

Chapter 3

SEQUENCE ALIGNMENTS

Michael S. Waterman

TABLE OF CONTENTS

I.	Introduction	54
II.	The Number of Alignments	55
III.	Dynamic Programming Alignment of Two Sequences	58
	A. Distance Alignment	58
	B. Similarity Alignment	61
	C. Fitting One Sequence into Another	64
	D. Identification of Similar Segments	65
	E. Near Optimal Alignments	69
IV.	Dynamic Programming Alignment of Multiple Sequences	70
	A. Line Geometries	70
	B. Generalization of the Two-Sequence Algorithm	71
V.	Consensus Alignment of Multiple Sequences	72
VI.	Data Base Searches	74
	A. A Regions Method	76
	B. Rapid Similarity Searches	76
VII.	The Statistical Distribution of Alignment Scores	77
	A. The Log(n) Law	77
	B. A Data Base Study	78
	C. Linear and Logarithmic Behavior	79
	Appendix	81
	Acknowledgments	90
	References	90

I. INTRODUCTION

When two or more sequences are displayed with one sequence written over another, the resulting configuration is known as an alignment of the set of sequences. These displays are very common in molecular biology as they communicate information about proposed common evolution or function of the nucleotide positions found in any given column of an alignment. Sequence alignments are frequently obtained in an ad hoc manner, and simply written down in a way that makes the result "look good". Sometimes this means that important but implicit criteria have been satisfied. No mathematical approach devoid of deep biological intuition could hope to compete with this method. What is to be hoped for is that, over a long period of cooperation, mathematical and biological scientists can make explicit some of the biological insights. However, alignments that are pleasing to the eye of the scientist may not have much merit beyond that. The commonly used phrases "aligned so as to maximize homology" and "gaps inserted in order to maximize homology" might conceal a fairly confused attempt to maximize matches and minimize gaps. It is often not clear what has been attempted and "homology" is left undefined. In these cases, an explicit optimization function is an advantage, both for checking whether the alignment is optimal and for discussing the desirability of the optimization function itself.

One of the first mathematical questions concerns how difficult the problem of sequence alignment really is. A naive approach to the problem is to systematically list all alignments, evaluating each one. Alas, even for the two sequence case, when gaps are allowed, there are a huge number of alignments, and this exhaustive approach succeeds on only the smallest of problems. Section II will give an account of what is known about the combinatorics of alignments. The message is that alignment is a difficult problem!

Next, Section III turns to the problem of aligning two sequences by dynamic programming methods. This problem has received more mathematical attention than any other in molecular sequence analysis and the length of Section III reflects this. Computer science has studied the same problem under the description of the string edit problem; the problem statement is to find the minimum number of changes (substitutions, insertions, deletions, inversions) to convert one string (sequence, file, or word) into another. This problem has often been studied in many different forms and a book¹ has appeared on dynamic programming approaches to the problems of sequence comparison. In the present chapter, Section III.D on most similar segments is probably the most useful dynamic programming algorithm for current problems in biology, and some new developments are reported. A program written in C for the similar segments algorithm appears in the Appendix. Both similarity and distance methods are given for aligning two sequences. Methods for producing the optimal and near optimal alignments are described.

Many very interesting problems involve more than two sequences. Section IV gives dynamic programming approaches to the several sequence problem. Line geometries, a recent method employing ideas from the geometry of geodesics, are described and illustrated. While the line geometry approach is practical, it is likely to fail when the sequences are processed using the most unrelated sequences first. Line geometries handle the sequences in a given order, not simultaneously. Another dynamic programming method is presented that does not suffer from this shortcoming although it is prohibitively expensive in terms of time storage.

It will have escaped few readers that alignment is closely related to consensus. While the methods for consensus patterns do not instantly solve the problem of aligning several sequences, Section V describes an important new and practical approach. Based on the consensus word algorithm,² these techniques can align many long sequences in reasonable time and storage. Many problems in biology involve more than two sequences.

One of the major tasks facing a sequence analyst is that of comparing new sequences with all, or a major portion of, a large data base. The most successful algorithms for that

problem are based on hashing the data base, hashing the new sequence, and then comparing hash tables. Section VI treats these issues.

Finally suppose that these or some other methods have produced an alignment of interest to the scientist. Section VII discusses how the alignment might be examined for statistical significance. Some mathematics has recently been worked out that solves several aspects of this problem. Extreme value theory has been employed to derive the so-called log(n) distribution,^{3,4} that gives the statistical distribution of the best matching segments between two or more long sequences.

Recently in Waterman,⁵ methods of sequence comparison were reviewed and organized. Many references of that paper do not appear in this chapter. Here, the goal is to more fully describe some major methods of sequence alignment and not to provide a review

II. THE NUMBER OF ALIGNMENTS

In this section, a combinatorial treatment of sequence alignments is given. Biology provides the motivation for aligning sequences and for considering how difficult alignment is. It is then a mathematical task to estimate the number of sequence alignments. The results are applicable to biology in a negative sense; they assure one that a huge number of possible alignments exist and that direct enumeration is hopeless. The two-sequence case is handled first.

Notation is important here. Let $a = a_1 a_2 \dots a_n$ and $b = b_1 b_2 \dots b_m$ be two sequences of length, n and m . One way to think of alignment is that an alignment is produced when null elements, ϕ , are inserted into the sequences; the new sequences must be of the same length, L . Then the two sequences are written, one over the other. $a = a_1 a_2 \dots a_n$ becomes, with the insertion of ϕ , $a^* = a_1^* a_2^* \dots a_n^*$ while $b = b_1 b_2 \dots b_m$ becomes $b^* = b_1^* b_2^* \dots b_m^*$. The subsequence of a^* or b^* whose elements are not equal to ϕ is the original sequence. The alignment is

$$\begin{array}{l} a_1^* a_2^* \dots a_n^* \\ b_1^* b_2^* \dots b_m^* \end{array}$$

To see this process, let $a = \text{ATAAGC}$ and $b = \text{AAAAACG}$. To obtain an alignment, one of many possibilities is to set $a^* = \phi \text{ATAAGC} \phi$ and $b^* = \text{AAAA} \phi \text{CG}$. For example $b_1^* = \text{A}$, $b_2^* = \phi$, and $b_3^* = \text{C}$ while $b_4^* = \text{A}$, $b_5^* = \text{C}$ and $b_6^* = \text{G}$. The alignment is written

$$\begin{array}{l} a^* = \phi \text{ATAAGC} \phi \\ b^* = \text{AAAA} \phi \text{CG} \end{array}$$

Here $b_1^* = \text{A}$ is said to be inserted into the first sequence or deleted from the second, depending on the point of view. $a_2^* = \text{A}$ matches $b_2^* = \text{A}$ while $a_1^* = \text{T}$ and $b_1^* = \text{A}$ constitutes a mismatch.

The problem of this section is to find how many ways can a^* be written. No alignment terms ϕ are allowed as there is no point in matching two deletions. This makes it clear that $\max\{n, m\} \leq L \leq n + m$. The case $L = n + m$ comes, e.g., by first deleting all a , and then deleting all b :

$$\begin{array}{l} a_1 a_2 \dots a_n \phi \phi \dots \phi \\ \phi \phi \dots \phi b_1 b_2 \dots b_m \end{array}$$

Combinatorial insight comes by recognizing that alignments of these two sequences can end in exactly one of three ways

$$\begin{matrix} \dots & a_n & \dots & a_n & \dots & \phi \\ \dots & \phi & \dots & b_m & \dots & b_m \end{matrix}$$

where ϕ corresponds to an insertion/deletion of a_n , ϕ corresponds to a match or mismatch of a and b , and ϕ corresponds to an insertion/deletion of b_m . Note that the fate of the unseen bases (those not displayed) is not specified. Define

$$f(i, j) = \text{number of all possible alignments of one sequence of } i \text{ letters with another of } j \text{ letters}$$

The above three cases imply

$$f(n, m) = f(n - 1, m) + f(n - 1, m - 1) + f(n, m - 1)$$

This follows because, for example, the first term corresponding to deleting a_n can obtain with the number of alignments of a_1, \dots, a_{n-1} and b_1, \dots, b_m , which is $f(n - 1, m)$.

The author obtained the above recursion equation and with P. R. Stein determined that it specified the Stanton-Cowan numbers.⁶ Then Stein communicated the problem of asymptotics to H. T. Laquer who obtained the following theorem in 1981.⁷

Theorem 1 — Let $f(n, m)$ be defined as above. Then

$$f(n, n) \sim (1 + \sqrt{2})^{2n+1} \sqrt{n}$$

as $n \rightarrow \infty$, where $c(n) \sim d(n)$ means $\lim_{n \rightarrow \infty} c(n)/d(n) = 1$.

Two sequences of length 1000, then have

$$f(1000, 1000) \sim (1 + \sqrt{2})^{2001} \sqrt{1000} = 10^{784}$$

alignments! There are approximately 10^{24} elementary particles in the universe; Avogadro's number is on the order of 10^{23} .

If it is agreed not to count

$$\begin{matrix} C\phi & \text{and} & \phi C \\ \phi G & & G\phi \end{matrix}$$

as distinct, the situation improves (slightly). Let $g(n, m)$ denote this smaller number of alignments. g_n has three possibilities

$$\begin{matrix} \dots & a_n & \phi & \dots & \phi & \phi & \dots & a_n \phi \\ \dots & b_m & \phi b_m & \dots & b_m & \phi b_m & \dots & \phi b_m \end{matrix}$$

while ϕ has

$$\begin{matrix} \dots & a_n & \phi a_n & \dots & a_n & \phi a_n & \dots & \phi a_n \\ \dots & b_m & \phi & \dots & \phi & \phi & \dots & \phi b_m \end{matrix}$$

The new version of the recursion equation is

$$g(n, m) = g(n - 1, m) + g(n, m - 1) + g(n - 1, m - 1) - g(n - 1, m - 1)$$

subtracting the double count. The result is given in the next theorem, where the asymptotics are derived from Stirling's formula.⁸

Table 1
BEHAVIOR OF $h(b, n) \sim \gamma_b n^{-1/2} D_b^n$,
THE NUMBER OF ALIGNMENTS OF
TWO SEQUENCES OF LENGTH n
WITH MATCHED BLOCKS OF
LENGTH AT LEAST b

b	D_b	γ_b
1	5.8284	0.57268
2	4.5189	0.53206
3	4.1489	0.54290
4	4.0400	0.55520
5	4.0103	0.56109
10	4.0001	0.56183
∞	4.0000	0.56419

Theorem 2 — If $g(n, m)$ is defined as above, $g(0, 0) = g(0, 1) = g(1, 0) = 1$, and $g(n, m) = \binom{n+m}{n}$. If $n = m$,

$$g(n, n) \sim \binom{2n}{n} \sim 2^{2n} (4\sqrt{n\pi})^{-1}, \text{ as } n \rightarrow \infty$$

Two sequences $n = m = 1000$ have $g(1000, 1000) \sim 10^{800}$ alignments so that direct search is still impossible.

