PART I

Exact String Matching: The Fundamental String Problem

Exact matching: what's the problem?

Given a string P called the *pattern* and a longer string T called the *text*, the **exact** matching problem is to find all occurrences, if any, of pattern P in text T.

For example, if P = aba and T = bbabaxababay then P occurs in T starting at locations 3, 7, and 9. Note that two occurrences of P may overlap, as illustrated by the occurrences of P at locations 7 and 9.

Importance of the exact matching problem

The practical importance of the exact matching problem should be obvious to anyone who uses a computer. The problem arises in widely varying applications, too numerous to even list completely. Some of the more common applications are in word processors; in utilities such as grep on Unix; in textual information retrieval programs such as Medline, Lexis, or Nexis; in library catalog searching programs that have replaced physical card catalogs in most large libraries; in internet browsers and crawlers, which sift through massive amounts of text available on the internet for material containing specific keywords;¹ in internet news readers that can search the articles for topics of interest; in the giant digital libraries that are being planned for the near future; in electronic journals that are already being "published" on-line; in telephone directory assistance; in on-line encyclopedias and other educational CD-ROM applications; in on-line dictionaries and thesauri, especially those with crossreferencing features (the Oxford English Dictionary project has created an electronic on-line version of the OED containing 50 million words); and in numerous specialized databases. In molecular biology there are several hundred specialized databases holding raw DNA, RNA, and amino acid strings, or processed patterns (called motifs) derived from the raw string data. Some of these databases will be discussed in Chapter 15.

Although the practical importance of the exact matching problem is not in doubt, one might ask whether the problem is still of any research or educational interest. Hasn't exact matching been so well solved that it can be put in a black box and taken for granted? Right now, for example, I am editing a ninety-page file using an "ancient" shareware word processor and a PC clone (486), and every exact match command that I've issued executes faster than I can blink. That's rather depressing for someone writing a book containing a large section on exact matching algorithms. So is there anything left to do on this problem?

The answer is that for typical word-processing applications there probably is little left to do. The exact matching problem is solved for those applications (although other more so-phisticated string tools might be useful in word processors). But the story changes radically

¹ I just visited the Alta Vista web page maintained by the Digital Equipment Corporation. The Alta Vista database contains over 21 billion words collected from over 10 million web sites. A search for all web sites that mention "Mark Twain" took a couple of seconds and reported that twenty thousand sites satisfy the query. For another example see [392].

for other applications. Users of Melvyl, the on-line catalog of the University of California library system, often experience long, frustrating delays even for fairly simple matching requests. Even grepping through a large directory can demonstrate that exact matching is not yet trivial. Recently we used GCG (a very popular interface to search DNA and protein databanks) to search Genbank (the major U.S. DNA database) for a thirty-character string, which is a small string in typical uses of Genbank. The search took over four hours (on a local machine using a local copy of the database) to find that the string was not there.² And Genbank today is only a fraction of the size it will be when the various genome programs go into full production mode, cranking out massive quantities of sequenced DNA. Certainly there are faster, common database searching programs (for example, BLAST), and there are faster machines one can use (for example, an e-mail server is available for exact and inexact database matching running on a 4,000 processor MasPar computer). But the point is that the exact matching problem is not so effectively and universally solved that it needs no further attention. It will remain a problem of interest as the size of the databases grow and also because exact matching will continue to be a subtask needed for more complex searches that will be devised. Many of these will be illustrated in this book.

But perhaps the most important reason to study *exact* matching in detail is to understand the various ideas developed for it. Even assuming that the exact matching problem itself is sufficiently solved, the entire field of string algorithms remains vital and open, and the education one gets from studying exact matching may be crucial for solving less understood problems. That education takes three forms: specific algorithms, general algorithmic styles, and analysis and proof techniques. All three are covered in this book, but style and proof technique get the major emphasis.

Overview of Part I

In Chapter 1 we present naive solutions to the exact matching problem and develop the fundamental tools needed to obtain more efficient methods. Although the classical solutions to the problem will not be presented until Chapter 2, we will show at the end of Chapter 1 that the use of fundamental tools alone gives a simple linear-time algorithm for exact matching. Chapter 2 develops several classical methods for exact matching, using the fundamental tools developed in Chapter 1. Chapter 3 looks more deeply at those methods and extensions of them. Chapter 4 moves in a very different direction, exploring methods for exact matching based on arithmetic-like operations rather than character comparisons.

Although exact matching is the focus of Part I, some aspects of inexact matching and the use of wild cards are also discussed. The exact matching problem will be discussed again in Part II, where it (and extensions) will be solved using suffix trees.

Basic string definitions

We will introduce most definitions at the point where they are first used, but several definitions are so fundamental that we introduce them now.

Definition A string S is an ordered list of characters written contiguously from left to right. For any string S, S[i..j] is the (contiguous) substring of S that starts at position

² We later repeated the test using the Boyer–Moore algorithm on our own raw copy of Genbank. The search took less than ten minutes, most of which was devoted to movement of text between the disk and the computer, with less than one minute used by the actual text search.

i and ends at position *j* of *S*. In particular, S[1..i] is the *prefix* of string *S* that ends at position *i*, and S[i..|S|] is the *suffix* of string *S* that begins at position *i*, where |S| denotes the number of characters in string *S*.

Definition S[i...j] is the empty string if i > j.

For example, *california* is a string, *lifo* is a substring, *cal* is a prefix, and *ornia* is a suffix.

Definition A *proper* prefix, suffix, or substring of S is, respectively, a prefix, suffix, or substring that is not the entire string S, nor the empty string.

Definition For any string S, S(i) denotes the *i*th character of S.

We will usually use the symbol S to refer to an arbitrary fixed string that has no additional assumed features or roles. However, when a string is known to play the role of a pattern or the role of a text, we will refer to the string as P or T respectively. We will use lower case Greek characters (α , β , γ , δ) to refer to variable strings and use lower case roman characters to refer to single variable characters.

Definition When comparing two characters, we say that the characters *match* if they are equal; otherwise we say they *mismatch*.

Terminology confusion

The words "string" and "word" are often used synonymously in the computer science literature, but for clarity in this book we will never use "word" when "string" is meant. (However, we do use "word" when its colloquial English meaning is intended.)

More confusing, the words "string" and "sequence" are often used synonymously, particularly in the biological literature. This can be the source of much confusion because "substrings" and "subsequences" are very different objects and because algorithms for substring problems are usually very different than algorithms for the analogous subsequence problems. The characters in a substring of S must occur contiguously in S, whereas characters in a subsequence might be interspersed with characters not in the subsequence. Worse, in the biological literature one often sees the word "sequence" used in place of "subsequence". Therefore, for clarity, in this book we will always maintain a distinction between "subsequence" and "substring" and never use "sequence" for "subsequence". We will generally use "string" when pure computer science issues are discussed and use "sequence" or "string" interchangeably in the context of biological applications. Of course, we will also use "sequence" when its standard mathematical meaning is intended.

The first two parts of this book primarily concern problems on strings and substrings. Problems on subsequences are considered in Parts III and IV.

Exact Matching: Fundamental Preprocessing and First Algorithms

1.1. The naive method

Almost all discussions of exact matching begin with the *naive method*, and we follow this tradition. The naive method aligns the left end of P with the left end of T and then compares the characters of P and T left to right until either two unequal characters are found or until P is exhausted, in which case an occurrence of P is reported. In either case, P is then shifted one place to the right, and the comparisons are restarted from the left end of P. This process repeats until the right end of P shifts past the right end of T.

Using *n* to denote the length of *P* and *m* to denote the length of *T*, the worst-case number of comparisons made by this method is $\Theta(nm)$. In particular, if both *P* and *T* consist of the same repeated character, then there is an occurrence of *P* at each of the first m - n + 1 positions of *T* and the method performs exactly n(m - n + 1) comparisons. For example, if P = aaa and T = aaaaaaaaaaa then n = 3, m = 10, and 24 comparisons are made.

The naive method is certainly simple to understand and program, but its worst-case running time of $\Theta(nm)$ may be unsatisfactory and can be improved. Even the practical running time of the naive method may be too slow for larger texts and patterns. Early on, there were several related ideas to improve the naive method, both in practice and in worst case. The result is that the $\Theta(n \times m)$ worst-case bound can be reduced to O(n + m). Changing " \times " to "+" in the bound is extremely significant (try n = 1000 and m = 10,000,000, which are realistic numbers in some applications).

1.1.1. Early ideas for speeding up the naive method

The first ideas for speeding up the naive method all try to shift P by more than one character when a mismatch occurs, but never shift it so far as to miss an occurrence of P in T. Shifting by more than one position saves comparisons since it moves P through T more rapidly. In addition to shifting by larger amounts, some methods try to reduce comparisons by skipping over parts of the pattern after the shift. We will examine many of these ideas in detail.

Figure 1.1 gives a flavor of these ideas, using P = abxyabxz and T = xabxyabxyabxz. Note that an occurrence of P begins at location 6 of T. The naive algorithm first aligns P at the left end of T, immediately finds a mismatch, and shifts P by one position. It then finds that the next seven comparisons are matches and that the succeeding comparison (the ninth overall) is a mismatch. It then shifts P by one place, finds a mismatch, and repeats this cycle two additional times, until the left end of P is aligned with character 6 of T. At that point it finds eight matches and concludes that P occurs in T starting at position 6. In this example, a total of twenty comparisons are made by the naive algorithm.

A smarter algorithm might realize, after the ninth comparison, that the next three

EXACT MATCHING





comparisons of the naive algorithm will be mismatches. This smarter algorithm skips over the next three shift/compares, immediately moving the left end of P to align with position 6 of T, thus saving three comparisons. How can a smarter algorithm do this? After the ninth comparison, the algorithm knows that the first seven characters of P match characters 2 through 8 of T. If it also knows that the first character of P (namely a) does not occur again in P until position 5 of P, it has enough information to conclude that character a does not occur again in T until position 6 of T. Hence it has enough information to conclude that there can be no matches between P and T until the left end of P is aligned with position 6 of T. Reasoning of this sort is the key to shifting by more than one character. In addition to shifting by larger amounts, we will see that certain aligned characters do not need to be compared.

An even smarter algorithm knows the next occurrence in P of the first three characters of P (namely abx) begin at position 5. Then since the first seven characters of P were found to match characters 2 through 8 of T, this smarter algorithm has enough information to conclude that when the left end of P is aligned with position 6 of T, the next three comparisons must be matches. This smarter algorithm avoids making those three comparisons. Instead, after the left end of P is moved to align with position 6 of T, the algorithm compares character 4 of P against character 9 of T. This smarter algorithm therefore saves a total of six comparisons over the naive algorithm.

