Identifying protein-binding sites from unaligned DNA fragments

(specificity/regulatory sites/pattern recognition/information theory)

GARY D. STORMO AND GEORGE W. HARTZELL III

Department of Molecular, Cellular and Developmental Biology, University of Colorado, Boulder, CO 80309

Communicated by William B. Wood, November 14, 1988

ABSTRACT The ability to determine important features within DNA sequences from the sequences alone is becoming essential as large-scale sequencing projects are being undertaken. We present a method that can be applied to the problem of identifying the recognition pattern for a DNA-binding protein given only a collection of sequenced DNA fragments, each known to contain somewhere within it a binding site for that protein. Information about the position or orientation of the binding sites within those fragments is not needed. The method compares the "information content" of a large number of possible binding site alignments to arrive at a matrix representation of the binding site pattern. The specificity of the protein is represented as a matrix, rather than a consensus sequence, allowing patterns that are typical of regulatory protein-binding sites to be identified. The reliability of the method improves as the number of sequences increases, but the time required increases only linearly with the number of sequences. An example, using known cAMP receptor proteinbinding sites, illustrates the method.

Gene expression is often controlled by protein factors that interact with DNA regions to affect transcription. Understanding the regulation of the expression requires knowing both the protein factors and the DNA sites at which they act. The sites have traditionally been determined by isolating cis-acting mutations that affect expression and then determining the changes in the DNA that accompany the mutant phenotypes (1). More recently, regulatory proteins have been used to affinity purify the DNA regions to which they bind, and "footprinting" techniques have further delimited the binding sites (2). Each of these methods is time consuming and gives only partial information about the binding site. The final determination of the binding site pattern usually includes a comparison of many example sites. A method to determine the binding sites from the DNA sequences alone could greatly facilitate the process. Current sequencing technology is rapid enough that the most efficient means of determining the binding specificity of a protein may be to sequence a collection of regions known to contain binding sites. These may be a set of restriction fragments that are shown to bind the protein or a collection of DNA segments to which binding sites have been mapped. Since each of the fragments contains a binding site, the pattern of bases recognized by the protein should be discernible as the most significant pattern in the collection.

Regulatory Patterns

The difficulty arises that binding site patterns are not usually simple strings of bases. For example, *Escherichia coli* promoter sequences have two highly conserved parts, called the -35 and -10 regions (3, 4). The consensus sequences for those are TTGACA and TATAAT, respectively. The con-

sensus spacing between those regions is 17 bases, but other spacings are also allowed. An individual promoter may match the consensus at only a few positions and, while some positions are more conserved than others, no position is absolutely conserved. The most conserved bases in the -10region are TAnnnT, but only $\approx 65\%$ of all promoters even match this limited criterion (4). This means that methods of identifying the binding pattern that rely on common substrings, or "words," will likely fail.

A better representation of a protein's binding specificity than a consensus sequence is a matrix that has an element for each possible base at each position of the site (5). The matrix elements represent contributions of the individual bases to the protein–DNA interaction. The affinity of the protein for any site depends on the sum of all the interactions between the DNA and the protein. The individual interactions may or may not be independent; the simplest representation assumes independence. If the sequences of several binding sites are known, they can be used to construct a matrix representation of the protein that will give scores to individual sites that correlate well with the relative binding affinities of those sites (5-8).

The sequences of 23 identified binding sites of cAMP receptor protein (CRP) are shown in Fig. 1A (8, 9). Fig. 1B is a matrix whose elements are the frequency that each base occurs at each position within the CRP-binding sites. Fig. 1C is the matrix representation of CRP specificity, based on the information at each position of the site (5, 7, 8). The matrix is a representation of the specificity of the binding protein and can be used to search for new sites (5, 11-13). It can also be used to rank the affinities of different sites, and matrices of this type usually do well as predictors of quantitative activity (5-8). Fig. 1D shows the "information content" at each position of the site (10), derived from the formula

$$I_{\text{seq}} = \sum_{b=A}^{I} f_{b} \log_2 \left(\frac{f_{b}}{p_{b}} \right), \quad [1]$$

where f_b is the observed frequency of each base in the collection of sites and p_b is the fraction of each base in the genome. Note that the f_b terms are the elements of the matrix of Fig. 1*B*, and the log₂(f_b/p_b) terms are the elements of the matrix of Fig. 1*C*. Therefore, the "information content" plotted in Fig. 1*D* is the dot product of the position (column) vectors from those two matrices. The "information content" is a measure of how constrained the choice of bases is at each position in the binding sites (5, 7, 8, 10).