It is possible to further reduce the number of alignments by requiring matches and mismatches to occur in blocks of length at least b without interruptions by deletions. The motivation for this is that biologists sometimes reject alignments with small groups of matches. The counting scheme of Theorem 1 is readopted with this new requirement. The following theorem appears in Griggs et al.⁹ where it is derived via generating functions.

Theorem 3 — Let $h(b, n)$ be the number of alignments of two sequences of length n where matches must occur in blocks of length at least $b \geq 1$. Define $\Phi(x) = (1 - x)^2 - 4x(x^b - x + 1)^2$, and let ρ be the smallest real root of $\Phi(0) = 0$. Then

$$h(b, n) \sim (\gamma_b n^{-1/2}) D_b^n, \text{ as } n \rightarrow \infty$$

where $D_b = \rho^{-1}$ and

$$\gamma_b = (\rho^b - \rho + 1)(-\pi\rho\Phi'(\rho))^{-1/2}$$

Notice that $b = 1$ has $h(1, n) = f(n, n)$. Table 1 shows the behavior with b . When $b = 2$, for example,

$$h(2, n) \sim (0.53206)n^{-1/2}(4.5189)^n$$

More than two sequences is bound to make the problem more complicated, greatly increasing the number of alignments. The question is to determine how many more alignments. Recently, Griggs et al.¹⁰ have provided an answer which is given in the next theorem. The methods of proof are the most difficult of any required so far, using a saddle point technique.

Theorem 4 — Let $l_k(n)$ be the number of alignments of k sequences of length n . For a fixed $k \geq 2$,

$$\lim_{n \rightarrow \infty} \frac{\ell_n \ell_n(n)}{n} = \ell n c_k$$

where

$$c_k = (2^{1/k} - 1)^{-1} = \frac{1}{\sqrt{2}} \left[\frac{k}{\ell n(2)} \right]$$

For our use notice that

$$\ell_n - c_k^n = (2^{1/k} - 1)^{-kn} \\ = 2^{-n/k} \left(\frac{k}{\ell n(2)} \right)^n$$

For $n = 1000$ and $k = 3$, $\ell_n(1000) \approx 10^{1731}$.

III. DYNAMIC PROGRAMMING ALIGNMENT OF TWO SEQUENCES

Needleman and Wunsch¹¹ wrote a paper titled "A general method applicable to the search for similarities in the amino acid sequence of two proteins". It was surely unknown to the authors that their method fit into a broad class of algorithms introduced by Richard Bellman under the name dynamic programming. Their paper has had a great deal of influence in biological sequence alignment. Its great advantage is that an explicit criterion for optimality of alignment is stated, as well as an efficient method of solution given. Insertions, deletions, mismatches (negative similarity), and matches (positive similarity) were allowed in the alignments.

During early 1970s, Stan Ulam and some other mathematicians became interested in defining a distance $D(a,b)$ on sequences. The minimum distance alignment was defined to be an alignment with the smallest weighted sum of mismatches, insertions, and deletions. The advantage of a distance was the construction of a metric space on the space of sequences:

1. $D(a,b) = 0$ if and only if $a = b$.
2. $D(a,b) = D(b,a)$ (symmetry)
3. $D(a,b) \leq D(a,c) + D(c,b)$ for any c (triangle inequality).

The emphasis on sequence metrics came from the fact that a matrix of sequence distances was often used to construct an evolutionary tree. P.H. Sellers¹² gave a dynamic programming algorithm, very similar to that of Needleman and Wunsch, to calculate the distance.

The historical order is reversed here. Distance methods are described in Section III.A, with similarity methods in Section III.B. As mentioned in the introduction, we find similarity to be the most satisfactory. All problems known to be solvable with distance methods can be solved with similarity methods. However, in Section III.D, a similarity solution is given that has no distance counterpart. Still, the metric space associated with a distance makes it worthwhile to present distance methods. Section III.C shows several simple modifications to solve related problems such as best fit of a short sequence into a long one. Section III.D studies the important problem of locating segments of two sequences which are unexpectedly similar, although the full sequences might not have a good alignment. New results are given here for this problem. Section III.E closes with a recent modification of the dynamic programming algorithms that allows all alignments near the optimal to be produced.

A. Distance Alignment

The sequences $a = a_1 a_2 \dots a_n$ and $b = b_1 b_2 \dots b_m$ are written over the alphabet

{A,C,G,T}. Any finite alphabet will of course work here. In particular the 20-letter, amino acid, alphabet of proteins, or the purine/pyrimidine alphabet for DNA can be used. Let $d(a,b)$ be a distance on the alphabet and let $g(a)$ be the positive cost of a gap of the letter "a". The distance $d(a,b)$ represents the cost of a mutation of a into b. If $d(a,b)$ is extended so that $d(a,\phi) = d(\phi,a) = g(a)$, then define

$$D(a,b) = \min \sum_{i=1}^l d(a'_i, b'_i)$$

where the minimum is extended over all alignments of a with b . Seller's result¹² can be summarized in the next Theorem.

Theorem 1 — If $a = a_1 a_2 \dots a_n$ and $b = b_1 b_2 \dots b_m$, define $D_{i,j} = D(a_1 \dots a_i, b_1 \dots b_j)$. Also set

$$D_{i0} = 0, D_{0j} = \sum_{k=1}^j d(\phi, b_k), \text{ and } D_{i,n} = \sum_{k=1}^i d(a_i, \phi)$$

Then

$$D_{i,j} = \min\{D_{i,j-1} + d(a_i, \phi), D_{i,j-1} + d(a_i, b_j), D_{i-1,j} + d(\phi, b_j)\} \quad (1)$$

If $d(\dots)$ is a metric on the alphabet, then $D(\dots)$ is a metric on the set of finite sequences.

Proof — We verify Equation 1 with reasoning similar to that for verifying the recursion equation for $\ell(n,m)$ in Section II. The alignment of $a_1 \dots a_i$ and $b_1 \dots b_j$ can end in one of three ways

$$\begin{array}{cccc} \dots & a_i & \dots & a_i & \dots & \phi \\ \dots & \phi & \dots & b_j & \dots & b_j \end{array}$$

If the optimal alignment ends in ϕ , the cost must be $D_{i,j-1} + d(a_i, \phi)$ since the initial part of the alignment must itself be optimal and align $a_1 \dots a_{i-1}$ with $b_1 \dots b_m$.

If the optimal alignment ends in b_j , the cost must be

$$D_{i,j-1} + d(a_i, b_j) \text{ since } a_1 \dots a_{i-1} \text{ and } b_1 \dots b_{j-1} \text{ must be optimally aligned}$$

The case ϕ is identical in reasoning with the case ϕ .

The optimal alignment has least cost of these three possibilities and Equation 1 is proven. Another statement of Equation 1 is

$$D_{i,j} = \min\{D_{i,j-1} + g, D_{i,j-1} + d(a_i, b_j), D_{i-1,j} + g\}$$

when g is a constant gap cost, $g = d(\phi, a) = d(a, \phi)$. These algorithms have computation cost proportional to nm , $O(nm)$.

To illustrate the algorithm, we align two *Escherichia coli* tRNA sequences, threonine tRNA (sequence a; GenBank name ECOTRTACU) and valine tRNA (sequence b; GenBank name ECOTRV1). The algorithm has $d(a,b) = 2$, if $a \neq b$, and $g = 2.5$. Entries in the matrix in Table 2 are multiplied by 10 to allow use of integer arithmetic. Table 2 shows the matrix for the 5' (left) ends of the sequences. There are 72 optimal alignments, one of which is shown next. Portions of the alignment common to all 72 alignments are boxed.

swered by Smith and Waterman¹⁶ and Fitch et al.¹⁷ When full sequences are aligned by distance (similarity), there is a similarity (distance) algorithm that gives the same set of optimal alignments. That is, finding similarity and distance alignments are dual problems.

Theorem 4 — Let a similarity measure be given with $s(a,b)$ and gap penalties \hat{g}_k and a distance measure be given with $d(a,b)$ and gap weight g_k . Assume there is a constant c , $0 \leq c \leq \max_{a,b} d(a',b')$ such that $s(a,b) = c - d(a,b)$ and $\hat{g}_k = g_k - (kc)/2$. Then an alignment is similarity optimal if and only if it is distance optimal.

Proof — For simplicity, the proof uses the single gap case ($g_k = \infty$ for $k \geq 2$). Now by elementary counting

$$n + m = 2\# \text{matches} + \# \text{gaps}$$

obviously holds. Using this simple equation,

$$\begin{aligned} D(a,b) &= \min \left\{ \sum_{\text{matches}} d(a,b) + g_k \# \text{gaps} \right\} \\ &= \min \left\{ \sum_{\text{matches}} c - \sum_{\text{matches}} s(a,b) + g_k \# \text{gaps} \right\} \\ &= \min \left\{ c(n + m)/2 - \sum_{\text{matches}} s(a,b) + (g_1 - c/2) \# \text{gaps} \right\} \\ &= c(n + m)/2 - \max_{\text{matches}} \left\{ \sum_{\text{matches}} s(a,b) - (g_1 - c/2) \# \text{gaps} \right\} \end{aligned}$$

Notice that

$$D(a,b) + S(a,b) = c(n + m)/2$$

so "large distance" is "small similarity". After seeing this equivalence, it is surprising that there are problems with a simple similarity algorithm for which no equivalent simple distance algorithm exists. This situation arises in Section D.

C. Fitting One Sequence into Another

Next the algorithms are modified to solve a new problem: the best fit of a "short" sequence into a "larger" sequence. An example of when this might be of interest is in locating a regulatory pattern in a nucleotide sequence, such as TATAAT in a bacterial promoter. The algorithm finds where the short pattern approximately appears in the longer sequence.