The above example illustrates the kinds of ideas that allow some comparisons to be skipped, although it should still be unclear how an algorithm can efficiently implement these ideas. Efficient implementations have been devised for a number of algorithms such as the Knuth-Morris-Pratt algorithm, a real-time extension of it, the Boyer-Moore algorithm, and the Apostolico-Giancarlo version of it. All of these algorithms have been implemented to run in linear time (O(n + m) time). The details will be discussed in the next two chapters.

1.2. The preprocessing approach

Many string matching and analysis algorithms are able to efficiently skip comparisons by first spending "modest" time learning about the internal structure of either the pattern P or the text T. During that time, the other string may not even be known to the algorithm. This part of the overall algorithm is called the *preprocessing* stage. Preprocessing is followed by a *search* stage, where the information found during the preprocessing stage is used to reduce the work done while searching for occurrences of P in T. In the above example, the

6

smarter method was assumed to know that character a did not occur again until position 5, and the even smarter method was assumed to know that the pattern abx was repeated again starting at position 5. This assumed knowledge is obtained in the preprocessing stage.

For the exact matching problem, all of the algorithms mentioned in the previous section preprocess pattern P. (The opposite approach of preprocessing text T is used in other algorithms, such as those based on suffix trees. Those methods will be explained later in the book.) These preprocessing methods, as originally developed, are similar in spirit but often quite different in detail and conceptual difficulty. In this book we take a different approach and do not initially explain the originally developed preprocessing methods. Rather, we highlight the similarity of the preprocessing *tasks* needed for several different matching algorithms, by first defining a *fundamental preprocessing* of P that is independent of any particular matching algorithm. Then we show how each specific matching algorithm uses the information computed by the fundamental preprocessing of P. The result is a simpler more uniform exposition of the preprocessing needed by several classical matching methods and a simple linear time algorithm for exact matching based only on this preprocessing (discussed in Section 1.5). This approach to linear-time pattern matching was developed in [202].

1.3. Fundamental preprocessing of the pattern

Fundamental preprocessing will be described for a general string denoted by S. In specific applications of fundamental preprocessing, S will often be the pattern P, but here we use S instead of P because fundamental preprocessing will also be applied to strings other than P.

The following definition gives the key values computed during the fundamental preprocessing of a string.

Definition Given a string S and a position i > 1, let $Z_i(S)$ be the *length* of the longest substring of S that *starts* at i and matches a prefix of S.

In other words, $Z_i(S)$ is the length of the longest *prefix* of S[i..|S|] that matches a prefix of S. For example, when S = aabcaabxaaz then

$$Z_5(S) = 3 \ (aabc...aabx...),$$

 $Z_6(S) = 1 \ (aa...ab...),$
 $Z_7(S) = Z_8(S) = 0,$
 $Z_9(S) = 2 \ (aab...aaz).$

When S is clear by context, we will use Z_i in place of $Z_i(S)$.

To introduce the next concept, consider the boxes drawn in Figure 1.2. Each box starts at some position j > 1 such that Z_j is greater than zero. The length of the box starting at j is meant to represent Z_j . Therefore, each box in the figure represents a maximal-length



Figure 1.2: Each solid box represents a substring of *S* that matches a prefix of *S* and that starts between positions 2 and *i*. Each box is called a *Z*-box. We use r_i to denote the *right-most* end of any *Z*-box that begins at or to the left of position *i* and α to denote the substring in the *Z*-box ending at r_i . Then l_i denotes the left end of α . The copy of α that occurs as a prefix of *S* is also shown in the figure.

substring of S that matches a prefix of S and that does not start at position one. Each such box is called a Z-box. More formally, we have:

Definition For any position i > 1 where Z_i is greater than zero, the Z-box at i is defined as the interval starting at i and ending at position $i + Z_i - 1$.

Definition For every i > 1, r_i is the right-most endpoint of the Z-boxes that begin at or before position *i*. Another way to state this is: r_i is the largest value of $j + Z_j - 1$ over all $1 < j \le i$ such that $Z_j > 0$. (See Figure 1.2.)

We use the term l_i for the value of j specified in the above definition. That is, l_i is the position of the *left end* of the Z-box that ends at r_i . In case there is more than one Z-box ending at r_i , then l_i can be chosen to be the left end of any of those Z-boxes. As an example, suppose S = aabaabcaxaabaabcy; then $Z_{10} = 7$, $r_{15} = 16$, and $l_{15} = 10$.

The linear time computation of Z values from S is the *fundamental* preprocessing task that we will use in all the classical linear-time matching algorithms that preprocess P. But before detailing those uses, we show how to do the fundamental preprocessing in linear time.

1.4. Fundamental preprocessing in linear time

The task of this section is to show how to compute all the Z_i values for S in linear time (i.e., in O(|S|) time). A direct approach based on the definition would take $\Theta(|S|^2)$ time. The method we will present was developed in [307] for a different purpose.

The preprocessing algorithm computes Z_i , r_i , and l_i for each successive position i, starting from i = 2. All the Z values computed will be kept by the algorithm, but in any iteration i, the algorithm only needs the r_j and l_j values for j = i - 1. No earlier r or l values are needed. Hence the algorithm only uses a single variable, r, to refer to the most recently computed r_j value; similarly, it only uses a single variable l. Therefore, in each iteration i, if the algorithm discovers a new Z-box (starting at i), variable r will be incremented to the end of that Z-box, which is the right-most position of any Z-box discovered so far.

To begin, the algorithm finds Z_2 by explicitly comparing, left to right, the characters of S[2..|S|] and S[1..|S|] until a mismatch is found. Z_2 is the length of the matching string. If $Z_2 > 0$, then $r = r_2$ is set to $Z_2 + 1$ and $l = l_2$ is set to 2. Otherwise r and l are set to zero. Now assume inductively that the algorithm has correctly computed Z_i for i up to k - 1 > 1, and assume that the algorithm knows the current $r = r_{k-1}$ and $l = l_{k-1}$. The algorithm next computes Z_k , $r = r_k$, and $l = l_k$.

The main idea is to use the already computed Z values to accelerate the computation of Z_k . In fact, in some cases, Z_k can be deduced from the previous Z values without doing any additional character comparisons. As a concrete example, suppose k = 121, all the values Z_2 through Z_{120} have already been computed, and $r_{120} = 130$ and $l_{120} = 100$. That means that there is a substring of length 31 starting at position 100 and matching a prefix of S (of length 31). It follows that the substring of length 10 starting at position 121 must match the substring of length 10 starting at position 22 of S, and so Z_{22} may be very helpful in computing Z_{121} . As one case, if Z_{22} is three, say, then a little reasoning shows that Z_{121} must also be three. Thus in this illustration, Z_{121} can be deduced without any additional character comparisons. This case, along with the others, will be formalized and proven correct below.



Figure 1.3: String S[k..r] is labeled β and also occurs starting at position k' of S.



Figure 1.4: Case 2a. The longest string starting at k' that matches a prefix of S is shorter than $|\beta|$. In this case, $Z_k = Z_{k'}$.



Figure 1.5: Case 2b. The longest string starting at k' that matches a prefix of S is at least $|\beta|$.

The Z algorithm

Given Z_i for all $1 < i \le k - 1$ and the current values of r and l, Z_k and the updated r and l are computed as follows:

Begin

- 1. If k > r, then find Z_k by explicitly comparing the characters starting at position k to the characters starting at position 1 of S, until a mismatch is found. The length of the match is Z_k . If $Z_k > 0$, then set r to $k + Z_k 1$ and set l to k.
- 2. If $k \le r$, then position k is contained in a Z-box, and hence S(k) is contained in substring S[l..r] (call it α) such that l > 1 and α matches a prefix of S. Therefore, character S(k) also appears in position k' = k l + 1 of S. By the same reasoning, substring S[k..r] (call it β) must match substring $S[k'..Z_l]$. It follows that the substring beginning at position k must match a prefix of S of length at least the *minimum* of $Z_{k'}$ and $|\beta|$ (which is r k + 1). See Figure 1.3.

We consider two subcases based on the value of that minimum.

- **2a.** If $Z_{k'} < |\beta|$ then $Z_k = Z_{k'}$ and r, l remain unchanged (see Figure 1.4).
- **2b.** If $Z_{k'} \ge |\beta|$ then the entire substring S[k..r] must be a prefix of S and $Z_k \ge |\beta| = r k + 1$. However, Z_k might be strictly larger than $|\beta|$, so compare the characters starting at position r + 1 of S to the characters starting a position $|\beta| + 1$ of S until a mismatch occurs. Say the mismatch occurs at character $q \ge r + 1$. Then Z_k is set to q k, r is set to q 1, and l is set to k (see Figure 1.5).

End

Theorem 1.4.1. Using Algorithm Z, value Z_k is correctly computed and variables r and l are correctly updated.

PROOF In Case 1, Z_k is set correctly since it is computed by explicit comparisons. Also (since k > r in Case 1), before Z_k is computed, no Z-box has been found that starts

between positions 2 and k - 1 and that ends at or after position k. Therefore, when $Z_k > 0$ in Case 1, the algorithm does find a new Z-box ending at or after k, and it is correct to change r to $k + Z_k - 1$. Hence the algorithm works correctly in Case 1.

In Case 2a, the substring beginning at position k can match a prefix of S only for length $Z_{k'} < |\beta|$. If not, then the next character to the right, character $k + Z_{k'}$, must match character $1 + Z_{k'}$. But character $k + Z_{k'}$ matches character $k' + Z_{k'}$ (since $Z_{k'} < |\beta|$), so character $k' + Z_{k'}$ must match character $1 + Z_{k'}$. However, that would be a contradiction to the definition of $Z_{k'}$, for it would establish a substring longer than $Z_{k'}$ that starts at k'and matches a prefix of S. Hence $Z_k = Z_{k'}$ in this case. Further, $k + Z_k - 1 < r$, so r and *l* remain correctly unchanged.

In Case 2b, β must be a prefix of S (as argued in the body of the algorithm) and since any extension of this match is explicitly verified by comparing characters beyond r to characters beyond the prefix β , the full extent of the match is correctly computed. Hence Z_k is correctly obtained in this case. Furthermore, since $k + Z_k - 1 \ge r$, the algorithm correctly changes r and l. \Box

Corollary 1.4.1. Repeating Algorithm Z for each position i > 2 correctly yields all the Z_i values.

Theorem 1.4.2. All the $Z_i(S)$ values are computed by the algorithm in O(|S|) time.

PROOF The time is proportional to the number of iterations, |S|, plus the number of character comparisons. Each comparison results in either a match or a mismatch, so we next bound the number of matches and mismatches that can occur.

