

Lecture I

OUTLINE OF ALGORITHMICS

We assume the student is familiar with computer programming and has a basic course in data structures. Problems solved using computers can be roughly classified into problems-in-the-large and problems-in-the-small. The former is associated with large software systems such as an airlines reservation system, compilers or text editors. The latter¹ is identified with mathematically well-defined problems such as sorting, multiplying two matrices or solving a linear program. The methodology for studying the large and small problems are quite distinct: Algorithmics is the study of the small problems and their algorithmic solution. In this introductory lecture, we presents an outline of this enterprise. Throughout this book, **computational problems** (or simply “problems”) refer to problems-in-the-small. It is the only kind of problem we address.

READING GUIDE: This chapter is mostly informal and depends on some prior understanding of algorithms. The rest of this book has no dependency on this chapter, save the definitions in §8 concerning asymptotic notations. Hence a light reading may be sufficient. We recommend re-reading this chapter after finishing the rest of the book, when many of the remarks here may take on more concrete meaning.

§1. What is Algorithmics?

Algorithmics is the systematic study of efficient algorithms for computational problems; it includes techniques of algorithm design, data structures, and mathematical tools for analyzing algorithms.

Why is algorithmics important? Because algorithms is at the core of all applications of computers. These algorithms are the “computational engines” that drive larger software systems. Hence it is important to learn how to construct algorithms and to analyze them. Although algorithmics provide the building blocks for large application systems, the construction of such systems usually require additional non-algorithmic techniques (e.g., database theory) which are outside our scope.

We can classify algorithmics according to its applications in subfields of the sciences and mathematics: thus we have computational geometry, computational topology, computational number theory, computer algebra, computational statistics, computational finance, computational physics and computational biology.

Another way to classify algorithmics is to look at the generic tools and techniques that are largely independent of subject matter. Along this line, we identify four basic themes:

- (a) data-structures (e.g, linked lists, stacks, search trees)
- (b) algorithmic techniques (e.g., divide-and-conquer, dynamic programming)
- (c) basic computational problems (e.g., sorting, graph-search, point location)
- (d) analysis techniques (e.g., recurrences, amortization, randomized analysis)

These themes interplay with each other. For instance, some data-structures naturally suggest certain algorithmic techniques. Or, an algorithmic technique may entail certain analysis methods (e.g., divide-and-conquer algorithms require recurrence solving). Complexity theory provides some unifying concepts for algorithmics; but complexity theory is too abstract to capture many finer distinctions we wish to make. Thus algorithmics often makes domain-dependent assumptions. For example, in the subfield of computer

¹If problems-in-the-large is macro-economics, then the latter is micro-economics.

algebra, the complexity model takes each algebraic operation as a primitive while in the subfield of computational number theory, these algebraic operations are reduced to some bit-complexity model primitives. In this sense, algorithmics is, say, more like combinatorics (which is eclectic) than group theory (which starts out from a unified framework).

§2. What are Computational Problems?

Despite its name, the starting point for algorithmics is **computational problems**, not algorithms. But what are computational problems? We mention three main categories.

(A) Input-output problems. Here is the simplest formulation: A **computational problem** is a precise specification of input and output formats, and for each input instance I , a description of the set of possible output instances $O = O(I)$.

The word “formats” emphasizes the fact the input and output representation is part and parcel of the problem. In practice, standard representations may be taken for granted (e.g., numbers are assumed to be in binary and set elements are arbitrarily listed without repetition). Note that the input-output relationship need not be functional: a given input may have several acceptable outputs.

■ **Example:** SORTING PROBLEM: Input is a sequence of numbers (a_1, \dots, a_n) and output is a rearrangement of these numbers (a'_1, \dots, a'_n) in non-decreasing order. An input instance is $(2, 5, 2, 1, 7)$, with corresponding output instance $(1, 2, 2, 5, 7)$.

■ **Example:** PRIMALITY TESTING: Input is a natural number n and output is either YES (if n is prime) or NO (if n is composite). This is an example of **decision or recognition problem**, where the outputs have only two possible answers (YES/NO, 0/1, Accept/Reject). Such problems play an important role in complexity theory. One can generalize this to problems whose output comes from a finite set. For instance, in computational geometry, the decision problems tend to have three possible answers: Positive/Negative/Zero or IN/OUT/ON. For instance, the **point classification problem** is where we are given a point and some geometric object such as a triangle or a cell. The point is either inside the cell, outside the cell or on the boundary of the cell.

(B) Static preprocessing problems. A generalization of input-output problems is what we call **preprocessing problem**: *given a set S of objects, construct a data structure $D(S)$ such that for an arbitrary ‘query’ (of a suitable type) about S , we can use $D(S)$ to efficiently answer the query.* This problem is a “static” preprocessing problem because the members of the set S does not change under the querying.

■ **Example:** RANKING PROBLEM: preprocessing input is a set S of numbers. A query on S is a number q for which we like to determine its rank in S . The rank of q in S is the number of items in S that are smaller than or equal to q . A standard solution to this problem is the *binary search tree* data structure $D(S)$ and the binary search algorithm on $D(S)$.

■ **Example:** POST OFFICE PROBLEM: Many problems in computational geometry and database search are the preprocessing type. The following is a geometric-database illustration: given a set S of points in the plane, find a data structure $D(S)$ such that for any query point p , we find an element in S that is closest to p . (Think of S as a set of post offices and we want to know the nearest post office to any position p). Note that the 1-dimensional version of this problem is closely allied to the ranking problem.

Two algorithms are needed to solve a preprocessing problem: one to construct $D(S)$ and another to answer queries. They correspond to the two stages of computation: an initial **preprocessing stage** to construct $D(S)$, and a subsequent **querying stage** in which the data structure $D(S)$ is used. There may be a tradeoff between the **preprocessing complexity** and the **query complexity**: $D_1(S)$ may be faster to construct than an alternative $D_2(S)$, but answering queries using $D_1(S)$ is less efficient than $D_2(S)$. But our general attitude to prefer $D_2(S)$ over $D_1(S)$ in this case: we prefer data structures $D(S)$ that support the fastest possible query complexity. Our attitude is often justified because the preprocessing complexity is a one-time cost.

Preprocessing problems can be seen as a special case of **partial evaluation problems**. In such problems, we construct partial answers or intermediate structures based on part of the inputs; these partial answers or intermediate structures must anticipate all possible extensions of the partial inputs.

(C) Dynamization and Online problems. Now assume the input S is a set, or more generally some kind of aggregate object. If S can be modified under queries, then we have a **dynamization problem**: with S and $D(S)$ as above, we must now design our solution with an eye to the possibility of modifying S (and hence $D(S)$). Typically, we want to insert and delete elements in S while at the same time, answer queries on $D(S)$ as before. A set S whose members can vary over time is called a **dynamic set** and hence the name for this class of problems.

Here is another formulation: *we are given a sequence (r_1, r_2, \dots, r_n) of **requests**, where a request is one of two types: either an **update** or a **query**. We want to ‘preprocess’ the requests in an online fashion, while maintaining a time-varying data structure D : for each update request, we modify D and for each query request, we use D to compute and retrieve an answer (D may be modified as a result).*

In the simplest case, updates are either “insert an object” or “delete an object” while queries are “is object x in S ?”. This is sometimes called the **set maintenance problem**. Preprocessing problems can be viewed as a set maintenance problem in which we first process a sequence of insertions (to build up the set S), followed by a sequence of queries.

■ **Example:** **DYNAMIC RANKING PROBLEM:** Any preprocessing problem can be systematically converted into a set maintenance problem. For instance, the ranking problem turns into the **dynamic ranking problem** in which we dynamically maintain the set S subject to intermittent rank queries. The data structures in solutions to this problem are usually called **dynamic search trees**.

■ **Example:** **GRAPH MAINTENANCE PROBLEMS:** Dynamization problems on graphs are more complicated than set maintenance problems (though one can still view it as maintaining a set of edges). One such problem is the **dynamic connected component problem**: updates are insertion or deletion of edges and/or vertices. Queries are pairs of vertices in the current graph, and we want to know if they are in the same component. The graphs can be directed or undirected.

■ **Example:** **ONLINE SCHEDULING PROBLEMS:** In this case, each request r_i ($i = 1, \dots, n$) is a task that needs to be serviced. The requirement of the online problem is that the service decision s_i (to service it or not, when to serve it, etc) must be made before the next request is seen. We can view this as an optimization problem, where we are trying to maximize some function $f(s_1, s_2, \dots, s_n)$ of the service decisions. There is a growing literature for a class of such algorithms that are “competitive” in the sense $f(s_1, \dots, s_n)$ is at least some constant factor of the optimal solution (s_1^*, \dots, s_n^*) .

(D) Pseudo-problems. Let us illustrate what we regard to be a pseudo-problem from the viewpoint of our subject. Suppose your boss asks your IT department to “build an integrated accounting system-cum-

employee database”. This may be a real world scenario but it is not a legitimate topic for algorithmics because part of the task is to figure out what the input and output of the system should be, and there are probably other implicit non-quantifiable criteria (such as available technology and economic realities).

§3. Computational Model: How do we solve problems?

Once we agree on the computational problem to be solved, we must choose the tools for solving it. This is given by the **computational model**. Any conventional programming languages such as C or Java (suitably abstracted, so that it does not have finite space bounds, etc) can be regarded as a computational model. A computational model is specified by

- (a) the kind of data objects that it deals with
- (b) the primitive operations to operate on these objects
- (c) rules for composing primitive operations into larger units called **programs**.

Programs can be viewed as individual instances of a computational model. For instance, the Turing model of computation is an important model in complexity theory and the programs here are called Turing machines.

Models for Sorting. To illustrate computational models, we consider the problem of sorting. The sorting problem has been extensively studied under several computational models. We briefly describe just three of them: the **comparison-tree model**, the **comparator circuit model**, and the **tape model**. In each models, the data objects are elements from a linear order.

The first model, comparison-trees, has only one primitive operation, viz., comparing the two elements x, y resulting in one of two outcomes $x < y$ or $x \geq y$. Such a comparison is usually denoted “ $x : y$ ”. We compose these primitive comparisons into a **tree program** by putting them at the internal nodes of binary tree. Tree programs represent flow of control and are more generally called **decision trees**. Figure 1(a) illustrates a comparison-tree on inputs x, y, z . The output of the decision tree is specified at each leaf. For instance, if the tree is used for sorting, we would want to write the sorted order of the input elements in each leaf. If the tree is used to find the maximum element of the input set, then each leaf would specify the maximum element.