First consider the problem of fitting $a = a_1 a_2 \dots a_n$ into $b = b_1 b_2 \dots b_m$. For the purpose of visualizing the problem think of n as much smaller than m . (The relative sizes of n and m are irrelevant to the mathematics.) The problem is to find i and j ($i \leq j$) such that

$$S(a, b_{i+1} \dots b_j) = \max_{i,j} S(a, b_{i+1} \dots b_j)$$

This problem is simply solved. We just use the similarity algorithm of Theorem 1, Section III.B. with $S_{0,j} = 0$ for all j . Then the required j is found by

$$S(a, b_{i+1} \dots b_j) = \max_{1 \leq i \leq n, 1 \leq j \leq m} S(i,j)$$

and i is found by tracing back beginning at (n,j) (instead of n,m). The same procedure works with distance, as first shown by Sellers.^{18,19}

To illustrate this algorithm, we take as sequence b the *E. coli* promoter sequence of IacI.²⁰ In *E. coli* promoter sequences, the -10 signal or pattern TATAAT is well known to have functional significance. The designation -10 refers to the distance to the left (5') of the mRNA start. We take $a = \text{TATAAT}$. As above $s(a,a) = 1$, $s(a,b) = -1$ if $a \neq b$, and $\hat{g} = 2$. The matrix $(S(i,j))$ is shown in Table 4. Searching the last row of the matrix gives two solutions of $\max_{1 \leq i \leq n, 1 \leq j \leq m} S(i,j) = 2$, at $(6,13)$ and $(6,43)$. The pattern at $(6,43)$ has the alignment

TATAAT
CATGAT

and is CATGAT in the promoter sequence, the canonical -10 pattern. The pattern at $(6,13)$ has the alignment

TATAAT
TCGAAT

with TCGAAT in the promoter sequence, an equally good fit.

This illustrates the utility of the algorithm, in that it locates the putative -10 signal CATGAT in IacI. It also emphasizes the difficulty of promoter signal analysis by finding an equally good pattern TCGAAT 30 bases 5' of the -10 pattern.

D. Identification of Similar Segments

Surprising relationships have been discovered between sequences that overall have little similarity. See Weiss,²¹ Doolittle et al.²² and Naharro²³ for accounts of some unexpected long matching segments between viral and host DNA. The subject of this subsection is a dynamic programming algorithm to find these similar segments. This is probably the most useful dynamic programming algorithm for current problems. For a mathematical statement of the problem, it is necessary to assume a similarity function $s(a,b)$. The object is to find

$$\max_{1 \leq i \leq n, 1 \leq j \leq m} S(a, a_{i+1} \dots a_j, b_{i+1} \dots b_j)$$

This amounts to $(\frac{n}{2})(\frac{m}{2})$ sequence alignment problems, and a new algorithm must be devised.

While Sellers^{18,19} began the work on problems of this type, his problem formulations were based on distance functions and his algorithms involved forward and backward recursions, each recursion requiring a matrix. Although later²⁴ similarity functions are recommended by Sellers, those algorithms still involve intersection of path graphs and are quite complex. The similarity formulation given above was presented by Smith and Waterman and is solved in a straightforward way. Define $H_{i,j}$ to be the maximum similarity of two segments ending at a_i and b_j :

$$H_{i,j} = \max(0; S(a_{i-1} \dots a_x, b_{j-1} \dots b_y); \quad 1 \leq x \leq i, 1 \leq y \leq j)$$

A recursion similar to those given for the similarity problems discussed above is obtained for $H_{i,j}$.²⁵

Theorem 1 — Set $H_{i,0} = H_{0,j} = 0$ for $1 \leq i \leq n$ and $1 \leq j \leq m$. Then

	A	G	T	C	G	A	G	G	C	T	A	C	T	C	T	A	C	T	G	A	A	C		
T	0	0	0	10	10	0	0	0	10	0	10	10	10	0	10	0	0	0	0	0	0	0	10	
C	0	0	0	10	20	1	0	0	0	10	1	10	1	10	1	0	10	1	0	0	0	0	10	
A	10	0	0	0	1	11	11	0	0	0	0	1	11	0	1	0	1	11	0	1	0	10	10	
A	10	1	0	0	0	0	21	2	0	0	0	11	2	0	0	0	11	2	0	0	10	20	1	
T	0	1	11	0	0	0	1	12	0	0	0	10	0	2	12	0	10	0	2	12	0	0	1	
C	0	0	0	21	10	0	0	0	1	0	10	0	0	20	0	0	20	2	32	12	0	20	0	0
T	0	0	10	1	12	1	0	0	0	0	20	0	0	0	20	0	0	2	32	12	0	20	0	0
A	10	0	0	1	0	3	11	0	0	0	0	30	10	0	11	12	42	22	2	11	10	10	0	
C	0	1	0	10	11	0	0	2	0	0	0	10	40	20	10	2	22	52	32	12	2	1	20	
T	0	0	11	0	1	2	0	0	0	0	0	20	50	30	20	2	32	62	42	22	2	0	0	
A	10	0	0	2	0	0	12	0	0	0	0	30	10	30	41	21	30	12	42	53	52	32	12	
C	0	1	0	10	12	0	0	1	0	0	10	0	10	40	20	40	32	40	22	33	44	43	42	
T	0	0	11	0	1	3	0	0	0	0	0	20	0	20	50	30	50	30	20	50	10	24	35	34
G	0	10	0	2	0	11	0	10	10	10	0	11	0	30	41	30	41	21	30	40	40	20	24	
C	0	0	1	10	12	0	2	0	1	1	20	0	0	21	10	40	32	21	51	31	40	51	31	30
T	0	0	10	0	1	3	0	0	0	0	0	30	10	1	31	20	50	30	31	41	41	31	42	22
T	0	0	10	1	0	0	0	0	0	0	0	10	21	1	11	22	30	41	21	41	52	32	22	31
G	0	10	0	1	0	10	0	10	10	0	0	1	12	0	2	11	21	32	21	51	43	23	13	
A	10	0	0	0	1	2	10	0	0	0	0	11	10	0	2	0	1	10	11	22	14	41	52	32
G	0	20	0	0	11	0	20	10	10	0	0	2	1	0	0	0	0	1	2	32	21	32	43	
T	0	0	30	10	0	0	2	0	11	1	1	10	0	0	11	0	10	0	11	12	23	12	23	
A	10	0	10	21	1	0	10	0	0	2	0	20	0	0	2	0	20	0	0	2	22	33	13	
C	0	1	0	20	11	11	0	1	0	0	12	0	0	30	10	10	0	0	30	10	0	2	13	43

A

	A	G	T	C	G	A	G	G	C	T	A	C	T	C	T	A	C	T	G	A	A	C				
T	0	0	0	10	10	0	0	0	0	0	0	0	10	0	10	0	0	0	0	0	0	0	10			
C	0	0	0	10	20	1	0	0	0	0	0	0	10	1	10	1	0	10	1	0	0	0	10			
A	10	0	0	0	1	11	11	0	0	0	0	1	0	1	11	0	1	0	1	0	10	10	0			
A	10	1	0	0	0	0	21	2	0	0	0	11	0	0	0	11	2	0	0	10	20	1	0			
T	0	1	11	0	0	0	1	12	0	0	0	10	0	2	0	0	10	0	2	12	0	0	1			
C	0	0	0	21	10	0	0	0	1	0	10	0	1	10	0	0	0	0	1	10	0	0	1			
T	0	0	10	1	12	1	0	0	0	0	20	0	0	0	0	0	0	0	0	20	0	0	0			
A	10	0	0	1	0	3	11	0	0	0	0	0	30	10	0	11	0	0	0	11	10	10	0			
C	0	1	0	10	11	0	0	2	0	0	0	10	40	20	10	2	0	0	0	0	2	1	20			
T	0	0	11	0	1	2	0	0	0	0	0	20	0	20	50	30	20	0	0	0	0	0	0			
A	10	0	0	2	0	0	12	0	0	0	0	30	10	30	41	21	30	10	0	0	0	0	0			
C	0	1	0	10	12	0	0	1	0	0	10	0	10	40	20	40	32	12	40	20	0	0	1	20		
T	0	0	11	0	1	3	0	0	0	0	0	20	0	20	50	30	50	30	20	50	10	10	0	0		
G	0	10	0	2	0	11	0	10	10	10	0	11	0	30	41	30	41	21	30	40	40	20	0	0		
C	0	0	1	10	12	0	2	0	1	1	20	0	0	21	10	40	32	21	51	31	40	51	31	30		
T	0	0	10	0	1	3	0	0	0	0	0	30	10	1	31	20	50	30	31	41	41	31	42	22		
T	0	0	10	1	0	0	0	0	0	0	10	21	1	11	22	30	41	21	41	52	32	22	31			
G	0	10	0	1	0	10	0	10	10	0	0	1	12	0	2	11	21	32	21	51	43	23	13			
C	0	0	1	10	11	0	1	0	1	1	20	0	0	11	10	0	2	0	1	10	11	22	14	41	52	32
A	10	0	0	0	1	2	10	0	0	0	0	11	10	0	2	0	1	10	11	22	14	41	52	32		
G	0	20	0	0	11	0	20	10	10	0	0	2	1	0	0	0	0	1	2	32	21	32	43			
T	0	0	30	10	0	0	2	0	11	1	1	10	0	0	11	0	10	0	11	12	23	12	23			
A	10	0	10	21	1	0	10	0	0	2	0	20	0	0	2	0	20	0	0	2	22	33	13			
C	0	1	0	20	11	11	0	1	0	0	12	0	0	30	10	10	0	0	30	10	0	2	13	43		

B

FIGURE 1. Maximum similarity segment analysis of two sequences. Here $s(a,a) = 1$, $s(a,b) = -1$ if $a \neq b$, and $g = 2$. In Figure 1A, the initial matrix is shown. Tracing back from the maximum entry 6.2 produces the alignment given in the text. In Figure 1B, the matrix is shown with recomputed entries with * to the right. Now the maximum entry is 6.1.

A computer program written in C for this algorithm is given in the Appendix. The program is also available on tape or disk from the author. To illustrate the algorithm, we set $s(a,a) = 1$, $s(a,b) = -0.9$, if $a \neq b$, and $g = 2$. Two sequences are compared and Figure 1A gives the matrix H ($\times 10$), where the best matching segments are

CCAATCTACT
CTACTCTACT

with score 6.2. The matrix N ($\times 10$) is shown in Figure 1B where the recomputed entries are shown with * to the right. For this step, the best matching segments are

CTACTACTGCT
CTACTCTACT

with score 6.1.

E. Near Optimal Alignments

The optimal alignments depend on the input sequences and the algorithm parameters. The weights assigned to mismatches and gaps are determined by experience. An effort is made to use biological data to infer meaningful values. Of course, in addition to assigning weights, there are sometimes unknown constraints on the sequences that cause the correct alignment to differ from the optimal alignment given by an algorithm. Hence, it is of some interest to produce all alignments with score (distance or similarity) within a specified distance of the optimum score. Recently, Waterman²⁷ and Byers and Waterman,²⁸ presented a new algorithm which accomplishes this. The algorithm has been previously presented for the distance algorithm. In this chapter, it is given for the similarity algorithm.