Each iteration that performs any character comparisons at all ends the first time it finds a mismatch; hence there are at most |S| mismatches during the entire algorithm. To bound the number of matches, note first that $r_k \ge r_{k-1}$ for every iteration k. Now, let k be an iteration where q > 0 matches occur. Then r_k is set to $r_{k-1} + q$ at least. Finally, $r_k \le |S|$, so the total number of matches that occur during any execution of the algorithm is at most |S|. \Box

1.5. The simplest linear-time exact matching algorithm

Before discussing the more complex (classical) exact matching methods, we show that fundamental preprocessing alone provides a simple linear-time exact matching algorithm. This is the simplest linear-time matching algorithm we know of.

Let S = P\$T be the string consisting of P followed by the symbol "\$" followed by T, where "\$" is a character appearing in neither P nor T. Recall that P has length n and T has length m, and $n \le m$. So, S = P\$T has length n + m + 1 = O(m). Compute $Z_i(S)$ for i from 2 to n + m + 1. Because "\$" does not appear in P or T, $Z_i \le n$ for every i > 1. Any value of i > n + 1 such that $Z_i(S) = n$ identifies an occurrence of P in T starting at position i - (n + 1) of T. Conversely, if P occurs in T starting at position j of T, then $Z_{(n+1)+j}$ must be equal to n. Since all the $Z_i(S)$ values can be computed in O(n + m) = O(m) time, this approach identifies all the occurrences of P in T in O(m) time.

The method can be implemented to use only O(n) space (in addition to the space needed for pattern and text) independent of the size of the alphabet. Since $Z_i \leq n$ for all *i*, position k' (determined in step 2) will always fall inside *P*. Therefore, there is no need to record the *Z* values for characters in *T*. Instead, we only need to record the *Z* values

for the *n* characters in *P* and also maintain the current *l* and *r*. Those values are sufficient to compute (but not store) the *Z* value of each character in *T* and hence to identify and output any position *i* where $Z_i = n$.

There is another characteristic of this method worth introducing here: The method is considered an *alphabet-independent* linear-time method. That is, we never had to assume that the alphabet size was finite or that we knew the alphabet ahead of time – a character comparison only determines whether the two characters match or mismatch; it needs no further information about the alphabet. We will see that this characteristic is also true of the Knuth-Morris-Pratt and Boyer–Moore algorithms, but not of the Aho–Corasick algorithm or methods based on suffix trees.

1.5.1. Why continue?

Since function Z_i can be computed for the pattern in linear time and can be used directly to solve the exact matching problem in O(m) time (with only O(n) additional space), why continue? In what way are more complex methods (Knuth-Morris-Pratt, Boyer-Moore, real-time matching, Apostolico-Giancarlo, Aho-Corasick, suffix tree methods, etc.) deserving of attention?

For the exact matching problem, the Knuth-Morris-Pratt algorithm has only a marginal advantage over the direct use of Z_i . However, it has historical importance and has been generalized, in the Aho-Corasick algorithm, to solve the problem of searching for a set of patterns in a text in time linear in the size of the text. That problem is not nicely solved using Z_i values alone. The real-time extension of Knuth-Morris-Pratt has an advantage in situations when text is input on-line and one has to be sure that the algorithm will be ready for each character as it arrives. The Boyer-Moore method is valuable because (with the proper implementation) it also runs in linear worst-case time but typically runs in sublinear time, examining only a fraction of the characters of T. Hence it is the preferred method in most cases. The Apostolico-Giancarlo method is valuable because it has all the advantages of the Boyer-Moore method and yet allows a relatively simple proof of linear worst-case running time. Methods based on suffix trees typically preprocess the text rather than the pattern and then lead to algorithms in which the search time is proportional to the size of the pattern rather than the size of the text. This is an extremely desirable feature. Moreover, suffix trees can be used to solve much more complex problems than exact matching, including problems that are not easily solved by direct application of the fundamental preprocessing.

1.6. Exercises

The first four exercises use the fact that fundamental processing can be done in linear time and that all occurrences of P in T can be found in linear time.

- 1. Use the existence of a linear-time exact matching algorithm to solve the following problem in linear time. Given two strings α and β , determine if α is a circular (or cyclic) rotation of β , that is, if α and β have the same length and α consists of a suffix of β followed by a prefix of β . For example, *defabc* is a circular rotation of *abcdef*. This is a classic problem with a very elegant solution.
- **2.** Similar to Exercise 1, give a linear-time algorithm to determine whether a linear string α is a substring of a *circular string* β . A circular string of length *n* is a string in which character *n* is considered to precede character 1 (see Figure 1.6). Another way to think about this



Figure 1.6: A circular string β . The linear string $\overline{\beta}$ derived from it is *accatggc*.

problem is the following. Let $\overline{\beta}$ be the *linear* string obtained from β starting at character 1 and ending at character *n*. Then α is a substring of circular string β if and only if α is a substring of some circular rotation of $\overline{\beta}$.

A digression on circular strings in DNA

The above two problems are mostly exercises in using the existence of a linear-time exact matching algorithm, and we don't know any critical biological problems that they address. However, we want to point out that *circular DNA* is common and important. Bacterial and mitochondrial DNA is typically circular, both in its genomic DNA and in additional small double-stranded circular DNA molecules called *plasmids*, and even some true eukaryotes (higher organisms whose cells contain a nucleus) such as yeast contain plasmid DNA in addition to their nuclear DNA. Consequently, tools for handling circular strings may someday be of use in those organisms. Viral DNA is not always circular, but even when it is linear some virus genomes exhibit circular properties. For example, in some viral populations the linear order of the DNA in one individual will be a circular rotation of the order in another individual [450]. Nucleotide mutations, in addition to rotations, occur rapidly in viruses, and a plausible problem is to determine if the DNA of two individual viruses have mutated away from each other only by a circular rotation, rather than additional mutations.

It is very interesting to note that the problems addressed in the exercises are actually "solved" in nature. Consider the special case of Exercise 2 when string α has length *n*. Then the problem becomes: Is α a circular rotation of $\overline{\beta}$? This problem is solved in linear time as in Exercise 1. Precisely this matching problem arises and is "solved" in *E. coli* replication under the certain experimental conditions described in [475]. In that experiment, an enzyme (RecA) and ATP molecules (for energy) are added to *E. coli* containing a single strand of one of its plasmids, called string β , and a double-stranded linear DNA molecule, one strand of which is called string α . If α is a circular rotation of $\overline{\beta}$ then the strand opposite to α (which has the DNA sequence complementary to α) hybridizes with β creating a proper double-stranded plasmid, leaving α as a single strand. This transfer of DNA may be a step in the replication of the plasmid. Thus the problem of determining whether α is a circular rotation of $\overline{\beta}$ is solved by this natural system.

Other experiments in [475] can be described as *substring* matching problems relating to circular and linear DNA in *E. coli*. Interestingly, these natural systems solve their matching problems faster than can be explained by kinetic analysis, and the molecular mechanisms used for such rapid matching remain undetermined. These experiments demonstrate the role of enzyme RecA in *E. coli* replication, but do not suggest immediate important computational problems. They do, however, provide indirect motivation for developing computational tools for handling circular strings as well as linear strings. Several other uses of circular strings will be discussed in Sections 7.13 and 16.17 of the book.

3. Suffix-prefix matching. Give an algorithm that takes in two strings α and β , of lengths n

and *m*, and finds the longest suffix of α that exactly matches a prefix of β . The algorithm should run in O(n + m) time.

4. Tandem arrays. A substring α contained in string *S* is called a *tandem array* of β (called the base) if α consists of more than one consecutive copy of β . For example, if S = xyzabcabcabcabcabcqq, then $\alpha = abcabcabcabcabc$ is a tandem array of $\beta = abc$. Note that *S* also contains a tandem array of abcabc (i.e., a tandem array with a longer base). A *maximal* tandem array is a tandem array that cannot be extended either left or right. Given the base β , a tandem array of β in *S* can be described by two numbers (*s*, *k*), giving its starting location in *S* and the number of times β is repeated. A tandem array is an example of a repeated substring (see Section 7.11.1).

Suppose *S* has length *n*. Give an example to show that two maximal tandem arrays of a given base β can overlap.

Now give an O(n)-time algorithm that takes S and β as input, finds every maximal tandem array of β , and outputs the pair (s, k) for each occurrence. Since maximal tandem arrays of a given base can overlap, a naive algorithm would establish only an $O(n^2)$ -time bound.

- 5. If the Z algorithm finds that $Z_2 = q > 0$, all the values $Z_3, \ldots, Z_{q+1}, Z_{q+2}$ can then be obtained immediately without additional character comparisons and without executing the main body of Algorithm Z. Flesh out and justify the details of this claim.
- **6.** In Case 2b of the Z algorithm, when $Z_{k'} \ge |\beta|$, the algorithm does explicit comparisons until it finds a mismatch. This is a reasonable way to organize the algorithm, but in fact Case 2b can be refined so as to eliminate an unneeded character comparison. Argue that when $Z_{k'} > |\beta|$ then $Z_k = |\beta|$ and hence no character comparisons are needed. Therefore, explicit character comparisons are needed only in the case that $Z_{k'} = |\beta|$.
- 7. If Case 2b of the Z algorithm is split into two cases, one for $Z_{k'} > |\beta|$ and one for $Z_{k'} = |\beta|$, would this result in an overall speedup of the algorithm? You must consider all operations, not just character comparisons.
- 8. Baker [43] introduced the following matching problem and applied it to a problem of software maintenance: "The application is to track down duplication in a large software system. We want to find not only exact matches between sections of code, but parameterized matches, where a parameterized match between two sections of code means that one section can be transformed into the other by replacing the parameter names (e.g., identifiers and constants) of one section by the parameter names of the other via a one-to-one function".

Now we present the formal definition. Let Σ and Π be two alphabets containing no symbols in common. Each symbol in Σ is called a *token* and each symbol in Π is called a *parameter*. A string can consist of any combinations of tokens and parameters from Σ and Π . For example, if Σ is the upper case English alphabet and Π is the lower case alphabet then *XYabCaCXZddW* is a legal string over Σ and Π . Two strings S_1 and S_2 are said to *p*-match if and only if

- **a.** Each token in S_1 (or S_2) is opposite a matching token in S_2 (or S_1).
- **b.** Each parameter in S_1 (or S_2) is opposite a parameter in S_2 (or S_1).
- **c.** For any parameter x, if one occurrence of x in S_1 (S_2) is opposite a parameter y in S_2 (S_1), then every occurrence of x in S_1 (S_2) must be opposite an occurrence of y in S_2 (S_1). In other words, the alignment of parameters in S_1 and S_2 defines a one-one correspondence between parameter names in S_1 and parameter names in S_2 .

For example, $S_1 = XYabCaCXZddbW$ *p*-matches $S_2 = XYdxCdCXZccxW$. Notice that parameter *a* in S_1 maps to parameter *d* in S_2 , while parameter *d* in S_1 maps to *c* in S_2 . This does not violate the definition of *p*-matching.

