

Chapter 1

Problem Solving Patterns

It's not that I'm so smart, it's just that I
stay with problems longer.

— A. EINSTEIN

Developing problem solving skills is like learning to play a musical instrument—books and teachers can point you in the right direction, but only your hard work will take you there. Just as a musician, you need to know underlying concepts, but theory is no substitute for practice.

Great problem solvers have skills that cannot be rigorously formalized. Still, when faced with a challenging programming problem, it is helpful to have a small set of “patterns”—general reusable solutions to commonly occurring problems—that may be applicable.

We now introduce several patterns and illustrate them with examples. We have classified these patterns into three categories:

- data structure patterns,
- algorithm design patterns, and
- abstract analysis patterns.

These patterns are summarized in Table 1.1 on the next page, Table 1.2 on Page 12, and Table 1.3 on Page 21, respectively.

The notion of patterns is very general; in particular, many patterns arise in the context of software design—the builder pattern, composition, publish-subscribe, etc. These are more suitable to large-scale systems, and as such are outside the scope of EPI, which is focused on smaller programs that can be solved in an interview.

Data structure patterns

A data structure is a particular way of storing and organizing related data items so that they can be manipulated efficiently. Usually the correct selection of data structures is key to designing a good algorithm. Different data structures are suited to different applications; some are highly specialized. For example, heaps are par-

Table 1.1: Data structure patterns.

Data structure	Key points
Primitive types	Know how <code>int</code> , <code>char</code> , <code>double</code> , etc. are represented in memory and the primitive operations on them.
Arrays	Fast access for element at an index, slow lookups (unless sorted) and insertions. Be comfortable with notions of iteration, resizing, partitioning, merging, etc.
Lists	Understand trade-offs with respect to arrays. Be comfortable with iteration, insertion, and deletion within singly and doubly linked lists. Know how to implement a list with dynamic allocation, and with arrays.
Stacks and queues	Understand insertion and deletion. Know array and linked list implementations.
Binary trees	Use for representing hierarchical data. Know about depth, height, leaves, search path, traversal sequences, successor/predecessor operations.
Heaps	Key benefit: $O(1)$ lookup find-min, $O(\log n)$ insertion, and $O(\log n)$ deletion of min. Node and array representations. Max-heap variant.
Hash tables	Key benefit: $O(1)$ insertions, deletions and lookups. Key disadvantages: not suitable for order-related queries; need for resizing; poor worst-case performance. Understand implementation using array of buckets and collision chains. Know hash functions for integers, strings, objects. Understand importance of equals function. Variants such as Bloom filters.
Binary search trees	Key benefit: $O(\log n)$ insertions, deletions, lookups, find-min, find-max, successor, predecessor when tree is balanced. Understand implementation using nodes and pointers. Be familiar with notion of balance, and operations maintaining balance. Know how to augment a binary search tree, e.g., interval trees and dynamic order statistics.

ticularly well-suited for algorithms that merge sorted data streams, while compiler implementations usually use hash tables to look up identifiers.

Solutions often require a combination of data structures. Our solution to the problem of tracking the most visited pages on a website (Solution 9.12 on Page 78) involves a combination of a heap, a queue, a binary search tree, and a hash table.

PRIMITIVE TYPES

You should be comfortable with the basic types (chars, integers, doubles, etc.), their variants (unsigned, long, etc.), and operations on them (bitwise operators, comparison, etc.). Don't forget that the basic types differ among programming languages. For example, Java has no unsigned integers, and the number of bits in an integer is compiler- and machine-dependent in C.

A common problem related to basic types is computing the number of bits set to 1 in an integer-valued variable x . To solve this problem you need to know how to manipulate individual bits in an integer. One straightforward approach is to iteratively test individual bits using an unsigned integer variable m initialized to 1. Iteratively identify bits of x that are set to 1 by examining the bitwise AND of m with x , shifting m left one bit at a time. The overall complexity is $O(n)$ where n is the length of the integer.

Another approach, which may run faster on some inputs, is based on computing $y = x \& !(x - 1)$, where $\&$ is the bitwise AND operator. This is 1 at exactly the rightmost bit of x . Consequently, this bit may be removed from x by computing $x \oplus y$. The time complexity is $O(s)$, where s is the number of bits set to 1 in x .

In practice if the computation is done repeatedly, the most efficient approach would be to create a lookup table. In this case, we could use a 256 entry integer-valued array P such that $P[i]$ is the number of bits set to 1 in i . If x is 32 bits, the result can be computed by decomposing x into 4 disjoint bytes, $b3, b2, b1$, and $b0$. The bytes are computed using bitmasks and shifting, e.g., $b1$ is $(x \& 0xff00) \gg 8$. The final result is $P[b3] + P[b2] + P[b1] + P[b0]$. Computing the parity of an integer is closely related to counting the number of bits set to 1, and we present a detailed analysis of the parity problem in Solution 2.1 on Page 167.

ARRAYS

Conceptually, an array maps integers in the range $[0, n - 1]$ to objects of a given type, where n is the number of objects in this array. Array lookup and insertion are fast, making arrays suitable for a variety of applications. Reading past the last element of an array is a common error, invariably with catastrophic consequences.

The following problem arises when optimizing quicksort: given an array A whose elements are comparable, and an index i , reorder the elements of A so that the initial elements are all less than $A[i]$, and are followed by elements equal to $A[i]$, which in turn are followed by elements greater than $A[i]$, using $O(1)$ space.

The key to the solution is to maintain two regions on opposite sides of the array that meet the requirements, and expand these regions one element at a time. Details are given in Solution 3.1 on Page 176.

LISTS

An abstract data type (ADT) is a mathematical model for a class of data structures that have similar functionality. Strictly speaking, a list is an ADT, and not a data structure. It implements an ordered collection of values, which may be repeated. In the context of this book we view a list as a sequence of nodes where each node has a link to the next node in the sequence. In a doubly linked list each node additionally has a link to the prior node.

A list is similar to an array in that it contains objects in a linear order. The key differences are that inserting and deleting elements in a list has time complexity $O(1)$. On the other hand, obtaining the k -th element in a list is expensive, having $O(n)$

time complexity. Lists are usually building blocks of more complex data structures. However, they can be the subject of tricky problems in their own right, as illustrated by the following:

Given a singly linked list $\langle l_0, l_1, l_2, \dots, l_{n-1} \rangle$, define the “zip” of the list to be $\langle l_0, l_{n-1}, l_1, l_{n-2}, \dots \rangle$. Suppose you were asked to write a function that computes the zip of a list, with the constraint that it uses $O(1)$ space. The operation of this function is illustrated in Figure 1.1.

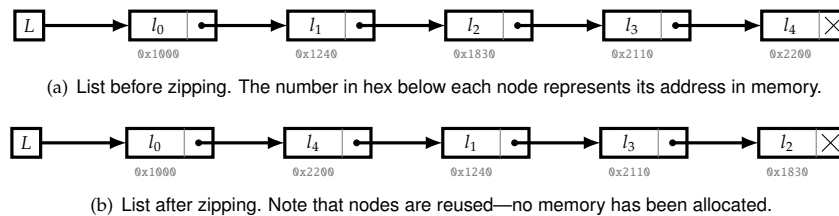


Figure 1.1: Zipping a list.