To be explicit, let $S = (S_{ij})$ be the single gap similarity matrix with

$$S_{ij} = \max\{S_{i-1,j-1} + s(a_i,b_j), S_{i-1,j} - g, S_{i,j-1} - g\}$$

The task is to find all alignments with score within $\epsilon > 0$ of the optimum value $S_{n,m}$. All optimum alignments are included.

At position (i,j) assume a traceback from (n,m) to $(0,0)$ is being performed that can result in an alignment with score greater than or equal to $S_{n,m} - \epsilon$. The score from (n,m) to be not including (i,j) is T_{ij} . T_{ij} is the sum of the possibly nonoptimal alignment weights to reach (i,j) . From (i,j) , as usual, three steps are possible: $(i-1,j)$, $(i-1,j-1)$, and $(i,j-1)$. Each step is in a desired alignment if and only if

$$T_{ij} - g + S_{i-1,j} \geq S_{n,m} - \epsilon$$

$$T_{ij} + s(a_i,b_j) + S_{i-1,j-1} \geq S_{n,m} - \epsilon$$

$$T_{ij} - g + S_{i,j-1} \geq S_{n,m} - \epsilon$$

respectively. Multiple near-optimal alignments can be produced by stacking unexplored directions. Of course, multiple insertions and deletions can be included.

A study of sequence alignment sensitivity to weights and multiple insertions or deletions has been carried out by Fitch and Smith.²⁹ The sequences displayed below are chick hemoglobin mRNA sequences, nucleotides 115-171 from the β chain (upper sequence) and 118-156 from the α chain (lower sequence):

UUUGCGUCCUUUGGGAACCCUCCAGCCCCACUCCAUCCUUGUCACACGGCAACCCCAUUGGUC
UUUCCCCACUUCG AUUUUUGUCACAC GGCCUCCGCUAAAUC

This alignment is presumed correct from the analysis of the many known amino acid sequences for which such RNA sequences code.

With a distance function, using a mismatch weight of 1 and a multiple insertion or deletion function $x_k = 2.5 + k$, where k is the length of the insertion or deletion, the correct alignment is found among the 14 optimal alignments. (This is region Q of the Fitch and Smith paper.) To indicate the size of neighborhoods in this example, there are 14 alignments within 0% of the optimum, 14 within 1%, 35 within 2%, 157 within 3%, 579 within 4% and 1317 within 5%.

A mismatch weight of 1 and a multiple insertion or deletion function $2.5 + 0.5k$ is region P of Fitch and Smith; accordingly, the correct alignment is not in the list of the 14

optimal alignments. This example illustrates the sensitivity of alignment to weighting functions.

IV. DYNAMIC PROGRAMMING ALIGNMENT OF MULTIPLE SEQUENCES

Many sequence alignment problems involve more than two sequences. The straightforward generalization of our two-sequence alignment algorithm to R sequences has complexity to $O(2^n n^R)$ when comparing length n sequences. For $R = 3$ this is frequently not practical and for $R > 3$, n must be very small. In Section IV.B these algorithms are discussed with some recent improvements. In the next section, IV.A, a geometrical approach, also based on dynamic programming, to these problems is presented.

A. Line Geometries

For this section a new, simple representation of sequences is required. It arises from consideration of a column in an alignment of several sequences. For example, the column

```
seq1 ... A ...
seq2 ... A ...
seq3 ... T ...
seq4 ... A ...
seq5 ...  $\phi$  ...
```

is frequently summarized as the "consensus" letter A. Much information is lost in this summary, and our representation of the column is with the "letter" $a = (p_A, p_C, p_G, p_T, p_\phi)$. For example, p_A represents the proportion of A in the given column. In this example $p = (3/5, 0, 0, 1/5, 1/5) = (0.6, 0, 0, 0.2, 0.2)$. Of course, usual sequences can be represented in this format by using $C = (0, 1, 0, 0, 0)$ etc. To compare two "letters" $a = (p_A, p_C, p_G, p_T, p_\phi)$ and $b = (q_A, q_C, q_G, q_T, q_\phi)$, the well-known metric

$$d(a,b) = \left(\sum_i w_i |p_i - q_i|^\alpha \right)^{1/\alpha}$$

is used, where w_i are weighting factors and $\alpha \geq 1$ is a constant. We have found in examples that w_ϕ needs to be larger than the other w_i ; certainly ϕ plays a quite different role from A, C, G, or T.

The distance $D(a,b)$ or similarly $S(a,b)$ between two sequences $a = a_1 a_2 \dots a_n$ and $b = b_1 b_2 \dots b_n$ is computed by dynamic programming as above even though the sequence elements are themselves vectors. Associated with $D(a,b)$ is an optimal alignment.

```
a_1^* a_2^* ... a_n^*
b_1^* b_2^* ... b_n^*
```

Now, for this alignment, define $c(\lambda) = \lambda a \oplus (1-\lambda)b$ by $c_i(\lambda) = \lambda a_i^* + (1-\lambda)b_i^*$. The last " \oplus " is simple vector addition. It is not surprising that $c(1/2)$ is midway between a and b . In fact, much more is true.

Theorem 1 — Let $c(\lambda) = \lambda a \oplus (1-\lambda)b$ where $0 \leq \lambda \leq 1$. Then

$$D(a,b) = D(a,c(\lambda)) + D(c(\lambda),b)$$

$$D(a,c(\lambda)) = (1-\lambda)D(a,b),$$

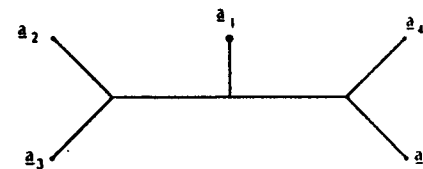


FIGURE 2. Tree relating a_1 (*E. coli*), a_2 (*B. stearo*), a_3 (*H. volcanii*), a_4 (*D. discoideum*), and a_5 (*S. cerevisiae*)

and

$$D(c(\lambda_1), c(\lambda_2)) = |\lambda_1 - \lambda_2| D(a,b)$$

This theorem is proved in Waterman and Perlwitz.¹⁰

This technique was devised to align several sequences, but it does not always behave well on sets of sequences. Methods of finding "center of gravity" sequences were studied and did not always converge rapidly. However, if the sequences are related by a given phylogenetic tree, useful alignments can be produced by use of " \oplus ".

As an example of this alignment algorithm, consider the tree of Figure 2 with the associated 5 sequences. Alignment is performed by using the tree to suggest the order of sequence alignment:

$$b = \frac{1}{5} a, \oplus \frac{4}{5} \left[\frac{1}{2} \left(\frac{1}{2} a_1 \oplus \frac{1}{2} a_2 \right) \oplus \frac{1}{2} \left(\frac{1}{2} a_3 \oplus \frac{1}{2} a_4 \right) \right]$$

The deletion term needs a heavier penalty than A, C, G, T:

$$d(a,b) = \sum_{i \in \{A,C,G,T\}} |p_i - q_i| + 4|p_\phi - q_\phi|$$

The resulting b is used to obtain an overall alignment by aligning each of a_1, \dots, a_5 with b . The result given below is identical with Woese et al.,¹¹ who obtain it in Table 31 of their analysis of 16S-like rRNAs.

<i>E. coli</i> (a_1)	CAACCCUUAUCC ϕ UUUGUAACUCAAGGAGGAAUGUUGGGA
<i>B. steato</i> (a_2)	CAACCCUCGCCU ϕ CUAGUCACUCUAGAGGGGAAGGUGGGGA
<i>H. volcanii</i> (a_3)	AGACCCGCACUU ϕ CUAAUUACAUUAGAAGGGGAAGGAACGGG
<i>D. discoideum</i> (a_4)	AGACCCUGACCCUGCUAACCUUCUUAGAGGGGAAGUCCGAGG
<i>S. cerevisiae</i> (a_5)	AGACCUAAACCUAUAACUUCUUAGAGGGGAAGUUGAGG

B. Generalization of the Two-Sequence Algorithm

For a straightforward extension of the two-sequence algorithm to more sequences, some notation is first set. There are R sequences, $a = a_1 a_2 \dots a_n$, $b = b_1 b_2 \dots b_n, \dots, r = r_1 r_2 \dots r_n$. The algorithm comes from considering an alignment of $a_1, \dots, a_n; b_1, \dots, b_n; r_1, \dots, r_n$. The last column of the alignment will appear as

```
 $\epsilon_1 a,$ 
 $\epsilon_2 b,$ 
 $\dots$ 
 $\epsilon_1 r,$ 
```

where $\epsilon_i = 0$ or 1 and $0 \cdot a = \phi$. If $\epsilon = (\epsilon_1, \dots, \epsilon_n) \neq 0$, it is required that there are $2^n - 1$ such columns. Recall that for $R = 2$ there are $2^2 - 1 = 3$ terms to optimize. Here, if a "distance" function is defined on R variables, the formula extends to

$$D_{i,j} = \min_{\epsilon \neq 0} \{D_{i,j}(\epsilon_1, \dots, \epsilon_n) + d(\epsilon_1 a_1, \epsilon_2 b_1, \dots, \epsilon_n r_1)\}$$

This formula appeared in Waterman et al.,¹³ and has computational complexity proportional to $O(2^n n^2)$ where the sequences are assumed to have length n . Storage takes $O(n^2)$. Carrillo and Lipman¹⁴ have improved the complexity of these algorithms with a branch and bound technique.

A deeper analysis by Sankoff¹⁵ appeared earlier. Sankoff assumes that the sequences are related by a given phylogenetic tree. His algorithm constructs sequences for each interior node of the tree and produces an alignment of the R sequences at the exterior nodes along with the (N) sequences (X_1, X_2, \dots, X_N) constructed for the interior nodes. Sankoff and Cedergren¹⁶ give an excellent account of this procedure. The dynamic programming step is given by

$$D_{i,j} = \min_{\epsilon \neq 0} \{D_{i,j}(\epsilon_1 a_1, \epsilon_2 a_2, \dots, \epsilon_n a_n) + \min_{X_1, \dots, X_N} \epsilon_n r_n\} \quad (55)$$

The second minimum is to indicate that X_1, \dots, X_N have been chosen for the interior nodes in such a way as to minimize number of mutations along the tree. A generalization of Fitch's parsimony method¹⁷ is used here and takes N steps. The computation requires $O(2^n n^2 N)$ steps. $R = 3$ is nearly the largest practical R and Sankoff has devised an iterative method, which works with groups of three sequences, to build up a solution for larger problems. Altschul and Lipman¹⁸ have extended Carrillo and Lipman¹⁴ to align multiple sequences with a given tree.