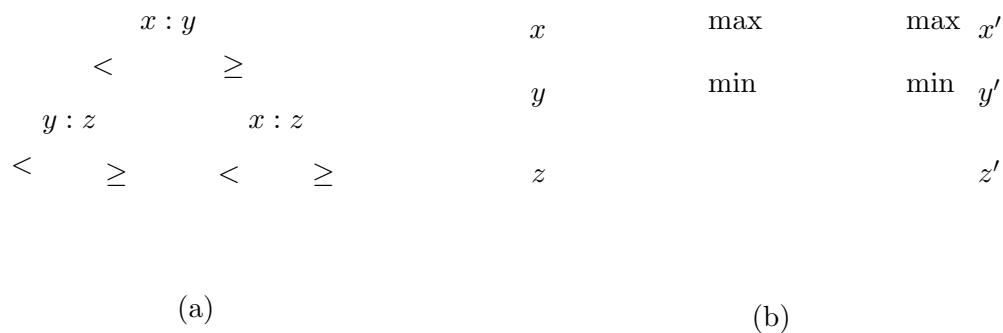


Figure 1: (a) A comparison-tree and (b) a comparator circuit

In the comparator circuit model, we also have one primitive operation which takes two input elements x, y and returns two outputs: one output is $\max\{x, y\}$, the other $\min\{x, y\}$. These are composed into **circuits**

which are directed acyclic graphs with n input nodes (in-degree 0) and n output nodes (out-degree 0) and some number of comparator nodes (in-degree and out-degree 2). In contrast to tree programs, the edges (called **wires**) in such circuits represent actual data movement. Figure 1(b) shows a comparator circuit on inputs x, y, z . Depending on the problem, the output of the comparator circuit may be the set of all output lines (x', y', z' in Figure 1(b)) or perhaps some subset of these lines.

A third model for sorting is the tape model. Assume one or more sequential tapes which can store the elements to be sorted. Each tape is either in the state of being read or being written. The key property of such tapes is that they allow only **sequential access**. Hence the tape positions can be changed in two ways: advance to the next element or reset (meaning go to the first element). Initially, the input has n elements stored in one tape. Finally, the elements are written out, in sorted order, on one of the tapes. We have a main memory that can only store m of these elements. Typically, $m \ll n$ and m may even be fixed. The instantaneous state of the algorithm will specify, for each tape, whether it is being read or being written, and the position of its the tape head. The reading and writing on each tape can proceed independently. We can detect the condition of reading past the last element in a tape. The sorting algorithm dictates the input and output behavior on each tape, including which tape(s) to be reset for reading or writing.

The tape model was important in the early days of computing when main memory was expensive and physical tapes is the standard medium for storing large databases. Interestingly, with the advent of the web-age, a variant of this model called **streaming data model** is coming back. Now we are faced with hugh amounts of real time data, and instead of sorting, we often need to compute some function of the data. Because of the large volume of data, we do not want to store this information but allow only one pass over the data. For more information about the tape model, we refer to Knuth [3, Chap. 5.4], under external sorting.

Algorithms versus programs. To use a computational model to solve a given problem, we must make sure there is a match between the data objects in the problem specification and the data objects handled by the computational model. If not, we must specify some suitable encoding of the former objects by the latter. Similarly, the input and output formats of the problem must be represented in some way. After making explicit such encoding conventions, we may call A an **algorithm for P** if, if the program A indeed computes a correct output for every legal input of P . Thus the term algorithm is a semantical concept, signifying a program in its relation to some problem. In contrast, programs are purely syntactic objects. E.g., the programs in figure 1(a,b) are both algorithms to compute the maximum of x, y, z . But what is the output convention for these two algorithms?

Uniform versus Non-uniform Computational Models. While problems generally admit inputs of arbitrarily large sizes (see discussion of size below), some computational models define programs that admit inputs of a fixed size only. This is true of the decision tree and circuit models of computation. In order to solve problems of infinite sizes, we must take a sequence of programs $P = (P_1, P_2, P_3, \dots)$ where P_i admits inputs of size i . We call such a program P a **non-uniform program** since we have no *a priori* connections between the different P_i 's. For this reason, we call the models whose programs admit only finite size inputs **non-uniform models**. The next section will introduce a **uniform model** called the RAM Model. Pointer machines (see Chapter 6) and Turing machines are other examples of uniform models. The relationship between complexity in uniform models and in non-uniform models is studied in complexity theory.

Problem of Merging Two Sorted Lists. Let us illustrate the difference between uniform and non-uniform algorithms. A subproblem that arises in sorting is the **merge problem** where we are given two sorted lists (x_1, x_2, \dots, x_m) and (y_1, y_2, \dots, y_n) and we want to produce a sorted list $(z_1, z_2, \dots, z_{m+n})$ where $\{z_1, \dots, z_{m+n}\} = \{x_1, \dots, x_m, y_1, \dots, y_n\}$. Assume these sorted lists are non-decreasing. Here is an algorithm for this problem, written in a generic conventional programming language:

```

MERGE ALGORITHM
Initialize
   $i \leftarrow 1, j \leftarrow 1.$ 
Loop
  while ( $i \leq m$  and  $j \leq n$ )
    if  $x_i \leq y_j$ 
      Output( $x_i$ ) and  $i \leftarrow i + 1.$ 
    else
      Output( $y_j$ ) and  $j \leftarrow j + 1.$ 
Terminate
  if  $i > m$ 
    Output the rest of the  $y$ 's.
  else
    Output the rest of the  $x$ 's.

```

This is a uniform algorithm for merging two lists. For each m, n , this algorithm can be “unwounded” into a comparison-tree $T_{m,n}$ for merging two sorted lists of sizes m and n (Exercise). Hence family $\{T_{m,n} : m, n \in \mathbb{N}\}$ is a non-uniform algorithm for merging two lists.

Program Correctness. This has to do with the relationship between an program and a computational problem. *A program that is correct relative to a problem is, by definition, an algorithm for that problem.* It is usual to divide correctness into two parts: partial correctness and halting. Partial correctness says that the algorithm gives the correct output provided it halts. In some algorithms, correctness may be trivial but this is not always true.

EXERCISES

Exercise 3.1: What problems do the programs in Figure 1(a) and (b) solve, respectively? You have some leeway in giving them suitable interpretations. \diamond

Exercise 3.2: (a) Extend the program in Figure 1(a) so that it sorts three input elements $\{x, y, z\}$.
 (b) In general, define what it means to say that a comparison-tree program sorts a set $\{x_1, \dots, x_n\}$ of elements. \diamond

Exercise 3.3: Design a tree program to merge two sorted lists (x, y, z) and (a, b, c, d) . The height of your tree should be 5 (the optimum). \diamond

END EXERCISES

§4. The RAM Model

The computational models such as decision-tree model are rather special, and suitable only for sorting-like problems. We now present a more general model which is suitable for most problems we shall encounter.

This is a uniform computational model called the **random-access memory model** (RAM model). It is basically an abstract assembly language. The RAM model can be regarded as a generalization of an even simpler model called the **Register model**. We describe a register machine first.

(a) Data objects. We assume infinitely many (**storage**) **registers**, each indexed by an integer $0, \pm 1, \pm 2, \pm 3, \dots$. Register 0 is special and is known as the **accumulator**. Each register can store an integer. Note that the integer is arbitrarily large (in contrast to real computers whose integer are limited to 32- or 63-bit values).

(b) Primitive Operations. Each operation has the 4-field format

$$[\langle LABEL \rangle] \langle OPERATOR \rangle \langle ARG \rangle [\langle COMMENTS \rangle]$$

where $\langle LABEL \rangle$ and $\langle COMMENTS \rangle$ are arbitrary non-empty alphabetic strings – the square brackets indicate that these are optional fields. There is one argument $\langle ARG \rangle$ whose nature is determined by the $\langle OPERATOR \rangle$. The argument $\langle ARG \rangle$ is either an integer (denoted n below) or a label (denoted ℓ below). The contents of register n is a number, denoted $c(n)$. Let $c : \mathbb{Z} \rightarrow \mathbb{Z}$ be the **contents function**. The operators, their arguments and actions are specified in the following table:

Operator	Argument	Semantics
GET	n	$c(0) \leftarrow c(n)$.
PUT	n	$c(n) \leftarrow c(0)$.
ZERO	n	$c(n) \leftarrow 0$.
INC	n	$c(n) \leftarrow c(n) + 1$.
ADD	n	$c(0) \leftarrow c(0) + c(n)$.
SUB	n	$c(0) \leftarrow c(0) - c(n)$.
MUL	n	$c(0) \leftarrow c(0) \times c(n)$.
DIV	n	$c(0) \leftarrow c(0) \div c(n)$, error if $c(n) = 0$.
JUMP	ℓ	Go to label ℓ .
JPOS	ℓ	If $c(0) > 0$ then go to ℓ .
JNEG	ℓ	If $c(0) < 0$ then go to ℓ .
HALT		The computation halts.

(c) Semantics. A **register program** is any finite sequence of such primitive operations. Associated with the program is a process called a **computation**, which is a sequence of transformations of the contents function. The first transformation is defined by executing the first instruction of the program. The transformation corresponding to each instruction is given by the last column of the above table. For instance, “GET 4” will put update the value of $c(0)$ to the value $c(4)$, but no other register contents are changed. Subsequently, transformations are defined by successive instructions of the program, unless there is a successful jump to some label ℓ . In that case, the next instruction with label ℓ is executed (multiply defined labels is no problem). The computation halts on encountering the HALT operation, or when there is no “next” instruction, or jumping to some non-existent label.

(d) Input/Output conventions. We assume that a finite number of registers is initialized with an encoding of the input, while the rest of the registers are initially zero. The convention for the output is some simple function of the final contents function. E.g., the output is defined to be $c(0)$. As an exercise, the reader may write a RAM program to compute the maximum of three numbers (convention: the input numbers are in the registers 1, 2, 3 and the maximum value is to be output in register 0).

Random Access Feature. Notice that a register program can only access a fixed set of registers. In order to allow a program to access an arbitrary number of registers, we introduce “indirect addressing”, a

concept from computer architecture. We allow another form of integer argument, denoted $@n$ where $n \in \mathbb{N}$. This means we are using the value $c(c(n))$ instead of $c(n)$ as the actual argument for the operation. Thus $c(n)$ is interpreted as the address of the actual argument. We call $@n$ an “indirect argument”. For instance, “GET @4” results in the assignment $c(0) \leftarrow c(c(4))$. If register 4 contains 256, and register 256 contains -1 , this means $c(0)$ is assigned the value -1 . Similarly, “PUT @4” will place the value $c(0)$ into register 256. Thus, a **RAM program** is basically a register program in which we allow indirect addressing.

Variants. By design, the instruction set of our RAM is parsimonious and rather primitive. There are obvious variants where we enrich the instruction set which amount to combinations of these primitive instructions. We can also allow another kind of integer argument denoted $=n$. In this case the value n itself is being used. This is called a “literal argument”. Literal arguments are useful for GET and the arithmetic instructions, but meaningless for PUT instruction.

The above model may be called an **integer RAM model**. This can be generalized to the **real RAM model** where where the registers can store an arbitrary real number, and possibly augmenting the primitive operations with other real functions (such as computing square roots). When we use a register n for indirect addressing ($@n$), we need to have some convention for handling the case where its contents $c(n)$ is not an integer.

Higher Level Languages. In practice, we may write our program using a more abstract language or a “higher level” language. Thus, we may use well-known constructs such as for-loops and if-then-else, and even allow recursion. Nevertheless, such extensions of the model can be translated into a standard RAM program.