The solution is based on an appropriate iteration combined with “pointer swapping”, i.e., updating next field for each node. Refer to Solution 4.11 on Page 208 for details.

STACKS AND QUEUES

Stacks support last-in, first-out semantics for inserts and deletes, whereas queues are first-in, first-out. Both are ADTs, and are commonly implemented using linked lists or arrays. Similar to lists, stacks and queues are usually building blocks in a solution to a complex problem, but can make for interesting problems in their own right.

As an example consider the problem of evaluating Reverse Polish notation expressions, i.e., expressions of the form “3,4,×,1,2,+,”, “1,1,+,-2,×”, or “4,6,/,2,/”. A stack is ideal for this purpose—operands are pushed on the stack, and popped as operators are processed, with intermediate results being pushed back onto the stack. Details are given in Solution 5.2 on Page 211.

BINARY TREES

A binary tree is a data structure that is used to represent hierarchical relationships. Binary trees most commonly occur in the context of binary search trees, wherein keys are stored in a sorted fashion. However, there are many other applications of binary trees. Consider a set of resources organized as nodes in a binary tree. Processes need to lock resource nodes. A node may be locked if and only if none of its descendants and ancestors are locked. Your task is to design and implement an application programming interface (API) for locking.

A reasonable API is one with `isLock()`, `lock()`, and `unlock()` methods. Naïvely implemented the time complexity for these methods is $O(n)$, where n is the number of nodes. However these can be made to run in time $O(1)$, $O(h)$, and $O(h)$, respectively,

where h is the height of the tree, if nodes have a parent field. Details are given in Solution 6.4 on Page 226.

HEAPS

A heap is a data structure based on a binary tree. It efficiently implements an ADT called a priority queue. A priority queue resembles a queue, with one difference: each element has a “priority” associated with it, and deletion removes the element with the highest priority.

Suppose you are given a set of files, each containing stock trade information. Each trade appears as a separate line containing information about that trade. Lines begin with an integer-valued timestamp, and lines within a file are sorted in increasing order of timestamp. Suppose you were asked to design an algorithm that combines the set of files into a single file R in which trades are sorted by timestamp.

This problem can be solved by a multistage merge process, but there is a trivial solution based on a min-heap data structure. Entries are trade-file pairs and are ordered by the timestamp of the trade. Initially the min-heap contains the first trade from each file. Iteratively delete the minimum entry $e = (t, f)$ from the min-heap, write t to R , and add in the next entry in the file f . Details are given in Solution 7.1 on Page 237.

HASH TABLES

A hash table is a data structure used to store keys, optionally with corresponding values. Inserts, deletes and lookups run in $O(1)$ time on average. One caveat is that these operations require a good hash function—a mapping from the set of all possible keys to the integers which is similar to a uniform random assignment. Another is that if the number of keys that is to be stored is not known in advance then the hash table needs to be periodically resized, which depending on how the resizing is implemented, can lead to some updates having $\Theta(n)$ complexity.

Suppose you were asked to write an application that compares n programs for plagiarism. Specifically, your application is to break every program into overlapping character strings, each of length 100, and report on the number of strings that appear in each pair of programs. A hash table can be used to perform this check very efficiently if the right hash function is used. Details are given in Solution 9.14 on Page 277.

BINARY SEARCH TREES

Binary search trees (BSTs) are used to store objects that are comparable. The underlying idea is to organize the objects in a binary tree in which the nodes satisfy the BST property: the key stored at any node is greater than or equal to the keys stored in its left subtree and less than or equal to the keys stored in its right subtree. Insertion and deletion can be implemented so that the height of the BST is $O(\log n)$, leading to fast ($O(\log n)$) lookup and update times. AVL trees and red-black trees are BST implementations that support this form of insertion and deletion.

BSTs are a workhorse of data structures and can be used to solve almost every data structures problem reasonably efficiently. It is common to augment the BST to make it possible to manipulate more complicated data, e.g., intervals, and efficiently support more complex queries, e.g., the number of elements in a range.

As an example application of BSTs, consider the following problem. You are given a set of line segments. Each segment is a closed interval $[l_i, r_i]$ of the x -axis, a color, and a height. For simplicity assume no two segments whose intervals overlap have the same height. When the x -axis is viewed from above the color at point x on the x -axis is the color of the highest segment that includes x . (If no segment contains x , the color is blank.) You are to implement a function that computes the sequence of colors as seen from the top.

The key idea is to sort the endpoints of the line segments and do a sweep from left-to-right. As we do the sweep, we maintain a list of line segments that intersect the current position as well as the highest line and its color. To quickly lookup the highest line in a set of intersecting lines we keep the current set in a BST, with the interval's height as its key. Details are given in Solution 11.15 on Page 310.

Other data structures

The data structures described above are the ones commonly used. Examples of other data structures that have more specialized applications include:

- Skip lists, which store a set of comparable items using a hierarchy of sorted linked lists. Lists higher in the hierarchy consist of increasingly smaller subsequences of the items. Skip lists implement the same functionality as balanced BSTs, but are simpler to code and faster, especially when used in a concurrent context.
- Treaps, which are a combination of a BST and a heap. When an element is inserted into a treap it is assigned a random key that is used in the heap organization. The advantage of a treap is that it is height balanced with high probability and the insert and delete operations are considerably simpler than for deterministic height balanced trees such as AVL and red-black trees.
- Fibonacci heaps, which consist of a series of trees. Insert, find minimum, decrease key, and merge (union) run in amortized constant time; delete and delete-minimum take $O(\log n)$ time. In particular Fibonacci heaps can be used to reduce the time complexity of Dijkstra's shortest path algorithm from $O((|E| + |V|) \log |V|)$ to $O(|E| + |V| \log |V|)$.
- Disjoint-set data structures, which are used to manipulate subsets. The basic operations are union (form the union of two subsets), and find (determine which set an element belongs to). These are used in a number of algorithms, notably in tracking connected components in an undirected graph and Kruskal's algorithm for the minimum spanning tree (MST). We use the disjoint-set data structure to solve the offline minimum problem (Solution 3.7 on Page 182).
- Tries, which are a tree-based data structure used to store strings. Unlike BSTs, nodes do not store keys; instead, the node's position in the tree determines the key it is associated with. Tries can have performance advantages with respect

Table 1.2: Algorithm design patterns.

Technique	Key points
Sorting	Uncover some structure by sorting the input.
Recursion	If the structure of the input is defined in a recursive manner, design a recursive algorithm that follows the input definition.
Divide and conquer	Divide the problem into two or more smaller independent subproblems and solve the original problem using solutions to the subproblems.
Dynamic programming	Compute solutions for smaller instances of a given problem and use these solutions to construct a solution to the problem.
The greedy method	Compute a solution in stages, making choices that are locally optimum at step; these choices are never undone.
Incremental improvement	Quickly build a feasible solution and improve its quality with small, local updates.
Elimination	Identify and rule out potential solutions that are sub-optimal or dominated by other solutions.
Parallelism	Decompose the problem into subproblems that can be solved independently on different machines.
Caching	Store computation and later look it up to save work.
Randomization	Use randomization within the algorithm to reduce complexity.
Approximation	Efficiently compute a suboptimum solution that is of acceptable quality.
State	Identify an appropriate notion of state.

to BSTs and hash tables; they can also be used to solve the longest matching prefix problem (Solution 16.3 on Page 393).