The Algorithm

The problem of identifying the binding sites from a collection of unaligned sequences is equivalent to finding an alignment that maximizes the "information content," at least within a local "window" that is the width of the binding site. That is, when the sequences are aligned by their binding sites there will be a peak of "information content," as in Fig. 1D, that

The publication costs of this article were defrayed in part by page charge payment. This article must therefore be hereby marked "*advertisement*" in accordance with 18 U.S.C. §1734 solely to indicate this fact.

Abbreviation: CRP, cAMP receptor protein.

Col E1 s Col E1 s ara site : bgl R m crp cya deo P2 si gal ilv B lac site 1 lac site 2 mal E mal K mal T comp A tna A uxu AB pBR P4 cat site 1 tdc	ite 2 ite 1 it it it te 2 ite 1	דד אדא G איז אד אדא דד א G איז אד אד א A A A G T C א א א	T T A T A T G T A A A A A T T A T A G G C A T	T T A A C A G A T A A A T C C T G T T G C A T	TT & TT TT TT TT CT T TT T CT T TT T	0 T 0 T 0 C 0 T 0 T 0 C 0 C 0 C 0 C 0 C	ŦŦŦŦŦĊŦŦŦŦŦŦŦŦŦŦŦŦŦŦŦŦŦ	00000 < 100 < 000 < 0000 < 0000 o 0000	G A T C A A A A A A T A A A A A A A A A A	UTUAGAAATTTGGUAUUTTAUUG	A C H C C G A C G C C H C A C A G H G A G G H	Τ G A G A G T C T C A T G G T C G C T T G T G	C T T G T A T A G A A A G A A A A G G A A T G	G T A C G C G G T T C G A G A G G A G C A G T	G T A G G G A A A G C C T A A T T T T C G A C	G T T T T T T T T T T C T A T C G T T T T G A T G	000000000000000000000000000000000000000	G	<pre><pre><pre><pre><pre><pre><pre><pre></pre></pre></pre></pre></pre></pre></pre></pre>	G A G A A A G A A A A T A A G A A A C A C C A	A A G C T T T T G C A C T C G T C T C G T A T	***	T A A T T A T A G T T T T T A C C T T A T C G T
В																							
6	A C G T	0.48 0.04 0.09 0.39	0.48 0.00 0.13 0.39	0.39 0.13 0.13 0.35	0.04 0.09 0.00 0.87	0.00 0.04 0.78 0.17	0.04 0.04 0.00 0.91	0.13 0.00 0.83 0.04	0.83 0.04 0.04 0.09	0.26 0.30 0.17 0.26	0.22 0.35 0.26 0.17	0.13 0.17 0.35 0.35	0.48 0.04 0.26 0.22	0.22 0.17 0.44 0.17	0.31 0.17 0.26 0.26	0.09 0.09 0.17 0.65	0.09 0.87 0.00 0.04	0.65 0.09 0.22 0.04	0.26 0.65 0.04 0.04	0.65 0.13 0.17 0.04	0.17 0.26 0.17 0.39	0.30 0.09 0.00 0.61	0.26 0.13 0.09 0.52
C																							
	A C G T	0.94 -2.64 -1.47 0.64	0.94 -2.75 -0.94 0.64	0.64 -0.94 -0.94 0.49	-2.64 -1.47 -2.75 1.80	-2.75 -2.63 1.66 -0.54	-2.63 -2.63 -2.75 1.88	-0.94 -2.75 1.73 -2.64	1.73 -2.64 -2.64 -1.47	0.07 0.28 -0.54 0.07	-0.18 0.49 0.06 -0.56	-0.94 -0.56 0.49 0.49	0.94 -2.64 0.06 -0.18	-0.18 -0.56 0.82 -0.56	0.31 -0.56 0.06 0.06	-1.47 -1.47 -0.56 1.38	-1.47 1.80 -2.75 -2.64	1.38 -1.47 -0.18 -2.64	0.07 1.39 -2.63 -2.63	1.39 -0.93 -0.54 -2.63	-0.54 0.07 -0.54 0.66	0.26 -1.47 -2.75 1.29	0.06 -0.94 -1.47 1.06
D																							
Iseq	1.4 1.2 1.0 0.8 0.6 0.4 0.2 0.0		~	\checkmark		✓ 	<u> </u>	`					<u> </u>		4	/	<u> </u>	\	<u> </u>	~~	\bigvee	\bigwedge	<u> </u>
		1	2	3	4	5	6	7	8	9	10	11	12	13	14	15	16	17	18	19	20	21	22
												Pos	ition										