 EXERCISES

Exercise 4.1: Recall the Merge Algorithm described in §3. Please convert it into a RAM program. \diamond

Exercise 4.2: Write RAM algorithms for the following problems:

- (a) Sort a sequence of input numbers.
 - (b) Compute the GCD of two integers.
 - (c) Solve the ranking problem (§2). You need two algorithms, for preprocessing and for answering queries.
 - (d) Solve the dynamic ranking problem (§2). You need four algorithms here: to initialize an empty data structure, to insert, to delete, and to answer rank queries.
- Be sure to state your input/output conventions. Also, describe any data structures you use. \diamond

Exercise 4.3: Reduce the number of primitive operators (listed in table above) for our RAM model to a minimum. Show that your minimal RAM model can simulate our original RAM model. \diamond

Exercise 4.4: Show how to implement higher level language constructs such as *for-loops*, *if-then-else*, *case-statements* in our RAM model. \diamond

 END EXERCISES

§5. Complexity Model: How to assess algorithms?

We now have a suitable computational model for solving our problem. What is a criteria to choose among different algorithms within a model? For this, we need to introduce a **complexity model**.

In most computational models, there are usually natural notions of **time** and **space**. These are two examples of **computational resources**. Naturally, resources are scarce and algorithms consume resources when they run. We want to choose algorithms that minimize the use of resources. In our discussions, we focus on only one resource at a time, usually time (occasionally space). So we avoid issues of trade-offs between two resources.

Next, for each primitive operation executing on a particular data, we need to know how much of the resource is consumed. For instance, in **Java**, we could define each execution of the addition operation on two numbers a, b to use time $\log(|a| + |b|)$. But it would be simpler to say that this operation takes unit time, independent of a, b . This simpler version is our choice throughout these lectures: *each primitive operation takes unit time, independent of the actual data.*

How is the running time for sorting 1000 elements related to the running time for sorting 10 elements? The answer lies in viewing running time as a function of the number of input elements, the “input size”. In general, problems usually have a natural notion of “input size” and this is the basis for understanding the complexity of algorithms.

So we want a notion of **size** on the input domain, and measure resource usage as a function of input size. The size $size(I)$ of an input instance I is a positive integer. We make a general assumption about the size function: *there are inputs of arbitrarily large size.*

For our running example of the sorting problem, it may seem natural to define the size of an input (a_1, \dots, a_n) to be n . But actually, this is only natural because we usually use computational models that compares a pair of numbers in unit time. For instance, if we must encode the input as binary strings (as in the Turing machine model), then input size is better taken to be $\sum_{i=1}^n (1 + \log(1 + |a_i|))$.

Suppose A is an algorithm for our problem P . For any input instance I , let $T_A(I)$ be the total amount of time used by A on input I . Naturally, $T_A(I) = \infty$ if A does not halt on I . Then we define the **worst case running time** of A to be the function $T_A(n)$ where

$$T_A(n) := \sup\{T_A(I) : size(I) = n\}$$

Using “sup” here is only one way to “aggregate” the set of numbers $\{T_A(I) : size(I) \leq n\}$. In general, we may apply a set-valued function G to the set,

$$T_A(n) = G(\{T_A(I) : size(I) \leq n\})$$

For instance, if G is the **average** function and we get **average time complexity**.

To summarize: a **complexity model** is a specification of

- (a) the computational resource,
- (b) the input size function,
- (c) the unit of resource, and
- (d) the method G of aggregating.

Once the complexity model is fixed, we can associate to each algorithm A a **complexity function** T_A .

■ **Example:** Consider the Comparison Tree Model for sorting. Let $T(n)$ be the worst case number of comparisons needed to sort n elements. Any tree program to sort n elements must have at least $n!$ leaves,

since we need at least one leaf for each possible sorting outcome. Since a binary tree with $n!$ leaves has height at least $\lceil \lg(n!) \rceil$.

LEMMA 1. *Every tree program for sorting n elements has height at least $\lceil \lg(n!) \rceil$, i.e., $T(n) \geq \lceil \lg(n!) \rceil$.*

This lower bound is called the **Information Theoretic Lower Bound** for sorting.

■ **Example:** In our RAM model (real or integer version), let the computational resource be time, where each primitive operation takes unit time. The input size function is the number of registers used for encoding the input. The aggregation method is the worst case (for any fixed input size). This is called the **unit time** complexity model.

Complexity of Merging. Define $M(m, n)$ to be the minimum height of any comparison tree for merging two sorted lists of sizes m and n , respectively. We can prove the following bounds

$$M(m, n) \leq m + n - 1$$

and

$$M(m, n) \geq 2 \min\{m, n\} - \delta(m, n)$$

where $\delta(m, n) = 1$ if $m = n$ and $\delta(m, n) = 0$ otherwise. The upper bound comes from the algorithm for merging described in §3: each comparison results in at least one output. But the last element can be output without any comparison. Hence we never make more comparisons than $m + n - 1$. The lower bound comes from the following example: assume the input is $x_1 < x_2 < \dots < x_m$ and $y_1 < \dots < y_n$ where $m \geq n$. We assume that

$$x_1 < y_1 < x_2 < y_2 < x_3 < \dots < x_n < y_n < x_{n+1} < \dots < x_m.$$

Note that each y_i must be compared to x_i and x_{i+1} (for $i = 1, \dots, n - 1$). Moreover, y_n must be compared to x_n , and in case $\delta(m, n) = 0$, y_n must also be compared to x_{n+1} . This proves $M(m, n) \geq 2n - \delta(m, n)$, where $n = \min\{m, n\}$. This method of proving lower bounds is called the **adversary argument**.

A corollary of the above upper and lower bounds are some exact bounds for the complexity of merging:

$$M(m, m) = 2m - 1$$

and

$$M(m, m + 1) = 2m.$$

Thus the uniform algorithm is optimal in these cases. More generally, $M(m, m + k) = 2m + k - 1$ for $k = 0, \dots, 4$ and $m \geq 6$ (see [3] and Exercise).

Now consider the other extreme where the two input lists are as disparate in lengths as possible: $M(m, 1)$. In this case, the information theoretic bound says that $M(m, 1) \geq \lceil \lg(m + 1) \rceil$ (why?). Also, by binary search, this lower bound is tight. Hence we now know another exact value:

$$M(m, 1) = \lceil \lg(m + 1) \rceil.$$

A non-trivial result from Hwang and Lin says

$$M(2, n) = \lceil \lg 7(n + 1)/12 \rceil + \lceil \lg 14(n + 1)/17 \rceil.$$

More generally, the **information-theoretic bound** says

$$M(m, n) \geq \lg \binom{m + n}{m}$$

since there are $\binom{m+n}{n}$ ways of merging the two sorted lists. To see this, imagine that we already have the sorted list of $m+n$ elements: but which of these elements come from the list of size m ? There are $\binom{m+n}{m}$ ways of choosing these elements.

Thus we have two distinct methods for proving lower bounds on $M(m, n)$: the adversary method is better when $|m - n|$ is small, and the information theoretic bound is better when this gap is large. The exact value of $M(m, n)$ is known for several other cases, but a complete description of this complexity function remains an open problem.

Other Complexity Measures. There are complexity models, especially in computational geometry, in which an output size function is taken into account. The complexity function would now take at least two arguments, $T(n, k)$ where n is the input size, but k is the output size. This is the **output-sensitive complexity model**.

Remarks:

1. Another kind of complexity measure is the **size** of a program. In the RAM model, this can be the number of primitive instructions. We can measure the complexity of a problem P in terms of the size $s(P)$ of the smallest program that solves P . This complexity measure assigns a single number $s(P)$, not a complexity function, to P . This **program size measure** is an instance of **static complexity measure**; in contrast, time and space are examples of **dynamic complexity measures**. Here “dynamic” (“static”) refers to fact that the measure depends (does not depend) on the running of a program. Complexity theory is mostly developed for dynamic complexity measures.
2. The comparison tree complexity model ignores all the other computational costs except comparisons. In most situations this is well-justified. But it is possible² to create conjure up ridiculous algorithms which minimize the comparison cost, at an exorbitant cost in other operations.

EXERCISES

Exercise 5.1: How many comparisons are required in the worst case to sort 10 elements? Give a lower bound in the comparison tree model. Note: to do the computation by hand, it is handy to know that $10! = 3,628,800$ and $2^{20} = 1,048,576$. ◇

Exercise 5.2: How good is the information theoretic lower bound for sorting 3 elements a sharp bound? In other words, can you find upper bounds that matches the information-theoretic lower bound? Repeat this exercise for 4 and 5 elements. ◇

Exercise 5.3: The following is a variant of the previous exercise. Is it always possible to sort n elements using a comparison tree with $n!$ leaves? Check this out for $n = 3, 4, 5$. ◇

Exercise 5.4: (a) Consider a variant of the unit time complexity model for the integer RAM model, called the **logarithmic time complexity model**. Each operand takes time that is logarithmic in the address of the register and logarithmic in the size of its operands. What is the relation between the logarithmic time and the unit time models?

(b) Is this model realistic in the presence of the arithmetic operators (ADD, SUB, MUL, DIV). Discuss. ◇

²My colleague, Professor Robert Dewar suggests the following example: given n numbers to be sorted, we first search for all potential comparison trees for sorting n elements. To make this search finite, we only evaluate comparison trees of height at most $n \lceil \lg n \rceil$. Among those trees that we have determined to be able to sort, we pick one of minimum height. Now we run this comparison tree on the given input.

Exercise 5.5: Describe suitable complexity models for the “space” resource in integer RAM models. Give two versions, analogous to the unit time and logarithmic time versions. What about real RAM models? \diamond

Exercise 5.6: Justify the claim that $M(m, 1) = \lceil \lg(m + 1) \rceil$. \diamond

Exercise 5.7: Using direct arguments, give your best upper and lower bounds for $M(2, 10)$. \diamond

Exercise 5.8: Prove that $M(m, m + i) = 2m + i - 1$ for $i = 2, 3, 4$ for $m \geq 6$. \diamond

Exercise 5.9: Prove that $M(k, m) \geq k \lg_2(m/k)$ for $k \leq m$. HINT: split the list of length m into three sublists of roughly equal sizes. \diamond

Exercise 5.10: Open problem: determine $M(m, 3)$ and $M(m, m + 5)$ for all m . \diamond

Exercise 5.11: With respect to the comparator circuit and tree program models in §3, describe suitable complexity models for each. \diamond

END EXERCISES

§6. Algorithmic Techniques: How to design algorithms

Now that we have some criteria to judge algorithms, we begin to design algorithms. There emerges some general paradigms of algorithms design: (i) Divide-and-conquer (e.g., merge sort)
(ii) Greedy method (e.g., Kruskal’s algorithm for minimum spanning tree)
(iii) Dynamic programming (e.g., multiplying a sequence of matrices)
(iv) Incremental method (e.g., insertion sort)

Let us briefly outline the merge sort algorithm to illustrate divide-and-conquer: Suppose you want to sort an array A of n elements. Assume n is a power of 2. Here is the Merge Sort algorithm on input A :

1. (Basis) If n is 1 simply return the array A .
2. (Divide) Divide the elements of A into two subarrays B and C of size $n/2$ each.
3. (Recurse) Recursively, call the Merge Sort algorithm on B . Do the same for C .
4. (Conquer) Merge the sorted arrays B and C and put the result back into array A .

There is only one non-trivial step, the merging of two sorted arrays. We leave this as an exercise.