Algorithm design patterns

SORTING

Certain problems become easier to understand, as well as solve, when the input is sorted. The solution to the calendar rendering problem (Problem 10.10 on Page 84) entails taking a set of intervals and computing the maximum number of intervals whose intersection is nonempty. Naïve strategies yield quadratic run times. However, once the interval endpoints have been sorted, it is easy to see that a point of maximum overlap can be determined by a linear time iteration through the endpoints.

Often it is not obvious what to sort on—for example, we could have sorted the intervals on starting points rather than endpoints. This sort sequence, which in some

respects is more natural, does not work. However, some experimentation with it will likely lead to the correct criterion.

Sorting is not appropriate when an $O(n)$ (or better) algorithm is possible, e.g., determining the k -th largest element (Problem 8.13 on Page 71). Furthermore, sorting can obfuscate the problem. For example, given an array A of numbers, if we are to determine the maximum of $A[i] - A[j]$, for $i < j$, sorting destroys the order and complicates the problem.

RECURSION

A recursive function consists of base cases, and calls to the same function with different arguments. A recursive algorithm is appropriate when the input is naturally expressed using recursive functions.

String matching exemplifies the use of recursion. Suppose you were asked to write a Boolean-valued function which takes a string and a matching expression, and returns true iff the string “matches” the matching expression. Specifically, the matching expression is itself a string, and could be

- x where x is a character, for simplicity assumed to be a lower-case letter (matches the string “ x ”).
- $.$ (matches any string of length 1).
- x^* (matches the string consisting of zero or more occurrences of the character x).
- $.*$ (matches the string consisting of zero or more of any characters).
- r_1r_2 where r_1 and r_2 are regular expressions of the given form (matches any string that is the concatenation of strings s_1 and s_2 , where s_1 matches r_1 and s_2 matches r_2).

This problem can be solved by checking a number of cases based on the first one or two characters of the matching expression, and recursively matching the rest of the string. Details are given in Solution 3.22 on Page 197.

DIVIDE AND CONQUER

A divide and conquer algorithm works by decomposing a problem into two or more smaller independent subproblems, until it gets to instances that are simple enough to be solved directly; the results from the subproblems are then combined. More details and examples are given in Chapter 12; we illustrate the basic idea below.

A triomino is formed by joining three unit-sized squares in an L-shape. A mutilated chessboard (henceforth 8×8 Mboard) is made up of 64 unit-sized squares arranged in an 8×8 square, minus the top-left square, as depicted in Figure 1.2(a) on the next page. Suppose you are asked to design an algorithm that computes a placement of 21 triominoes that covers the 8×8 Mboard. Since the 8×8 Mboard contains 63 squares, and we have 21 triominoes, a valid placement cannot have overlapping triominoes or triominoes which extend out of the 8×8 Mboard.

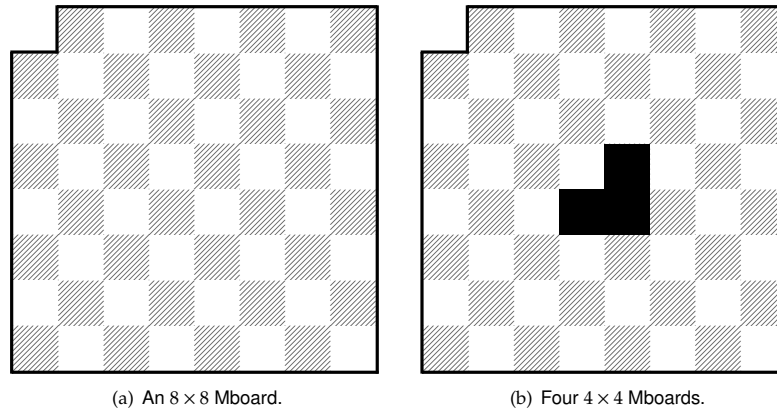


Figure 1.2: Mutilated chessboards.

Divide and conquer is a good strategy for this problem. Instead of the 8×8 Mboard, let's consider an $n \times n$ Mboard. A 2×2 Mboard can be covered with one triomino since it is of the same exact shape. You may hypothesize that a triomino placement for an $n \times n$ Mboard with the top-left square missing can be used to compute a placement for an $(n + 1) \times (n + 1)$ Mboard. However you will quickly see that this line of reasoning does not lead you anywhere.

Another hypothesis is that if a placement exists for an $n \times n$ Mboard, then one also exists for a $2n \times 2n$ Mboard. Now we can apply divide and conquer. work. Take four $n \times n$ Mboards and arrange them to form a $2n \times 2n$ square in such a way that three of the Mboards have their missing square set towards the center and one Mboard has its missing square outward to coincide with the missing corner of a $2n \times 2n$ Mboard, as shown in Figure 1.2(b). The gap in the center can be covered with a triomino and, by hypothesis, we can cover the four $n \times n$ Mboards with triominoes as well. Hence a placement exists for any n that is a power of 2. In particular, a placement exists for the $2^3 \times 2^3$ Mboard; the recursion used in the proof directly yields the placement.

Divide and conquer is usually implemented using recursion. However, the two concepts are not synonymous. Recursion is more general—subproblems do not have to be of the same form.

In addition to divide and conquer, we used the generalization principle above. The idea behind generalization is to find a problem that subsumes the given problem and is easier to solve. We used it to go from the 8×8 Mboard to the $2^n \times 2^n$ Mboard.

Other examples of divide and conquer include counting the number of pairs of elements in an array that are out of sorted order (Solution 12.2 on Page 315) and computing the closest pair of points in a set of points in the plane (Solution 12.3 on Page 316).

DYNAMIC PROGRAMMING

Dynamic programming (DP) is applicable when the problem has the “optimal sub-

structure” property, that is, it is possible to reconstruct a solution to the given instance from solutions to subinstances of smaller problems of the same kind. A key aspect of DP is maintaining a cache of solutions to subinstances. DP can be implemented recursively (in which case the cache is typically a dynamic data structure such as a hash table or a BST), or iteratively (in which case the cache is usually a one- or multi-dimensional array). It is most natural to design a DP algorithm using recursion. Usually, but not always, it is more efficient to implement it using iteration.

As an example of the power of DP, consider the problem of determining the number of combinations of 2, 3, and 7 point plays that can generate a score of 222. Let $C(s)$ be the number of combinations that can generate a score of s . Then $C(222) = C(222 - 7) + C(222 - 3) + C(222 - 2)$, since a combination ending with a 2 point play is different from one ending with a 3 point play, and a combination ending with a 3 point play is different from one ending with a 7 point play, etc.

The recursion breaks down for small scores, specifically, when (1.) $s < 0 \Rightarrow C(s) = 0$, and (2.) $s = 0 \Rightarrow C(s) = 1$.

Implementing the recursion naïvely results in multiple calls to the same subinstance. Let $C(a) \rightarrow C(b)$ indicate that a call to C with input a directly calls C with input b . Then $C(213)$ will be called in the order $C(222) \rightarrow C(222 - 7) \rightarrow C((222 - 7) - 2)$, as well as $C(222) \rightarrow C(222 - 3) \rightarrow C((222 - 3) - 3) \rightarrow C(((222 - 3) - 3) - 3)$.