V. CONSENSUS ALIGNMENT OF MULTIPLE SEQUENCES

As Section IV illustrates, dynamic programming methods have generally not been found practical for more than two sequences. It is natural to ask whether methods for finding consensus patterns can be applied to find alignments, which are in a real sense consensus patterns themselves. Usually, it will be reasonable to limit the amount of shifting one sequence can have, relative to the others. That is, the alignment

```

sequence 1 ... a,
sequence 2 ... b,
...
sequence r ... r,
    
```

is only allowed when $|i - j|, \dots, |j - x|, |j - x|, \dots$ are all less than, or equal to, some bound. It is feasible to set the bound equal to sequence length and, therefore, to allow unrestricted shifting. The object of this section is to use the algorithm for consensus words

to make a useful, practical, alignment algorithm (see the chapter on consensus patterns⁷). This algorithm is presented in a recent paper.¹⁹

To conform with the conventions of the previous chapter, the object is to find consensus words w , k letters in length, where the window width is W . This means the maximum shift between matched words is $W - k$. The bound on $|i - j|, \dots, |i - x|, |j - x|, \dots$ is thus $W - k$. The window at position i appears as:

	Position	
	$i + 1 \dots i + W$	
sequence 1	...	$a_{i+1}, a_{i+2}, \dots, a_i, W$
sequence 2	...	$b_{i+1}, b_{i+2}, \dots, b_i, W$
...		
sequence R	...	$r_{i+1}, r_{i+2}, \dots, r_i, W$
		window width = W

A neighborhood of, say, less than, or equal to, 2 mismatches is specified for matching words, and a score $s_w(v)$ is given to a word v in the neighborhood of w , where $0 \leq s_w(v) \leq 1$. Define the maximum scoring word in the window of sequence i to be v_i so that v_i best matches w . Then the score of w in the sequence set is

$$S(w) = \sum_{i=1}^R \max_v \{s_w(v)\} = \sum_{i=1}^R s_w(v_i)$$

A consensus word w^* is one satisfying

$$S(w^*) = \max_w S(w)$$

where w ranges over all k -letter words.

An illustration of the location of such a consensus word is given next, where only the window is displayed.

	Position
	i
sequence 1 v_1
sequence 2
sequence 3	v_3
...
sequence R v_R
	W

The words v_i are in the neighborhood of w . Here, sequence 2 fails to have a word v_2 in the neighborhood of w .

Define a partial order on consensus words by their location in the sequences as follows: $w^{(1)} < w^{(2)}$ if the occurrence of $w^{(1)}$ in sequence i is to the left of (and not overlapping) the occurrence of $w^{(2)}$ in sequence i , for all $i = 1, 2, \dots, R$. If $v_i^{(1)}, v_i^{(2)}$ mark the individual sequence patterns, respectively, in sequence i , then the order appears as

..... $v_1^{(1)}$ $v_1^{(2)}$
 $v_2^{(1)}$ $v_2^{(2)}$
 $v_3^{(1)}$
 $v_n^{(1)}$ $v_n^{(2)}$

An optimal alignment A , is one that satisfies

$$S(A) = \max_{A'} S(A') = \max_{A'} \sum_{w_1, \dots, w_j} S(w_i)$$

In general, it is not known how to find $S(A)$. However, one version of the problem can be solved exactly.

Let $w_i | w_j$ mean that consensus words w_i and w_j can be found in nonoverlapping windows:

..... $v_1^{(1)}$ $v_1^{(2)}$
 $v_2^{(1)}$
 $v_n^{(1)}$ $v_n^{(2)}$

The modified optimization problem is to find the alignment A , satisfying

$$R(A) = \max_{A'} R(A') = \max_{A'} \sum_{w_i | w_j} R(w_i)$$

A simple recursion for R is easily given which is linear in sequence length. Define R_i to be the value of R for the set of sequences ending with column i . Then

$$R_i = \max\{R_{j-1} + S(w_j); i = W + 1 \leq j \leq i\}$$

where w_j is the consensus word with the window from column j to column i .

If much shifting is taking place in the optimal alignment A associated with $S(A)$, then the alignment A^* in $R(A^*)$ is unsatisfactory and overlapping windows must be considered. To overcome this difficulty in a practical way, find

$$\hat{T}_i = \max\{\hat{T}_{j-1} + \hat{S}(w_j); i - W + 1 \leq j \leq i\}$$

where $\hat{S}(w_j)$ is the largest score for a consensus word in the window of width W , such that all occurrences of w_j are to the right of the matches in \hat{T}_{j-1} . While this algorithm is not guaranteed to find $S(A)$, it is much more useful than $R(A)$ for most problems. In addition, it is practical for the alignment of many sequences, while dynamic programming is not.

Another hierarchical algorithm is to find the maximum scoring word w in the entire sequence, $\max_{w_1, \dots, w_n} S(w') = S(w)$, and to realign the sequences on w . This breaks the set into two shorter blocks. The same procedure can be reapplied, recursively, to produce an alignment.

VI. DATA BASE SEARCHES

In summer 1988, there were approximately 20×10^6 bp of data in GenBank, a nucleic acid data base, containing approximately the same data which is in the EMBL data base in Heidelberg. The data represent nearly 5600 nucleotide sequences, at an average of approximately 900 to 1000 bp per sequence (see Chapter 1). When a new DNA sequence is determined, there are usually some sequence comparisons suggested by the nature of the sequence. If the new sequence is a coding region for a rat immunoglobulin (Ig), running the

new sequence against all known mammalian Ig sequences is a natural thing to do. However, in addition to these sorts of comparisons, there is increasing interest in seeing if any unexpected relationships show up when the sequence is run against all mammalian DNA (all eukaryotic DNA. The search need not stop even there. Ribosomal RNA structures have been preserved across all life forms.³¹ There seems to be some homology between ribosomal proteins from rat and *E. coli*³²). Therefore, it is not unreasonable to just run the new sequence against the entire data base and see if any unexpected matches show up.

If the search for a match between a new sequence of 10^3 bp and a data base of 6×10^7 bp were performed with the dynamic programming algorithm of Section III.D, then the running time is proportional to $6 \times 10^9 = (10^3)(6 \times 10^6)$. On a VAX 11/780, this might take a day or more of CPU time. With increasing computer capacity, this should not be too disturbing, but it indicates that such a venture is not to be lightly undertaken.

One response to this important problem has been to utilize the computer science technique of hashing to obtain information about possible regions of high matching. The first use of hashing in molecular biology was made by Dumas and Ninio.³³ The method was later utilized by Wilbur and Lipman³⁴ and Karlin et al.³⁵ in developing their own approaches to sequence comparisons.

The basic method begins with a choice of word size, say k letter words. Then a DNA sequence $a = a_1 a_2 \dots a_n$ is transformed to a sequence of integers by associating each triple, $a_i a_{i+1} \dots a_{i+k-1}$, with one of the integers $0, 1, \dots, 4^k - 1$. For example, let

$$\hat{a} = \begin{cases} 0, & \text{if } a = A \\ 1, & \text{if } a = C \\ 2, & \text{if } a = G \\ 3, & \text{if } a = T \end{cases}$$

Then define

$$x_i = \hat{a}_i \cdot 4^{i-1} + \hat{a}_{i+1} \cdot 4^{i-2} + \dots + \hat{a}_{i+k-1} \cdot 4^1 + \hat{a}_{i+k} \cdot 4^0$$

Notice that $x_{i+1} = x_i \cdot 4 + \hat{a}_{i+1} \pmod{4^k}$, so that these transformations can be rapidly made. To continue with a numerical example, let $a = \text{TAGAGCA}$. With $k = 2$, $4^k = 16$ and

$$\begin{aligned} x_1 &= 3 \cdot 4 + 0 = 12 \\ x_2 &= 12 \cdot 4 + 2 \pmod{16} = 2 \\ x_3 &= 2 \cdot 4 + 0 \pmod{16} = 8 \\ x_4 &= 8 \cdot 4 + 2 \pmod{16} = 2 \\ x_5 &= 2 \cdot 4 + 1 \pmod{16} = 9 \\ x_6 &= 9 \cdot 4 + 0 \pmod{16} = 4 \end{aligned}$$

and $a = (12)(2)(8)(2)(9)(4)$.

Next, lists are made of locations of all occurrences of the integers $0, 1, \dots, 4^k - 1$. In our numerical example, there are only five nonempty lists

0 : ϕ
 1 : ϕ

2: 2,4
 3: ϕ
 4: 6
 ...: ϕ
 8: 3
 9: 5
 ...: ϕ
 12: 1
 ...: ϕ

Another approach is to keep a list of positions of the first occurrences of 0, 1, ..., $4^k - 1$. In place of the x_i at those locations, are pointers to the next occurrences of the associated integer. In either case, storage is $O(n)$ and the hashing can be performed in time $O(n)$.

Next, two methods, both due to Wilbur and Lipman, which allow use of this information in sequence comparison.

A. A Regions Method

The algorithm described here appears in Wilbur and Lipman. Define a region r by $(v; i, j)$ where v is a word of length k which begins at position i in a and position j in b .

Define $r_1 = (v; i_1, j_1) < r_2 = (v; i_2, j_2)$ if $i_1 + k - 1 < i_2$ and $j_1 + k - 1 < j_2$. Set $r_0 = (\phi; 0, 0)$ as a least element and $r^* = (\phi; n, m)$ as a greatest element. $\Gamma = (r_1, r_2, \dots, r_k)$ is a path if $p < q$ implies $r_p < r_q$. The score of path Γ is given by

$$\text{score}(\Gamma) = \sum_{i=1}^k s(r_i) - \sum_{i=1}^{k-1} g(i_{i+1} - |w_i| - i_i - 1, j_{i+1} - |w_i| - j_i - 1)$$

where $s(\cdot)$ is a similarity score for region r_i , like $|w_i|$, and $g(\cdot, \cdot)$ is a gap penalty. Then

$$\text{score}(a, b) = \max\{\text{score}(\Gamma) : \Gamma \text{ is a path from } r_0 \text{ to } r^*\}.$$

The algorithm makes two lists of regions: L_1 , ordered by $<$ and L_2 , ordered by the usual order $<<$ of best scores from r_0 to the region listed in L_1 .

$$\gamma = \text{score}(r_q) - g(i_q - |w_q| - i_{q-1} - 1, j_q - |w_q| - j_{q-1} - 1) + s(r_q) + \text{score}(r_q)$$

If there is no region r_q in L_2 below r_q under $<$ with a score greater than $\text{score}(r_q) - s(r_q)$, or if the inequality cannot be satisfied, go to (C).

