

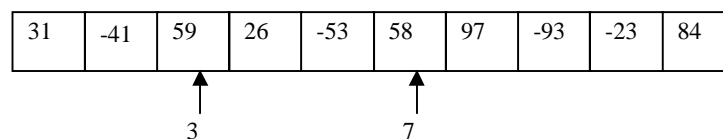
Column 7: Algorithm Design Techniques

Column 2 describes the "everyday" impact that algorithm design can have on programmers: an algorithmic view of a problem gives insights that can make a program simpler to understand and to write. In this column we'll study a contribution of the field that is less frequent but more impressive: sophisticated algorithmic methods sometimes lead to dramatic performance improvements.

This column is built around one small problem, with an emphasis on the algorithms that solve it and the techniques used to design the algorithms. Some of the algorithms are a little complicated, but with justification. While the first program we'll study takes thirty-nine days to solve a problem of size ten thousand, the final one solves the same problem in less than a second.

7.1 The Problem and a Simple Algorithm

The problem arose in one-dimensional pattern recognition; I'll describe its history later. The input is a vector \mathbf{X} of N real numbers; the output is the maximum sum found in any *contiguous* subvector of the input. For instance, if the input vector is



then the program returns the sum of $\mathbf{X}[3..7]$, or 187. The problem is easy when all the numbers are positive—the maximum subvector is the entire input vector. The rub comes when some of the numbers are negative: should we include a negative number in hopes that the positive numbers to its sides will compensate for its negative contribution? To complete the definition of the problem, we'll say that when all inputs are negative the maximum sum subvector is the empty vector, which has sum zero.

The obvious program for this task iterates over all pairs of integers L and U satisfying $1 \leq L \leq U \leq N$; for each pair it computes the sum of $\mathbf{X}[L..U]$ and checks whether that sum is greater than the maximum sum so far. The pseudocode for Algorithm 1 is

```
MaxSoFar := 0.0
for L := 1 to N do
  for U := L to N do
    Sum := 0.0
    for I := L to U do
      Sum := Sum + X[I]
    /* Sum now contains the sum of X[L..U] */
    MaxSoFar := max(MaxSoFar, Sum)
```

This code is short, straightforward, and easy to understand. Unfortunately, it has the severe disadvantage of being slow. On the computer I typically use, for instance, the code takes about an hour if N is 1000 and thirty-nine days if N is 10,000; we'll get to timing details in Section 7.5.

Those times are anecdotal; we get a different kind of feeling for the algorithm's efficiency using the "big-oh" notation described in Section 5.1. The statements in the outermost loop are executed exactly N times, and those in the middle loop are executed at most N times in each execution of the

outer loop. Multiplying those two factors of N shows that the four lines contained in the middle loop are executed $O(N^2)$ times. The loop in those four lines is never executed more than N times, so its cost is $O(N)$. Multiplying the cost per inner loop times its number of executions shows that the cost of the entire program is proportional to N cubed, so we'll refer to this as a cubic algorithm.

This example illustrates the technique of big-oh analysis of run time and many of its strengths and weaknesses. Its primary weakness is that we still don't really know the amount of time the program will take for any particular input; we just know that the number of steps it executes is $O(N^3)$. That weakness is often compensated for by two strong points of the method. Big-oh analyses are usually easy to perform (as above), and the asymptotic run time is often sufficient for a back-of-the-envelope calculation to decide whether or not a program is efficient enough for a given application.

The next several sections use asymptotic run time as the only measure of program efficiency. If that makes you uncomfortable, peek ahead to Section 7.5, which shows that for this problem such analyses are extremely informative. Before you read further, though, take a minute to try to find a faster algorithm.

7.2 Two Quadratic Algorithms

Most programmers have the same response to Algorithm 1: "There's an obvious way to make it a lot faster." There are two obvious ways, however, and if one is obvious to a given programmer then the other often isn't. Both algorithms are quadratic—they take $O(N^2)$ steps on an input of size N —and both achieve their run time by computing the sum of $X[L..U]$ in a constant number of steps rather than in the $U-L+1$ steps of Algorithm 1. But the two quadratic algorithms use very different methods to compute the sum in constant time.