This phenomenon results in the run time increasing exponentially with the size of the input. The solution is to store previously computed values of C in an array of length 223. Details are given in Solution 12.15 on Page 333.

Sometimes it is profitable to study the set of partial solutions. Specifically it may be possible to “prune” dominated solutions, i.e., solutions which cannot be better than previously explored solutions. The candidate solutions are referred to as the “efficient frontier” that is propagated through the computation.

For example, if we are to implement a stack that supports a max operation, which returns the largest element stored in the stack, we can record for each element in the stack what the largest value stored at or below that element is by comparing the value of that element with the value of the largest element stored below it. Details are given in Solution 5.1 on Page 210. The largest rectangle under the skyline (Problem 12.8 on Page 323) provides a more sophisticated example of the efficient frontier concept.

Another consideration is how the partial solutions are organized. In the solution to the longest nondecreasing subsequence problem 12.6 on Page 320, it is better to keep the efficient frontier sorted by length of each subsequence rather than its final index.

THE GREEDY METHOD

A greedy algorithm is one which makes decisions that are locally optimum and never changes them. This strategy does not always yield the optimum solution. Furthermore, there may be multiple greedy algorithms for a given problem, and only some of them are optimum.

For example, consider $2n$ cities on a line, half of which are white, and the other half are black. We want to map white to black cities in a one-to-one onto fashion so that

the total length of the road sections required to connect paired cities is minimized. Multiple pairs of cities may share a single section of road, e.g., if we have the pairing $(0, 4)$ and $(1, 2)$ then the section of road between Cities 0 and 4 can be used by Cities 1 and 2. The most straightforward greedy algorithm for this problem is to scan through the white cities, and, for each white city, pair it with the closest unpaired black city. It leads to suboptimum results: consider the 4 city case, where the white cities are at 0 and at 3 and the black cities are at 2 and at 5. If the straightforward greedy algorithm processes the white city at 3 first, it pairs it with 2, forcing the cities at 0 and 5 to pair up, leading to a road length of 5, whereas the pairing of cities at 0 and 2, and 3 and 5 leads to a road length of 4.

However, a slightly more sophisticated greedy algorithm does lead to optimum results: iterate through the white cities in left-to-right order, pairing each white city with the nearest unpaired black city. We can prove this leads to an optimum pairing by induction. The idea is that the pairing for the first city must be optimum, since if it were to be paired with any other city, we could always change its pairing to be with the nearest black city without adding any road.

INCREMENTAL IMPROVEMENT

When you are faced with the problem of computing an optimum solution, it is often straightforward to come up with a candidate solution, which may be a partial solution. This solution can be incrementally improved to make it optimum. This is especially true when a solution has to satisfy a set of constraints.

As an example consider a department with n graduate students and n professors. Each student begins with a rank ordered preference list of the professors based on how keen he is to work with each of them. Each professor has a similar preference list of students. Suppose you were asked to devise an algorithm which takes as input the preference lists and outputs a one-to-one pairing of students and advisers in which there are no student-adviser pairs (s_0, a_0) and (s_1, a_1) such that s_0 prefers a_1 to a_0 and a_1 prefers s_0 to s_1 .

Here is an algorithm for this problem in the spirit of incremental improvement. Each student who does not have an adviser “proposes” to the most-preferred professor to whom he has not yet proposed. Each professor then considers all the students who have proposed to him and says to the student in this set he most prefers “I accept you”; he says “no” to the rest. The professor is then provisionally matched to a student; this is the candidate solution. In each subsequent round, each student who does not have an adviser proposes to the professor to whom he has not yet proposed who is highest on his preference list. He does this regardless of whether the professor has already been matched with a student. The professor once again replies with a single accept, rejecting the rest. In particular, he may leave a student with whom he is currently paired. That this algorithm is correct is nontrivial—details are presented in Solution 18.18 on Page 433.

Many other algorithms are in this spirit: the standard algorithms for bipartite matching (Solution 18.19 on Page 435), maximum flow (Solution 18.21 on Page 436), and computing all pairs of shortest paths in a graph (Solutions 13.12 on Page 365

and 13.11 on Page 364) use incremental improvement. Other famous examples include the simplex algorithm for linear programming, and Euler's algorithm for computing a path in a graph which covers each edge once.

Sometimes it is easier to start with an infeasible solution that has a lower cost than the optimum solution, and incrementally update it to arrive at a feasible solution that is optimum. The standard algorithms for computing an MST (Solution 14.6 on Page 371) and shortest paths in a graph from a designated vertex (Solution 13.9 on Page 362) proceed in this fashion.

It is noteworthy that naïvely applying incremental improvement does not always work. For the professor-student pairing example above, if we begin with an arbitrary pairing of professors and students, and search for pairs p and s such that p prefers s to his current student, and s prefers p to his current professor and reassign such pairs, the procedure will not always converge.

Incremental improvement is often useful when designing heuristics, i.e., algorithms which are usually faster and/or simpler to implement than algorithms which compute an optimum result, but may return a suboptimal result. The algorithm we present for computing a tour for a traveling salesman (Solution 14.6 on Page 371) is in this spirit.

ELIMINATION

One common approach to designing an efficient algorithm is to use elimination—that is to identify and rule out potential solutions that are suboptimal or dominated by other solutions. Binary search, which is the subject of a number of problems in Chapter 8, uses elimination. Solution 8.9 on Page 257, where we use elimination to compute the square root of a real number, is especially instructive. Below we consider a fairly sophisticated application of elimination.

Suppose you have to build a distributed storage system. A large number, n , of users will share data on your system, which consists of m servers, numbered from 0 to $m-1$. One way to distribute users across servers is to assign the user with login id l to the server $h(l) \bmod m$, where $h()$ is a hash function. If the hash function does a good job, this approach distributes users uniformly across servers. However, if certain users require much more storage than others, some servers may be overloaded while others idle.

Let b_i be the number of bytes of storage required by user i . We will use values $k_0 < k_1 < \dots < k_{m-2}$ to partition users across the m servers—a user with hash code c gets assigned to the server with the lowest id i such that $c \leq k_i$, or to server $m-1$ if no such i exists. We would like to select k_0, k_1, \dots, k_{m-2} to minimize the maximum number of bytes stored at any server.

The optimum values for k_0, k_1, \dots, k_{m-2} can be computed via DP—the essence of the program is to add one server at a time. The straightforward formulation has an $O(nm^2)$ time complexity.

However, there is a much faster approach based on elimination. The search for values k_0, k_1, \dots, k_{m-2} such that no server stores more than b bytes can be performed

in $O(n)$ time by greedily selecting values for the k_i s. We can then perform binary search on b to get the minimum b and the corresponding values for k_0, k_1, \dots, k_{m-2} . The resulting time complexity is $O(n \log W)$, where $W = \sum_{i=0}^{m-1} b_i$.

For the case of 10000 users and 100 servers, the DP algorithm took over an hour; the approach using binary search for b with greedy assignment took 0.1 seconds. Details are given in Solution 12.24 on Page 343.

The takeaway is that there may be qualitatively different ways to search for a solution, and that it is important to look for ways in which to eliminate candidates. The efficient frontier concept, described on Page 15, has some commonalities with elimination.