- (B) Set $\text{score } r_q = \gamma$ and go to (A).
 (C) Remove r_q from L_1 and insert it in L_2 under $<<$.
 If $L_1 \neq \phi$, go to (A).

Changes can be made to this algorithm to give maximum scoring segments. A more practical algorithm is given next.

B. Rapid Similarity Searches

The method described here is due to Wilbur and Lipman¹⁹ and was extended, especially for protein data base searches, by Lipman and Pearson.⁴¹ Recently, Pearson and Lipman⁴²

have modified their earlier work to DNA sequences. Each of the 4^k words w_i of length k , can be located in each sequence, and the positions of the word w_i in a can be considered to match the b positions of w_i . Each such match, say position i in a and position j in b , has an associated offset $i - j$. The offsets or diagonals with a large number of matches are candidates for good matching regions. This technique is not so rigorous as the dynamic programming methods but it is a great deal faster. In fact, searches of the entire DNA data base can be accomplished on an IBM PC and several groups have programmed this algorithm.

VII. THE STATISTICAL DISTRIBUTION OF ALIGNMENT SCORES

Much has been written about the statistical distribution of distance and similarity scores. In this section, we focus attention on the statistical distribution of the scores of maximum similarity segments (see Section III.D), when the scores are computed for random sequences. The idea is that when the sequences satisfy some model of randomness, such as uniform and independent bases, there is a resulting distribution of maximum similarity segment scores. The scientist can use this distribution to ascertain whether the scores from real, biological sequences are, in the statistical sense, significantly larger than those from random sequences. Of course, there is little agreement about the appropriate model of randomness. Fortunately, as we discuss in Section VII.B, the distribution of maximum similarity segment scores for real, unrelated biological sequences coincides with that of independent, identically distributed sequences of the same composition.

Although simulations are often recommended to determine statistical significance, there are several drawbacks to this approach. First of all, it is expensive in terms of computer time. Also, it often requires more time from the scientist to set up and process the simulations. Frequently, the simulations fail to give the desired results. If statistical significance of a real sequence matching is to be estimated and that significance is α , then $1/\alpha$ simulations must, on the average, be run before seeing a result as extreme as in the real sequences. When $\alpha = 0.0001$, a not unreasonable case, we must do on the order of 10,000 runs. While simulations are often unavoidable, they are not ideal.

Recently, however, some new results in probability theory give very precise answers to certain problems that arise in practice. These are discussed in the subsections on the "log(n) law". What is less well-understood, is that the log(n) law and another equally special result give guidance on what to expect for the distribution of any maximum similarity algorithm. There are only two behaviors of expected similarity score with sequence length: either proportional to sequence length or to logarithm of sequence length. This result, along with a table of means and variances, is given in Section VII.C on expected behavior.

A. The Log(n) Law

Erdos and Renyi⁴³ proved in 1970 that the length R_n of the longest run of heads in independent coin tosses with $P(\text{Heads}) = p$ grows like $\log_{1/p}(n)$:

$$P\left(\lim_{n \rightarrow \infty} \frac{R_n}{\log_{1/p}(n)} = 1\right) = 1$$

If our sequences are perfectly aligned

$$\begin{array}{l} a_1 a_2 \dots a_n \\ b_1 b_2 \dots b_n \end{array}$$

this mathematical result provides the answer to questions about longest match. Simply write H if $a_i = b_i$, T if $a_i \neq b_i$, and calculate

$$p = p_a^2 + p_{kc}^2 + p_0^2 + p_t^2$$

The length of the longest match is the length of the longest head run. Notice that we assume the distributions of the sequences are identical. The results stated here hold if the sequences have base distributions that are "not too different".

When the sequences are allowed to shift relative to one another, the answer changes. Now matches can include offsets such as

$$\begin{array}{cccccccc} a_1 & a_2 & \dots & a_n & \dots & a_n & \dots & a_n \\ b_1 & b_2 & \dots & b_j & b_{j+1} & \dots & b_n & \end{array}$$

It has been shown³⁾ that shifting doubles the growth of the longest match length M_n :

$$P\left(\lim_{n \rightarrow \infty} \frac{M_n}{\log_{1/p}(n)} = 2\right) = 1$$

Much more precise results have been obtained. With $a = a_1 a_2 \dots a_n$ and $b = b_1 b_2 \dots b_m$, it is required below that $\log(m)/\log(n) \rightarrow 1$. The expectation (mean) of M , the longest match length including k mismatches, is approximately

$$E(M) \approx \log(qmn) + k \log(qmn) + k \log(q/p) - \log(k!) + \gamma \log(e) - 1/2$$

where $q = 1 - p$, $\log = \log_{1/p}$, and $\gamma = 0.577 \dots$ is the Euler-Mascheroni constant. The variance is approximately

$$\text{Var}[M(n,m)] \approx [\pi \log(e)]^{2/6} + 1/12$$

In the case of $k = 0$,

$$E(M) \approx \log(qmn) + \gamma \log(e) - 1/2$$

While it might seem natural to perform the normal approximation at this point, that is not correct! The tail behavior of the distribution is extreme value which has exponential tails. The correct procedure can be found in Theorem 2 of Arratia et al.⁴ Results closely related to these were first announced for repeats in a single sequence, with $k = 0$, and minor differences in constants, by Karlin et al.¹ The more general case stated here is from Arratia et al.⁴ As discussed in Chapter 6, these formulas can be extended to study the longest match common to R of N sequences. Instead, here we pursue generalizations that allow less exact matching.

B. A Data Base Study

In Smith et al.,⁴⁾ the following similarity and gap functions are used:

$$\begin{aligned} s(a,a) &= 1 \\ s(a,b) &= -0.9 \quad \text{if } a \neq b \\ w(1) &= 2 \\ w(k) &= \infty \quad \text{if } k \geq 2 \end{aligned}$$

to study maximum similarity segment scores from real, biological sequences. While the $\log(n)$ law appears useful for a fixed number of mismatches, what is the distribution of

$$S(a,b) = \max_{\substack{I \subseteq a \\ J \subseteq b}} S(I,J)$$

for a similarity function such as that specified above? Here, I and J denote contiguous segments of sequence.

In the study cited above, all pairwise comparisons were made with a set of 204 vertebrate and eukaryotic DNA sequences and complements for a total of $\binom{204}{2} = 20,706$ comparisons. The results display remarkable linear behavior with $\log_{1/p}(nm)$. The best fit of the data, after adjusting for outliers by techniques of robust statistics, is

$$S = 2.5 \log_{1/p}(nm) - 8.99$$

This fit of real sequence data is identical to the fit of simulated sequences reported below.

Therefore, allowing insertions and deletions does not change the $\log(nm)$ growth, although the slope is certainly changed. Note that $\log(nm) = \log(n^2) = 2 \log n$ if $n = m$.

C. Linear and Logarithmic Behavior

Let μ be the non-negative mismatch penalty and δ the non-negative deletion penalty, the maximum segment similarity; measure S and write $S = S(\mu, \delta)$ to indicate this dependence. The goal of this section is to give information about the probability distribution $S(\mu, \delta)$ for all $\mu \geq 0$, $\delta \geq 0$. It was only recently that anything was learned about these questions and much work remains for mathematicians to complete the theory.⁵⁾ Still, the broad outlines have been established and the results are both theoretically and practically of interest.

Now for two random sequences of length n ,

$$S(\infty, \infty) \approx 2 \log(n)$$

by the above discussion. In contrast, if no penalty is associated with unmatched bases, it is known³⁾ that

$$S(0,0) \approx c \cdot n$$

where c is a constant. While the precise distribution of $S(\infty, \infty)$ is known, even the value $c = c(0,0)$ has eluded probabilists for more than 10 years. It has recently been shown that this contrasting behavior between linear and logarithmic growth is general. There is a 2-dimensional curve through $[0, \infty]^2$, such that the behavior is linear on one side of the curve and logarithmic on the other.⁶⁾ Let $R(0)$ be the region containing $(0,0)$ and $R(\infty)$ the region containing (∞, ∞) . Then the theorem is that

$$S(\mu, \delta) \approx c(\mu, \delta) \cdot n \quad \text{if } (\mu, \delta) \in R(0)$$

and

$$S(\mu, \delta) \approx d(\mu, \delta) \cdot \log(n) \quad \text{if } (\mu, \delta) \in R(\infty)$$

Table 5 gives means and variances for a large simulation. In this simulation, 28 independent pairs of sequences were generated of lengths $2^1, 2^2, \dots, 2^{12}$ where $p_a = p_c = p_t = p_g = 1/4$. At each (μ, δ) of Table 5 and for each of the $28 \binom{6}{2} = 168$ pairs of sequences, $S(\mu, \delta)$ was calculated. In the region $R(0)$, (which contains all of column $\mu = 0.50$), a linear fit was made to the data for each pair (μ, δ) . In Table 5, the triple

Table 5
SLOPES, INTERCEPTS, AND STANDARD DEVIATIONS FOR $S(\mu, \delta)$ IN THE
LINEAR, $R(0)$, AND LOGARITHMIC, $R(\infty)$, REGIONS