In Baker's application, a token represents a part of the program that cannot be changed,

whereas a parameter represents a program's variable, which can be renamed as long as all occurrences of the variable are renamed consistently. Thus if S_1 and S_2 *p*-match, then the variable names in S_1 could be changed to the corresponding variable names in S_2 , making the two programs identical. If these two programs were part of a larger program, then they could both be replaced by a call to a single subroutine.

The most basic *p*-match problem is: Given a text *T* and a pattern *P*, each a string over Σ and Π , find all substrings of *T* that *p*-match *P*. Of course, one would like to find all those occurrences in O(|P| + |T|) time. Let function Z_i^p for a string *S* be the length of the longest string starting at position *i* in *S* that *p*-matches a prefix of *S*[1..*i*]. Show how to modify algorithm *Z* to compute all the Z_i^p values in O(|S|) time (the implementation details are slightly more involved than for function Z_i , but not too difficult). Then show how to use the modified algorithm *Z* to find all substrings of *T* that *p*-match *P*, in O(|P| + |T|) time.

In [43] and [239], more involved versions of the *p*-match problem are solved by more complex methods.

The following three problems can be solved without the Z algorithm or other fancy tools. They only require thought.

- **9.** You are given two strings of *n* characters each and an additional parameter *k*. In each string there are n k + 1 substrings of length *k*, and so there are $\Theta(n^2)$ pairs of substrings, where one substring is from one string and one is from the other. For a pair of substrings, we define the *match-count* as the number of opposing characters that match when the two substrings of length *k* are aligned. The problem is to compute the match-count for each of the $\Theta(n^2)$ pairs of substrings from the two strings. Clearly, the problem can be solved with $O(kn^2)$ operations (character comparisons plus arithmetic operations). But by better organizing the computations, the time can be reduced to $O(n^2)$ operations. (From Paul Horton.)
- **10.** A DNA molecule can be thought of as a string over an alphabet of four characters $\{a, t, c, g\}$ (nucleotides), while a protein can be thought of as a string over an alphabet of twenty characters (amino acids). A gene, which is physically embedded in a DNA molecule, typically encodes the amino acid sequence for a particular protein. This is done as follows. Starting at a particular point in the DNA string, every three consecutive DNA characters encode a single amino acid character in the protein string. That is, three DNA nucleotides specify one amino acid. Such a coding triple is called a *codon*, and the full association of codons to amino acids is called the *genetic code*. For example, the codon *ttt* codes for the amino acid Phenylalanine (abbreviated in the single character amino acid alphabet as *F*), and the codon *gtt* codes for the amino acids, there is a possibility that two or more triples form codons for the same amino acid and that some triples do not form codons. In fact, this is the case. For example, the amino acid Leucine is coded for by six different codons.

Problem: Suppose one is given a DNA string of *n* nucleotides, but you don't know the correct "reading frame". That is, you don't know if the correct decomposition of the string into codons begins with the first, second, or third nucleotide of the string. Each such "frameshift" potentially translates into a different amino acid string. (There are actually known genes where each of the three reading frames not only specifies a string in the amino acid alphabet, but each specifies a functional, yet different, protein.) The task is to produce, for each of the three reading frames, the associated amino acid string. For example, consider the string *atggacgga*. The first reading frame has three complete codons, *atg, gac,* and *gga,* which in the genetic code specify the amino acids *Met, Asp,* and *Gly.* The second reading frame has two complete codons, *tgg* and *acg,* coding for amino acids *Trp* and *Thr.* The third reading frame has two complete codons, *gga* and *cgg,* coding for amino acids *Gly* and *Arg.*

The goal is to produce the three translations, using the fewest number of character exami-

nations of the DNA string and the fewest number of indexing steps (when using the codons to look up amino acids in a table holding the genetic code). Clearly, the three translations can be done with 3n examinations of characters in the DNA and 3n indexing steps in the genetic code table. Find a method that does the three translations in at most n character examinations and n indexing steps.

Hint: If you are acquainted with this terminology, the notion of a finite-state transducer may be helpful, although it is not necessary.

11. Let *T* be a text string of length *m* and let *S* be a *multiset* of *n* characters. The problem is to find all substrings in *T* of length *n* that are formed by the characters of *S*. For example, let $S = \{a, a, b, c\}$ and T = abahgcabah. Then *caba* is a substring of *T* formed from the characters of *S*.

Give a solution to this problem that runs in O(m) time. The method should also be able to state, for each position *i*, the length of the longest substring in *T* starting at *i* that can be formed from S.

Fantasy protein sequencing. The above problem may become useful in sequencing protein from a particular organism after a large amount of the genome of that organism has been sequenced. This is most easily explained in prokaryotes, where the DNA is not interrupted by introns. In prokaryotes, the amino acid sequence for a given protein is encoded in a contiguous segment of DNA - one DNA codon for each amino acid in the protein. So assume we have the protein molecule but do not know its sequence or the location of the gene that codes for the protein. Presently, chemically determining the amino acid sequence of a protein is very slow, expensive, and somewhat unreliable. However, finding the multiset of amino acids that make up the protein is relatively easy. Now suppose that the whole DNA sequence for the genome of the organism is known. One can use that long DNA sequence to determine the amino acid sequence of a protein of interest. First, translate each codon in the DNA sequence into the amino acid alphabet (this may have to be done three times to get the proper frame) to form the string T; then chemically determine the multiset S of amino acids in the protein; then find all substrings in T of length |S| that are formed from the amino acids in S. Any such substrings are candidates for the amino acid sequence of the protein, although it is unlikely that there will be more than one candidate. The match also locates the gene for the protein in the long DNA string.

- 12. Consider the two-dimensional variant of the preceding problem. The input consists of two-dimensional text (say a filled-in crossword puzzle) and a multiset of characters. The problem is to find a *connected* two-dimensional substructure in the text that matches all the characters in the multiset. How can this be done? A simpler problem is to restrict the structure to be rectangular.
- 13. As mentioned in Exercise 10, there are organisms (some viruses for example) containing intervals of DNA encoding not just a single protein, but three viable proteins, each read in a different reading frame. So, if each protein contains n amino acids, then the DNA string encoding those three proteins is only n + 2 nucleotides (characters) long. That is a very compact encoding.

(**Challenging problem?**) Give an algorithm for the following problem: The input is a protein string S_1 (over the amino acid alphabet) of length n and another protein string of length m > n. Determine if there is a string specifying a DNA encoding for S_2 that contains a substring specifying a DNA encoding of S_1 . Allow the encoding of S_1 to begin at any point in the DNA string for S_2 (i.e., in any reading frame of that string). The problem is difficult because of the degeneracy of the genetic code and the ability to use any reading frame.

2.1. Introduction

This chapter develops a number of classical comparison-based matching algorithms for the exact matching problem. With suitable extensions, all of these algorithms can be implemented to run in linear worst-case time, and all achieve this performance by preprocessing pattern P. (Methods that preprocess T will be considered in Part II of the book.) The original preprocessing methods for these various algorithms are related in spirit but are quite different in conceptual difficulty. Some of the original preprocessing methods are quite difficult.¹ This chapter does not follow the original preprocessing methods but instead exploits fundamental preprocessing, developed in the previous chapter, to implement the needed preprocessing for each specific matching algorithm.

Also, in contrast to previous expositions, we emphasize the Boyer–Moore method over the Knuth-Morris-Pratt method, since Boyer–Moore is the practical method of choice for exact matching. Knuth-Morris-Pratt is nonetheless completely developed, partly for historical reasons, but mostly because it generalizes to problems such as real-time string matching and matching against a set of patterns more easily than Boyer–Moore does. These two topics will be described in this chapter and the next.

2.2. The Boyer–Moore Algorithm

As in the naive algorithm, the Boyer-Moore algorithm successively aligns P with T and then checks whether P matches the opposing characters of T. Further, after the check is complete, P is shifted right relative to T just as in the naive algorithm. However, the Boyer-Moore algorithm contains three clever ideas not contained in the naive algorithm: the right-to-left scan, the bad character shift rule, and the good suffix shift rule. Together, these ideas lead to a method that typically examines fewer than m + n characters (an expected sublinear-time method) and that (with a certain extension) runs in linear worstcase time. Our discussion of the Boyer-Moore algorithm, and extensions of it, concentrates on *provable* aspects of its behavior. Extensive experimental and practical studies of Boyer-Moore and variants have been reported in [229], [237], [409], [410], and [425].

2.2.1. Right-to-left scan

For any alignment of P with T the Boyer-Moore algorithm checks for an occurrence of P by scanning characters from *right to left* rather than from left to right as in the naive

¹ Sedgewick [401] writes "Both the Knuth-Morris-Pratt and the Boyer-Moore algorithms require some complicated preprocessing on the pattern that is difficult to understand and has limited the extent to which they are used". In agreement with Sedgewick, I still do not understand the original Boyer-Moore preprocessing method for the *strong* good suffix rule.

algorithm. For example, consider the alignment of P against T shown below:

To check whether P occurs in T at this position, the Boyer-Moore algorithm starts at the *right* end of P, first comparing T(9) with P(7). Finding a match, it then compares T(8) with P(6), etc., moving right to left until it finds a mismatch when comparing T(5) with P(3). At that point P is shifted *right* relative to T (the amount for the shift will be discussed below) and the comparisons begin again at the right end of P.

Clearly, if P is shifted right by one place after each mismatch, or after an occurrence of P is found, then the worst-case running time of this approach is O(nm) just as in the naive algorithm. So at this point it isn't clear why comparing characters from right to left is any better than checking from left to right. However, with two additional ideas (the *bad character* and the *good suffix* rules), shifts of more than one position often occur, and in typical situations large shifts are common. We next examine these two ideas.

2.2.2. Bad character rule

To get the idea of the bad character rule, suppose that the last (right-most) character of P is y and the character in T it aligns with is $x \neq y$. When this initial mismatch occurs, if we know the right-most position in P of character x, we can safely shift P to the right so that the right-most x in P is below the mismatched x in T. Any shorter shift would only result in an immediate mismatch. Thus, the longer shift is correct (i.e., it will not shift past any occurrence of P in T). Further, if x never occurs in P, then we can shift P completely past the point of mismatch in T. In these cases, some characters of T will never be examined and the method will actually run in "sublinear" time. This observation is formalized below.

Definition For each character x in the alphabet, let R(x) be the position of right-most occurrence of character x in P. R(x) is defined to be zero if x does not occur in P.

It is easy to preprocess P in O(n) time to collect the R(x) values, and we leave that as an exercise. Note that this preprocessing does not require the fundamental preprocessing discussed in Chapter 1 (that will be needed for the more complex shift rule, the good suffix rule).