PARALLELISM

In the context of interview questions, parallelism is useful when dealing with scale, i.e., when the problem is too large to fit on a single machine or would take an unacceptably long time on a single machine. The key insight you need to display is that you know how to decompose the problem so that

1. each subproblem can be solved relatively independently, and
2. the solution to the original problem can be efficiently constructed from solutions to the subproblems.

Efficiency is typically measured in terms of central processing unit (CPU) time, random access memory (RAM), network bandwidth, number of memory and database accesses, etc.

Consider the problem of sorting a petascale integer array. If we know the distribution of the numbers, the best approach would be to define equal-sized ranges of integers and send one range to one machine for sorting. The sorted numbers would just need to be concatenated in the correct order. If the distribution is not known then we can send equal-sized arbitrary subsets to each machine and then merge the sorted results, e.g., using a min-heap. Details are given in Solution 10.4 on Page 283.

CACHING

Caching is a great tool whenever computations are repeated. For example, the central idea behind dynamic programming is caching results from intermediate computations. Caching is also extremely useful when implementing a service that is expected to respond to many requests over time, and many requests are repeated. Workloads on web services exhibit this property. Solution 15.1 on Page 381 sketches the design of an online spell correction service; one of the key issues is performing cache updates in the presence of concurrent requests.

RANDOMIZATION

Suppose you were asked to write a routine that takes an array A of n elements and an integer k between 1 and n , and returns the k -th largest element in A .

This problem can be solved by first sorting the array, and returning the element at index k in the sorted array. The time complexity of this approach is $O(n \log n)$. However, sorting performs far more work than is needed. A better approach is to eliminate parts of the array. We could use the median to determine the $n/2$ largest elements of A ; if $n/2 \geq k$, the desired element is in this set, otherwise we search for the $(k - n/2)$ -th largest element in the $n/2$ smallest elements.

It is possible, though nontrivial, to compute the median in $O(n)$ time without using randomization. However, an approach that works well is to select an index r at random and reorder the array so that elements greater than or equal to $A[r]$ appear first, followed by $A[r]$, followed by elements less than or equal to $A[r]$. Let $A[r]$ be the k -th element in the reordered array A' . If $k = n/2$, $A'[k] = A[r]$ is the desired element. If $k > n/2$, we search for the $n/2$ -th largest element in $A'[0 : k - 1]$. Otherwise we search for the $(n/2 - k)$ -th largest element in $A'[k + 1 : n - 1]$. The closer $A[r]$ is the true median, the faster the algorithm runs. A formal analysis shows that the probability of the randomly selected element repeatedly being far from the desired element falls off exponentially with n . Details are given in Solution 8.13 on Page 260.

Randomization can also be used to create “signatures” to reduce the complexity of search, analogous to the use of hash functions. Consider the problem of determining whether an $m \times m$ array S of integers is a subarray of an $n \times n$ array T . Formally, we say S is a subarray of T iff there are p, q such that $S[i][j] = T[p + i][q + j]$, for all $0 \leq i, j \leq m - 1$. The brute-force approach to checking if S is a subarray of T has complexity $O(n^2 m^2)$ — $O(n^2)$ individual checks, each of complexity $O(m^2)$. We can improve the complexity to $O(n^2 m)$ by computing a hash code for S and then computing the hash codes for $m \times m$ subarrays of T . The latter hash codes can be computed incrementally in $O(m)$ time if the hash function is chosen appropriately. For example, if the hash code is simply the XOR of all the elements of the subarray, the hash code for a subarray shifted over by one column can be computed by XORing the new elements and the removed elements with the previous hash code. A similar approach works for more complex hash functions, specifically for those that are a polynomial.

APPROXIMATION

In the real-world it is routine to be given a problem that is difficult to solve exactly, either because of its intrinsic complexity, or the complexity of the code required. Developers need to recognize such problems, and be ready to discuss alternatives with the author of the problem. In practice a solution that is “close” to the optimum solution is usually perfectly acceptable.

Let $\{A_0, A_1, \dots, A_{n-1}\}$ be a set of n cities, as in Figure 1.3 on the next page. Suppose we need to choose a subset of A to locate warehouses. Specifically, we want to choose k cities in such a way that cities are close to the warehouses. Define the cost of a warehouse assignment to be the maximum distance of any city to a warehouse.

The problem of finding a warehouse assignment that has the minimum cost is known to be NP-complete. However, consider the following algorithm for com-

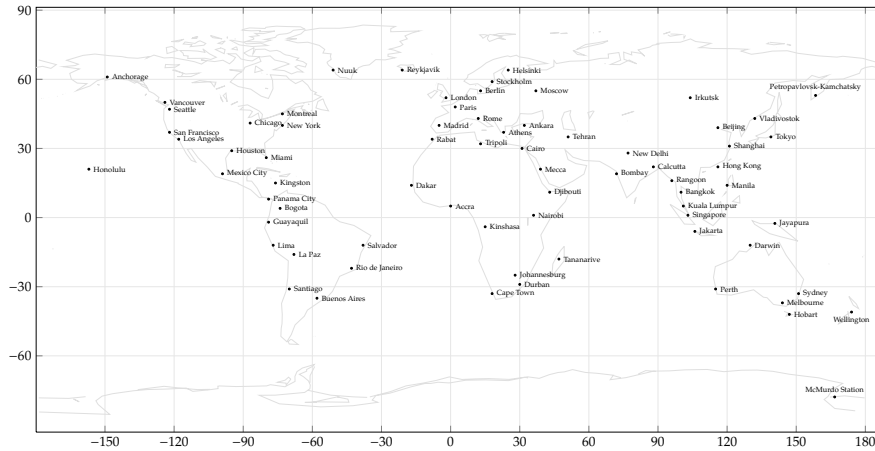


Figure 1.3: An instance of the warehouse location problem. The distance between cities at (p, q) and (r, s) is $\sqrt{(p - r)^2 + (q - s)^2}$.

putting k cities. We pick the first warehouse to be the city for which the cost is minimized—this takes $\Theta(n^2)$ time since we try each city one at a time and check its distance to every other city. Now let's say we have selected the first $i - 1$ warehouses $\{c_1, c_2, \dots, c_{i-1}\}$ and are trying to choose the i -th warehouse. A reasonable choice for c_i is the city that is the farthest from the $i - 1$ warehouses already chosen. This city can be computed in $O(ni)$ time. This greedy algorithm yields a solution whose cost is no more than $2 \times$ that of the optimum solution; some heuristic tweaks can be used to further improve the quality. Details are given in Solution 14.7 on Page 372.

As another example of approximation, consider the problem of determining the k most frequent elements of a very large array. The direct approach of maintaining counts for each element may not be feasible because of space constraints. A natural approach is to *sample* the set to determine a set of candidates, exact counts for which are then determined in a second pass. The size of the candidate set depends on the distribution of the elements.

STATE

Formally, the state of a system is information that is sufficient to determine how that system evolves as a function of future inputs. Identifying the right notion of state can be critical to coming up with an algorithm that is time and space efficient, as well as easy to implement and prove correct.

There may be multiple ways in which state can be defined, all of which lead to correct algorithms. When computing the max-difference (Problem 3.3 on Page 38), we could use the values of the elements at all prior indices as the state when we iterate through the array. Of course, this is inefficient, since all we really need is the minimum value.

Table 1.3: Abstract analysis techniques.