δ	0.25	0.50	0.75	1.00	1.25	1.50	2.00	2.50	3.00	3.50	4.00	4.50	5.00
0.25	-2.86 0.51 5.14	-2.82 0.48 5.73	-2.82 0.48 5.73	-2.82 0.48 5.73	-2.82 0.48 5.73	-2.82 0.48 5.73	-2.82 0.48 5.73	-2.82 0.48 5.73	-2.82 0.48 5.73	-2.82 0.48 5.73	-2.82 0.48 5.73	-2.82 0.48 5.73	-2.82 0.48 5.73
0.50	-2.05 0.42 5.64	-2.03 0.38 6.78	-1.73 0.34 6.78	-1.18 0.30 7.25	-1.18 0.30 7.25	-1.18 0.30 7.25	-1.18 0.30 7.25	-1.18 0.30 7.25	-1.18 0.30 7.25	-1.18 0.30 7.25	-1.18 0.30 7.25	-1.18 0.30 7.25	-1.18 0.30 7.25
0.75	-0.93 0.35 5.60	-0.47 0.29 6.24	0.03 0.24 6.66	1.01 0.20 7.29	2.47 0.16 7.82	4.53 0.13 8.63	4.53 0.13 8.63	4.53 0.13 8.63	4.53 0.13 8.63	4.53 0.13 8.63	4.53 0.13 8.63	4.53 0.13 8.63	4.53 0.13 8.63
1.00	-0.19 0.31 5.63	1.54 0.22 6.07	2.53 0.16 6.65	4.49 0.11 6.99	7.79 0.07 7.08	11.88 0.03 7.07							
1.25	0.77 0.28 5.79	2.92 0.18 6.10	6.45 0.10 6.53	9.98 0.04 6.63			-4.13 1.85 1.63	-3.08 1.64 1.49	-3.08 1.64 1.49	-3.08 1.64 1.49	-3.08 1.64 1.49	-3.08 1.64 1.49	-3.08 1.64 1.49
1.50	1.60 0.25 5.83	4.13 0.15 5.63	8.95 0.06 6.13		-6.48 2.21 1.75	-3.93 1.76 1.41	-2.34 1.46 1.33	-2.01 1.38 1.22	-2.03 1.36 1.20	-2.03 1.36 1.20	-2.03 1.36 1.20	-2.03 1.36 1.20	-2.03 1.36 1.20
2.00	2.75 0.21 5.19	6.70 0.09 5.99		-5.09 1.95 1.71	-2.14 1.44 1.35	-1.40 1.27 1.27	-1.60 1.23 1.20	-1.50 1.20 1.14	-1.42 1.18 1.12	-1.36 1.17 1.10	-1.29 1.16 1.10	-1.29 1.16 1.10	-1.29 1.16 1.10
2.50	3.98 0.19 5.25	9.52 0.06 5.66	-9.10 2.69 2.24	-2.94 1.58 1.44	-1.52 1.30 1.26	-1.19 1.20 1.22	-1.23 1.15 1.14	-1.21 1.14 1.09	-1.22 1.13 1.05	-1.18 1.12 1.04	-1.15 1.12 1.03	-1.14 1.11 1.02	-1.13 1.11 1.01
3.00	5.11 0.17 5.08	12.02 0.04 5.32	-6.20 2.21 2.00	-2.45 1.47 1.37	-1.46 1.26 1.23	-1.24 1.18 1.20	-1.11 1.12 1.14	-1.08 1.11 1.08	-1.05 1.09 1.05	-1.02 1.09 1.03	-1.00 1.08 1.02	-0.99 1.08 1.00	-0.97 1.08 1.00
3.50	6.23 0.16 4.82	14.25 0.02 5.17	-5.40 2.07 1.93	-2.00 1.40 1.36	-1.31 1.23 1.24	-1.21 1.17 1.20	-1.06 1.11 1.13	-1.02 1.10 1.06	-0.98 1.08 1.03	-0.94 1.07 1.01	-0.91 1.07 1.00	-0.90 1.07 1.00	-0.89 1.07 1.00
4.00	7.25 0.15 4.96		-5.28 2.03 1.91	-1.80 1.37 1.37	-1.24 1.22 1.23	-1.20 1.17 1.19	-1.05 1.11 1.12	-1.01 1.09 1.05	-0.98 1.08 1.02	-0.94 1.07 1.01	-0.91 1.07 1.00	-0.90 1.06 1.00	-0.89 1.06 1.00
4.50	8.06 0.14 5.20		-5.23 2.01 1.89	-1.74 1.36 1.31	-1.21 1.21 1.22	-1.19 1.16 1.17	-1.04 1.11 1.10	-0.99 1.09 1.04	-0.95 1.07 1.02	-0.92 1.07 1.00	-0.89 1.06 1.00	-0.89 1.06 1.00	-0.89 1.06 1.00
5.00	8.65 0.13 5.44		-5.26 2.01 1.87	-1.71 1.35 1.31	1.19 1.21 1.21	1.18 1.16 1.16	-1.02 1.10 1.09	0.97 1.07 1.03	0.92 1.07 1.02	0.90 1.07 1.01	-0.88 1.06 1.01	-0.88 1.06 1.00	-0.88 1.06 1.00

a
b
s

in region $R(0)$ means that

$$S(\mu, \delta) = a + b\mu$$

with standard deviation s ; while for (μ, δ) in region $R(\infty)$, the triple means that

$$S(\mu, \delta) = a + b \log_2(n)$$

with standard deviation s .

These means and variances can be used for estimates of statistical significance. Above, it was noted that $S(\infty, \infty)$ has exponential tails, not the tails of the normal distribution. We conjecture that all of region $R(\infty)$ has exponential tails, but that is by no means established mathematically. Moreover, convergence to extreme value distributions is usually very slow, so that except for $S(\infty, \infty)$, it is unclear what constants to use. To further complicate the situation, in $R(0)$, knowledge beyond linear behavior is completely lacking. Still, it is possible to use the present state of knowledge to obtain useful information.

For illustration, consider *E. coli* threonine tRNA and *E. coli* valine tRNA. Both sequences are 76 in length. To be consistent with the analysis of these sequences in Section III.A, take $\mu = 1$ and $\delta = 2$. A computer sequence comparison yields $S(1, 2) = 24 - (17) - 2(1) = 15$. This is the score of the longest boxed portion of the alignment in III.A. The pair (1, 2) is in $R(\infty)$ and Table 5 gives $a = -5.09$, $b = 1.95$, and $s = 1.71$. Therefore,

$$a + b \log_2(n) = -5.09 + 1.95 \log_2(78) = 7.0935$$

and

$$S(1, 2) = 15 = 7.1665 + 4.624(1.71)$$

so that $S(1, 2)$ is 4.624 standard deviations above the mean.

If additional information is required about the significance, Chebyshev's inequality can be used, as in Smith et al.⁴⁵ This conservative result, valid for all probability distributions, is that the probability of $|S(\mu, \delta) - \text{mean}| \geq \lambda$ is less than, or equal to, $(s/\lambda)^2$. In our example, $\lambda = (4.581)s$ and the statistical significance of our result is therefore no larger than $(1/4.581)^2 = 0.0477$.

APPENDIX

```

/* code for finding segments from two sequences with maximum similarity */
/* Michael S. Waterman and Mark Eggert */
/* see "A NEW ALGORITHM FOR BEST SUBSEQUENCE ALIGNMENTS
WITH APPLICATION TO TRNA-RRNA COMPARISONS
Michael S. Waterman and Mark Eggert. Journal of Molecular
Biology (1987) 197,723-728 */

#include "maxsegs.h"

/* sequence lengths */
int m, n;

/* sequences as arrays of characters */
char *x, /* x[1]...x[m] */
      *y; /* y[1]...y[n] */

```

```

/* score cells as arrays of arrays of short integers */
short
**F, /* F[0][0]...F[l][n] */
**S; /* S[0][0]...S[m][n] */

int maxantidiag; /* maximum antidiagonal of traceback end */
int maxS, maxSI, maxSJ; /* maximum score and its location */
int maxi; /* maximum l of traceback end */
int force; /* flag to force reporting maximum score */

/* weighting constants for SIMPLE match/mismatch score */
int match, delta; /* s(a,b)=match if a==b
                  delta if a!=b */

/* weighting constants for LINEAR insertion/deletion score */
int alpha, beta; /* w(k)=alpha+beta*(k-1) */

/* cutoff, tracebacks are not done from cells with scores below this value */
int cutoff;

/* structure for cell information */
typedef struct {int i, j;} cellinfo;

/* array of cell structures to be ordered */
cellinfo *list;

/* number of cell structures in array */
int listcount;

.....
/* maxsegs() is the entry point after all initialization is done */
/* BEFORE maxsegs() is called: */
/* match, delta, alpha, beta and cutoff must be assigned */
/* x and y must be loaded with the sequences to be compared */
/* m and n must be assigned the sequence lengths */
/* the score arrays must be initialized thus: */
/* S[i][0]=0 for all i */
/* F[0][j]=0 for all j */
/* S[0][j]=0 for all j */
.....

maxsegs()
{
    make_S();
    trace();
}

#define FAST_REG

/* load score cells with pre-traceback values */

make_S()
{
    register int i, j;
    /* register to hold x[i] */
    register char cx;
    /* registers to make array access easier */
    register short E, *Sp, *Fp=F[0], *bFp=F[1], *sFp;

    maxS=0;
    listcount=0;

    for (j=1; j<=n; ++j)
        Fp[j]=0;
    /* recursive S for LINEAR weighting */
    for (i=1; i<=m; ++i)
    {
        E=0;
        sFp=Fp;
        Fp=bFp;
        bFp=sFp;
        Sp=S[i];
        cx=x[i];
        for (j=1; j<=n; ++j)
        {
            /* F[i][j] = max { 0, max over k { S[i][j-k]-w(k) } */
            E=max(0,

```

```

            Sp[j]-alpha,
            E-beta);
        /* F[i][j] = max { 0, max over k { S[i][j-k]-w(k) } */
        Fp[j]=max(0,
            S[i-1][j]-alpha,
            bFp[j]-beta);

        Sp[j]=max(S[i-1][j-1]+s(cx, y[j]),
            (int)Fp[j],
            (int)E);

        if (Sp[j]>maxS)
        {
            maxS=Sp[j];
            maxSI=i;
            maxSJ=j;
        }
        if (Sp[j]>=cutoff)
        {
            list[listcount].i=i;
            list[listcount].j=j;
            ++listcount;
        }
    }
}

/* match/mismatch score */
int s(a, b) char a, b;
{
    if (a==b)
        return(match);
    else
        return(delta);
}

/* deletion/insertion score */
int w(k) int k;
{
    return(alpha+beta*(k-1));
}

/* recalculation variables */
int lastl, lastj;
/* If row j is being analysed
   all cells from cell[lastl][j-1] to cell[m][j-1] are unchanged. */
/* If column i is being analysed
   all cells from cell[i-1][lastj] to cell[i-1][n] are unchanged. */

/* Cellorder returns a value of type RELATION. The RELATION type, and its value:
   BEFORE, SAME_AS and AFTER are defined in max_segs_defs. The cellorder code
   compares two cell structures, item1 is "cellorder{item1, item2}" item2. Our
   program uses this, but its use is hidden in the sorting code */

RELATION cellorder(item1, item2) cellinfo *item1, *item2;
{
    /* order by cell's content */
    if (S[item1->i][item1->j]>S[item2->i][item2->j])
        return(BEFORE);
    if (S[item1->i][item1->j]<S[item2->i][item2->j])
        return(AFTER);

    /* order by distance of cell's antidiagonal from origin */
    if ((item1->i+item1->j)<(item2->i+item2->j))
        return(BEFORE);
    if ((item1->i+item1->j)>(item2->i+item2->j))
        return(AFTER);

    /* order by cell's position on cell's antidiagonal */
    if (item1->i<item2->i)
        return(BEFORE);
    if (item1->i>item2->i)
        return(AFTER);
}