There are many variations or refinements of these paradigms. E.g., Kirkpatrick and Seidel [2] introduced a form of divide-and-conquer (called “marriage-before-dividing”) that leads to an output-sensitive convex

hull algorithm. There may be domain specific versions of these methods. E.g., plane sweep is an incremental method suitable for problems on points in Euclidean space.

Closely allied with the choice of algorithmic technique is the choice of *data structures*. A data structure is a representation of a complex mathematical structure (such as sets, graphs or matrices), together with algorithms to support certain querying or updating operations. The following are some basic data structures.

- (a) **Linked lists:** each list stores a sequence of objects together with operations for (i) accessing the first object, (ii) accessing the next object, (iii) inserting a new object after a given object, and (iv) deleting any object.
- (b) **LIFO, FIFO queues:** each queue stores a set of objects under operations for insertion and deletion of objects. The queue discipline specifies which object is to be deleted. There are two³ basic disciplines: last-in first-out (LIFO) or first-in first-out (FIFO). Note that recursion is intimately related to LIFO.
- (c) **Binary search trees:** each tree stores a set of elements from a linear ordering together with the operations to determine the smallest element in the set larger than a given element. A dynamic binary search tree supports, in addition, the insertion and deletion of elements.
- (d) **Dictionaries:** each dictionary stores a set of elements and supports the operations of (i) inserting a new element into the set, (ii) deleting an element, and (iii) testing if a given element is a member of the set.
- (e) **Priority queues:** each queue stores a set of elements from a linear ordering together with the operations to (i) insert a new element, (ii) delete the minimum element, and (iii) return the minimum element (without removing it from the set).

EXERCISES

- Exercise 6.1:** (a) Give a pseudo-code description of $Merge(B, C, A)$ which, given two sorted arrays B and C of size n each, returns their merged (hence sorted) result into the array A of size $2n$.
 (b) Why did we assume n is a power of 2 in the description of merge sort? How can we justify this assumption in theoretical analysis? How can we handle this assumption in practice? \diamond

- Exercise 6.2:** Design an incremental sorting algorithm based on the following principle: assuming that the first m elements have been sorted, try to add (“insert”) the $m + 1$ st element into the first m elements to extend the inductive hypothesis. \diamond

END EXERCISES

§7. Analysis: How to estimate complexity

We have now a measure T_A of the complexity of our algorithm A , relative to some complexity model. Unfortunately, the function T_A is generally too complex to admit a simple description, or to be expressed in terms of familiar mathematical functions. Instead, we aim to give upper and lower bounds on T_A . This

³A discipline of a different sort is called GIGO, or, garbage-in garbage-out. This is really a law of nature.

constitutes the subject of **algorithmic analysis** which is a major part of this book. The tools for this analysis depends to a large extent on the algorithmic paradigm or data structure used by A . We give two examples.

■ **Example:** (Divide-and-conquer) If we use divide-and-conquer then it is likely we need to solve some recurrence equations. In our Merge Sort algorithm, assuming n is a power of 2, we obtain the following recurrence:

$$T(n) = 2T(n/2) + Cn$$

for $n \geq 2$ and $T(1) = 1$. Here $T(n)$ is the (worst case) number of comparisons needed by our algorithm to sort n elements. The solution is $T(n) = \Theta(n \log n)$. In the next chapter, we study the solutions of such equations.

■ **Example:** (Amortization) If we employ certain data-structures that might be described as “lazy” then amortization analysis might be needed. Let us illustrate this with the problem of maintaining a binary search tree under repeated insertion and deletion of elements. Ideally, we want the binary tree to have height $O(\log n)$ if there are n elements in the tree. There are a number of known solutions for this problem (see Chapter 3). Such a solution achieves the optimal logarithmic complexity for *each* insertion/deletion operation. But it may be advantageous to be lazy about maintaining this logarithmic depth property: such laziness may be rewarded by a simpler coding or programming effort. The price for laziness is that our complexity may be linear for individual operations, but we still logarithmic cost in the **amortized sense**. To illustrate this idea, suppose we allow the tree to grow to non-logarithmic depth as long as it does not cost us anything (*i.e.*, there are no queries on a leaf with big depth). But when we have to answer a query on a “deep leaf”, we take this opportunity to restructure the tree so that the depth of this leaf is now reduced (say halved). Thus repeated queries to this leaf will make it shallow. The cost of a single query could be linear time, but we hope that over a long sequence of such queries, the cost is amortized to something small (say logarithmic). This technique prevents an adversary from repeated querying of a “deep leaf”. Unfortunately, this is not enough because the very first query into a “deep leaf” has to be amortized as well (since there may be no subsequent queries). To anticipate this amortization cost, we “pre-charge” the requests (insertions) that lead to this inordinate depth. Using a financial paradigm, we put the pre-paid charges into some bank account. Then the “deep queries” can be paid off by withdrawing from this account. Amortization is both an algorithmic paradigm as well as an analysis technique. This will be treated in Chapter 6.

§8. Asymptotics: How robust is the model?

This section contains important definitions for the rest of the book.

We started with a problem, selected a computational model and an associated complexity model, designed an algorithm and managed to analyze its complexity. Looking back at this process, we are certain to find arbitrariness in our choices. For instance, would a simple change in the set of primitive operations change the complexity of your solution? Or what if we charge two units of time for some of the operations? Of course, there is no end to such revisionist afterthoughts. What we are really seeking is a certain robustness or invariance in our results.

What is a complexity function? In this book, we call a partial real function

$$f : \mathbb{R} \rightarrow \mathbb{R}$$

a **complexity function** (or simply, “function”). We use complexity functions to quantify the complexity of our algorithms. Why do we consider *partial* functions? For one thing, many functions of interest are only

defined on positive integers. For example, the running time $T_A(n)$ of an algorithm A that takes discrete inputs is a partial real function (normally defined only when n is a natural number). Of course, if the domain of T_A is taken to be \mathbb{N} , then $T_A(n)$ would be total. So why do we think of \mathbb{R} as the domain of $T_A(n)$? Again, we often use functions such $f(n) = n/2$ or $f(n) = \sqrt{n}$, to bound our complexity functions, and these are naturally defined on the real domain; all the tools of analysis and calculus becomes available to analyze such functions. Many common real functions such as $f(n) = 1/n$ or $f(n) = \log n$ are partial functions because $1/n$ is undefined at $n = 0$ and $\log n$ is undefined for $n \leq 0$. If $f(n)$ is not defined at n , we write $f(n) = \uparrow$, otherwise $f(n) = \downarrow$. Since complexity functions are partial, we have to be careful about operations such as functional composition.

Designated variable and Anonymous functions. In general, we will write “ n^2 ” and “ $\log x$ ” to refer to the functions $f(n) = n^2$ or $g(x) = \log x$, respectively. Thus, the functions denoted n^2 or $\log x$ are **anonymous** (or self-naming). This convention is very convenient, but it relies on an understanding that “ n ” in n^2 or “ x ” in $\log x$ is the **designated variable** in the expression. For instance, the anonymous complexity function $2^x n$ is a linear function if n is the designated variable, but an exponential function if x is the designated variable. *The designated variable in complexity functions, by definition, range over real numbers.* This may be a bit confusing when the designated variable is “ n ” since in mathematical literature, n is usually a natural number.

Robustness or Invariance issue. Let us return to the robustness issue which motivated this section. The motivation was to state complexity results that have general validity, or independent of many apparently arbitrary choices in the process of deriving our results. There are many ways to achieve this: for instance, we can specify complexity functions up to “polynomial smearing”. Two real functions f, g , are equivalent in this sense if for some $c > 0$, $f(n) \leq cg(n)^c$ and $g(n) \leq cf(n)^c$ for all n large enough. This is *extremely* robust but alas, too drastic for most purposes. The most widely accepted procedure is to take two smaller steps:

- Step 1: We are interested in the eventual behavior of functions (e.g., if $T(n) = 2^n$ for $n \leq 1000$ and $T(n) = n$ for $n > 1000$, then we want to regard $T(n)$ as a linear function).
- Step 2: We distinguish functions only up to multiplicative constants (e.g., $n/2$, n and $10n$ are indistinguishable),

These two decisions give us most of the robustness properties we desire, and are captured in the following language of asymptotics.

Eventuality. This is Step 1 in our search for invariance. Given two functions, we say “ $f \leq g$ **eventually**”, written

$$f \leq g \text{ (ev.)}, \tag{1}$$

if $f(x) \leq g(x)$ holds for all x large enough. More precisely, this means there is some x_0 such that the following statement is true:

$$(\forall x)[x \geq x_0 \Rightarrow f(x) \leq g(x)]. \tag{2}$$

We must be careful: *what does “ $f(x) \leq g(x)$ ” mean when either $f(x)$ or $g(x)$ may be undefined?* The answer depends on the quantifier that bounds (i.e., that controls) x , whether x is bounded by an existential quantifier or by a universal quantifier. If a universal quantifier (as in (1)), we declare the predicate “ $f(x) \leq g(x)$ ” to be true if either $f(x)$ or $g(x)$ is undefined. If an existential quantifier, we declare the predicate “ $f(x) \leq g(x)$ ” to be false if either $f(x)$ or $g(x)$ is undefined. So (2) can be expanded into:

$$(\forall x)[x \geq x_0 \wedge f(x) = \downarrow \wedge g(x) = \downarrow \Rightarrow f(x) \leq g(x)].$$

REMARK: We can generalize this treatment of quantification over any **partial predicate** $P(x)$ on a real variable x . By partial predicate, we mean $P(x)$ may be undefined for some values of x . In the previous example, $P(x)$ is just “ $f(x) \leq g(x)$ ”. The universally quantified statement “ $(\forall x)[P(x)]$ ” should be interpreted as saying “for all $x \in \mathbb{R}$, if $P(x)$ is defined then $P(x)$ is true”. Similarly, the existentially quantified statement “ $(\exists x)[P(x)]$ ” says “there exists some $x \in \mathbb{R}$ such that $P(x)$ is defined and $P(x)$ is true”. Note that the definition of “holding eventually” involves both kinds of quantifiers. To show the role of the x variable, we may also write (1) as

$$f(x) \leq g(x) \text{ (ev. } x\text{)}.$$

Clearly, this is a transitive relation.

The “eventually” terminology is quite general: if a predicate $R(x)$ is parametrized by x in some real domain $D \subseteq \mathbb{R}$, and $R(x)$ holds for all $x \in D$ larger than some x_0 , then we say $R(x)$ **holds eventually** (abbreviated, ev.). We can also extend this to predicates $R(x, y, z)$ on several variables. A related notion is this: if $R(x)$ holds for infinitely many values of $x \in D$, we say $R(x)$ **holds infinitely often** (abbreviated, i.o.).

If $g \leq f$ (ev.) and $f \leq g$ (ev.), then clearly

$$g = f \text{ (ev.)}.$$

Thus means $f(x) = g(x)$ for sufficiently large x , whenever both sides are defined. Most natural functions f in complexity satisfies $f \geq 0$ (ev.) and are non-decreasing eventually.