We use the *R* values in the following way, called the *bad character shift rule*:

Suppose for a particular alignment of P against T, the right-most n - i characters of P match their counterparts in T, but the next character to the left, P(i), mismatches with its counterpart, say in position k of T. The bad character rule says that P should be shifted right by max [1, i - R(T(k))] places. That is, if the right-most occurrence in P of character T(k) is in position j < i (including the possibility that j = 0), then shift P so that character j of P is below character k of T. Otherwise, shift P by one position.

The point of this shift rule is to shift P by more than one character when possible. In the above example, T(5) = t mismatches with P(3) and R(t) = 1, so P can be shifted right by two positions. After the shift, the comparison of P and T begins again at the right end of P.

Extended bad character rule

The bad character rule is a useful heuristic for mismatches near the right end of P, but it has no effect if the mismatching character from T occurs in P to the right of the mismatch point. This may be common when the alphabet is small and the text contains many similar, but not exact, substrings. That situation is typical of DNA, which has an alphabet of size four, and even protein, which has an alphabet of size twenty, often contains different regions of high similarity. In such cases, the following *extended bad character rule* is more robust:

When a mismatch occurs at position i of P and the mismatched character in T is x, then shift P to the right so that the closest x to the left of position i in P is below the mismatched x in T.

Because the extended rule gives larger shifts, the only reason to prefer the simpler rule is to avoid the added implementation expense of the extended rule. The simpler rule uses only $O(|\Sigma|)$ space (Σ is the alphabet) for array R, and one table lookup for each mismatch. As we will see, the extended rule can be implemented to take only O(n) space and at most one extra step per character comparison. That amount of added space is not often a critical issue, but it is an empirical question whether the longer shifts make up for the added time used by the extended rule. The original Boyer-Moore algorithm only uses the simpler bad character rule.

Implementing the extended bad character rule

We preprocess P so that the extended bad character rule can be implemented efficiently in both time and space. The preprocessing should discover, for each position i in P and for each character x in the alphabet, the position of the closest occurrence of x in P to the left of i. The obvious approach is to use a two-dimensional array of size n by $|\Sigma|$ to store this information. Then, when a mismatch occurs at position i of P and the mismatching character in T is x, we look up the (i, x) entry in the array. The lookup is fast, but the size of the array, and the time to build it, may be excessive. A better compromise, below, is possible.

During preprocessing, scan P from right to left collecting, for each character x in the alphabet, a list of the positions where x occurs in P. Since the scan is right to left, each list will be in decreasing order. For example, if P = abacbabc then the list for character a is 6, 3, 1. These lists are accumulated in O(n) time and of course take only O(n) space. During the search stage of the Boyer-Moore algorithm if there is a mismatch at position i of P and the mismatching character in T is x, scan x's list from the top until we reach the first number less than i or discover there is none. If there is none then there is no occurrence of x before i, and all of P is shifted past the x in T. Otherwise, the found entry gives the desired position of x.

After a mismatch at position i of P the time to scan the list is at most n - i, which is roughly the number of characters that matched. So in worst case, this approach at most doubles the running time of the Boyer-Moore algorithm. However, in most problem settings the added time will be vastly less than double. One could also do binary search on the list in circumstances that warrant it.

2.2.3. The (strong) good suffix rule

The bad character rule by itself is reputed to be highly effective in practice, particularly for English text [229], but proves less effective for small alphabets and it does not lead to a linear worst-case running time. For that, we introduce another rule called the *strong*



Figure 2.1: Good suffix shift rule, where character x of T mismatches with character y of P. Characters y and z of P are guaranteed to be distinct by the good suffix rule, so z has a chance of matching x.

good suffix rule. The original preprocessing method [278] for the strong good suffix rule is generally considered quite difficult and somewhat mysterious (although a weaker version of it is easy to understand). In fact, the preprocessing for the strong rule was given incorrectly in [278] and corrected, without much explanation, in [384]. Code based on [384] is given without real explanation in the text by Baase [32], but there are no published sources that try to fully explain the method.² Pascal code for strong preprocessing, based on an outline by Richard Cole [107], is shown in Exercise 24 at the end of this chapter.

In contrast, the fundamental preprocessing of P discussed in Chapter 1 makes the needed preprocessing very simple. That is the approach we take here. The strong good suffix rule is:

Suppose for a given alignment of P and T, a substring t of T matches a suffix of P, but a mismatch occurs at the next comparison to the left. Then find, if it exists, the right-most copy t' of t in P such that t' is not a suffix of P and the character to the left of t' in P differs from the character to the left of t in P. Shift P to the right so that substring t' in P is below substring t in T (see Figure 2.1). If t' does not exist, then shift the left end of P past the left end of t in T. If no such shift is possible, then shift P by n places to the right. If an occurrence of P is found, then shift P by the least amount so that a proper prefix of the shifted P matches a suffix of the shifted P matches a suffix of the shift P by n places, that is, shift P past t in T.

For a specific example consider the alignment of P and T given below:

When the mismatch occurs at position 8 of P and position 10 of T, t = ab and t' occurs in P starting at position 3. Hence P is shifted right by six places, resulting in the following alignment:

² A recent plea appeared on the internet newsgroup comp. theory:

I am looking for an elegant (easily understandable) proof of correctness for a part of the Boyer-Moore string matching algorithm. The difficult-to-prove part here is the algorithm that computes the dd_2 (good-suffix) table. I didn't find much of an understandable proof yet, so I'd much appreciate any help!

| | 0 | 1 | |
|----|---------|------------|----|
| | 1234567 | 890123456 | 78 |
| т: | prstabs | tubabvqxrs | st |
| P: | | qcabdabda | ab |

Note that the extended bad character rule would have shifted P by only one place in this example.

Theorem 2.2.1. The use of the good suffix rule never shifts P past an occurrence in T.

PROOF Suppose the right end of P is aligned with character k of T before the shift, and suppose that the good suffix rule shifts P so its right end aligns with character k' > k. Any occurrence of P ending at a position l strictly between k and k' would immediately violate the selection rule for k', since it would imply either that a closer copy of t occurs in P or that a longer prefix of P matches a suffix of t. \Box

The original published Boyer-Moore algorithm [75] uses a simpler, weaker, version of the good suffix rule. That version just requires that the shifted P agree with the t and does not specify that the next characters to the left of those occurrences of t be different. An explicit statement of the weaker rule can be obtained by deleting the italics phrase in the first paragraph of the statement of the strong good suffix rule. In the previous example, the weaker shift rule shifts P by three places rather than six. When we need to distinguish the two rules, we will call the simpler rule the *weak* good suffix rule and the rule stated above the *strong* good suffix rule. For the purpose of proving that the search part of Boyer-Moore runs in linear worst-case time, the weak rule is not sufficient, and in this book the strong version is assumed unless stated otherwise.

2.2.4. Preprocessing for the good suffix rule

We now formalize the preprocessing needed for the Boyer-Moore algorithm.

Definition For each *i*, L(i) is the largest position less than *n* such that string P[i..n] matches a suffix of P[1..L(i)]. L(i) is defined to be zero if there is no position satisfying the conditions. For each *i*, L'(i) is the largest position less than *n* such that string P[i..n] matches a suffix of P[1..L'(i)] and such that the character preceding that suffix is not equal to P(i-1). L'(i) is defined to be zero if there is no position satisfying the conditions.

For example, if P = cabdabdab, then L(8) = 6 and L'(8) = 3.

L(i) gives the right end-position of the right-most copy of P[i..n] that is not a suffix of P, whereas L'(i) gives the right end-position of the right-most copy of P[i..n] that is not a suffix of P, with the stronger, added condition that its preceding character is unequal to P(i - 1). So, in the strong-shift version of the Boyer-Moore algorithm, if character i - 1 of P is involved in a mismatch and L'(i) > 0, then P is shifted right by n - L'(i) positions. The result is that if the right end of P was aligned with position k of T before the shift, then position L'(i) is now aligned with position k.

During the preprocessing stage of the Boyer-Moore algorithm L'(i) (and L(i), if desired) will be computed for each position *i* in *P*. This is done in O(n) time via the following definition and theorem.

Definition For string P, $N_j(P)$ is the length of the longest suffix of the substring P[1..j] that is also a suffix of the full string P.

For example, if P = cabdabdab, then $N_3(P) = 2$ and $N_6(P) = 5$.

Recall that $Z_i(S)$ is the length of the longest substring of S that starts at *i* and matches a prefix of S. Clearly, N is the reverse of Z, that is, if P^r denotes the string obtained by reversing P, then $N_j(P) = Z_{n-j+1}(P^r)$. Hence the $N_j(P)$ values can be obtained in O(n)time by using Algorithm Z on P^r . The following theorem is then immediate.

Theorem 2.2.2. L(i) is the largest index j less than n such that $N_j(P) \ge |P[i..n]|$ (which is n-i+1). L'(i) is the largest index j less than n such that $N_j(P) = |P[i..n]| = (n-i+1)$.

Given Theorem 2.2.2, it follows immediately that all the L'(i) values can be accumulated in linear time from the N values using the following algorithm:

Z-based Boyer-Moore

for i := 1 to n do L'(i) := 0; for j := 1 to n - 1 do begin $i := n - N_j(P) + 1$; L'(i) := j; end;

The L(i) values (if desired) can be obtained by adding the following lines to the above pseudocode:

L(2) := L'(2);for i := 3 to n do $L(i) := \max[L(i - 1), L'(i)];$

Theorem 2.2.3. The above method correctly computes the L values.

PROOF L(i) marks the right end-position of the right-most substring of P that matches P[i..n] and is not a suffix of P[1..n]. Therefore, that substring begins at position L(i)-n+i, which we will denote by j. We will prove that $L(i) = \max[L(i-1), L'(i)]$ by considering what character j - 1 is. First, if j = 1 then character j - 1 doesn't exist, so L(i - 1) = 0 and L'(i) = 1. So suppose that j > 1. If character j - 1 equals character i - 1 then L(i) = L(i-1). If character j - 1 does not equal character i - 1 then L(i) = L'(i). Thus, in all cases, L(i) must either be L'(i) or L(i - 1).

However, L(i) must certainly be greater than or equal to both L'(i) and L(i - 1). In summary, L(i) must either be L'(i) or L(i - 1), and yet it must be greater or equal to both of them; hence L(i) must be the maximum of L'(i) and L(i - 1). \Box

Final preprocessing detail

The preprocessing stage must also prepare for the case when L'(i) = 0 or when an occurrence of P is found. The following definition and theorem accomplish that.

Definition Let l'(i) denote the length of the largest suffix of P[i..n] that is also a prefix of P, if one exists. If none exists, then let l'(i) be zero.

Theorem 2.2.4. l'(i) equals the largest $j \leq |P[i..n]|$, which is n - i + 1, such that $N_j(P) = j$.

