at position $0 \leq i < \text{index}$ are left unchanged.
2. The newly created node (at position `index`) is pointed to by the node at position `index-1`, or by `head` if `index = 0`.
3. All nodes previously at positions $\text{index} \leq i < N_{\text{old}}$ are now at position $i + 1$.
4. The node at position $N_{\text{new}} - 1$ has an empty `next`.

The code has three cases that we'll look at individually.

Case 1: Initially empty list: Going down the list of things we need to prove:

1. There are no other nodes, so nothing changes.
2. The newly created node is at position 0, and so is correctly pointed to by `head` after the code runs.
3. There are no other nodes, so nothing needs to change.
4. Since $N_{\text{new}} = 1$, we need the node at position $N_{\text{new}} - 1 = 0$ to have an empty `next`. Since this is the node we just created, this is true.

Case 2: Insertion at head: Going down the list of things we need to prove:

1. Since `index = 0`, there can be no nodes at positions before `index`, so nothing changes.
2. The newly created node is at position 0, and so is correctly pointed to by `head` after the code runs.

For the remaining two rules, consider the fact that the (newly created) 0th node's `next` field now points at the node that used to be the head.

- By our rules for linked lists, this means that the old 0th element is now the 1st element (correct).
- The new 1st element is pointing at what used to be the 1st element, but is now the 2nd element (also correct).
- The new 2nd element is pointing at what used to be the 2nd element, but is now the 3rd element (also correct).

- The new 3rd element is pointing at what used to be the 3rd element, but is now the 4th element (also correct).

This pattern continues: The new i th element is pointing at what used to be the i th element, but is now the $i + 1$ th element (still correct), up until we get to the $N_{\text{new}} - 1$ th element, which used to be the $N_{\text{old}} - 1$ th element, which still (correctly) has an empty next.

Case 3: Insertion anywhere else: Going down the list of things we need to prove:

1. The code finds the node at position $\text{index}-1$, skipping over all preceding nodes. These are unchanged.
2. The newly created node is pointed to by the node at position $\text{index}-1$, putting it at position index .

For the remaining two rules, we want to consider what happens to the node previously at position index , but also need to consider the possibility that there was no node at this index to begin with (if $\text{index} = N_{\text{old}}$). If the node exists, then as in **Case 2**, the node at position index moves to position $\text{index} + 1$, likewise repositioning every following node, including the final one. If the node doesn't exist, then the newly created node's next field is assigned an empty value, (correctly) making it the last node in the list.

4.4.2.4. Aside: Invariants and Rule Preservation

It's worth taking a step back at this point and reviewing what we just did and why. We defined the expected state of the linked list (in terms of the behavior of `get` and our rules for a correct linked list) *after* the update, but we did so relatively to its state *before* the update. Specifically, we showed that, if the list was correct before we called `update`, it would still be correct (for the new list) after we called `update`, or in other words, we showed that `update` **preserves the invariants** we defined for the linked list.

Invariants are an extremely useful debugging technique for data structures (and other code). If you can precisely define one or more rules (invariants) for your data structure, you can check to see whether your code follows the rules: If you get a correct data structure as input to your function, is it always the case that you get a correct data structure as an output (keeping in mind that the function may change the definition of correctness).

4.4.2.5. Linked List `add` runtime

The runtime of `add` follows a pattern very similar to that of `update`:

```
public void add(int index, T element)
{
    if( /* Theta(1) */ ){ /* Theta(1) */ }
    /* Theta(1) */
    if( /* Theta(1) */ ){ /* Case 1: Initially empty list */
        /* Theta(1) */
    } else if( /* Theta(1) */ ){ /* Case 2: Insertion at head */
        /* Theta(1) */
    } else { /* Case 3: Insertion elsewhere */
        /* Theta(1) */
        for(i = 0; i < index-1; i++){
            /* Theta(1) */
        }
        if( /* Theta(1) */ ){ /* Theta(1) */ }
```

```

        /* Theta(1) */
    }
}

```

Based on the the outer if statement, we have three cases:

$$T_{\text{add}}(\text{index}) = \begin{cases} \Theta(1) & \text{if case 1} \\ \Theta(1) & \text{if case 2} \\ \Theta(1) + \sum_{i=0}^{\text{index}-1} (\Theta(1)) & \text{if case 3} \end{cases}$$

Simplifying the third case, we get:

- $\Theta(1) + \sum_{i=0}^{\text{index}-1} (\Theta(1))$
- $\Theta(1) + (\text{index} - 1 + 0 + 1)(\Theta(1))$
- $\Theta(1) + \text{index} \cdot \Theta(1)$
- $\Theta(1) + \Theta(\text{index})$
- $\Theta(\text{index})$

$$\text{So... } T_{\text{add}}(\text{index}) = \begin{cases} \Theta(1) & \text{if case 1} \\ \Theta(1) & \text{if case 2} \\ \Theta(\text{index}) & \text{if case 3} \end{cases}$$