The first quadratic algorithm computes the sum quickly by noticing that the sum of $X[L..U]$ has an intimate relationship to the sum previously computed, that of $X[L..U-1]$. Exploiting that relationship leads to Algorithm 2.

```

MaxSoFar := 0.0
for L := 1 to N do
  Sum := 0.0
  for U := L to N do
    Sum := Sum + X[U]
    /* Sum now contains the sum of X[L..U] */
    MaxSoFar := max(MaxSoFar, Sum)

```

The statements inside the first loop are executed N times, and those inside the second loop are executed at most N times on each execution of the outer loop, so the total run time is $O(N^2)$.

An alternative quadratic algorithm computes the sum in the inner loop by accessing a data structure built before the outer loop is ever executed. The I^{th} element of **CumArray** contains the cumulative sum of the values in $X[1..I]$, so the sum of the values in $X[L..U]$ can be found by computing $\text{CumArray}[U] - \text{CumArray}[L-1]$. This results in the following code for Algorithm 2b.

```

CumArray[0] := 0.0
for I := 1 to N do
  CumArray[I] := CumArray[I-1] + X[I]

```

```

MaxSoFar := 0.0
for L := 1 to N do
  for U := L to N do
    Sum := CumArray[U] - CumArray[L-1]
    /* Sum now contains the sum of X[L..U] */
    MaxSoFar := max(MaxSoFar, Sum)

```

This code takes $O(N^2)$ time; the analysis is exactly the same as the analysis of Algorithm 2.

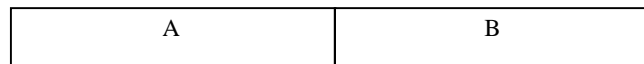
The algorithms we've seen so far inspect all possible pairs of starting and ending values of subvectors and consider the sum of the numbers in that subvector. Because there are $O(N^2)$ subvectors, any algorithm that inspects all such values must take at least quadratic time. Can you think of a way to sidestep this problem and achieve an algorithm that runs in less time?

7.3 A Divide-and-Conquer Algorithm

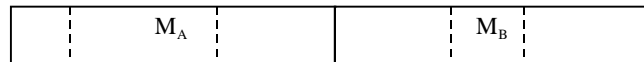
Our first subquadratic algorithm is complicated; if you get bogged down in its details, you won't lose much by skipping to the next section. It is based on the following divide-and-conquer schema:

To solve a problem of size N , recursively solve two subproblems of size approximately $N/2$, and combine their solutions to yield a solution to the complete problem.

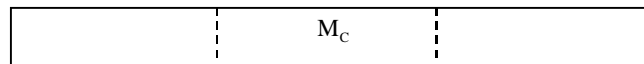
In this case the original problem deals with a vector of size N , so the most natural way to divide it into subproblems is to create two subvectors of approximately equal size, which we'll call **A** and **B**.



We then recursively find the maximum subvectors in **A** and **B**, which we'll call M_A and M_B .



It is tempting to think that we have solved the problem because the maximum sum subvector of the entire vector must be either M_A or M_B , and that is almost right. In fact, the maximum is either entirely in **A**, entirely in **B**, or it crosses the border between **A** and **B**; we'll call that M_C for the maximum *crossing* the border.



Thus our divide-and-conquer algorithm will compute M_A and M_B recursively, compute M_C by some other means, and then return the maximum of the three.

That description is almost enough to write code. All we have left to describe is how we'll handle small vectors and how we'll compute M_C . The former is easy: the maximum of a one-element vector is the only value in the vector or zero if that number is negative, and the maximum of a zero-element vector was previously defined to be zero. To compute M_C we observe that its component in **A** is the largest subvector starting at the boundary and reaching into **A**, and similarly for its

component in **B**. Putting these facts together leads to the following code for Algorithm 3, which is originally invoked by the procedure call

```

Answer := MaxSum(1, N)

recursive function MaxSum(L, U)
  if L > U then      /* Zero-element vector */
    return 0.0
  if L = U then     /* One-element vector */
    return max(0.0, X[L])

  M := (L+U)/2      /* A is X[L..M], B is X[M+1..U] */
  /* Find max crossing to left */
  Sum := 0.0; MaxToLeft := 0.0
  for I := M downto L do
    Sum := Sum + X[I]
    MaxToLeft := max(MaxToLeft, Sum)
  /* Find max crossing to right */
  Sum := 0.0; MaxToRight := 0.0
  for I := M+1 to U do
    Sum := Sum + X[I]
    MaxToRight := max(MaxToRight, Sum)
  MaxCrossing := MaxToLeft + MaxToRight
  MaxInA := MaxSum(L,M)
  MaxInB := MaxSum(M+1,U)
  return max(MaxCrossing, MaxInA, MaxInB)