Biochemistry: Stormo and Hartzell ۸

1184

FIG. 1. CRP binding sites. (A) Twenty-three sites identified as binding to CRP (8, 9). (B) The frequency with which each base occurs at each position in the CRP binding sites. (C) The specificity matrix for the protein, based on the binding sites, which is calculated as $\log_2(f_b/p_b)$, where f_b is the observed frequency of each base (from the matrix above) and p_b is the *a priori* probability of obtaining base *b*. In this example $p_b = b_b = b_b + b_b$ 0.25 for all b, approximating the E. coli genomic composition. At positions for which $f_b = 0$, an estimated frequency of 0.5/23 is used in the calculation. (D) The "information content" at each position of the CRP binding sites is plotted (10). The sum of all positions is 13.06 bits.

is greater (statistically more significant) than other peaks that occur by chance. A rigorous approach to identifying the binding sites would be to calculate, for every possible alignment of the fragments, the "information content." However, if there are N sequences, each of length L, then there are L^N possible alignments to consider, impractical for any interesting case. Methods have been developed for finding significant patterns in approximately aligned sequences (14, 15), but these are not useful when no alignment information is available. An alternative approach is to assume that the specificity of the protein can be represented as a matrix and to search for a matrix that gives at least one high-scoring site on each fragment. We have developed an algorithm that adds sequences to the analysis one at a time, at each step keeping many matrices (23). As more sequences are added to the analysis, the matrix representing the binding sites of the protein emerges from the background of possible alignments as the one with the highest "information content.'

The basic algorithm can be outlined as follows:

(i) Each of the k-words (k-long substrings, with k a variable chosen by the user) of the first sequence are used as an initial set of "interesting" matrices. These initial matrices contain one element of value 1.00 in each column, corresponding to the base at that position of the sequence, and all other elements are of value 0. Each of these words is a potential binding site, and with no further information they are equally likely.

(ii) The next sequence on the list is added to the analysis. Each of the set of interesting matrices is compared with each position of the new sequence.

(iii) A new set of interesting matrices is generated by updating some of the previous matrices by the addition of sites from the new sequence.

(iv) Steps 2 and 3 are repeated until all sequences have been included in the analysis.

The set of new interesting matrices may be selected in several ways. For instance, a threshold of "information content" could be used, and all new matrices that exceed that threshold would be kept for the next step. Alternatively, if computer memory is adequate to store X matrices, then the X-best might be kept for the next step. We have used a simpler criterion in the example below. Each existing matrix is updated with its best scoring site on the new sequence. When there are ties, which occur when two (or more) sites score highest from a particular matrix, each site is kept. This can lead to a slight increase in the number of matrices at each successive step, but the total number is approximately equal to the length of the first sequence.

The choice of the matrix width is not especially critical; the size of the matrix need not match the binding-site size for the method to work. When the matrix is larger than the binding



FIG. 2. The example data set. These are 105-base regions surrounding the sites in Fig. 1A (shown capitalized). They were obtained from GenBank Release 55, with the LOCUS name shown at the left. (The *tdc* gene is not in that release; the sequence was obtained from ref. 16.)

site, several matrices will be found that each include the binding sites within. When the matrix is smaller than the binding site, several overlapping matrices will be found that combine to give the entire site. The important consideration is that the matrix be large enough for the information in the binding sites to stand out against the background of other matrices. We have chosen to start with a matrix width of 20 bases; this is a fairly typical size for a binding site, at least for prokaryotic regulatory proteins (10).

An Example

Fig. 2 shows the data used in this example. Each sequence is 105 bases long and contains at least one CRP-binding site. In each case, the DNA strand shown is the one that appears in GenBank, and the order of sequences is alphabetical by GenBank LOCUS name. The regulated genes can occur at either end of these sequences. This example is typical of the type of data one might have from which to deduce binding sites for a protein. Were the orientation of the sites not known and no assumption of symmetry possible, both DNA strands would have to be compared.

The 86 20-long words of the first sequence constitute the initial set of matrices. As described above, each of these words is compared with each of the 20-long words of the next sequence and the best match for each matrix kept as a two-sequence matrix. As there was one tie, 87 matrices are kept at this step. Each of those is then compared with the 20-long words of the next sequence, and the best match to each matrix is kept again. This procedure is followed until all 18 sequences have been included; the total number of matrices at that point is 94. The "information content" of each matrix is calculated by Eq. 1. Fig. 3A shows the distribution of "information content" of these 94 matrices. Three matrices stand out as clearly significant above the others. These are overlapping matrices that collectively cover 22 bases, essentially as in Fig. 1. The other cluster of 13 matrices below the most significant, but still above the main distribution, are variations on the binding sites. In some cases they overlap parts of the binding site but not all of it. In other cases they overlap it completely but have some additional sites that decrease the total "information content."