Analysis principle	Key points
Case analysis	Split the input/execution into a number of cases and solve each case in isolation.
Small examples	Find a solution to small concrete instances of the problem and then build a solution that can be generalized to arbitrary instances.
Iterative refinement	Most problems can be solved using a brute-force approach. Find such a solution and improve upon it.
Reduction	Use a well known solution to some other problem as a subroutine.
Graph modeling	Describe the problem using a graph and solve it using an existing algorithm.
Write an equation	Express relationships in the problem in the form of equations (or inequalities).
Variation	Solve a slightly different (possibly more general) problem and map its solution to the given problem.
Invariants	Find a function of the state of the given system that remains constant in the presence of (possibly restricted) updates to the state. Use this function to design an algorithm, prove correctness, or show an impossibility result.

One solution to computing the Levenshtein distance between two strings (Problem 12.11 on Page 101) entails creating a 2D array whose dimensions are $(m+1) \times (n+1)$, where m and n are the lengths of the strings being compared. For large strings this size may be unacceptably large. The algorithm iteratively fills rows of the array, and reads values from the current row and the previous row. This observation can be used to reduce the memory needed to two rows. A more careful implementation can reduce the memory required to just one row.

More generally, the efficient frontier concept on Page 15 demonstrates how an algorithm can be made to run faster and with less memory if state is chosen carefully. Other examples illustrating the benefits of careful state selection include string matching (Problem 3.19 on Page 43) and lazy initialization (Problem 3.2 on Page 37).

Abstract analysis patterns

CASE ANALYSIS

In case analysis a problem is divided into a number of separate cases, and analyzing each such case individually suffices to solve the initial problem. Cases do not have to be mutually exclusive; however, they must be exhaustive, that is cover all possibilities. For example, to prove that for all n , $n^3 \bmod 3$ is 0, 1, or 8, we can consider the cases $n = 3m$, $n = 3m + 1$, and $n = 3m + 2$. These cases are individually

easy to prove, and are exhaustive. Case analysis is commonly used in mathematics and games of strategy. Here we consider an application of case analysis to algorithm design.

Suppose you are given a set S of 25 distinct integers and a CPU that has a special instruction, SORT5, that can sort five integers in one cycle. Your task is to identify the largest, second-largest, and third-largest integers in S using SORT5 to compare and sort subsets of S ; furthermore, you must minimize the number of calls to SORT5.

If all we had to compute was the largest integer in the set, the optimum approach would be to form five disjoint subsets S_1, \dots, S_5 of S , sort each subset, and then sort $\{\max S_1, \dots, \max S_5\}$. This takes six calls to SORT5 but leaves ambiguity about the second and third largest integers.

It may seem like many additional calls to SORT5 are still needed. However if you do a careful case analysis and eliminate all $x \in S$ for which there are at least three integers in S larger than x , only five integers remain and hence just one more call to SORT5 is needed to compute the result. Details are given in Solution 18.2 on Page 423.

SMALL EXAMPLES

Problems that seem difficult to solve in the abstract can become much more tractable when you examine small concrete instances. For instance, consider the following problem. Five hundred closed doors along a corridor are numbered from 1 to 500. A person walks through the corridor and opens each door. Another person walks through the corridor and closes every alternate door. Continuing in this manner, the i -th person comes and toggles the state (open or closed) of every i -th door starting from Door i . You must determine exactly how many doors are open after the 500-th person has walked through the corridor.

It is difficult to solve this problem using an abstract approach, e.g., introducing Boolean variables for the state of each door and a state update function. However if you try the same problem with 1, 2, 3, 4, 10, and 20 doors, it takes a short time to see that the doors that remain open are 1, 4, 9, 16 . . . , regardless of the total number of doors. The 10 doors case is illustrated in Figure 1.4 on the next page. Now the pattern is obvious—the doors that remain open are those corresponding to the perfect squares. Once you make this connection, it is easy to prove it for the general case. Hence the total number of open doors is $\lfloor \sqrt{500} \rfloor = 22$. Solution 18.1 on Page 422 develops this analysis in more detail.

Optimally selecting a red card (Problem 17.18 on Page 143) and avoiding losing at the alternating coin pickup game (Problem 18.20 on Page 151) are other problems that benefit from use of the “small example” principle.

ITERATIVE REFINEMENT OF A BRUTE-FORCE SOLUTION

Many problems can be solved optimally by a simple algorithm that has a high time/space complexity—this is sometimes referred to as a brute-force solution. Other

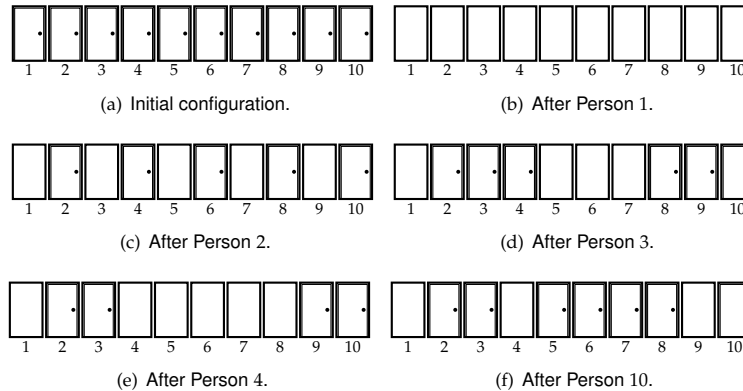


Figure 1.4: Progressive updates to 10 doors.

terms are *exhaustive search* and *generate-and-test*. Often this algorithm can be refined to one that is faster. At the very least it may offer hints into the nature of the problem.

As an example, suppose you were asked to write a function that takes an array A of n numbers, and rearranges A 's elements to get a new array B having the property that $B[0] \leq B[1] \geq B[2] \leq B[3] \geq B[4] \leq B[5] \geq \dots$.

One straightforward solution is to sort A and interleave the bottom and top halves of the sorted array. Alternately, we could sort A and then swap the elements at the pairs $(A[1], A[2]), (A[3], A[4]), \dots$. Both these approaches have the same time complexity as sorting, namely $O(n \log n)$.

You will soon realize that it is not necessary to sort A to achieve the desired configuration—you could simply rearrange the elements around the median, and then perform the interleaving. Median finding can be performed in time $O(n)$ (Solution 8.13 on Page 260), which is the overall time complexity of this approach.

Finally, you may notice that the desired ordering is very local, and realize that it is not necessary to find the median. Iterating through the array and swapping $A[i]$ and $A[i + 1]$ when i is even and $A[i] > A[i + 1]$ or i is odd and $A[i] < A[i + 1]$, achieves the desired configuration. In code:

```

1 template <typename T>
2 void rearrange(vector<T> &A) {
3     for (int i = 0; i < A.size() - 1; ++i) {
4         if ((i & 1 && A[i] < A[i + 1]) || ((i & 1) == 0 && A[i] > A[i + 1])) {
5             swap(A[i], A[i + 1]);
6         }
7     }
8 }

```

This approach has time complexity $O(n)$, which is the same as the approach based on median finding. However it is much easier to implement, and operates in an online fashion, i.e., it never needs to store more than two elements in memory or read a previous element.

As another example of iterative refinement, consider the problem of string search (Problem 3.19 on Page 43): given two strings s (search string) and t (text), find all occurrences of s in t . Since s can occur at any offset in t , the brute-force solution is to test for a match at every offset. This algorithm is perfectly correct; its time complexity is $O(nm)$, where n and m are the lengths of s and t .