```

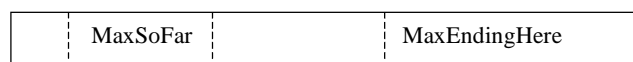
The code is complicated and easy to get wrong, but it solves the problem in $O(N \log N)$ time. There are a number of ways to prove this fact. An informal argument observes that the algorithm does $O(N)$ work on each of $O(\log N)$ levels of recursion. The argument can be made more precise by the use of recurrence relations. If $T(N)$ denotes the time to solve a problem of size N , then $T(1)=O(1)$ and

$$T(N) = 2T(N/2) + O(N).$$

Problem 11 shows that this recurrence has the solution $T(N) = O(N \log N)$.

7.4 A Scanning Algorithm

We'll now use the simplest kind of algorithm that operates on arrays: it starts at the left end (element $X[1]$) and scans through to the right end (element $X[N]$), keeping track of the maximum sum subvector seen so far. The maximum is initially zero. Suppose that we've solved the problem for $X[1..I-1]$; how can we extend that to a solution for the first I elements? We use reasoning similar to that of the divide-and-conquer algorithm: the maximum sum in the first I elements is either the maximum sum in the first $I-1$ elements (which we'll call **MaxSoFar**), or it is that of a subvector that ends in position I (which we'll call **MaxEndingHere**).



Recomputing **MaxEndingHere** from scratch using code like that in Algorithm 3 yields yet another quadratic algorithm. We can get around this by using the technique that led to Algorithm 2: instead of computing the maximum subvector ending in position **I** from scratch, we'll use the maximum subvector that ends in position **I-1**. This results in Algorithm 4.

```

MaxSoFar := 0.0
MaxEndingHere := 0.0
for I := 1 to N do
  /* Invariant: MaxEndingHere and MaxSoFar are accurate for
  X[1..I-1] */
  MaxEndingHere := max(MaxEndingHere+X[I], 0.0)
  MaxSoFar := max(MaxSoFar, MaxEndingHere)

```

The key to understanding this program is the variable **MaxEndingHere**. Before the first assignment statement in the loop, **MaxEndingHere** contains the value of the maximum subvector ending in position **I-1**; the assignment statement modifies it to contain the value of the maximum subvector ending in position **I**. The statement increases it by the value **X**[**I**] so long as doing so keeps it positive; when it goes negative, it is reset to zero because the maximum subvector ending at **I** is the empty vector. Although the code is subtle, it is short and fast: its run time is $O(N)$, so we'll refer to it as a linear algorithm. David Gries systematically derives and verifies this algorithm in his paper "A Note on the Standard Strategy for Developing Loop Invariants and Loops" in the journal *Science of Computer Programming* 2, pp. 207-214.

7.5 What Does It Matter?

So far I've played fast and loose with "big-ohs"; it's time for me to come clean and tell about the run times of the programs. I implemented the four primary algorithms (all except Algorithm 2b) in the C language on a VAX11/750, timed them, and extrapolated the observed run times to achieve the following table.