Domination. We now take Step 2 towards invariance. We say g **dominates** f , written

$$f \preceq g,$$

if there exists $C > 0$ such that $f \leq C \cdot g$ (ev.). Equivalently, $f \preceq g$ is written as $g \succeq f$. This notation naturally suggests the transitivity property: $f \preceq g$ and $g \preceq h$ implies $f \preceq h$. Of course, the reflexivity property holds: $f \preceq f$. If $f \preceq g$ and $g \preceq f$ then we write

$$f \asymp g.$$

Clearly \asymp is an equivalence relation. The equivalence classes of f is called the Θ -**order** of f ; more on this below. If $f \preceq g$ but not $g \preceq f$ then we write

$$f \prec g.$$

E.g., $1 + \frac{1}{n} \prec n \prec n^2$.

The big-Oh notation. We write

$$\mathcal{O}(f)$$

(and read **order of** f or **big-Oh of** f) to denote the set of all complexity functions g such that

$$0 \preceq g \preceq f.$$

Note that each function in $\mathcal{O}(f)$ dominates 0. In other words, if $g = \mathcal{O}(f)$ then there is some $C > 0$ and x_0 such that for all $x \geq x_0$, if $g(x) = \downarrow$ and $f(x) = \downarrow$ then $0 \leq g(x) \leq Cf(x)$.

E.g., The set $\mathcal{O}(1)$ is the set of functions f that is bounded. The function $1 + \frac{1}{n}$ is a member of $\mathcal{O}(1)$.

The simplest usage of this \mathcal{O} -notation is as follows: we write

$$g = \mathcal{O}(f)$$

(and read ‘ g is **big-Oh** of f ’ or ‘ g is **order of** f ’) to mean g is a member of the set $\mathcal{O}(f)$. The equality symbol ‘ $=$ ’ here is “uni-directional”: $g = \mathcal{O}(f)$ does not mean the same thing as $\mathcal{O}(f) = g$. Below, we will see how to interpret the latter expression. The equality symbol in this context is called a **one-way equality**. Why not just use ‘ \in ’ for the one-way equality? A partial explanation is that one common use of the equality symbol has a uni-directional flavor where we transform a formula from an unknown form into a known form, separated by an equality symbol. Our one-way equality symbol for \mathcal{O} -expressions lends itself to a similar manipulation. For example, the following sequence of one-way equalities

$$f(n) = \sum_{i=1}^n \left(i + \frac{n}{i}\right) = \left(\sum_{i=1}^n i\right) + \left(\sum_{i=1}^n \frac{n}{i}\right) = \mathcal{O}(n^2) + \mathcal{O}(n \log n) = \mathcal{O}(n^2)$$

may be viewed as a derivation to show f is at most quadratic.

Big-Oh expressions. The expression ‘ $\mathcal{O}(f(n))$ ’ is an example of an \mathcal{O} -expression, which we now define. In any \mathcal{O} -expression, there is a **designated variable** which is the real variable that goes⁴ to infinity. For instance, the \mathcal{O} -expression $\mathcal{O}(n^k)$ would be ambiguous were it not for the tacit convention that ‘ n ’ is normally the designated variable. Hence k is assumed to be constant. We shall define **\mathcal{O} -expressions** as follows:

(Basis) If f is the symbol for a function, then f is an \mathcal{O} -expression. If n is the designated variable for \mathcal{O} -expressions and c a real constant, then both ‘ n ’ and ‘ c ’ are also \mathcal{O} -expressions.

(Induction) If E, F are \mathcal{O} -expressions and f is a symbol denoting a complexity function then the following are \mathcal{O} -expressions:

$$\mathcal{O}(E), \quad f(E), \quad E + F, \quad EF, \quad -F, \quad 1/F, \quad E^F.$$

Each \mathcal{O} -expression E denotes a set \tilde{E} of partial real functions in the obvious manner: in the basis case, a function symbol f denotes the singleton set $\tilde{f} = \{f\}$. Inductively, the expression $E + F$ (for instance) denotes the set $\tilde{E} + \tilde{F}$ of all functions $f + g$ where $f \in \tilde{E}$ and $g \in \tilde{F}$.

Examples of \mathcal{O} -expressions:

$$2^n - \mathcal{O}(n^2 \log n), \quad n^{n+\mathcal{O}(\log n)}, \quad f(1 + \mathcal{O}(1/n)) - g(n).$$

If E, F are two \mathcal{O} -expressions, we may write

$$E = F$$

to denote $\tilde{E} \subseteq \tilde{F}$, *i.e.*, the equality symbol stands for set inclusion! This generalizes our earlier “ $f = \mathcal{O}(g)$ ” interpretation. Some examples of this usage:

$$\mathcal{O}(n^2) - 5^{\mathcal{O}(\log n)} = \mathcal{O}(n^{\log n}), \quad n + (\log n)\mathcal{O}(\sqrt{n}) = n^{\log \log n}, \quad 2^n = \mathcal{O}(1)^{n-\mathcal{O}(1)}.$$

An ambiguity arises from the fact that if \mathcal{O} does not occur in an \mathcal{O} -expression, it is indistinguishable from an ordinary expression. We must be explicit about our intention, or else rely on the context in such cases. Normally, at least one side of the one-sided equation ‘ $E = F$ ’ contains an occurrence of ‘ \mathcal{O} ’, in which case, the other side is automatically assumed to be an \mathcal{O} -expression. Some common \mathcal{O} -expressions are:

- $\mathcal{O}(1)$, the bounded functions.

⁴More generally, we can consider x approaching some other limit, such as 0.

- $1 \pm \mathcal{O}(1/n)$, a set of functions that tends to 1^\pm .
- $\mathcal{O}(n)$, the linearly bounded functions.
- $n^{\mathcal{O}(1)}$, the functions bounded by polynomials.
- $\mathcal{O}(1)^n$ or $2^{\mathcal{O}(n)}$, the functions bounded by simple exponentials.
- $\mathcal{O}(\log n)$, the functions bounded by some multiple of the logarithm.

Extended big-Oh notations. We introduce two simple extensions of the \mathcal{O} -notation:

1) **Inequality interpretation:** For \mathcal{O} -expressions E, F , we may write $E \neq F$ to mean that the set of functions denoted by E is not contained in the set denoted by F . For instance, $f(n) \neq \mathcal{O}(n^2)$ means that for all $C > 0$, there are infinitely many n such that $f(n) > Cn^2$.

2) **Subscripting convention:** We can subscript the big-Oh's in an \mathcal{O} -expression. For example,

$$O_A(n), \quad O_1(n^2) + O_2(n \log n).$$

The intent is that each subscript ($A, 1, 2$) picks out a specific but anonymous function in (the set denoted by) the unsubscripted \mathcal{O} -notation. Furthermore, within a given context, two occurrences of an identically subscripted \mathcal{O} -notation are meant to refer to the same function.

For instance, if A is a linear time algorithm, we may say that “ A runs in time $O_A(n)$ ” to indicate that the choice of the function $O_A(n)$ depends on A . Further, all occurrences of “ $O_A(n)$ ” in the same discussion will refer to the same anonymous function. Again, we may write

$$n2^k = O_k(n), \quad n2^k = O_n(2^k)$$

depending on one's viewpoint. Especially useful is the ability to do “in-line calculations”. As an example, we may write

$$g(n) = O_1(n \log n) = O_2(n^2)$$

where, it should be noted, the equalities here are true equalities of functions.

Related Asymptotic Notations: The above discussion extends in a natural way to several other related notations.

Big-Omega notation: $\Omega(f)$ is the set of all complexity functions g such that for some constant $C > 0$,

$$C \cdot g \geq f \geq 0 \text{ (ev.)}.$$

Note that $\Omega(f)$ is empty unless it is eventually non-negative. Clearly, big-Omega is just the reverse of the big-Oh relation: g is in $\Omega(f)$ iff $f = \mathcal{O}(g)$.

Theta notation: $\Theta(f)$ is the intersection of the sets $\mathcal{O}(f)$ and $\Omega(f)$. So g is in $\Theta(f)$ iff $g \asymp f$.

Small-oh notation: $o(f)$ is the set of all complexity functions g such that for all $C > 0$,

$$C \cdot f \geq g \geq 0 \text{ (ev.)}.$$

Thus g is in $o(f)$ implies $g(n)/f(n) \rightarrow 0$ as $n \rightarrow \infty$. Also, $o(f) \subseteq \mathcal{O}(f)$. A related notation is this: we say

$$f \sim g$$

if $f = g \pm o(g)$ or $f(x) = g(x)[1 \pm o(1)]$.

Small-omega notation: $\omega(f)$ is the set of all functions g such that for all $C > 0$,

$$C \cdot g \geq f \geq 0 \text{ (ev.)}.$$

Thus g is in $\omega(f)$ implies $g(n)/f(n) \rightarrow \infty$ as $n \rightarrow \infty$. Clearly $\omega(f) \subseteq \Omega(f)$.

For each of these notations, we again define the \circ -expressions ($\circ \in \{\Omega, \Theta, o, \omega\}$), use the one-way inequality instead of set-membership or set-inclusion, and employ the subscripting convention. Thus, we write “ $g = \Omega(f)$ ” instead of saying “ g is in $\Omega(f)$ ”. We call the set $\circ(f)$ the \circ -**order** of f . Here are some immediate relationships among these notations:

- $f = \mathcal{O}(g)$ iff $g = \Omega(f)$.
- $f = \Theta(g)$ iff $f = \mathcal{O}(g)$ and $f = \Omega(g)$.
- $f = \mathcal{O}(f)$ and $\mathcal{O}(\mathcal{O}(f)) = \mathcal{O}(f)$.
- $f + o(f) = \Theta(f)$.
- $o(f) \subseteq \mathcal{O}(f)$.
- $g = \omega(f)$ iff $f = o(g)$.

Lower Bounds. Based on our asymptotic notations, we can specify lower bounds on a complexity function $f(n)$ in one of three form:

- $f(n) = \Omega(g(n))$.
- $f(n) \neq \mathcal{O}(g(n))$.
- $f(n) \neq o(g(n))$.

It is not hard to see that each form is less stringent than the previous. See Exercise for how these are used in practice.

Discussion. There is some debate over the best way to define the asymptotic concepts. There is considerable divergence in the literature on the details. Here we note just two alternatives:

1. Perhaps the most common definition follows Knuth [4, p. 104] who defines “ $g = \mathcal{O}(f)$ ” to mean there is some $C > 0$ such that $|f(x)|$ dominates $C|g(x)|$. Using this definition, both $\mathcal{O}(-f)$ and $-\mathcal{O}(f)$ would mean the same thing as $\mathcal{O}(f)$. Our definition, on the contrary, allows us to distinguish⁵ between $1 + \mathcal{O}(1/n)$ and $1 - \mathcal{O}(1/n)$.

2. Again, we could have defined $\mathcal{O}(f)$ more simply to comprise those g such that for some $C > 0$, $g \leq C \cdot f$ (ev.). That is, we omit the requirement $g \geq 0$ (ev.) from our original definition. This definition is attractive because of its simplicity. But with this “simplified definition”, $\mathcal{O}(f)$ contains arbitrarily negative functions. Thus, the expression $1 - \mathcal{O}(1/n)$ is useful as an upper and lower bound under our official notation. But with the simplified definition, the expression $1 - \mathcal{O}(1/n)$ has no value as an upper bound. Our official definition opted for something that is intermediate between this simplified version and Knuth’s.