After trying some examples, you may see that there are several ways to improve the time complexity of the brute-force algorithm. As an example, if the character $t[i]$ is not present in s you can advance the matching by n characters. Furthermore, this skipping works better if we match the search string from its end and work backwards. These refinements will make the algorithm very fast (linear time) on random text and search strings; however, the worst case complexity remains $O(nm)$.

You can make the additional observation that a partial match of s that does not result in a full match implies other offsets that cannot lead to full matches. If $s = ab\dabcabc$ and if, starting backwards, we have a partial match up to $abcabc$ that does not result in a full match, we know that the next possible matching offset has to be at least three positions ahead (where we can match the second abc from the partial match).

By putting together these refinements you will have arrived at the famous Boyer-Moore string search algorithm—its worst-case time complexity is $O(n + m)$ (which is the best possible from a theoretical perspective); it is also one of the fastest string search algorithms in practice.

Many other sophisticated algorithms can be developed in this fashion. As another example, the brute-force solution to computing the maximum subarray sum for an integer array of length n is to compute the sum of all subarrays, which has $O(n^3)$ time complexity. This can be improved to $O(n^2)$ by precomputing the sums of all the prefixes of the given arrays; this allows the sum of a subarray to be computed in $O(1)$ time. The natural divide and conquer algorithm has an $O(n \log n)$ time complexity. Finally, one can observe that a maximum subarray must end at one of n indices, and the maximum subarray sum for a subarray ending at index i can be computed from previous maximum subarray sums, which leads to an $O(n)$ algorithm. Details are presented on Page 98.

REDUCTION

Consider the problem of finding if one string is a rotation of the other, e.g., “car” and “arc” are rotations of each other. A natural approach may be to rotate the first string by every possible offset and then compare it with the second string. This algorithm would have quadratic time complexity.

You may notice that this problem is quite similar to string search, which can be done in linear time, albeit using a somewhat complex algorithm. Therefore it is natural to try to reduce this problem to string search. Indeed, if we concatenate the second string with itself and search for the first string in the resulting string, we will find a match iff the two original strings are rotations of each other. This reduction yields a linear time algorithm for our problem.

The reduction principle is also illustrated in the problem of checking whether a road network is resilient in the presence of blockages (Problem 13.4 on Page 114) and the problem of finding the minimum number of pictures needed to photograph a set of teams (Problem 13.7 on Page 116).

Usually you try to reduce the given problem to an easier problem. Sometimes, however, you need to reduce a problem known to be difficult to the given problem. This shows that the given problem is difficult, which justifies heuristics and approximate solutions. Such scenarios are described in more detail in Chapter 14.

GRAPH MODELING

Drawing pictures is a great way to brainstorm for a potential solution. If the relationships in a given problem can be represented using a graph, quite often the problem can be reduced to a well-known graph problem. For example, suppose you are given a set of exchange rates among currencies and you want to determine if an arbitrage exists, i.e., there is a way by which you can start with one unit of some currency C and perform a series of barter which results in having more than one unit of C .

Table 1.4 shows a representative example. An arbitrage is possible for this set of exchange rates: $1 \text{ USD} \rightarrow 1 \times 0.8123 = 0.8123 \text{ EUR} \rightarrow 0.8123 \times 1.2010 = 0.9755723 \text{ CHF} \rightarrow 0.9755723 \times 80.39 = 78.426257197 \text{ JPY} \rightarrow 78.426257197 \times 0.0128 = 1.00385609212 \text{ USD}$.

Table 1.4: Exchange rates for seven major currencies.

Symbol	USD	EUR	GBP	JPY	CHF	CAD	AUD
USD	1	0.8148	0.6404	78.125	0.9784	0.9924	0.9465
EUR	1.2275	1	0.7860	96.55	1.2010	1.2182	1.1616
GBP	1.5617	1.2724	1	122.83	1.5280	1.5498	1.4778
JPY	0.0128	0.0104	0.0081	1	1.2442	0.0126	0.0120
CHF	1.0219	0.8327	0.6546	80.39	1	1.0142	0.9672
CAD	1.0076	0.8206	0.6453	79.26	0.9859	1	0.9535
AUD	1.0567	0.8609	0.6767	83.12	1.0339	1.0487	1

We can model the problem with a graph where currencies correspond to vertices, exchanges correspond to edges, and the edge weight is set to the logarithm of the exchange rate. If we can find a cycle in the graph with a positive weight, we would have found such a series of exchanges. Such a cycle can be solved using the Bellman-Ford algorithm, as described in Solution 13.12 on Page 365.

WRITE AN EQUATION

Some problems can be solved by expressing them in the language of mathematics. Suppose you were asked to write an algorithm that computes binomial coefficients, $\binom{n}{k} = \frac{n!}{k!(n-k)!}$.

The problem with computing the binomial coefficient directly from the definition is that the factorial function grows quickly and can overflow an integer variable. If

we use floating point representations for numbers, we lose precision and the problem of overflow does not go away. These problems potentially exist even if the final value of $\binom{n}{k}$ is small. One can try to factor the numerator and denominator and try to cancel out common terms but factorization is itself a hard problem.

The binomial coefficients satisfy the *addition formula*:

$$\binom{n}{k} = \binom{n-1}{k} + \binom{n-1}{k-1}$$

This identity leads to a straightforward recursion for computing $\binom{n}{k}$ which avoids the problems described above. DP has to be used to achieve good time complexity—details are in Solution 12.14 on Page 332.

VARIATION

The idea of the variation pattern is to solve a slightly different (possibly more general) problem and map its solution to your problem.

Suppose we were asked to design an algorithm which takes as input an undirected graph and produces as output a black or white coloring of the vertices such that for every vertex at least half of its neighbors differ in color from it.

We could try to solve this problem by assigning arbitrary colors to vertices and then flipping colors wherever constraints are not met. However this approach may lead to increasing the number of vertices that do not satisfy the constraint.

It turns out we can define a slightly different problem whose solution will yield the desired coloring. Define an edge to be *diverse* if its ends have different colors. It is straightforward to verify that a coloring that maximizes the number of diverse edges also satisfies the constraint of the original problem, so there always exists a coloring satisfying the constraint.

It is not necessary to find a coloring that maximizes the number of diverse edges. All that is needed is a coloring in which the set of diverse edges is maximal with respect to single vertex flips. Such a coloring can be computed efficiently; details are given in Problem 12.29 on Page 109.

INVARIANTS

The following problem was popular at interviews in the early 1990s. You are given an 8×8 square with two unit sized squares at the opposite ends of a diagonal removed, leaving 62 squares, as illustrated in Figure 1.5(a) on the facing page. You are given 31 rectangular dominoes. Each can cover exactly two squares. How would you cover all the 62 squares with the dominoes?

You can spend hours trying unsuccessfully to find such a covering. This experience will teach you that a problem may be intentionally worded to mislead you into following a futile path.