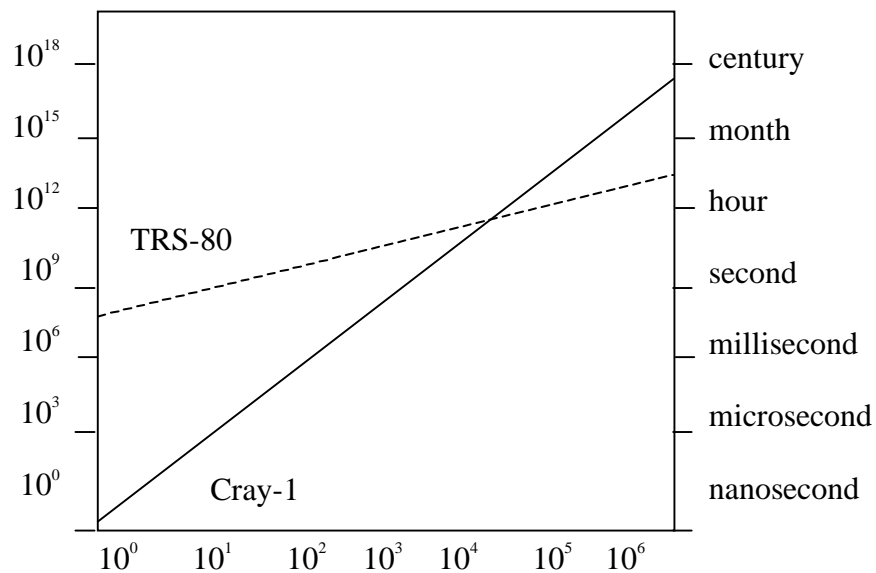
ALGORITHM		1	2	3	4
Lines of C Code		8	7	14	7
Run time in microseconds		$3.4N^3$	$13N^2$	$46N \log_2 N$	$33N$
Time to solve a problem of size	10^2	3.4 secs	.13 secs	.03 secs	.003 secs
	10^3	.94 hrs	13 secs	.45 secs	.033 secs
	10^4	39 days	22 mins	6.1 secs	.33 secs
	10^5	108 yrs	1.5 days	1.3 mins	3.3 secs
	10^6	108 mill	5 mos	15 mins	33 secs
Max size problem solved in one	sec	67	280	2000	30,000
	min	260	2200	82,000	2,000,000
	hr	1000	17,000	3,500,000	120,000,000
	day	3000	81,000	73,000,000	2,800,000,000
If N multiplies by 10, time multiplies by		1000	100	10+	10
If time multiplies by 10, N multiplies by		2.15	3.16	10-	10

This table makes a number of points. The most important is that proper algorithm design can make a big difference in run time; that point is underscored by the middle rows. The last two rows show how increases in problem size are related to increases in run time.

Another important point is that when we're comparing cubic, quadratic, and linear algorithms with one another, the constant factors of the programs don't matter much. (The discussion of the $O(N!)$ algorithm in Section 2.4 shows that constant factors matter even less in functions that grow faster than polynomially.) To underscore this point, I conducted an experiment in which I tried to make the constant factors of two algorithms differ by as much as possible. To achieve a huge constant factor I implemented Algorithm 4 on a BASIC interpreter on a Radio Shack TRS-80 Model III microcomputer. For the other end of the spectrum, Eric Grosse and I implemented Algorithm 1 in fine-tuned FORTRAN on a Cray-1 supercomputer. We got the disparity we wanted: the run time of the cubic algorithm was measured as $3.0N^3$ nanoseconds, while the run time of the linear algorithm was $19.5N$ milliseconds, or $19,500,000N$ nanoseconds. This table shows how those expressions translate to times for various problem sizes.

N	CRAY-1, FORTRAN, CUBIC ALGORITHM	TRS-80, BASIC, LINEAR ALGORITHM
10	3.0 microsecs	200 milliseecs
100	3.0 milliseecs	2.0 secs
1000	3.0 secs	20 secs
10,000	49 mins	3.2 mins
100,000	35 days	32 mins
1,000,000	95 yrs	5.4 hrs

The difference in constant factors of six and a half million allowed the cubic algorithm to start off faster, but the linear algorithm was bound to catch up. The break-even point for the two algorithms is around 2,500, where each takes about fifty seconds.



7.6 Principles

The history of the problem sheds light on the algorithm design techniques. The problem arose in a pattern-matching procedure designed by Ulf Grenander of Brown University in the two-dimensional form described in Problem 7. In that form, the maximum sum subarray was the maximum likelihood estimator of a certain kind of pattern in a digitized picture. Because the two-dimensional problem required too much time to solve, Grenander simplified it to one dimension to gain insight into its structure.

Grenander observed that the cubic time of Algorithm 1 was prohibitively slow, and derived Algorithm 2. In 1977 he described the problem to Michael Shamos of UNILOGIC, Ltd. (then of Carnegie-Mellon University) who overnight designed Algorithm 3. When Shamos showed me the problem shortly thereafter, we thought that it was probably the best possible; researchers had just shown that several similar problems require time proportional to $N \log N$. A few days later Shamos described the problem and its history at a Carnegie-Mellon seminar attended by statistician Jay Kadane, who designed Algorithm 4 within a minute. Fortunately, we know that there is no faster algorithm: any correct algorithm must take $O(N)$ time.