⁵On the other hand, there is no easy way to recover Knuth’s definition using our definitions. It may be useful to retain Knuth’s definition using the special notation $|\mathcal{O}(f(n))|$, etc.

We are following Cormen et al [1] in restricting the elements of $\mathcal{O}(f)$ to complexity functions that dominate 0. This approach has its own burden: thus whenever we say “ $g = \mathcal{O}(f)$ ”, we have to check that g dominates 0 (cf. exercise 1 below). In practice, this requirement is not much of a burden, and is silently passed over. If $|f - g| = \mathcal{O}(1)$, we cannot simply say “ $f = g + \mathcal{O}(1)$ ” without further analysis, because the correct statement might be $f = g - \mathcal{O}(1)$.

A common abuse is to use big-Oh notations in conjunction with the less-than or greater-than symbol: it is very tempting to write “ $f(n) \leq \mathcal{O}(g)$ ” instead of “ $f(n) = \mathcal{O}(g)$ ”. At best, this is redundant. The problem is that, once this is admitted, one may in the course of a long derivation eventually write “ $f(n) \geq E$ ” where E is an \mathcal{O} -expression. The latter is not very meaningful. Hence we regard any use of \leq or \geq symbols in \mathcal{O} -notations as illegitimate.

Perhaps most confusion (and abuse) in the literature arises from the variant definitions of the Ω -notation. For instance, one may have only shown a lower bound of the form $g(n) \neq \mathcal{O}(f(n))$ but this is claimed as a $g(n) = \Omega(f(n))$ result. In other words, the expression “ $g = \Omega(f)$ ” is interpreted to mean that there exists (or for all) $C > 0$ such that for infinitely many x , $g(x) \geq Cf(x)$.

Evidently, these asymptotic notations can be intermixed. E.g., $o(n^{\mathcal{O}(\log n)}) - \Omega(n)$. However, they can be tricky to understand and there seems to be little need for them.

 EXERCISES

Exercise 8.1: Assume $f(n) > 1$ (ev.).

- (a) Show that $f(n) = n^{\mathcal{O}(1)}$ iff there exists $k > 0$ such that $f(n) = \mathcal{O}(n^k)$. This is mainly an exercise in unraveling our notations!
- (b) Show a counter example to (a) in case $f(n) > 1$ (ev.) is false. ◇

Exercise 8.2: Prove or disprove: $f = \mathcal{O}(1)^n$ iff $f = 2^{\mathcal{O}(n)}$. ◇

Exercise 8.3: Unravel the meaning of the \mathcal{O} -expression: $1 - \mathcal{O}(1/n) + \mathcal{O}(1/n^2) - \mathcal{O}(1/n^3)$. Does the \mathcal{O} -expression have any meaning if we extend this into an infinite expression with alternating signs? ◇

Exercise 8.4: For basic properties of the logarithm and exponential functions, see the appendix in the next lecture. Show the following (remember that n is the designated variable). In each case, you must explicitly specify the constants n_0, C , etc, implicit in the asymptotic notations.

- (a) $(n + c)^k = \Theta(n^k)$. Note that c, k can be negative.
- (b) $\log(n!) = \Theta(n \log n)$.
- (c) $n! = o(n^n)$.
- (d) $\lceil \log n \rceil! = \Omega(n^k)$ for any $k > 0$.
- (e) $\lceil \log \log n \rceil! \leq n$ (ev.). ◇

Exercise 8.5: Provide either a counter-example when false or a proof when true. The base b of logarithms is arbitrary but fixed, and $b > 1$. Assume the functions f, g are arbitrary (do not assume that f and g are ≥ 0 eventually).

- (a) $f = \mathcal{O}(g)$ implies $g = \mathcal{O}(f)$.
- (b) $\max\{f, g\} = \Theta(f + g)$.
- (c) If $g > 1$ and $f = \mathcal{O}(g)$ then $\ln f = \mathcal{O}(\ln g)$. HINT: careful!

- (d) $f = \mathcal{O}(g)$ implies $f \circ \log = \mathcal{O}(g \circ \log)$. Assume that $g \circ \log$ and $f \circ \log$ are complexity functions.
 (e) $f = \mathcal{O}(g)$ implies $2^f = \mathcal{O}(2^g)$.
 (f) $f = o(g)$ implies $2^f = \mathcal{O}(2^g)$.
 (g) $f = \mathcal{O}(f^2)$.
 (h) $f(n) = \Theta(f(n/2))$.

◇

Exercise 8.6: Re-solve the previous exercise, assuming that f, g are unbounded and ≥ 0 eventually. ◇

Exercise 8.7: Suppose $T_A(n)$ is the running time of an algorithm A .

- (a) Suppose you have constructed an infinite sequence of inputs I_1, I_2, \dots of sizes $n_1 < n_2 < \dots$ such that A on I_i takes time more than $f(n_i)$. How can you express this lower bound result using our asymptotic notations?
 (b) In the spirit of (a), what would it take to prove a lower bound of the form $T_A(n) \neq \mathcal{O}(f(n))$? What must you show about of your constructed inputs I_1, I_2, \dots ?
 (c) Again, what does it take to prove a lower bound of the form $T_A(n) = \Omega(f(n))$?

◇

Exercise 8.8: Show some examples where you might want to use “mixed” asymptotic expressions. ◇

Exercise 8.9: Discuss the meaning of the expressions $n - \mathcal{O}(\log n)$ and $n + \mathcal{O}(\log n)$ under (1) our definition, (2) Knuth’s definition and (3) the “simplified definition” in the discussion. ◇

 END EXERCISES

§9. Two Dictums of Algorithmics

We state two principles in algorithmics. They will justify many of our procedures and motivate some of the fundamental questions we ask.

(A) *Complexity functions are determined only up to Θ -order.* This simply formalizes our motivation for introducing asymptotic notations, namely, concern for robust complexity results. For instance, we might prove a theorem that the running time $T(n)$ of an algorithm is “linear time” meaning $T(n) = \Theta(n)$. Then simple and local modifications to the algorithm should generally not affect the validity of this theorem.

There are of course several caveats: A consequence of this dictum is that a “new” algorithm is not considered significant unless its asymptotic order is less than previous known algorithms. This attitude could be counter-productive if it is abused. For instance, an asymptotically superior algorithm may actually be inferior when compared to another slower algorithm on all inputs of realistic sizes. For special problems, we might be interested in constant multiplicative factors. For instance, let $c > 0$ be the least constant such the median of n numbers can be found with $\leq cn$ comparisons. It is known that $c \geq 2$ [John W. John (1985)] and $c \leq 2.95$ [Zwick and Zwick (1994)].

(B) *Problems with complexity that are polynomial-bounded are feasible. Moreover, there is an unbridgeable gap between polynomial-bounded problems and those that are not polynomial-bounded.* This principle goes back to Cobham and Edmonds in the late sixties and relates to the P versus NP question. Hence, the first question we ask concerning any problem is whether it is polynomially-bounded. The answer may depend on

the particular complexity model. E.g., a problem may be polynomial-bounded in space-resource but not in time-resource, although at this moment it is unknown if this possibility can arise. Of course, polynomial-bounded complexity may not exactly be practical. In practice, a polynomial-time complexity of $T(n) = n^c$ is considered impractical if (say) c is larger than 5. So the “practically feasible class” is a rather small slice of P .

Despite the caveats, these two dictums turn out to be extremely useful. The landscape of computational problems is thereby simplified and made “understandable”. The quest for asymptotically good algorithms helps us understand the nature of the problem. Often, after a complicated but asymptotically good algorithm has been discovered, we find ways to achieve the same asymptotic result in a simpler (practical) way.

§10. Inherent Complexity: Lower Bounds

Let P be a computational problem. Suppose that we have designed an algorithm A for P and showed that it beats every known algorithm. But how do we know it would not be beaten in the future? This anxiety, common to all champions, may be abated somewhat if we could prove that A is somehow optimal.

Another way to look at this is that we want to define the **inherent complexity** of P . With any problem P , we want to associate a complexity function T_P . This function T_P will be defined in terms of the functions T_A where A range over all algorithms for P in our computational model. We can say A is **optimal** for P if A is an algorithm for P , and for every other algorithm B for P , we have $T_A = \mathcal{O}(T_B)$. In this case, we may define the inherent complexity T_P of P to be T_A . Clearly T_P is defined up to Θ -order.

Trouble is, the inherent complexity T_P may not exist. It certainly exist in non-uniform complexity models. For the sorting problem in the comparison model, it is known that $T_P(n) = n \log n$ and is achieved by several known algorithms. Such provably optimal non-linear time algorithms are extremely rare. Of course, linear time algorithms are usually optimal, but for trivial reasons.

Connection to complexity theory. Algorithmics is sometimes viewed as a subarea of complexity theory (algorithmics used to be called “concrete complexity” although this label has fallen out of use, perhaps because it sounds too much like a subfield of civil engineering). One view of complexity theory is that it is the study of inherent complexity. But in fact inherent complexity is just as important in algorithmics. The true distinction seems to be that algorithmics is the study of the inherent complexity of *individual problems* while complexity theory seeks to study properties of **classes of problems** (cf. [5, Preface]). The use of complexity classes leads to a way out of the dilemma of the missing optimal algorithms [5]. Because of its emphasis on individual problems, different machine models becomes interesting for algorithmics as they may be appropriate for different problems. But complexity theory prefers to use general-purpose machine models (Turing machines for example). Thus the comparison-tree model is extremely specialized but it is important for studying the complexity of problems on linear orders. A more general-purpose model would blur many distinctions in such problems which are of interest in algorithmics.

§A. APPENDIX: General Notations

We gather some general notations used throughout these lectures. Use this as reference. If there is a notation you do not understand from elsewhere in the book, this is a first place to look.

Definition. We write $X := \dots Y \dots$ when defining a term X in terms of $\dots Y \dots$

Numbers. Denote the set of natural numbers⁶ by $\mathbb{N} = \{0, 1, 2, \dots\}$, integers by $\mathbb{Z} = \{0, \pm 1, \pm 2, \dots\}$, rational numbers by $\mathbb{Q} = \{p/q : p, q \in \mathbb{Z}, q \neq 0\}$, the reals \mathbb{R} and complex numbers \mathbb{C} . The positive and non-negative reals are denoted $\mathbb{R}_{>0}$ and $\mathbb{R}_{\geq 0}$, respectively. The set of integers $\{i, i+1, \dots, j-1, j\}$ where $i, j \in \mathbb{N}$ is denoted $[i..j]$. So the size of $[i..j]$ is $\max\{0, j-i+1\}$. If r is a real number, let its **ceiling** $\lceil r \rceil$ be the smallest integer greater than or equal to r . Similarly, its **floor** $\lfloor r \rfloor$ is the largest integer less than or equal to r . Clearly, $\lfloor r \rfloor \leq r \leq \lceil r \rceil$. For instance, $\lfloor 0.5 \rfloor = 0$, $\lfloor -0.5 \rfloor = -1$ and $\lceil -2.3 \rceil = -2$.