We leave the proof, as well as the problem of how to accumulate the l'(i) values in linear time, as a simple exercise. (Exercise 9 of this chapter)

2.2.5. The good suffix rule in the search stage of Boyer–Moore

Having computed L'(i) and l'(i) for each position *i* in *P*, these preprocessed values are used during the search stage of the algorithm to achieve larger shifts. If, during the search stage, a mismatch occurs at position i - 1 of *P* and L'(i) > 0, then the good suffix rule shifts *P* by n - L'(i) places to the right, so that the L'(i)-length prefix of the shifted *P* aligns with the L'(i)-length suffix of the unshifted *P*. In the case that L'(i) = 0, the good suffix rule shifts *P* by n - l'(i) places. When an occurrence of *P* is found, then the rule shifts *P* by n - l'(2) places. Note that the rules work correctly even when l'(i) = 0.

One special case remains. When the first comparison is a mismatch (i.e., P(n) mismatches) then P should be shifted one place to the right.

2.2.6. The complete Boyer–Moore algorithm

We have argued that neither the good suffix rule nor the bad character rule shift P so far as to miss any occurrence of P. So the Boyer–Moore algorithm shifts by the largest amount given by either of the rules. We can now present the complete algorithm.

The Boyer-Moore algorithm

```
{Preprocessing stage}
     Given the pattern P,
    Compute L'(i) and l'(i) for each position i of P,
     and compute R(x) for each character x \in \Sigma.
{Search stage}
    k := n;
     while k \leq m do
         begin
         i := n;
          h := k;
          while i > 0 and P(i) = T(h) do
               begin
               i := i - 1;
               h := h - 1;
               end;
          if i = 0 then
               begin
               report an occurrence of P in T ending at position k.
               k := k + n - l'(2);
               end
          else
               shift P (increase k) by the maximum amount determined by the
               (extended) bad character rule and the good suffix rule.
          end;
```

Note that although we have always talked about "shifting P", and given rules to determine by how much P should be "shifted", there is no shifting in the actual implementation. Rather, the index k is increased to the point where the right end of P would be "shifted". Hence, each act of shifting P takes constant time.

We will later show, in Section 3.2, that by using the strong good suffix rule alone, the

Boyer-Moore method has a worst-case running time of O(m) provided that the pattern does not appear in the text. This was first proved by Knuth, Morris, and Pratt [278], and an alternate proof was given by Guibas and Odlyzko [196]. Both of these proofs were quite difficult and established worst-case time bounds no better than 5m comparisons. Later, Richard Cole gave a much simpler proof [108] establishing a bound of 4m comparisons and also gave a difficult proof establishing a tight bound of 3m comparisons. We will present Cole's proof of 4m comparisons in Section 3.2.

When the pattern does appear in the text then the original Boyer-Moore method runs in $\Theta(nm)$ worst-case time. However, several simple modifications to the method correct this problem, yielding an O(m) time bound in all cases. The first of these modifications was due to Galil [168]. After discussing Cole's proof, in Section 3.2, for the case that P doesn't occur in T, we use a variant of Galil's idea to achieve the linear time bound in all cases.

At the other extreme, if we only use the bad character shift rule, then the worst-case running time is O(nm), but assuming randomly generated strings, the expected running time is sublinear. Moreover, in typical string matching applications involving natural language text, a sublinear running time is almost always observed in practice. We won't discuss random string analysis in this book but refer the reader to [184].

Although Cole's proof for the linear worst case is vastly simpler than earlier proofs, and is important in order to complete the full story of Boyer–Moore, it is not trivial. However, a fairly simple extension of the Boyer–Moore algorithm, due to Apostolico and Giancarlo [26], gives a "Boyer–Moore–like" algorithm that allows a fairly direct proof of a 2m worst-case bound on the number of comparisons. The Apostolico–Giancarlo variant of Boyer–Moore is discussed in Section 3.1.

2.3. The Knuth-Morris-Pratt algorithm

The best known linear-time algorithm for the exact matching problem is due to Knuth, Morris, and Pratt [278]. Although it is rarely the method of choice, and is often much inferior in practice to the Boyer–Moore method (and others), it can be simply explained, and its linear time bound is (fairly) easily proved. The algorithm also forms the basis of the well-known Aho–Corasick algorithm, which efficiently finds all occurrences in a text of any pattern from a *set* of patterns.³

2.3.1. The Knuth-Morris-Pratt shift idea

For a given alignment of P with T, suppose the naive algorithm matches the first i characters of P against their counterparts in T and then mismatches on the next comparison. The naive algorithm would shift P by just *one* place and begin comparing again from the left end of P. But a larger shift may often be possible. For example, if P = abcxabcde and, in the present alignment of P with T, the mismatch occurs in position 8 of P, then it is easily deduced (and we will prove below) that P can be shifted by four places without passing over any occurrences of P in T. Notice that this can be deduced without even knowing what string T is or exactly how P is aligned with T. Only the location of the mismatch in P must be known. The Knuth-Morris-Pratt algorithm is based on this kind of reasoning to make larger shifts than the naive algorithm makes. We now formalize this idea.

³ We will present several solutions to that set problem including the Aho-Corasick method in Section 3.4. For those reasons, and for its historical role in the field, we fully develop the Knuth-Morris-Pratt method here.

Definition For each position *i* in pattern *P*, define $sp_i(P)$ to be the *length* of the longest proper *suffix* of P[1..i] that matches a prefix of *P*.

Stated differently, $sp_i(P)$ is the length of the longest proper substring of P[1..i] that ends at *i* and that matches a prefix of *P*. When the string is clear by context we will use sp_i in place of the full notation.

For example, if P = abcaeabcabd, then $sp_2 = sp_3 = 0$, $sp_4 = 1$, $sp_8 = 3$, and $sp_{10} = 2$. Note that by definition, $sp_1 = 0$ for any string.

An optimized version of the Knuth-Morris-Pratt algorithm uses the following values.

Definition For each position *i* in pattern *P*, define $sp'_i(P)$ to be the length of the longest proper suffix of P[1..i] that matches a prefix of *P*, with the added condition that characters P(i + 1) and $P(sp'_i + 1)$ are unequal.

Clearly, $sp'_i(P) \leq sp_i(P)$ for all positions *i* and any string *P*. As an example, if P = bbccaebbcabd, then $sp_8 = 2$ because string *bb* occurs both as a proper prefix of P[1..8] and as a suffix of P[1..8]. However, both copies of the string are followed by the same character *c*, and so $sp'_8 < 2$. In fact, $sp'_8 = 1$ since the single character *b* occurs as both the first and last character of P[1..8] and is followed by character *b* in position 2 and by character *c* in position 9.

The Knuth-Morris-Pratt shift rule

We will describe the algorithm in terms of the sp' values, and leave it to the reader to modify the algorithm if only the weaker sp values are used.⁴ The Knuth-Morris-Pratt algorithm aligns P with T and then compares the aligned characters from *left to right*, as the naive algorithm does.

For any alignment of P and T, if the first mismatch (comparing from left to right) occurs in position i + 1 of P and position k of T, then shift P to the right (relative to T) so that $P[1..sp'_i]$ aligns with $T[k - sp'_i..k - 1]$. In other words, shift P exactly $i + 1 - (sp'_i + 1) = i - sp'_i$ places to the right, so that character $sp'_i + 1$ of P will align with character k of T. In the case that an occurrence of P has been found (no mismatch), shift P by $n - sp'_n$ places.

The shift rule guarantees that the prefix $P[1..sp'_i]$ of the shifted P matches its opposing substring in T. The next comparison is then made between characters T(k) and $P[sp'_i+1]$. The use of the stronger shift rule based on sp'_i guarantees that the same mismatch will not occur again in the new alignment, but it does not guarantee that $T(k) = P[sp'_i + 1]$.

In the above example, where P = abcxabcde and $sp'_7 = 3$, if character 8 of P mismatches then P will be shifted by 7-3 = 4 places. This is true even without knowing T or how P is positioned with T.

The advantage of the shift rule is twofold. First, it often shifts P by more than just a single character. Second, after a shift, the left-most sp'_i characters of P are guaranteed to match their counterparts in T. Thus, to determine whether the newly shifted P matches its counterpart in T, the algorithm can start comparing P and T at position $sp'_i + 1$ of P (and position k of T). For example, suppose P = abcxabcde as above, T = xyabcxabcxadcdqfeg, and the left end of P is aligned with character 3 of T. Then P and T will match for 7 characters but mismatch on character 8 of P, and P will be shifted

⁴ The reader should be alerted that traditionally the Knuth-Morris-Pratt algorithm has been described in terms of *failure functions*, which are related to the sp_i values. Failure functions will be explicitly defined in Section 2.3.3.



Figure 2.2: Assumed missed occurrence used in correctness proof for Knuth-Morris-Pratt.

by 4 places as shown below:

1 2 123456789012345678 xyabcxabcxadcdqfeg abcxabcde abcxabcde

As guaranteed, the first 3 characters of the shifted P match their counterparts in T (and their counterparts in the unshifted P).

Summarizing, we have

Theorem 2.3.1. After a mismatch at position i + 1 of P and a shift of $i - sp'_i$ places to the right, the left-most sp'_i characters of P are guaranteed to match their counterparts in T.

Theorem 2.3.1 partially establishes the correctness of the Knuth-Morris-Pratt algorithm, but to fully prove correctness we have to show that the shift rule never shifts too far. That is, using the shift rule no occurrence of P will ever be overlooked.

Theorem 2.3.2. For any alignment of P with T, if characters 1 through i of P match the opposing characters of T but character i + 1 mismatches T(k), then P can be shifted by $i - sp'_i$ places to the right without passing any occurrence of P in T.

PROOF Suppose not, so that there is an occurrence of P starting strictly to the left of the shifted P (see Figure 2.2), and let α and β be the substrings shown in the figure. In particular, β is the prefix of P of length sp'_i , shown relative to the shifted position of P. The unshifted P matches T up through position i of P and position k - 1 of T, and all characters in the (assumed) missed occurrence of P match their counterparts in T. Both of these matched regions contain the substrings α and β , so the unshifted P and the assumed occurrence of P match on the entire substring $\alpha\beta$. Hence $\alpha\beta$ is a suffix of P[1..i] that matches a proper prefix of P. Now let $l = |\alpha\beta| + 1$ so that position l in the "missed occurrence" of P is opposite position k in T. Character P(l) cannot be equal to P(i + 1) since P(l) is assumed to match T(k) and P(i + 1) does not match T(k). Thus $\alpha\beta$ is a proper suffix of P[1..i] that matches a prefix of P, and the next character is unequal to P(i + 1). But $|\alpha| > 0$ due to the assumption that an occurrence of P starts strictly before the shifted P, so $|\alpha\beta| > |\beta| = sp'_i$, contradicting the definition of sp'_i . Hence the theorem is proved. \Box

Theorem 2.3.2 says that the Knuth-Morris-Pratt shift rule does not miss any occurrence of P in T, and so the Knuth-Morris-Pratt algorithm will correctly find all occurrences of P in T. The time analysis is equally simple.