Even though the one-dimensional problem is completely solved, Grenander's original two-dimensional problem remained open eight years after it was posed, as this book went to press. Because of the computational expense of all known algorithms, Grenander had to abandon that approach to the pattern-matching problem. Readers who feel that the linear-time algorithm for the one-dimensional problem is "obvious" are therefore urged to find an "obvious" algorithm for Problem 7!

The algorithms in this story were never incorporated into a system, but they illustrate important algorithm design techniques that have had substantial impact on many systems (see Section 7.9).

Save state to avoid recomputation. This simple form of dynamic programming arose in Algorithms 2 and 4. By using space to store results, we avoid using time to recompute them.

Preprocess information into data structures. The **CumArray** structure in Algorithm 2b allowed the sum of a subvector to be computed in just a couple of operations.

Divide-and-conquer algorithms. Algorithm 3 uses a simple form of divide-and-conquer; textbooks on algorithm design describe more advanced forms.

Scanning algorithms. Problems on arrays can often be solved by asking "how can I extend a solution for $X[1..I-1]$ to a solution for $X[1..I]$?" Algorithm 4 stores both the old answer and some auxiliary data to compute the new answer.

Cumulatives. Algorithm 2b uses a cumulative table in which the I^{th} element contains the sum of the first I values of X ; such tables are common when dealing with ranges. In business data processing applications, for instance, one finds the sales from March to October by subtracting the February year-to-date sales from the October year-to-date sales.

Lower bounds. Algorithm designers sleep peacefully only when they know their algorithms are the best possible; for this assurance, they must prove a matching lower bound. The linear lower bound for this problem is the subject of Problem 9; more complex lower bounds can be quite difficult.

7.7 Problems

- Algorithms 3 and 4 use subtle code that is easy to get wrong. Use the program verification techniques of Column 4 to argue the correctness of the code; specify the loop invariants carefully.
- Our analysis of the four algorithms was done only at the "big-oh" level of detail. Analyze the number of **max** functions used by each algorithm as exactly as possible; does this exercise give any insight into the running times of the programs? How much space does each algorithm require?
- We defined the maximum subvector of an array of negative numbers to be zero, the sum of the empty subvector. Suppose that we had instead defined the maximum subvector to be the value of the largest element; how would you change the programs?
- Suppose that we wished to find the subvector with the sum closest to zero rather than that with maximum sum. What is the most efficient algorithm you can design for this task? What algorithm design techniques are applicable? What if we wished to find the subvector with the sum closest to a given real number **T**?
- A turnpike consists of **N-1** stretches of road between **N** toll stations; each stretch has an associated cost of travel. It is trivial to tell the cost of going between any two stations in $O(N)$ time using only an array of the costs or in constant time using a table with $O(N^2)$ entries. Describe a data structure that requires $O(N)$ space but allows the cost of any route to be computed in constant time.
- After the array $X[1..N]$ is initialized to zero, **N** of the following operations are performed

```

for I := L to U do
    X[I] := X[I] + V

```

where **L**, **U** and **V** are parameters of each operation (**L** and **U** are integers satisfying $1 \leq L \leq U \leq N$ and **V** is a real). After the **N** operations, the values of $X[1]$ through $X[N]$ are reported in order. The method just sketched requires $O(N^2)$ time. Can you find a faster algorithm?

- In the maximum subarray problem we are given an $N \times N$ array of reals, and we must find the maximum sum contained in any rectangular subarray. What is the complexity of this problem?
- Modify Algorithm 3 (the divide-and-conquer algorithm) to run in linear worst-case time.
- Prove that any correct algorithm for computing maximum subvectors must inspect all **N** inputs. (Algorithms for some problems may correctly ignore some inputs; consider Saxe's algorithm in Solution 2.2 and Boyer and Moore's substring searching algorithm in the October 1977 CACM.)
- Given integers **M** and **N** and the real vector $X[1..N]$, find the integer **I** ($I \leq N-M$) such that the sum $X[I] + \dots + X[I+M]$ is nearest zero.
- What is the solution of the recurrence $T(N) = 2T(N/2) + CN$ when $T(1) = 0$ and **N** is a power of two? Prove your result by mathematical induction. What if $T(1) = C$?