Examination of the three best matrices reveals two things. (i) The overlap shows that the site is better represented as a 22-long matrix. (ii) The most significant part of each matrix is contained in a central region of only 16 bases, positions 4–19 of Fig. 1. The entire analysis was repeated using each of these alternative matrix widths. In each case a single best matrix was found that stands out above the distribution (Fig. 3 B and C) and represents the known binding sites. The best scoring 16-wide matrix is shown in Fig. 4. This "information content" is very similar to that obtained from the collection of known binding sites (Fig. 1). The "information content" in Fig. 4C is actually greater and somewhat more symmetric than from positions 4–19 of Fig. 1D. This difference is due to the inclusion in Fig. 1A of several second, and presumably weaker, sites that lower the total information and that are least conserved on the 3' side. The matrices of Fig. 4 come from only one site per sequence, thereby selecting for the most highly conserved collection of sites.

The matrices of Fig. 4 have the highest "information content" of all the 16-wide matrices obtained by our algorithm. To test whether the matrix of Fig. 4B represents the specificity of CRP, every position of each sequence shown in Fig. 2 was evaluated with the matrix (as described in refs. 5 and 11–13). All the identified CRP-binding sites were among the highest scoring sites. With two exceptions, the highest scoring site on each sequence was an identified CRP-binding site. These two exceptions are interesting to examine. The most highly conserved bases in the binding sites constitute a symmetric consensus sequence of TGTGAnnnnnTCACA (Fig. 1 and ref. 8). The identified CRP-binding site for the malK gene has six matches to that, five on the 5' side (Fig.



FIG. 3. Distribution of "information content" of the matrices at the end of each analysis, described in the text. Count is the number of matrices with "information content" in the interval shown. (A) The 94 20-wide matrices. (B) The 93 22-wide matrices. (C) The 93 16-wide matrices.



FIG. 4. The best 16-wide matrix. The positions are numbered 4-19, corresponding to the central positions of Fig. 1. (A) The frequency of each base for the sites included in the best matrix, as in Fig. 1B. (B) The specificity matrix determined from the frequency matrix, as in Fig. 1C. In this case the *a priori* values were determined from the data set shown in Fig. 2: $p_A = 0.30$; $p_C = 0.18$; $p_G = 0.21$; and $p_T = 0.31$. Analyses were also performed using $p_b = 0.25$ for all b. In that case, essentially the same matrix remained the best, but the distribution had more high-scoring matrices due to the high probability of adenine and thymine matches. The specificity matrix values for positions with $f_b = 0$ were estimated using $f_b = 0.5/18$ from the 18 sequences in the data set. (C) The "information content" at each position of the matrix is plotted. The sum from all positions is 12.15 bits.

1). The site with the highest score by the matrix of Fig. 4B occurs 3' of the one identified, with seven matches to the consensus, cGTGAtgttgcTtgCA (Fig. 2). It is conceivable that the CRP protein binds to this site instead of the one identified or that it binds to both (17).

The other exception is in the cat regulatory region. Two binding sites have been identified (Figs. 1 and 2) and shown to interact with the CRP protein in footprinting experiments (18). Site 1 has six matches with the consensus, and site 2 has five matches. Site 1 demonstrated the tightest binding of the two sites. Although both of these sites score high with the matrix of Fig. 4, the highest scoring site on the fragment is between them, gGTGtccctgtTgAtA, also having six matches with the consensus. Binding to this site has not been reported. Site 2 is especially interesting because it has nine matches with the consensus when an additional base is allowed between the two highly conserved regions, TGT-GAcggccgcTCACt. The fact that site 1 binds more tightly than site 2 suggests that either the variable spacing is not allowed or it causes enough strain in the interaction to limit the gain from additional consensus-base contacts. Variable spacing between the conserved regions of regulatory sites, excluding promoters, are not a common feature, although one has been verified (19). Our method could be amended to allow for a small number of gaps, perhaps one per site, but the complexity of the algorithm would increase substantially.

Conclusions

The complexity of the algorithm used in this study is one of its most appealing features. The memory required is independent of the number of sequences and linearly dependent on their lengths. The time required is linearly dependent on the number of sequences and the square of their lengths. This is in sharp contrast to rigorous methods that require comparing all possible alignments that require $O(L^N)^*$ time and space. The savings comes from selecting only a small subset of possible alignments to consider those that are "interesting," in that they are the best matches on each sequence to a set of possible matrix representations of the protein's specificity. The number of matrices that are considered at each step could be quite large, up to the limit of the computer memory available, although in the CRP example we needed to keep only one for each possible word in the first sequence. Reliability of the best matrix representation increases with the number of sequences. In the CRP example, in which we know the pattern of the binding sites, one pass through the data was sufficient to determine that pattern and identify the sites. If we hadn't known the answer, we would randomize the order of sequences several times and repeat the analysis to see whether we always got a similar answer or, if not, determine the variation among the most significant patterns.