Since $\text{index} = 0$ in cases 1 and 2, we can further simplify to

$$T_{\text{add}}(\text{index}) = \begin{cases} \Theta(\text{index}) & \text{if case 1} \\ \Theta(\text{index}) & \text{if case 2} = \Theta(\text{index}) \\ \Theta(\text{index}) & \text{if case 3} \end{cases}$$

As before, $\Theta(\text{index}) \in O(N)$ and $\Theta(\text{index}) \in \Omega(1)$

4.4.2.6. Linked List size (Take 1)

The size of the list is the number of elements. Thinking back to our rules, the last node is the node with an empty next pointer, so we need to figure out where this node is located

```

public int size()
{
    Option<Node> currentNode = head;
    int count = 0;
    while( ! currentNode.isEmpty() ){
        currentNode = currentNode.get().next
        count += 1
    }
    return count
}

```

Let's see if we can prove that this code is correct. This code is a bit harder to think about than `get` and `update`, since in both of those cases we had a nice handy for loop to keep track of how many 'steps' we take through the list. Here, the while loop just keeps going until `currentNode.isEmpty()`. When trying to tackle a tricky proof like this, it can often help to break down the proof into cases.

Let's start with the simplest case: What happens if the list is empty ($N = 0$)?

1. `count` starts off at 0.
2. If the list is empty, then by our rules for linked lists `head` is empty, so we start off with `currentNode` empty as well.
3. Since `currentNode` is empty from the start, we skip the while loop's body entirely.

4. `count` is never modified, and we correctly return 0.

Let's tackle the next simplest case next: What happens if the list only has 1 element ($N = 1$)

1. `count` starts off at 0.
2. `head` points to the 0th element of the list, so `currentNode` starts off pointing to the 0th node.
3. Since `currentNode` is not empty, we enter the body of the while loop.
4. We increment `count` to 1, and update `currentNode` from the 0th node to its next reference.
5. By our rules for linked lists, since the list has only 1 element, the 0th node's next reference is empty, so `currentNode` is now empty and we exit the while loop.
6. We correctly return `count = 1`,

Moving on, let's see what happens if the list has 2 elements ($N = 2$)

1. `count` starts off at 0.
2. `head` points to the 0th element of the list, so `currentNode` starts off pointing to the 0th node.
3. Since `currentNode` is not empty, we enter the body of the while loop, increment `count`, and update `currentNode` to its next reference.
4. We increment `count` to 1, and update `currentNode` from the 0th node to its next reference, the 1st node.
5. We increment `count` to 2, and update `currentNode` from the 1st node to its next reference.
6. By our rules for linked lists, since the list has only 2 elements, the 1st node's next reference is empty, so `currentNode` is now empty and we exit the while loop.
7. We correctly return `count = 2`,

Moving on, let's try a list with 3 elements ($N = 3$)

1. `count` starts off at 0.
2. `head` points to the 0th element of the list, so `currentNode` starts off pointing to the 0th node.
3. Since `currentNode` is not empty, we enter the body of the while loop, increment `count`, and update `currentNode` to its next reference.
4. We increment `count` to 1, and update `currentNode` from the 0th node to its next reference, the 1st node.
5. We increment `count` to 2, and update `currentNode` from the 1st node to its next reference, the 2nd node.
6. We increment `count` to 3, and update `currentNode` from the 2nd node to its next reference.
7. By our rules for linked lists, since the list has only 3 elements, the 2nd node's next reference is empty, so `currentNode` is now empty and we exit the while loop.
8. We correctly return `count = 3`,

Although the proof is different for each case, you might notice a pattern forming. For a list of *any* size, we can come up with a similar proof by adding a bunch of lines into the middle, all from the template:

We increment `count` to i , and update `currentNode` from the $i - 1$ th node to its next reference, the i th node.

In other words, it looks like the code ties the values of `count` and `currentNode` together. We can use this intuition to define a rule for ourselves: If `currentNode` is not empty, it always points to the `count` node of the list. This is true at the start of the loop, since `count = 0` and by our rules for linked lists, `currentNode` points at the 0th node. Now, if we know that `currentNode` (if non-empty) points to the

count node of the linked list when we **start** the loop body, we can show that it's still true at the end of the loop body:

1. If `currentNode` points to the i th node of the list, by our rules for linked lists, its next element points to the $i + 1$ th node.
2. If `count` is i to start, then incrementing it sets it to $i + 1$
3. From lines 1 and 2, after the loop body, `count = $i + 1$` and `currentNode` is the $i + 1$ th node, and we're still following the rule.

Since the rule holds at the **start** of the loop, and since the loop body preserves the rule, we know that the rule is also followed at the **end** of the loop. We call a rule that is satisfied at the start of a loop and preserved by the loop body, a **loop invariant**.