7.8 Further Reading

Only extensive study can put algorithm design techniques at your fingertips; most programmers will get this only from a textbook on algorithms. *Data Structures and Algorithms* by Aho, Hopcroft and Ullman (published by Addison-Wesley in 1983) is an excellent undergraduate text. Chapter 10 on "Algorithm Design Techniques" is especially relevant to this column.

7.9 The Impact of Algorithms [Sidebar]

Although the problem studied in this column illustrates several important techniques, it's really a toy—it was never incorporated into a system. We'll now survey a few real problems in which algorithm design techniques proved their worth.

Numerical Analysis. The standard example of the power of algorithm design is the discrete Fast Fourier Transform (FFT). Its divide-and-conquer structure reduced the time required for Fourier analysis from $O(N^2)$ to $O(N \log N)$. Because problems in signal processing and time series analysis frequently process inputs of size $N = 1000$ or greater, the algorithm speeds up programs by factors of more than one hundred.

In Section 10.3.C of his *Numerical Methods, Software, and Analysis* (published in 1983 by McGraw-Hill), John Rice chronicles the algorithmic history of three-dimensional elliptic partial differential equations. Such problems arise in simulating VLSI devices, oil wells, nuclear reactors, and airfoils. A small part of that history (mostly but not entirely from his book) is given in the following table. The run time gives the number of floating point operations required to solve the problem on an $N \times N \times N$ grid.

METHOD	YEAR	RUN TIME
Gaussian Elimination	1945	N^7
SOR Iteration (Suboptimal Paramcters)	1954	$8N^5$
SOR Iteration (Optimal Parameters)	1960	$8N^4 \log_2 N$
Cyclic Rcdution	1970	$8N^3 \log_2 N$
Multigrid	1978	$60N^3$

SOR stands for "successive over-relaxation". The $O(N^3)$ time of Multigrid is within a constant factor of optimal because the problem has that many inputs. For typical problem sizes ($N=64$), the speedup is a factor of a quarter million. Pages 1090-1091 of "Programming Pearls" in the November 1984 *Communications of the ACM* present data to support Rice's argument that the algorithmic speedup from 1945 to 1970 exceeds the hardware speedup during that period.

Graph Algorithms. In a common method of building integrated circuitry, the designer describes an electrical circuit as a graph that is later transformed into a chip design. A popular approach to laying out the circuit uses the "graph partitioning" problem to divide the entire electrical circuit into subcomponents. Heuristic algorithms for graph partitioning developed in the early 1970's used $O(N^2)$ time to partition a circuit with a total of N components and wires. Fiduccia and Mattheyses describe "A linear-time heuristic for improving network partition" in the *19th Design Automation*

Conference. Because typical problems involve a few thousand components, their method reduces layout time from a few hours to a few minutes.

Geometric Algorithms. Late in their design, integrated circuits are specified as geometric “artwork” that is eventually etched onto chips. Design systems process the artwork to perform tasks such as extracting the electrical circuit it describes, which is then compared to the circuit the designer specified. In the days when integrated circuits had $N= 1000$ geometric figures that specified 100 transistors, algorithms that compared all pairs of geometric figures in $O(N^2)$ time could perform the task in a few minutes. Now that VLSI chips contain millions of geometric components, quadratic algorithms would take months. “Plane sweep” or “scan line” algorithms have reduced the run time to $O(N \log N)$, so the designs can now be processed in a few hours. Szymanski and Van Wyk’s “Space efficient algorithms for VLSI artwork analysis” in the *20th Design Automation Conference* describes efficient algorithms for such tasks that use only $O(\sqrt{N})$ primary memory (a later version of the paper appears in the June 1985 *IEEE Design and Test*).

Appel’s program described in Section 5.1 uses a tree data structure to represent points in 3-space and thereby reduces an $O(N^2)$ algorithm to $O(N \log N)$ time. That was the first step in reducing the run time of the complete program from a year to a day.

<p><u>Programming Pearls</u> by Jon Bentley Addison-Wesley Publishing Company Reading, Massachusetts April, 1986</p>
--