As the amount of determined DNA sequence increases, methods that identify the important features of those sequences using only the information within the sequences will become increasingly important. Projects like the Human Genome Initiative (20) will generate such an enormous amount of sequence data that efficient methods of pattern identification will be essential to elucidate those features. We present a method to help in one such task—that of identifying regulatory protein-binding sites. Because the method relies on an information measure of similarity among sequences, which has been shown to give reliable representations of protein specificity (5, 7, 8), we expect the method to work for any sequence-specific DNA-binding proteins. Variations on the basic method should also be applicable to similar prob-

^{*}The notation $O(L^N)$ is read "order L^N " and means that the calculation time and space (computer memory) required is $aL^N + b$, where a and b are constants. The algorithm described in this paper is O(L) in space and $O(L^2N)$ in time.

lems in which functionally analogous sites can be identified by information within linear sequences. Functional domains of some proteins may be of this type (21, 22), and related methods may be useful in identifying them.

We thank Drs. Gerald Hertz and Charles E. Lawrence for stimulating discussions and helpful suggestions. We are grateful for comments on the manuscript from Drs. Susan Dutcher, Matt Scott, Bill Wood, and Mike Yarus. Sequences were extracted from Gen-Bank using the EuGene software from the Molecular Biology Information Resource at Baylor College of Medicine (Houston). This work was supported by Grant GM28755 from the National Institutes of Health.

- 1. Gilbert, W., Gralla, J., Majors, J. & Maxam, A. (1975) in *Protein-Ligand Interactions*, eds. Sund, H. & Blauer, G. (de Gruyter, Berlin), pp. 193-206.
- Kadonaga, J. T., Jones, K. A. & Tjian, R. (1986) Trends Biochem. Sci. 11, 20-23.
- 3. Hawley, D. K. & McClure, W. R. (1983) Nucleic Acids Res. 11, 2237-2255.
- 4. Harley, C. B. & Reynolds, R. P. (1987) Nucleic Acids Res. 15, 2343–2361.
- 5. Stormo, G. D. (1988) Annu. Rev. Biophys. Biophys. Chem. 17, 241-263.
- Mulligan, M. E., Hawley, D. K., Entriken, R. & McClure, W. R. (1984) Nucleic Acids Res. 12, 789-800.
- 7. Berg, O. G. & von Hippel, P. H. (1987) J. Mol. Biol. 193, 723-750.

- 8. Berg, O. G. & von Hippel, P. H. (1988) J. Mol. Biol. 200, 209-223.
- 9. de Crombrugghe, B., Busby, S. & Buc, H. (1984) Science 224, 831-838.
- Schneider, T. D., Stormo, G. D., Gold, L. & Ehrenfeucht, A. (1986) J. Mol. Biol. 188, 415-431.
- 11. Staden, R. (1984) Nucleic Acids Res. 12, 505-519.
- 12. Staden, R. (1988) Comput. Appl. Biosci. 4, 53-60.
- 13. Stormo, G. D. (1987) in Nucleic Acid and Protein Sequence Analysis: A Practical Approach, eds. Bishop, M. J. & Rawlings, C. J. (IRL, Oxford), pp. 231-258.
- Galas, D. J., Waterman, M. S. & Eggert, M. (1985) J. Mol. Biol. 186, 117-128.
- 15. Mengeritsky, G. & Smith, T. F. (1987) Comput. Appl. Biosci. 3, 323-328.
- 16. Goss, T. J. & Datta, P. (1985) Mol. Gen. Genet. 201, 308-314.
- Le Grice, S., Matzura, H., Marcoli, R., Iida, S. & Bickle, T. (1982) J. Bacteriol. 150, 312-318.
- Sadler, J. R., Sasmor, H. & Betz, J. L. (1983) Proc. Natl. Acad. Sci. USA 80, 6785–6789.
- 20. Lewin, R. (1988) Science 240, 602-604.
- Gribskov, M., McLachlan, A. D. & Eisenberg, D. (1987) Proc. Natl. Acad. Sci. USA 84, 4355-4358.
- 22. Gribskov, M., Homyak, M., Edenfield, J. & Eisenberg, D. (1988) Comput. Appl. Biosci. 4, 61-66.
- 23. Bacon, D. J. & Anderson, W. F. (1986) J. Mol. Biol. 191, 153-161.