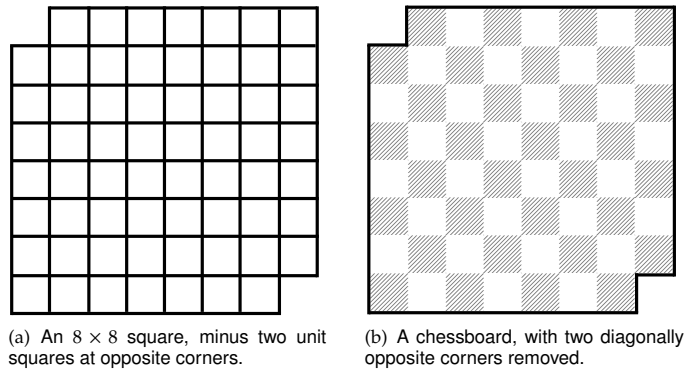


Figure 1.5: Invariant analysis exploiting auxiliary elements.

Here is a simple argument that no covering exists. Think of the 8×8 square as a chessboard as shown in Figure 1.5(b). Then the two removed squares will always have the same color, so there will be either 30 black and 32 white squares to be covered, or 32 black and 30 white squares to be covered. Each domino will cover one black and one white square, so the number of black and white squares covered as you successively put down the dominoes is equal. Hence it is impossible to cover the given chessboard.

This proof of impossibility is an example of invariant analysis. An invariant is a function of the state of a system being analyzed that remains constant in the presence of (possibly restricted) updates to the state. Invariant analysis is particularly powerful at proving impossibility results as we just saw with the chessboard tiling problem. The challenge is finding a simple invariant.

The argument above also used the auxiliary elements pattern, in which we added a new element to our problem to get closer to a solution. The original problem did not talk about the colors of individual squares; adding these colors made proving impossibility much easier.

It is possible to prove impossibility without appealing to square colors. Specifically, orient the board with the missing pieces on the lower right and upper left. An impossibility proof exists based on a case-analysis for each column on the height of the highest domino that is parallel to the base. However, the proof given above is much simpler.

Invariant analysis can be used to design algorithms, as well as prove impossibility results. In the coin selection problem, sixteen coins are arranged in a line, as in Figure 1.6 on the following page. Two players, F and S , take turns at choosing one coin each—they can only choose from the two coins at the ends of the line. Player F goes first. The game ends when all the coins have been picked up. The player whose coins have the higher total value wins.

The optimum strategy for F can be computed using DP (Solution 12.19 on Page 337). However, if F 's goal is simply to ensure he does not do worse than

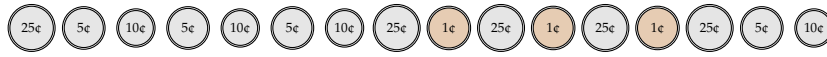


Figure 1.6: Coins in a row.

S , he can achieve this goal with much less computation. Specifically, he can number the coins from 1 to 16 from left-to-right, and compute the sum of the even-index coins and the sum of the odd-index coins. Suppose the odd-index sum is larger. Then F can force S to always select an even-index coin by selecting the odd-index coins when it is his own turn, ensuring that S cannot win. The same principle holds when the even-index sum is larger, or the sums are equal. Details are given in Solution 18.5 on Page 424.

Invariant analysis can be used with symmetry to solve very difficult problems, sometimes in less than intuitive ways. This is illustrated by the game known as “chomp” in which Player F and Player S alternately take bites from a chocolate bar. The chocolate bar is an $n \times n$ rectangle; a bite must remove a square and all squares above and to the right in the chocolate bar. The first player to eat the lower leftmost square, which is poisoned, loses. Player F can force a win by first selecting the square immediately above and to the right of the poisoned square, leaving the bar shaped like an L , with equal vertical and horizontal sides. Now whatever move S makes, F can play a symmetric move about the line bisecting the chocolate bar through the poisoned square to recreate the L shape (this is the invariant), which forces S to be the first to consume the poisoned square. Details are given in Solution 18.6 on Page 424.

Algorithm design using invariants is also illustrated in Solution 9.7 on Page 270 (can the characters in a string be permuted to form a palindrome?) and in Solution 10.14 on Page 291 (are there three elements in an array that sum to a given number?).

Complexity Analysis

The run time of an algorithm depends on the size of its input. One common approach to capture the run time dependency is by expressing asymptotic bounds on the worst-case run time as a function of the input size. Specifically, the run time of an algorithm on an input of size n is $O(f(n))$ if, for sufficiently large n , the run time is not more than $f(n)$ times a constant. The big- O notation simply indicates an upper bound; if the run time is asymptotically proportional to $f(n)$, the complexity is written as $\Theta(f(n))$. (Note that the big- O notation is widely used where sometimes Θ is more appropriate.) The notation $\Omega(f(n))$ is used to denote an asymptotic lower bound of $f(n)$ on the time complexity of an algorithm.

As an example, searching an unsorted array of integers of length n , for a given integer, has an asymptotic complexity of $\Theta(n)$ since in the worst-case, the given integer may not be present. Similarly, consider the naïve algorithm for testing primality that tries all numbers from 2 to the square root of the input number n . What is its complexity? In the best case, n is divisible by 2. However in the worst-case the

input may be a prime, so the algorithm performs \sqrt{n} iterations. Furthermore, since the number n requires $\lg n$ bits to encode, this algorithm's complexity is actually exponential in the size of the input. The big-Omega notation is illustrated by the $\Omega(n \log n)$ information-theoretic lower bound on any comparison-based array sorting algorithm.

Generally speaking, if an algorithm has a run time that is a polynomial, i.e., $O(n^k)$ for some fixed k , where n is the size of the input, it is considered to be efficient; otherwise it is inefficient. Notable exceptions exist—for example, the simplex algorithm for linear programming is not polynomial but works very well in practice. On the other hand, the AKS primality testing algorithm has polynomial run time but the degree of the polynomial is too high for it to be competitive with randomized algorithms for primality testing.

Complexity theory is applied in a similar way when analyzing the space requirements of an algorithm. Usually, the space needed to read in an instance is not included; otherwise, every algorithm would have $\Omega(n)$ space complexity.

Several of our problems call for an algorithm that uses $O(1)$ space. Conceptually, the memory used by such an algorithm should not depend on the size of the input instance. Specifically, it should be possible to implement the algorithm without dynamic memory allocation (explicitly, or indirectly, e.g., through library routines). Furthermore, the maximum depth of the function call stack should also be a constant, independent of the input. The standard algorithm for depth-first search of a graph is an example of an algorithm that does not perform any dynamic allocation, but uses the function call stack for implicit storage—its space complexity is not $O(1)$.

A streaming algorithm is one in which the input is presented as a sequence of items and is examined in only a few passes (typically just one). These algorithms have limited memory available to them (much less than the input size) and also limited processing time per item. Algorithms for computing summary statistics on log file data often fall into this category.

As a rule, algorithms should be designed with the goal of reducing the worst-case complexity rather than average-case complexity for several reasons:

1. It is very difficult to define meaningful distributions on the inputs.
2. Pathological inputs are more likely than statistical models may predict. A worst-case input for a naïve implementation of quicksort is one where all entries are the same, which is not unlikely in a practical setting.
3. Malicious users may exploit bad worst-case performance to create denial-of-service attacks.

Conclusion

In addition to developing intuition for which patterns may apply to a given problem, it is also important to know when your approach is not working. In an interview setting, even if you do not end up solving the problem entirely you will get credit for approaching problems in a systematic way and clearly communicating your approach to the problem. We cover nontechnical aspects of problem solving in Chapter 20.