Sets. The **size** or **cardinality** of a set S is the number of elements in S and denoted $|S|$. The empty set is \emptyset . A set of size one is called a **singleton**. The disjoint union of two sets is denoted $X \uplus Y$. Thus, $X = X_1 \uplus X_2 \uplus \dots \uplus X_n$ to denote a partition of X into n subsets. If X is a set, then 2^X denotes the set of all subsets of X . The **Cartesian product** $X_1 \times \dots \times X_n$ of the sets X_1, \dots, X_n is the set of all n -tuples of the form (x_1, \dots, x_n) where $x_i \in X_i$. If $X_1 = \dots = X_n$ then we simply write this as X^n . If $n \in \mathbb{N}$ then a n -set refers to one with cardinality n , and $\binom{X}{n}$ denotes the set of n -subsets of X .

Functions. If $f : X \rightarrow Y$ is a partial function, then write $f(x) \uparrow$ if $f(x)$ is undefined and $f(x) \downarrow$ otherwise. Function composition will be denoted $f \circ g : X \rightarrow Z$ where $g : X \rightarrow Y$ and $f : Y \rightarrow Z$. Thus $(f \circ g)(x) = f(g(x))$. We say a total function f is **injective** or 1-1 if $f(x) = f(y)$ implies $x = y$; it is **surjective** or **onto** if $f(X) = Y$; it is **bijective** if it is both injective and surjective.

The special functions of exponentiation $\exp_b(x)$ and logarithm $\log_b(x)$ to base $b > 0$ are more fully described in the Appendix of Chapter 2. Although these functions can be viewed as complex functions, we will exclusively treat them as real functions in this book. In particular, it means $\log_b(x)$ is undefined for $x \leq 0$. When the base b is not explicitly specified, it is assumed to be some constant $b > 1$. Two special bases deserve their own notations: $\lg x$ and $\ln x$ refer to logarithms to base $b = 2$ and base $b = e = 2.718\dots$, respectively. For any real i , we write $\log^i x$ as short hand for $(\log x)^i$. E.g., $\log^2 x = (\log x)^2$. But if i is a natural number then $\log^{(i)} x$ will denote the i -fold application of the log-function. E.g., $\log^{(2)} x = \log(\log x) = \log \log x$ and $\log^{(0)} x = x$. In fact, this notation can be extended to any integer i , where $i < 0$ indicates the $|i|$ -fold application of exp.

Programs. We write programs in a pseudo-language with standard programming constructs such as if-then-else and while statements. However, to reduce clutter, we indicate the structure of these programming constructs by indentation and newlines.

Assignment to programming variables is written $x \leftarrow \dots$, but we may also write $\dots \rightarrow x$ when convenient. We write $x++$ (resp., $++x$) to indicate the increment of a numerical variable x in a program. The value of this expression is the value of x before (resp., after) incrementing. There is an analogous notation for decrementing, $x--$ and $--x$.

⁶Zero is considered natural here, although the ancients do not consider it so. \mathbb{Z} comes from the German 'zahlen', to count.

Comments in a program are indicated in two ways: \triangleright *This is a forward comment* and \triangleleft *This is a backward comment*. These comments either precede (in case of forward comment) or follows (in case of backward comment) the code that it describes. Here is a recursive program to compute the Factorial function:

```

F(n)
  Input natural number n.
  Output n!
   $\triangleright$  Base Case
  1. if  $n \leq 1$  return(1).
   $\triangleright$  General Case
  2. return( $n \cdot F(n - 1)$ ).       $\triangleleft$  This is a recursive call

```

Note that we use indentation to indicate the block structure of our code.

Logic, Proofs, Induction. The student should know basic propositional (or Boolean) logic. But mathematical facts goes beyond propositional logic. Here is an example⁷ of a mathematical assertion $P(x, y)$ where x, y are real variables:

$$P(x, y) : \text{There exists a real } z \text{ such that if } x < y \text{ then } x < z < y. \quad (3)$$

The student should know how to parse such assertions. The assertion $P(x, y)$ happens to be true. This is logically equivalent to

$$(\forall x, y \in \mathbb{R})[P(x, y)]. \quad (4)$$

All mathematical assertions are of this nature. It is said that mathematical truths are universal: truthhood does not allow exceptions. If an assertion $P(x, y)$ has exceptions, and we can explicitly characterize the exceptions $E(x, y)$, then the pair $P(x, y) \vee E(x, y)$ constitute a true assertion.

Assertions contain variables: for example, $P(x, y)$ in (3) contains x, y, z . Each variable has an implied or explicit range (x, y, z range over “real numbers”), and each variable is either **quantified** (either by “for all” or “there exists”) or **unquantified**. Alternatively, they are either **bounded** or **free**. In our example $P(x, y)$, z is bounded while x, y are free. It is conventional to display the free variables as functional parameters of an assertion. The symbol \forall stands for “for all” and is called the **universal quantifier**. Likewise, the symbol \exists stands for “there exists” and is called the **existential quantifier**. Assertions with no free variables are called **statements**. We can always convert an assertion into a statement by adding some prefix to quantify each of the free variables. Thus, $P(x, y)$ can be converted into statements such as in (4) or as in $(\exists x \in \mathbb{R})(\forall y \in \mathbb{R})[P(x, y)]$.

Constructing proofs or providing counter examples to mathematical statements is a basic skill to cultivate. Three kinds of proofs are widely used: (i) case analysis, (ii) induction, and (iii) contradiction.

A proof by case analysis is often a matter of patience. But sometimes a straightforward enumeration of the possibilities will yield too many cases; clever insights may be needed to compress the argument. Induction is sometimes mechanical as well but very complicated inductions can also arise (Chapter 2 treats induction). Proofs by contradiction usually has a creative element: you need to find an assertion to be contradicted!

In proofs by contradiction, you will need to routinely negate a logical statement. Let us first consider the simple case of propositional logic. Here, you basically apply what is called De Morgan’s Law: if A are

⁷When we formalize the logical language of discussion, what is called “assertion” here is often called “formula”.

B are truth values, then $\neg(A \vee B) = (\neg A) \wedge (\neg B)$ and $\neg(A \wedge B) = (\neg A) \vee (\neg B)$. For instance suppose you want to contradict the proposition $A \Rightarrow B$. You need to first know that $A \Rightarrow B$ is the same as $(\neg A) \vee B$. Negating this by de Morgan's law gives us $A \wedge (\neg B)$.

Next consider the case of quantified logic. De Morgan's law becomes the following: $\neg((\forall x)P)$ is equivalent to $(\exists x)(\neg P)$; $\neg((\exists x)P)$ is equivalent to $(\forall x)(\neg P)$. A useful place to exercise these rules is to do some proofs involving the asymptotic notation (big-Oh, big-Omega, etc). See Exercise.

Formal languages. An **alphabet** is a finite set Σ of symbols. A finite sequence $w = x_1x_2 \cdots x_n$ of symbols from Σ is called a **word** or **string** over Σ ; the **length** of this string is n and denoted⁸ $|w|$. When $n = 0$, this is called the **empty string** or **word** and denoted with the special symbol ϵ . The set of all strings over Σ is denoted Σ^* . A **language** over Σ is a subset of Σ^* .

Graphs. A **hypergraph** is a pair $G = (V, E)$ where V is any set and $E \subseteq 2^V$. We call elements of V **vertices** and elements of E **hyper-edges**. In case $E \subseteq \binom{V}{k}$, we call G a k -graph. The case $k = 2$ is important and is called a **bigraph** (or more commonly, **undirected graph**). A **digraph** or **directed graph** is $G = (V, E)$ where $E \subseteq V^2 = V \times V$. For any digraph $G = (V, E)$, its **reverse** is the digraph (V, E') where $(u, v) \in E$ iff $(v, u) \in E'$. In this book, the word “graph” shall refer to a bigraph or digraph; the context should make the intent clear. The edges of graphs are often written ‘ (u, v) ’ or ‘ uv ’ where u, v are vertices. We will prefer⁹ to denote edge-hood by the notation $u-v$. Of course, in the case of bigraphs, $u-v = v-u$.

Often a graph $G = (V, E)$ comes with auxiliary data, say d_1, d_2 , etc. In this case we denote the graph by

$$G = (V, E; d_1, d_2, \dots)$$

using the semi-colon to mark the presence of auxiliary data. For example:

- (i) Often one or two vertices in V are distinguished. If $s, t \in V$ are distinguished, we might write $G = (V, E; s, t)$. This notation might be used in shortest path problems where s is the source and t is the target for the class of paths under consideration.
- (ii) A “weight” function $W : V \rightarrow \mathbb{R}$, and we denote the corresponding weighted graph by $G = (V, E; W)$.
- (iii) Another kind of auxiliary data is **vertex coloring** of G , i.e., a function $C : V \rightarrow S$ where S is any set. Then $C(v)$ is called the **color** of $v \in V$. If $|S| = k$, we call C a k -coloring. The **chromatic graph** is therefore given by the triple $G = (V, E; C)$. An **edge coloring** is similarly defined, $C : E \rightarrow S$.

We introduce terminology for some special graphs: If V is the empty set, A graph $G = (V, E)$ is called the **empty graph**. If E is the empty set, $G = (V, E)$ is called the **trivial graph**. Hence empty graphs are necessarily trivial but not vice-versa. $K_n = (V, \binom{V}{2})$ denotes the **complete graph** on $n = |V|$ vertices. A **bipartite graph** $G = (V, E)$ is a digraph such that $V = V_1 \uplus V_2$ and $E \subseteq V_1 \times V_2$. It is common to write $G = (V_1, V_2, E)$ in this case. Thus, $K_{m,n} = (V_1, V_2, V_1 \times V_2)$ denotes the **complete bipartite graph** where $m = |V_1|$ and $n = |V_2|$.

Two graphs $G = (V, E), G' = (V', E')$ are **isomorphic** if there is some bijection $\phi : V \rightarrow V'$ such that $\phi(E) = E'$ (the notation $\phi(E)$ has the obvious meaning).

If $G = (V, E), G' = (V', E')$ where $V' \subseteq V$ and $E' \subseteq E$ then we call G' a **subgraph** of G . In case E' is the restriction of E to the edges in V' , then we call G' the **V' -induced subgraph** of G or G' is the **restriction** of G to V' . We may write $G|V'$ for G' .

⁸This notation should not be confused with the absolute value of a number or the size of a set. The context will make this clear.

⁹When we write $u-v$, it is really an assertion that the (u, v) is an edge. So it is redundant to say “ $u-v$ is an edge”.

A **path** (from v_1 to v_k) is a sequence (v_1, v_2, \dots, v_k) of vertices such that (v_i, v_{i+1}) is an edge. Thus, we may also denote this path as $(v_1 - v_2 - \dots - v_k)$. A digraph path is a **cycle** if $v_1 = v_k$ and $k > 1$. In the case of bigraphs, the path is a **cycle** if $v_1 = v_k$, $k > 2$ and for all $1 < i < k$, $v_{i-1} \neq v_{i+1}$. A graph is **acyclic** if it has no cycles. Sometimes acyclic bigraphs are called **forests**, and acyclic digraph are called **dags** (“directed acyclic graph”).

Two vertices u, v are **connected** if there is a path from u to v , and a path from v to u . (Note that in the case of bigraphs, there is a path from u to v iff there is a path from v to u .) We shall say v is **adjacent to** u if $u-v$. Clearly, connectivity and adjacency are symmetric binary relation. It is easily seen that connectivity is also reflexive and transitive. This relation partitions the set of vertices into **connected components**.