Theorem 2.3.3. In the Knuth-Morris-Pratt method, the number of character comparisons is at most 2m.

PROOF Divide the algorithm into compare/shift phases, where a single phase consists of the comparisons done between successive shifts. After any shift, the comparisons in the phase go left to right and start either with the last character of T compared in the previous phase or with the character to its right. Since P is never shifted left, in any phase at most one comparison involves a character of T that was previously compared. Thus, the total number of character comparisons is bounded by m + s, where s is the number of shifts done in the algorithm. But s < m since after m shifts the right end of P is certainly to the right of the right end of T, so the number of comparisons done is bounded by 2m. \Box

2.3.2. Preprocessing for Knuth-Morris-Pratt

The key to the speed up of the Knuth-Morris-Pratt algorithm over the naive algorithm is the use of sp' (or sp) values. It is easy to see how to compute all the sp' and sp values from the Z values obtained during the fundamental preprocessing of P. We verify this below.

Definition Position j > 1 maps to i if $i = j + Z_j(P) - 1$. That is, j maps to i if i is the right end of a Z-box starting at j.

Theorem 2.3.4. For any i > 1, $sp'_i(P) = Z_j = i - j + 1$, where j > 1 is the smallest position that maps to *i*. If there is no such *j* then $sp'_i(P) = 0$. For any i > 1, $sp_i(P) = i - j + 1$, where *j* is the smallest position in the range $1 < j \le i$ that maps to *i* or beyond. If there is no such *j*, then $sp_i(P) = 0$.

PROOF If $sp'_i(P)$ is greater than zero, then there is a proper suffix α of P[1..i] that matches a prefix of P, such that P[i+1] does not match $P[|\alpha|+1]$. Therefore, letting j denote the start of α , $Z_j = |\alpha| = sp'_i(P)$ and j maps to i. Hence, if there is no j in the range $1 < j \le i$ that maps to i, then $sp'_i(P)$ must be zero.

Now suppose $sp'_i(P) > 0$ and let j be as defined above. We claim that j is the smallest position in the range 2 to i that maps to i. Suppose not, and let j^* be a position in the range $1 < j^* < j$ that maps to i. Then $P[j^*..i]$ would be a proper suffix of P[1..i] that matches a prefix (call it β) of P. Moreover, by the definition of mapping, $P(i + 1) \neq P(|\beta|)$, so $sp'_i(P) \ge |\beta| > |\alpha|$, contradicting the assumption that $sp'_i = \alpha$.

The proofs of the claims for $sp_i(P)$ are similar and are left as exercises. \Box

Given Theorem 2.3.4, all the sp' and sp values can be computed in linear time using the Z_i values as follows:

Z-based Knuth-Morris-Pratt

for i := 1 to n do $sp'_i := 0;$ for j := n downto 2 do begin $i := j + Z_j(P) - 1;$ $sp'_i := Z_j;$ end;

The sp values are obtained by adding the following:

 $sp_n(P) := sp'_n(P);$ for i := n - 1 downto 2 do $sp_i(P) := \max[sp_{i+1}(P) - 1, sp'_i(P)]$

2.3.3. A full implementation of Knuth-Morris-Pratt

We have described the Knuth-Morris-Pratt algorithm in terms of shifting P, but we never accounted for time needed to implement shifts. The reason is that shifting is only conceptual and P is never explicitly shifted. Rather, as in the case of Boyer-Moore, pointers to P and T are incremented. We use pointer p to point into P and one pointer c (for "current" character) to point into T.

Definition For each position *i* from 1 to n + 1, define the failure function F'(i) to be $sp'_{i-1} + 1$ (and define $F(i) = sp_{i-1} + 1$), where sp'_0 and sp_0 are defined to be zero.

We will only use the (stronger) failure function F'(i) in this discussion but will refer to F(i) later.

After a mismatch in position i + 1 > 1 of P, the Knuth-Morris-Pratt algorithm "shifts" P so that the next comparison is between the character in position c of T and the character in position $sp'_i + 1$ of P. But $sp'_i + 1 = F'(i+1)$, so a general "shift" can be implemented in constant time by just setting p to F'(i+1). Two special cases remain. When the mismatch occurs in position 1 of P, then p is set to F'(1) = 1 and c is incremented by one. When an occurrence of P is found, then P is shifted right by $n - sp'_n$ places. This is implemented by setting F'(n+1) to $sp'_n + 1$.

Putting all the pieces together gives the full Knuth-Morris-Pratt algorithm.

Knuth-Morris-Pratt algorithm

```
begin
Preprocess P to find F'(k) = sp'_{k-1} + 1 for k from 1 to n + 1.
    c := 1;
     p := 1;
     While c + (n - p) \le m
    do begin
          While P(p) = T(c) and p \le n
         do begin
              p := p + 1;
              c := c + 1;
         end;
    if p = n + 1 then
         report an occurrence of P starting at position c - n of T.
    if p := 1 then c := c + 1
     p := F'(p);
    end;
```

end.

2.4. Real-time string matching

In the search stage of the Knuth-Morris-Pratt algorithm, P is aligned against a substring of T and the two strings are compared left to right until either all of P is exhausted (in which

case an occurrence of P in T has been found) or until a mismatch occurs at some positions i+1 of P and k of T. In the latter case, if $sp'_i > 0$, then P is shifted right by $i - sp'_i$ positions, guaranteeing that the prefix $P[1..sp'_i]$ of the shifted pattern matches its opposing substring in T. No explicit comparison of those substrings is needed, and the next comparison is between characters T(k) and $P(sp'_i + 1)$. Although the shift based on sp'_i guarantees that P(i + 1) differs from $P(sp'_i + 1)$, it does not guarantee that $T(k) = P(sp'_i + 1)$. Hence T(k) might be compared several times (perhaps $\Omega(|P|)$ times) with differing characters in P. For that reason, the Knuth-Morris-Pratt method is not a *real-time* method.

To be real time, a method must do at most a *constant* amount of work between the time it first examines any position in T and the time it last examines that position. In the Knuth-Morris-Pratt method, if a position of T is involved in a match, it is never examined again (this is easy to verify) but, as indicated above, this is not true when the position is involved in a mismatch. Note that the definition of real time only concerns the search stage of the algorithm. Preprocessing of P need not be real time. Note also that if the search stage is real time it certainly is also linear time.

The utility of a real-time matcher is two fold. First, in certain applications, such as when the characters of the text are being sent to a small memory machine, one might need to guarantee that each character can be fully processed before the next one is due to arrive. If the processing time for each character is constant, independent of the length of the string, then such a guarantee may be possible. Second, in this particular real-time matcher, the shifts of P may be longer but never shorter than in the original Knuth-Morris-Pratt algorithm. Hence, the real-time matcher may run faster in certain problem instances.

Admittedly, arguments in favor of real-time matching algorithms over linear-time methods are somewhat tortured, and the real-time matching is more a theoretical issue than a practical one. Still, it seems worthwhile to spend a little time discussing real-time matching.

2.4.1. Converting Knuth-Morris-Pratt to a real-time method

We will use the Z values obtained during fundamental preprocessing of P to convert the Knuth-Morris-Pratt method into a real-time method. The required preprocessing of P is quite similar to the preprocessing done in Section 2.3.2 for the Knuth-Morris-Pratt algorithm. For historical reasons, the resulting real-time method is generally referred to as a *deterministic finite-state string matcher* and is often represented with a finite state machine diagram. We will not use this terminology here and instead represent the method in pseudo code.

Definition Let x denote a character of the alphabet. For each position i in pattern P, define $sp'_{(i,x)}(P)$ to be the length of the longest proper suffix of P[1..i] that matches a prefix of P, with the added condition that character $P(sp'_i + 1)$ is x.

Knowing the $sp'_{(i,x)}$ values for each character x in the alphabet allows a shift rule that converts the Knuth-Morris-Pratt method into a real-time algorithm. Suppose P is compared against a substring of T and a mismatch occurs at characters T(k) = x and P(i+1). Then P should be shifted right by $i - sp'_{(i,x)}$ places. This shift guarantees that the prefix $P[1..sp'_{(i,x)}]$ matches the opposing substring in T and that T(k) matches the next character in P. Hence, the comparison between T(k) and $P(sp'_{(i,x)} + 1)$ can be skipped. The next needed comparison is between characters $P(sp'_{(i,x)} + 2)$ and T(k + 1). With this shift rule, the method becomes real time because it still never reexamines a position in T involved in a match (a feature inherited from the Knuth-Morris-Pratt algorithm), and it now also never reexamines a position involved in a mismatch. So, the search stage of this algorithm never examines a character in T more than once. It follows that the search is done in real time. Below we show how to find all the $sp'_{(i,x)}$ values in linear time. Together, this gives an algorithm that does linear preprocessing of P and real-time search of T.

It is easy to establish that the algorithm finds all occurrences of P in T, and we leave that as an exercise.

2.4.2. Preprocessing for real-time string matching

Theorem 2.4.1. For $P[i+1] \neq x$, $sp'_{(i,x)}(P) = i - j + 1$, where j is the smallest position such that j maps to i and $P(Z_j + 1) = x$. If there is no such j then $sp'_{(i,x)}(P) = 0$.

The proof of this theorem is almost identical to the proof of Theorem 2.3.4 (page 26) and is left to the reader. Assuming (as usual) that the alphabet is finite, the following minor modification of the preprocessing given earlier for Knuth-Morris-Pratt (Section 2.3.2) yields the needed $sp'_{(i,x)}$ values in linear time:

Z-based real-time matching

```
for i := 1 to n do

sp'_{(i,x)} := 0 for every character x;

for j := n downto 2 do

begin

i := j + Z_j(P) - 1;

x := P(Z_j + 1);

sp'_{(i,x)} := Z_j;

end;
```

Note that the linear time (and space) bound for this method require that the alphabet Σ be finite. This allows us to do $|\Sigma|$ comparisons in constant time. If the size of the alphabet is explicitly included in the time and space bounds, then the preprocessing time and space needed for the algorithm is $O(|\Sigma|n)$.