The loop ends at the very first point where `currentNode` becomes empty. This happens in one of two cases:

- `head` is empty.
- The next element of `currentNode` is empty.

By our rules for linked lists, the first case happens only when the list is empty, and we've already shown that `size` is correct in this case. Similarly, by our rules for linked lists, the second case happens only when `currentNode` points to the last (i.e., $N - 1$ th) node of the linked list. From that, and our loop invariant, it must be the case that `count` is $N - 1$. Since `count` gets incremented one last time on the 2nd line of the loop body, when the loop ends `count = N` , which is also correct.

4.4.2.7. Aside: Loop Invariants

Once again, let's take a quick moment to review what we just did, because it's very useful trick for debugging loops in your code²³. We started by trying to prove that our code was correct for a couple of simple example cases. From those simple example cases, we found a pattern in the proof: a rule that our code seemed to obey throughout our proof. We then set up a **loop invariant** based on that pattern and showed that (i) it held at the start of the loop, and (ii) each line of the loop preserved the rule. Because we could show that the invariant held at the start of the loop, and each loop iteration preserved the invariant, we inferred that the invariant had to be true at the end of the loop as well. We used the fact that the invariant was true at the end of the loop to prove that `size` was correct.

When debugging code with loops, try tracing through the loop manually for a few iterations. You'll often start noticing patterns and relationships between elements of the code. Try to pin down **exactly** what that pattern is (like `currentNode` always points to the `count` node in our example), and write down the loop invariant. Then, try to prove to yourself that (i) the pattern holds at the start of the loop, (ii) the pattern is preserved by the loop body.

4.4.2.8. Linked List size (Take 1) runtime

As usual, we can replace all primitive operations with $\Theta(1)$

```
public int size()
{
    /* Theta(1) */
    while( /* Theta(1) */ ){
```

²³Although it's a bit too early to use the term in the main body of the text, loop invariants are a specific example of a proof trick called "induction". We'll come back to induction as a more general technique next chapter, when we talk about recursion.

```

    /* Theta(1) */
}
return /* Theta(1) */
}

```

Once again, the while loop makes our job a little harder. However, just as before, we can use the loop invariant to help ourselves out:

1. count starts at 0.
2. Each loop iteration increments count by 1.
3. The loop ends when count = N .

Combining these facts, we can definitively state that the loop goes through N iterations. So, we get:

$$\Theta(1) + \underbrace{\Theta(1) + \dots + \Theta(1)}_{N \text{ times}} + \Theta(1) = (N + 2) \cdot \Theta(1) = \Theta(N + 2) = \Theta(N)$$

4.4.2.9. Linked List size (Take 2)

It's not great when a function has a $\Theta(N)$ runtime, and it's really bad when it's something as fundamental as size. Many algorithms need to know the size of a collection, so it would be useful to have a way to cut the runtime to something more practical. In the case of size, there's one simple trick that lets us cut the runtime down to $\Theta(1)$: Precomputing the size. Let's add a N field to the linked list and modify the size method to just return it.

```

public class SinglyLinkedList<T> implements List<T>
{
    class Node
    {
        T value;
        Optional<Node> next = Optional.empty();
        public this(T value) { this.value = value; }
    }
    /** The first element of the list, or None if empty */
    Optional<Node> head = Option.none();

    /** The precomputed size of the list */
    private int N = 0;

    public int size() { return this.N; }

    /* method implementations */
}

```

The list starts off empty ($N = 0$), so initially this method is correct. However, we need to make sure that it stays correct as we change the structure. There are two operations, add and remove, that modify the size of a list, incrementing and decrementing the list's size. We need to modify these operations to properly maintain N :

```

public void add(int index, T element)
{
    /* ... as before ... */
    this.N += 1;
}

```

```
public void remove(int index)
{
    /* ... as before ... */
    this.N -= 1;
}
```

Let's review the trade-off we just made: Both changes add an additional $\Theta(1)$ runtime cost to add and remove. This doesn't change the overall code complexity of either ($\Theta(\text{index})$, or equivalently $O(N)$)²⁴. In exchange, the asymptotic runtime complexity of `size` drops from $\Theta(N)$ to $\Theta(1)$. To summarize, the asymptotic runtime complexity of all the other operations stays the same, and `size`'s asymptotic runtime complexity drops massively. This is almost always a worthwhile trade.

4.5. Iterators

4.6. Doubly Linked Lists

²⁴In practice, incrementing and decrementing an integer is a relatively inexpensive operation in most cases, so even outside of the magical land of asymptotics, this is usually a worthwhile trade. That said, there are *some* exceptions, primarily when dealing with concurrency (synchronized integers are expensive) or with random memory access (we'll talk about this later in the book).