In a digraph, **out-degree** and **in-degree** of a vertex is the number of edges issuing (respectively) from and into that vertex. The **out-degree** (resp., **in-degree**) of a digraph is the maximum of the out-degrees (resp., in-degrees) of its vertices. The vertices of out-degree 0 are called **sinks** and the vertices of in-degree 0 are called **sources**. The **degree** of a vertex in a bigraph is the number of adjacent vertices; the **degree** of a bigraph is the maximum of degrees of its vertices.

Trees. A connected acyclic bigraph is called a **free-tree**. A digraph such that there is a unique source vertex (called the **root**) and all the other vertices have in-degree 1, is called¹⁰ a **tree**. The sinks in a tree are called **leaves** or **external nodes** and non-leaves are called **internal nodes**. In general, we prefer a terminology in which the vertices of trees are called **nodes**. Thus there is a unique path from the root to each node in a tree. If u, v are nodes in T then u is a **descendent** of v if there is a path from v to u . Every node v is a descendent of itself, called the **improper descendent** of v . All other descendents of v are called **proper**. We may speak of the **child** or **grandchild** of any node in the obvious manner. The reverse of the descendent binary relation is the **ancestor** relation; thus we have **proper ancestors**, **parent** and **grandparent** of a node.

The **subtree** at any node u of T is the subgraph of T obtained by restricting to the descendents of u . The **depth** of a node u in a tree T is the length of the path from the root to u . So the root is the unique node of depth 0. The **depth of** T is the maximum depth of a node in T . The **height** of a node u is just the depth of the subtree at u ; alternatively, it is the length of the longest path from u to its descendents. Thus u has height 0 iff u is a leaf iff u has no children. The collection of all nodes at depth i is also called the **i th level** of the tree. Thus level zero is comprised of just the root. We normally draw a tree with the root at the top of the figure, and edges are implicitly direction from top to bottom.

Binary Trees. These trees are of major importance in computer science. A **binary tree** is a tree T of out-degree at most two and every edge is labeled either **left** or **right**, with the restriction that if a node has two outgoing edges, then these two edges have different labels. If there is a **left** edge from v to u then we call u a **left child** of v and u is denoted $v.\text{left}$. Similarly the **right child** of v is denoted $v.\text{right}$. Note that if a node has only one child, then that child could be either a left or a right child.

In conventional programming languages, the $u.\text{left}$ and $u.\text{right}$ are often implemented as pointers. These pointers are pre-allocated to the node u . Hence, even if u does not have a left child, the pointer $u.\text{left}$ is present but it is set to a special value called nil. In many applications, we further assume that each node has a pointer to its parent, and this is denoted $u.\text{parent}$.

Any node v in T can be viewed as the root of another binary tree called the **subtree at** v . The subtree

¹⁰One can also define trees in which the sense of the edges are reversed: the root is a sink and all the leaves are sources. We will often go back and forth between these two view points without much warning. E.g., we might speak of the “path from a node to the root”. While it is clear what is meant here, but to be technically correct, we ought to speak awkwardly of the path in the “reverse of the tree”.

at the left child of v is called the **left subtree of v** ; the left subtree of v is **empty** if v has no left child. Similarly for **right subtree of v** .

Like the admission of 0 to the counting numbers, it is mathematically pleasing to allow the the possibility of an **empty** binary tree. An empty binary tree has no nodes.

The following recursive definition of a binary tree is *the most convenient form to use* in proving facts about binary trees. A **binary tree** T with **root** $r(T)$ is defined to be a directed graph on a set V of nodes that satisfies:

BASE CASE: $V = \emptyset$ and $r(T) = \text{nil}$.

INDUCTION: $r(T) \in V$ and there are two binary trees T_L, T_R such that $v_0.\text{left} = r(T_L)$ and $v_0.\text{right} = r(T_R)$. Moreover, T_i is a binary tree on a set V_i ($i = L, R$) and $V = V_L \uplus \{v_0\} \uplus V_R$.

Viewed T as a graph, its edge set $E(T)$ is empty in the base case, and inductively given by

$$E(T) = \{(r(T), r(T_i)) : r(T_i) \neq \text{nil}, i = L, R\} \cup E(T_L) \cup E(T_R).$$

There are 5 paths associated with a node u in a tree T :

1. The unique path from the root of T to u . NOTE: it is often convenient to assume that each node has a pointer to its parent (in practice, most binary trees do have such pointers). Following these parent pointers, we can then talk about the unique path **from u to the root**.
2. The path $\pi = (u_0, u_1, \dots, u_\ell)$ where $u_0 = u$ and u_i is the left child of u_{i-1} ($i = 1, \dots, \ell$) and u_ℓ is a leaf. We call π the **left subpath** of u . We similarly define the **right subpath** of u .
3. The path $\pi = (u_0, u_1, \dots, u_\ell)$ where $u_0 = u$, u_1 is the left child of u_0 , and (u_1, \dots, u_ℓ) is the right subpath of u_1 . We call π the **left spine** of u . We similarly define the **right spine** of u .

In case u has no left child, its left subpath and left spines is the trivial path (u) of length 0. We often identify a binary tree T with its root u_0 . So when we speak of the left subtree, right spine, etc., of T , we mean the left subtree, right spine, etc., of u_0 .

A binary tree is **complete** if every leaf has the same depth. It is easy to see that if the depth is $d \geq 0$ then the tree has $2^{d+1} - 1$ nodes and 2^d leaves. A node in a binary tree is **full** if it has two children; a binary tree is said to be **full** if every internal node is full. A simple result about non-empty full binary trees is that it has exactly one fewer internal node than the number of leaves. Thus, if it $k \geq 1$ leaves iff it has $k - 1 \geq 0$ internal nodes. As corollary, a full binary tree has an odd number of nodes.

A basic fact about binary trees is that if it has n nodes, then its height is at least $\lceil \lg(1+n) \rceil - 1$ (see Exercise). Hence the height of a binary tree on n nodes is $\Omega(\lg n)$. A family of binary trees is **balanced** if for each n , there are trees in the family with n nodes and every such tree has height $\mathcal{O}(\log n)$. The implicit constant in the \mathcal{O} -notation depends on the particular family.

We caution the reader that there is an alternative view of binary trees which we call **extended binary trees**; this is described in Chapter III.

Binary Search Trees. The most important use of binary trees is for storing and searching items. A **binary search tree** is a binary tree that in addition stores a **key** at each node. The key at node u is usually denoted $u.\text{Key}$. Keys are elements from a totally ordered set so that we may compare them. The

keys in a binary search tree need not be distinct but they must satisfy the **binary search tree property**: all the keys in the left (resp., right) subtree of a node v are less (resp., greater) than or equal to the key at v .

The nodes of a binary tree can be listed in one of three standard ways: **pre-order**, **in-order** or **post-order**. These correspond to listing a node (respectively) (i) *before* any of its descendants, (ii) *after* all nodes in its left subtree but *before* any of the nodes in its right subtree, and (iii) *after* all of its descendants. In-order is also called **symmetric order**. Node listing is also called **tree traversal**.

EXERCISES

Exercise A.1: The following is a useful result about iterated floors and ceilings. It shows that $\lfloor \lfloor n/2 \rfloor / 2 \rfloor = \lfloor n/4 \rfloor$, for instance.

(a) Let n, b be positive integers. Let $N_0 := n$ and for $i \geq 0$, $N_{i+1} := \lceil N_i/b \rceil$. Show that $N_i = \lceil n/b^i \rceil$. Similarly for floors. HINT: use the fact that $N_{i+1} \leq N_i/b + (b-1)/b$.

(b) Let $u_0 = 1$ and $u_{i+1} = \lfloor 5u_i/2 \rfloor$ for $i \geq 0$. Show that for $i \geq 4$, $0.76(5/2)^i < u_i \leq 0.768(5/2)^i$. HINT: $r_i := u_i(2/5)^i$ is non-increasing; give a lower bound on r_i ($i \geq 4$) based on r_4 . \diamond

Exercise A.2: Let x, a, b be positive real numbers. Show that $\lfloor x/ab \rfloor \geq \lfloor \lfloor x/a \rfloor / b \rfloor$. Is it true that they are equal? \diamond

Exercise A.3: Suppose you want to prove that

$$f(n) \neq O(f(n/2))$$

where $f(n) = (\log n)^{\log n}$.

(a) Using de Morgan's law, show that this amounts to saying that for all $C > 0, n_0$ there exists n such that

$$(n \geq n_0) \wedge f(n) > Cf(n/2).$$

(b) Complete the proof by finding a suitable n for any given C, n_0 . \diamond

Exercise A.4: Basic facts about binary trees:

(a) A full binary tree on n leaves has $n-1$ internal nodes.

(b) Show that every binary tree on $n \geq 1$ nodes has height at least $\lceil \lg(1+n) \rceil - 1$. HINT: define $M(h)$ to be the maximum number of nodes in a binary tree of height h .

(c) Show that the bound in (b) is tight for each n .

(d) Show that a binary tree on $n \geq 1$ leaves has height $\lceil \lg n \rceil$. HINT: use a modified version of $M(h)$. \diamond

Exercise A.5: (Erdős-Rado) Show that in any 2-coloring of the edges of the complete graph K_n , there is a monochromatic spanning tree of K_n . HINT: use induction. \diamond

Exercise A.6: Let T be a binary tree on n nodes.

(a) What is the minimum possible number of leaves in T ?

(b) Show by strong induction on the structure of T that T has at most $\lfloor \frac{n+1}{2} \rfloor$ leaves. This is an exercise in case analysis, so proceed as follows: first let n be odd (say, $n = 2N+1$) and assume T has $k = 2K+1$ children in the left subtree. There are 3 other cases.

- (c) Give an alternative proof of part (b): show the result for n by a weaker induction on $n - 1$ and $n - 2$.
- (d) Show that the bound in part (b) is the best possible by describing a T with $\lfloor \frac{n+1}{2} \rfloor$ leaves. HINT: first show it when $n = 2^t - 1$. Alternatively, consider binary heaps. \diamond

Exercise A.7:

- (a) A binary tree with a key associated to each node is a binary search tree iff the in-order listing of these keys is in non-decreasing order.
- (b) Given *both* the post-order and in-order listing of the nodes of a binary tree, we can reconstruct the tree. \diamond

END EXERCISES

References

- [1] T. H. Corman, C. E. Leiserson, R. L. Rivest, and C. Stein. *Introduction to Algorithms*. The MIT Press and McGraw-Hill Book Company, Cambridge, Massachusetts and New York, second edition, 2001.
- [2] D. G. Kirkpatrick and R. Seidel. The ultimate planar convex hull algorithm? *SIAM J. Comput.*, 15:287–299, 1986.
- [3] D. E. Knuth. *The Art of Computer Programming: Sorting and Searching*, volume 3. Addison-Wesley, Boston, 1972.
- [4] D. E. Knuth. *The Art of Computer Programming: Fundamental Algorithms*, volume 1. Addison-Wesley, Boston, 2nd edition edition, 1975.
- [5] C. K. Yap. Introduction to the theory of complexity classes, 1987. Book Manuscript. Preliminary version (on ftp since 1990),
URL <ftp://cs.nyu.edu/pub/local/yap/complexity-bk>.