2.5. Exercises

- In "typical" applications of exact matching, such as when searching for an English word in a book, the simple bad character rule seems to be as effective as the extended bad character rule. Give a "hand-waving" explanation for this.
- 2. When searching for a single word or a small phrase in a large English text, brute force (the naive algorithm) is reported [184] to run faster than most other methods. Give a handwaving explanation for this. In general terms, how would you expect this observation to hold up with smaller alphabets (say in DNA with an alphabet size of four), as the size of the pattern grows, and when the text has many long sections of similar but not exact substrings?
- 3. "Common sense" and the $\Theta(nm)$ worst-case time bound of the Boyer–Moore algorithm (using only the bad character rule) both would suggest that empirical running times increase with increasing pattern length (assuming a fixed text). But when searching in actual English

EXACT MATCHING:CLASSICAL COMPARISON-BASED METHODS

30

texts, the Boyer–Moore algorithm runs faster in practice when given longer patterns. Thus, on an English text of about 300,000 characters, it took about five times as long to search for the word "Inter" as it did to search for "Interactively".

Give a hand-waving explanation for this. Consider now the case that the pattern length increases without bound. At what point would you expect the search times to stop decreasing? Would you expect search times to start increasing at some point?

- 4. Evaluate empirically the utility of the extended bad character rule compared to the original bad character rule. Perform the evaluation in combination with different choices for the two good-suffix rules. How much more is the average shift using the extended rule? Does the extra shift pay for the extra computation needed to implement it?
- 5. Evaluate empirically, using different assumptions about the sizes of P and T, the number of occurrences of P in T, and the size of the alphabet, the following idea for speeding up the Boyer-Moore method. Suppose that a phase ends with a mismatch and that the good suffix rule shifts P farther than the extended bad character rule. Let x and y denote the mismatching characters in T and P respectively, and let z denote the character in the shifted P below x. By the suffix rule, z will not be y, but there is no guarantee that it will be x. So rather than starting comparisons from the right of the shifted P, as the Boyer-Moore method would do, why not first compare x and z? If they are equal then a right-to-left comparison is begun from the right end of P, but if they are unequal then we apply the extended bad character rule from z in P. This will shift P again. At that point we must begin a right-to-left comparison of P against T.
- 6. The idea of the bad character rule in the Boyer–Moore algorithm can be generalized so that instead of examining characters in P from right to left, the algorithm compares characters in P in the order of how unlikely they are to be in T (most *unlikely* first). That is, it looks first at those characters in P that are least likely to be in T. Upon mismatching, the bad character rule or extended bad character rule is used as before. Evaluate the utility of this approach, either empirically on real data or by analysis assuming random strings.
- 7. Construct an example where fewer comparisons are made when the bad character rule is used alone, instead of combining it with the good suffix rule.
- 8. Evaluate empirically the effectiveness of the strong good suffix shift for Boyer–Moore versus the weak shift rule.
- **9.** Give a proof of Theorem 2.2.4. Then show how to accumulate all the l'(i) values in linear time.
- 10. If we use the weak good suffix rule in Boyer–Moore that shifts the closest copy of t under the matched suffix t, but doesn't require the next character to be different, then the preprocessing for Boyer–Moore can be based directly on spi values rather than on Z values. Explain this.
- 11. Prove that the Knuth-Morris-Pratt shift rules (either based on sp or sp') do not miss any occurrences of P in T.
- 12. It is possible to incorporate the bad character shift rule from the Boyer–Moore method to the Knuth-Morris-Pratt method or to the naive matching method itself. Show how to do that. Then evaluate how effective that rule is and explain why it is more effective when used in the Boyer–Moore algorithm.
- **13.** Recall the definition of l_i on page 8. It is natural to conjecture that $sp_i = i l_i$ for any index *i*, where $i \ge l_i$. Show by example that this conjecture is incorrect.
- **14.** Prove the claims in Theorem 2.3.4 concerning $sp'_i(P)$.
- **15.** Is it true that given only the *sp* values for a given string *P*, the *sp'* values are completely determined? Are the *sp* values determined from the sp' values alone?

Using sp values to compute Z values

In Section 2.3.2, we showed that one can compute all the *sp* values knowing only the *Z* values for string *S* (i.e., not knowing *S* itself). In the next five exercises we establish the converse, creating a linear-time algorithm to compute all the *Z* values from *sp* values alone. The first exercise suggests a natural method to accomplish this, and the following exercise exposes a hole in that method. The final three exercises develop a correct linear-time algorithm, detailed in [202]. We say that *sp_i* maps to *k* if $k = i - sp_i + 1$.

- **16.** Suppose there is a position *i* such that sp_i maps to *k*, and let *i* be the largest such position. Prove that $Z_k = i - k + 1 = sp_i$ and that $r_k = i$.
- 17. Given the answer to the previous exercise, it is natural to conjecture that Z_k always equals sp_i , where *i* is the largest position such that sp_i maps to *k*. Show that this is not true. Given an example using at least three distinct characters.

Stated another way, give an example to show that Z_k can be greater than zero even when there is *no* position *i* such that sp_i maps to *k*.

- **18.** Recall that r_{k-1} is known at the start of iteration k of the Z algorithm (when Z_k is computed), but r_k is known only at the end of iteration k. Suppose, however, that r_k is known (somehow) at the start of iteration k. Show how the Z algorithm can then be modified to compute Z_k using *no* character comparisons. Hence this modified algorithm need not even know the string S.
- **19.** Prove that if Z_k is greater than zero, then r_k equals the largest position *i* such that $k \ge i sp_i$. Conclude that r_k can be deduced from the *sp* values for every position *k* where Z_k is not zero.
- 20. Combine the answers to the previous two exercises to create a linear-time algorithm that computes all the Z values for a string S given only the sp values for S and not the string S itself.

Explain in what way the method is a "simulation" of the Z algorithm.

- It may seem that I'(i) (needed for Boyer-Moore) should be spn for any i. Show why this is not true.
- 22. In Section 1.5 we showed that all the occurrences of P in T could be found in linear time by computing the Z values on the string S = P\$T. Explain how the method would change if we use S = PT, that is, we do not use a separator symbol between P and T. Now show how to find all occurrences of P in T in linear time using S = PT, but with sp values in place of Z values. (This is not as simple as it might at first appear.)
- 23. In Boyer–Moore and Boyer–Moore–like algorithms, the search moves right to left in the pattern, although the pattern moves left to right relative to the text. That makes it more difficult to explain the methods and to combine the preprocessing for Boyer–Moore with the preprocessing for Knuth-Morris-Pratt. However, a small change to Boyer–Moore would allow an easier exposition and more uniform preprocessing. First, place the pattern at the right end of the text, and conduct each search *left to right* in the pattern, shifting the pattern *left* after a mismatch. Work out the details of this approach, and show how it allows a more uniform exposition of the preprocessing needed for it and for Knuth-Morris-Pratt. Argue that on average this approach has the same behavior as the original Boyer–Moore method.
- 24. Below is working Pascal code (in Turbo Pascal) implementing Richard Cole's preprocessing, for the strong good suffix rule. It is different than the approach based on fundamental preprocessing and is closer to the original method in [278]. Examine the code to extract the algorithm behind the program. Then explain the idea of the algorithm, prove correctness of the algorithm, and analyze its running time. The point of the exercise is that it is difficult to convey an algorithmic idea using a program.

```
program gsmatch(input,output);
{This is an implementation of Richard Cole's
preprocessing for the strong good suffix rule}
type
tstring = string[200];
indexarray = array[1..100] of integer;
const
zero = 0;
var
p:tstring;
bmshift,matchshift:indexarray;
m,i:integer;
procedure readstring(var p:tstring; var m:integer);
begin
read(p);
m:=Length(p);
writeln('the length of the string is ', m);
end;
procedure gsshift(p:tstring; var
gs_shift:indexarray;m:integer);
var
i,j,j_old,k:integer;
kmp_shift:indexarray;
go_on:boolean;
begin {1}
  for j:= 1 to m do
    gs_shift[j] := m;
  kmp_shift[m]:=1;
{stage 1}
  j:=m;
  for k:=m-1 downto 1 do
    begin {2}
       go_on:=true;
       while (p[j] <> p[k]) and go_on do
          begin {3}
          if (gs_shift[j] > j-k) then gs_shift[j] := j-k;
          if (j < m) then j := j + kmp_shift[j+1]
          else go_on:=false;
          end; {3}
```

```
if (p[k] = p[j]) then
          begin {3}
          kmp_shift[k]:=j-k;
         j:=j-1;
          end {3}
       else
       kmp_shift[k]:=j-k+1;
    end; {2}
 {stage 2}
  j:=j+1;
  j_old:=1;
  while (j <= m) do
    begin {2}
      for i:=j_old to j-1 do
      if (gs_shift[i] > j-1) then gs_shift[i]:=j-1;
      j_old:=j;
      j:=j+kmp_shift[j];
    end; {2}
 end; {1}
begin {main}
writeln('input a string on a single line');
readstring(p,m);
gsshift(p,matchshift,m);
writeln('the value in cell i is the number of positions to shift');
writeln('after a mismatch occurring in position i of the pattern');
for i:= 1 to m do
write(matchshift[i]:3);
writeln;
end. {main}
```

- Prove that the shift rule used by the real-time string matcher does not miss any occurrences of P in T.
- 26. Prove Theorem 2.4.1.
- 27. In this chapter, we showed how to use Z values to compute both the sp'_i and sp_i values used in Knuth-Morris-Pratt and the $sp'_{i,x}$ values needed for its real-time extension. Instead of using Z values for the $sp'_{i,x}$ values, show how to obtain these values from the sp_i and/or sp'_i values in linear $[O(n|\Sigma|)]$ time, where n is the length of P and $|\Sigma|$ is the length of the alphabet.
- 28. Although we don't know how to simply convert the Boyer–Moore algorithm to be a real-time method the way Knuth-Morris-Pratt was converted, we can make similar changes to the strong shift rule to make the Boyer–Moore shift more effective. That is, when a mismatch occurs between *P*(*i*) and *T*(*h*) we can look for the right-most copy in *P* of *P*[*i* + 1..*n*] (other than *P*[*i* + 1..*n*] itself) such that the preceding character is *T*(*h*). Show how to modify



Figure 2.3: The pattern P = aqra labels two subpaths of paths starting at the root. Those paths start at the root, but the subpaths containing *aqra* do not. There is also another subpath in the tree labeled *aqra* (it starts above the character *z*), but it violates the requirement that it be a subpath of a path starting at the root. Note that an edge label is displayed from the top of the edge down towards the bottom of the edge. Thus in the figure, there is an edge labeled "qra", not "arq".

the Boyer-Moore preprocessing so that the needed information is collected in linear time, assuming a fixed size alphabet.

29. Suppose we are given a tree where each edge is labeled with one or more characters, and we are given a pattern *P*. The label of a subpath in the tree is the concatenation of the labels on the edges in the subpath. The problem is to find all subpaths of paths starting at the root that are labeled with pattern *P*. Note that although the subpath must be part of a path directed from the root, the subpath itself need not start at the root (see Figure 2.3). Give an algorithm for this problem that runs in time proportional to the total number of characters on the edges of the tree plus the length of *P*.