```



```

/* same content, same position on same antidiagonal */
return(SAME_AS);
}

/* make tracebacks */
trace()
{
    int i, j, k;
    int scorei, scorej, scoreS;

    /* enqueue qualifying cells */
    /* if the highest is too low */
    printf("maxS is %d\n", maxS);
    if (maxS < cutoff)
    {
        if (force)
        {
            list[listcount].i = maxSi;
            list[listcount].j = maxSj;
            ++listcount;
            cutoff = 0;
        }
    }
    else
    {
        /* sort array with an algorithm following cellorder's behavior: */
        /* sort(pointer to array, number of objects in array) */
        sort(list, listcount);
    }

    /* start traceback at each queued entry,
    taking entries from "top" of array */
    /* this allows us to easily sort remaining entries
    indexed 0 to listcount-1 */
    while (--listcount >= 0)
    {
        scorei = list[listcount].i;
        scorej = list[listcount].j;
        scoreS = S[scorei][scorej];
        /* check if entry is OK */
        if (scoreS >= cutoff)
        {
            maxantidiag = 0;
            maxi = m;
            /* finds "shortest" traceback */
            findnode(scorei, scorej);
            lasti = lastj = 0;
            /* if this traceback succeeds... */
            if (plot(scorei, scorej))
            {
                /* recalculate scores off end of traceback */
                redoscores(scorei, scorej);
                /* print the alignment in */
                /* a form the user likes */
                diagnostic_stuff(scorei, scorej, scoreS);
                /* resort list */
                sort(list, listcount);
            }
        }
        else
        {
            /* if the top entry is too low, all are too low */
            break;
        }
    }
}

BOOLEAN plot(i, j) int i, j;
{
    int k, l;

    /* stop at zeros, return TRUE if this is shortest */
    if (i < maxantidiag)
        return(FALSE);
    if (S[i][j] == 0)
    {
        if (i == j == maxantidiag)
            return(TRUE);
        else
            return(FALSE);
    }
}

```

```

return(FALSE);
}

/* check deletion branch */
for (k=1; k<i; ++k)
{
    if (S[i-k][j] == 0)
        break;
    else if (S[i][j] == S[i-k][j] - w(k))
    {
        if (plot(i-k, j))
        {
            for (l=k-1; l>=0; --l)
            {
                /* note that cell i-l, j
                /* is deleted
                /* output something like
                /* x[i-l]
                /* -
                /*
                diag(i-l, j, NDEL);
                redoone(i-l, j);
                redocolumn(i-l, j);
            }
            return(TRUE);
        }
    }
}

/* check match/mismatch branch */
if (S[i][j] == S[i-1][j-1] + s(x[i], y[j]))
{
    if (plot(i-1, j-1))
    {
        /* note that cell i, j has had match or mismatch
        /* produce alignment output something like:
        /* x[i]
        /* y[j]
        /*
        /* S is built "left to right"
        /* recursion is "right to left"
        /* alignment output occurs on recursion exit
        /* so it comes out "left to right"
        diag(i, j, MATCH);
        if (redoone(i, j))
        {
            redocolumn(i, j);
            redorow(i, j);
        }
        return(TRUE);
    }
}

/* check other deletion branch */
for (k=1; k<j; ++k)
{
    if (S[i][j-k] == 0)
        break;
    else if (S[i][j] == S[i][j-k] - w(k))
    {
        if (plot(i, j-k))
        {
            for (l=k-1; l>=0; --l)
            {
                /* note that cell i, j-l is deleted
                /* output something like
                /* -
                /* y[j-l]
                /*
                diag(i, j-l, YDEL);
                redoone(i, j-l);
                redorow(i, j-l);
            }
            return(TRUE);
        }
    }
}

return(FALSE);
}
}

```

```

}
/* find "shortest" traceback */
findplot(i, j) register int i, j;
{
    register int k, l;
    /* stop if already too low */
    if (i+j < maxantidiag)
        return;
    /* stop at zeros */
    if (S[i][j] == 0)
    {
        if (i > maxantidiag)
        {
            maxantidiag = i+1;
            maxi = i;
        }
        else if (i+j == maxantidiag)
        {
            if (i > maxi)
                maxi = i;
        }
        return;
    }
    /* check match/mismatch branch */
    if (S[i][j] == S[i-1][j-1] + s(x[i], y[j]))
        findplot(i-1, j-1);
    /* check deletion branch */
    for (k=1; k<i; ++k)
    {
        if (S[i-k][j] == 0)
            break;
        else if (S[i][j] == S[i-k][j] - w(k))
            findplot(i-k, j);
    }
    /* check other deletion branch */
    for (k=1; k<j; ++k)
    {
        if (S[i][j-k] == 0)
            break;
        else if (S[i][j] == S[i][j-k] - w(k))
            findplot(i, j-k);
    }
}

/* find cell common to all tracebacks with minimum i+j */
findnode(i, j) int i, j;
{
    int nodei=i, nodej=j, righti=i, rightj=j, lefti=i, leftj=j;
    int k, leftindex=lefti+leftj, rightindex=righti+rightj;
    /* move down left-most and right-most tracebacks
    until they meet for the last time */
    for (;;)
    {
        /* left side */
        do
        {
            if (S[lefti][leftj] == 0)
            {
                break;
            }
            /* check deletion branch */
            for (k=1; k<=lefti; ++k)
            {
                if (S[lefti-k][leftj] == 0)
                {
                    if (S[lefti][leftj] ==
                        S[lefti-1][leftj]-1) +
                        s(x[lefti], y[leftj]))
                    {
                        --lefti;
                        --leftj;
                        leftindex -= 2;
                    }
                }
            }
        }
    }
}

```

```

}
else
{
    for (k=1; k<=leftj; ++k)
    {
        if (S[lefti][leftj] ==
            S[lefti][leftj-k] - w(k))
        {
            leftj -= k;
            leftindex -= k;
            break;
        }
    }
    break;
}
else if (S[lefti][leftj] == S[lefti-k][leftj] - w(k))
{
    lefti -= k;
    leftindex -= k;
    break;
}
}
}
while (leftindex > rightindex);
if (S[lefti][leftj] == 0)
{
    if (leftindex > maxantidiag)
    {
        maxantidiag = leftindex;
        maxi = lefti;
    }
    else if (leftindex == maxantidiag)
    {
        if (lefti > maxi)
            maxi = lefti;
    }
    break;
}
/* right side */
do
{
    if (S[righti][rightj] == 0)
    {
        break;
    }
    for (k=1; k<=rightj; ++k)
    {
        if (S[righti][rightj-k] == 0)
        {
            if (S[righti][rightj] ==
                S[righti-1][rightj-1] +
                s(x[righti], y[rightj]))
            {
                --rightj;
                --righti;
                rightindex -= 2;
            }
        }
        else
        {
            for (k=1; k<=righti; ++k)
            {
                if (S[righti][rightj] ==
                    S[righti-k][rightj] - w(k))
                {
                    righti -= k;
                    rightindex -= k;
                    break;
                }
            }
        }
    }
    break;
}
else if (S[righti][rightj] ==
    S[righti][rightj-k] - w(k))
{
}
}

```

```

        rightj--k;
        rightindex--k;
        break;
    }
}
while (leftindex<rightindex):
if (S[righti][rightj]==0)
{
    if (rightindex>maxantidiag)
    {
        maxantidiag=rightindex;
        maxi=righti;
    }
    else if (rightindex==maxantidiag)
    {
        if (righti>maxi)
            maxi=righti;
    }
    break;
}
if (lefti==righti)
{
    nodei=lefti;
    nodej=leftj;
}
}
/* now, find the "shortest" path from this node */
findplot(nodei, nodej);
}
SE(i, j)
register int i, j;
register int k, m=0, s0;
register short *Sp=S[i];
for (k=1; k<j; ++k)
{
    s0=Sp[j-k];
    if (s0==0)
        break;
    s0--w(k);
    if (s0>m)
        m=s0;
}
return(m);
}
SF(i, j)
register int i, j;
register int k, m=0, s0;
for (k=1; k<i; ++k)
{
    s0=S[i-k][j];
    if (s0==0)
        break;
    s0--w(k);
    if (s0>m)
        m=s0;
}
return(m);
}
/* return TRUE if further cells will be affected */
NOCLEAN redocone(i, j) int i, j;
{
    int newS;
    /* recalculate matrix at one place */
    /* omit scores from pathways already taken */
    newS=max(readdiag(i, j, MATCH)?0:(S[i-1][j]-1)*s(x[i], y[j])),
            readdiag(i, j, XDEL)?0:SF(i, j),
            readdiag(i, j, YDEL)?0:SE(i, j));
}

```

```

/* note result */
if (S[i][j]!=newS)
{
    S[i][j]=newS;
    return(TRUE);
}
/* if further cells will be unaffected, report this */
if ((i>=lasti)&&(j>=lastj))
{
    return(FALSE);
}
/* report that further cells will be affected */
return(TRUE);
}
/* recalculate column of matrix */
redocolumn(i, j) int i, j;
{
    int newS;
    int firstj=0; /* j of cell with lowest j of contiguous unchanged cells */
    while (++j<=n)
    {
        /* recalculate S sans previous paths */
        /* readdiag=true if cell[i][j] was MATCHed or XDEled or YDEled */
        newS=max(readdiag(i, j, MATCH)?0:(S[i-1][j]-1)*s(x[i], y[j])),
                readdiag(i, j, XDEL)?0:SF(i, j),
                readdiag(i, j, YDEL)?0:SE(i, j));
        if (S[i][j]==newS)
        {
            /* remember lowest cell
            where all further cells don't change */
            if (firstj==0)
                firstj=j;
            /* when no further cells will change, we're done */
            if (j>=lastj)
            {
                /* remember position for next scan */
                lastj=firstj;
                break;
            }
        }
        else
        {
            /* note that cell has changed */
            /* clear j of lowest of contiguous unchanged cells */
            firstj=0;
            /* actually change S */
            S[i][j]=newS;
        }
    }
}
/* recalculate row */
redorow(i, j) int i, j;
{
    int newS;
    int firsti=0;
    while (++i<=m)
    {
        newS=max(readdiag(i, j, MATCH)?0:(S[i-1][j]-1)*s(x[i], y[j])),
                readdiag(i, j, XDEL)?0:SF(i, j),
                readdiag(i, j, YDEL)?0:SE(i, j));
        if (S[i][j]==newS)
        {
            if (firsti==0)
                firsti=i;
            if (i>=lasti)
            {
                /* the first shall be last */
                lasti=firsti;
                break;
            }
        }
    }
}

```

